

Breaking the sample complexity barrier to regret-optimal model-free reinforcement learning

GEN LI

Department of Statistics and Data Science, The Wharton School, University of Pennsylvania,
Philadelphia, PA 19104, USA

LAIXI SHI

Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh,
PA 15213, USA

YUXIN CHEN

Department of Statistics and Data Science, The Wharton School, University of Pennsylvania,
Philadelphia, PA 19104, USA

AND

YUEJIE CHI

Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh,
PA 15213, USA

[†]Corresponding author. Email: yuejiechi@cmu.edu —This paper was presented in part at the 2021 Conference on Neural Information Processing Systems (NeurIPS).

[Received on 10 October 2021; revised on 22 August 2022; accepted on 16 October 2022]

Abstract

Achieving sample efficiency in online episodic reinforcement learning (RL) requires optimally balancing exploration and exploitation. When it comes to a finite-horizon episodic Markov decision process with S states, A actions and horizon length H , substantial progress has been achieved toward characterizing the minimax-optimal regret, which scales on the order of $\sqrt{H^2SAT}$ (modulo log factors) with T the total number of samples. While several competing solution paradigms have been proposed to minimize regret, they are either memory-inefficient, or fall short of optimality unless the sample size exceeds an enormous threshold (e.g. $S^6A^4 \text{poly}(H)$ for existing model-free methods).

To overcome such a large sample size barrier to efficient RL, we design a novel model-free algorithm, with space complexity $O(SAH)$, that achieves near-optimal regret as soon as the sample size exceeds the order of $SA \text{poly}(H)$. In terms of this sample size requirement (also referred to the initial burn-in cost), our method improves—by at least a factor of S^5A^3 —upon any prior memory-efficient algorithm that is asymptotically regret-optimal. Leveraging the recently introduced variance reduction strategy (also called *reference-advantage decomposition*), the proposed algorithm employs an *early-settled* reference update rule, with the aid of two Q-learning sequences with upper and lower confidence bounds. The design

principle of our early-settled variance reduction method might be of independent interest to other RL settings that involve intricate exploration–exploitation trade-offs.

Keywords: model-free RL; memory efficiency; variance reduction; Q-learning; upper confidence bounds; lower confidence bounds;

2010 Math Subject Classification: 68T09; 68T37; 68W27; 90C40.

1. Introduction

Contemporary reinforcement learning (RL) has to deal with unknown environments with unprecedentedly large dimensionality. How to make the best use of samples in the face of high-dimensional state–action space lies at the core of modern RL practice. An ideal RL algorithm would learn to act favorably even when the number of available data samples scales sub-linearly in the ambient dimension of the model, i.e. the number of parameters needed to describe the transition dynamics of the environment. The challenge is further compounded when this task needs to be accomplished with limited memory.

Simultaneously achieving the desired sample and memory efficiency is particularly challenging when it comes to online episodic RL scenarios. In contrast to the simulator setting that permits sampling of any state–action pair, an agent in online episodic RL is only allowed to draw sample trajectories by executing a policy in the unknown Markov decision process (MDP), where the initial states are pre-assigned and might even be chosen by an adversary. Careful deliberation needs to be undertaken when deciding what policies to use to allow for effective interaction with the unknown environment, how to optimally balance exploitation and exploration and how to process and store the collected information intelligently without causing redundancy.

1.1 *Regret-optimal model-free RL? A sample size barrier*

In order to evaluate and compare the effectiveness of RL algorithms in high dimension, a recent body of works sought to develop a finite-sample theoretical framework to analyze the algorithms of interest, with the aim of delineating the dependency of algorithm performance on all salient problem parameters in a non-asymptotic fashion [14, 33]. Such finite-sample guarantees are brought to bear toward understanding and tackling the challenges in the sample-starved regime commonly encountered in practice. To facilitate discussion, let us take a moment to summarize the state-of-the-art theory for episodic finite-horizon MDPs with non-stationary transition kernels, focusing on minimizing cumulative regret—a metric that quantifies the performance difference between the learned policy and the true optimal policy—with the fewest number of samples. Here and throughout, we denote by S , A and H the size of the state space, the size of the action space and the horizon length of the MDP, respectively, and let T represent the sample size. In addition, the immediate reward gained at each time step is assumed to lie between 0 and 1.

Fundamental regret lower bound. Following the arguments in [3, 28], the recent works [15, 29] developed a fundamental lower bound¹ on the expected total regret for this setting. Specifically, this lower bound claims that no matter what algorithm to use, one can find an MDP such that the accumulated

¹ It is worth emphasizing that [15] adopts the notation T to represent the number of trajectories (with each trajectory containing H samples), while this paper employs K to denote the number of sample trajectories and $T = KH$ the total number of samples.

regret incurred by the algorithm necessarily exceeds the order of

$$(\text{lower bound}) \quad \sqrt{H^2 SAT}, \quad (1.1)$$

as long as $T \geq H^2 SA$.² This sublinear regret lower bound in turn imposes a sampling limit if one wants to achieve ε average regret.

Model-based RL. Moving beyond the lower bound, let us examine the effectiveness of model-based RL—an approach that can be decoupled into a model estimation stage (i.e. estimating the transition kernel using available data) and a subsequent stage of planning using the learned model [2, 7, 20, 28, 48]. In order to ensure a sufficient degree of exploration, [7] came up with an algorithm called UCB-VI that blends model-based learning and the optimism principle, which achieves a regret bound³ $\tilde{O}(\sqrt{H^2 SAT})$ that nearly attains the lower bound (1.1) as T tends to infinity. Caution needs to be exercised, however, that existing theory does not guarantee the near optimality of this algorithm unless the sample size T surpasses

$$T \geq S^3 AH^6,$$

a threshold that is significantly larger than the dimension of the underlying model. This threshold can also be understood as the initial *burn-in cost* of the algorithm, namely, a sampling burden needed for the algorithm to exhibit the desired performance. In addition, model-based algorithms typically require storing the estimated probability transition kernel, resulting in a space complexity that could be as high as $O(S^2 AH)$ [7].

Model-free RL. Another competing solution paradigm is the model-free approach, which circumvents the model estimation stage and attempts to learn the optimal values directly [8, 29, 56, 70]. In comparison with the model-based counterpart, the model-free approach holds the promise of low space complexity, as it eliminates the need of storing a full description of the model. In fact, in a number of previous works (e.g. [29, 56]), an algorithm is declared to be model-free only if its space complexity is $o(S^2 AH)$ regardless of the sample size T .

- *Memory-efficient model-free methods.* [29] proposed the first memory-efficient model-free algorithm—which is an optimistic variant of classical Q-learning—that achieves a regret bound proportional to \sqrt{T} with a space complexity $O(SAH)$. Compared with the lower bound (1.1), however, the regret bound in [29] is off by a factor of \sqrt{H} and hence suboptimal for problems with long horizon. This drawback has recently been overcome in [75] by leveraging the idea of variance reduction (or the so-called ‘reference-advantage decomposition’) for large enough T . While the resulting regret matches the information-theoretic limit asymptotically, its optimality in the non-asymptotic regime is not guaranteed unless the sample size T exceeds (see (75, Lemma 7))

$$T \geq S^6 A^4 H^{28},$$

Consequently, the lower bound developed in [15] for non-stationary finite-horizon MDPs reads $\Omega(\sqrt{H^3 SAK})$, or equivalently, $\Omega(\sqrt{H^2 SAT})$ using the notation adopted herein.

² Given that a trivial upper bound on the regret is T , one needs to impose a lower bound $T \geq H^2 SA$ in order for (1.1) to be meaningful.

³ Here and throughout, we use the standard notation $f(n) = O(g(n))$ to indicate that $f(n)/g(n)$ is bounded above by a constant as n grows. The notation $\tilde{O}(\cdot)$ resembles $O(\cdot)$ except that it hides any logarithmic scaling. The notation $f(n) = o(g(n))$ means that $\lim_{n \rightarrow \infty} f(n)/g(n) = 0$.

TABLE 1 Comparisons between prior art and our results for non-stationary episodic MDPs when $T \geq H^2 SA$. The table includes the order of the regret bound, the range of sample sizes that achieve the optimal regret $\tilde{O}(\sqrt{H^2 SAT})$, and the memory complexity, with all logarithmic factors omitted for simplicity of presentation. The red text highlights the suboptimal part of the respective algorithms.

Algorithm	Regret	Range of sample sizes T that attain optimal regret	Space complexity
UCB-VI [7]	$\sqrt{H^2 SAT} + H^4 S^2 A$	$[S^3 A H^6, \infty)$	$S^2 A H$
UCB-Q-Hoeffding [29]	$\sqrt{H^4 SAT}$	never	SAH
UCB-Q-Bernstein [29]	$\sqrt{H^3 SAT} + \sqrt{H^9 S^3 A^3}$	never	SAH
UCB2-Q-Bernstein [8]	$\sqrt{H^3 SAT} + \sqrt{H^9 S^3 A^3}$	never	SAH
UCB-Q-Advantage [75]	$\sqrt{H^2 SAT} + H^8 S^2 A^{\frac{3}{2}} T^{\frac{1}{4}}$	$[S^6 A^4 H^{28}, \infty)$	SAH
UCB-M-Q [44]	$\sqrt{H^2 SAT} + H^4 SA$	$[SAH^6, \infty)$	$S^2 A H$
Q-EarlySettled-Advantage (this work)	$\sqrt{H^2 SAT} + H^6 SA$	$[SAH^{10}, \infty)$	SAH
Lower bound [15]	$\sqrt{H^2 SAT}$	n/a	n/a

a requirement that is even far more stringent than the burn-in cost imposed by [7].

- A *memory-inefficient ‘model-free’ variant*. The recent work [44] put forward a novel sample-efficient variant of Q-learning called UCB-M-Q, which relies on a carefully chosen momentum term for bias reduction. This algorithm is guaranteed to yield near-optimal regret $\tilde{O}(\sqrt{H^2 SAT})$ as soon as the sample size exceeds $T \geq SA\text{poly}(H)$, which is a remarkable improvement vis-à-vis previous regret-optimal methods [7, 75]. Nevertheless, akin to the model-based approach, it comes at a price in terms of the space and computation complexities, as the space required to store all bias-value function is $O(S^2 A H)$ and the computation required is $O(ST)$, both of which are larger by a factor of S than other model-free algorithms like [75]. In view of this memory inefficiency, UCB-M-Q falls short of fulfilling the definition of model-free algorithms in [29, 56]. See (44, Section 3.3) for more detailed discussions.

A more complete summary of prior results can be found in Table 1.

1.2 A glimpse of our contributions

In brief, while it is encouraging to see that both model-based and model-free approaches allow for near-minimal regret as T tends to infinity, they are either memory-inefficient, or require the sample size to exceed a threshold substantially larger than the model dimensionality. In fact, no prior algorithms have been shown to be *simultaneously regret-optimal and memory-efficient* unless

$$T \geq S^6 A^4 \text{poly}(H),$$

which constitutes a stringent sample size barrier constraining their utility in the sample-starved and memory-limited regime. The presence of this sample complexity barrier motivates one to pose a natural question:

Is it possible to design an algorithm that accommodates a significantly broader sample size range without compromising regret optimality and memory efficiency?

In this paper, we answer this question affirmatively, by designing a new model-free algorithm, dubbed as **Q-EarlySettled-Advantage**, that enjoys the following performance guarantee.

THEOREM 1.1. The proposed **Q-EarlySettled-Advantage** algorithm, which has a space complexity $O(SAH)$, achieves near-optimal regret $\tilde{O}(\sqrt{H^2SAT})$ as soon as the sample size exceeds $T \geq SA \text{poly}(H)$.

As can be seen from Table 1, the space complexity of the proposed algorithm is $O(SAH)$, which is far more memory-efficient than both the model-based approach in [7] and the **UCB-M-Q** algorithm in [44] (both of these prior algorithms require S^2AH units of space). In addition, the sample size requirement $T \geq SA \text{poly}(H)$ of our algorithm improves—by a factor of at least S^5A^3 —upon that of any prior algorithm that is simultaneously regret-optimal and memory-efficient. In fact, this requirement is nearly sharp in terms of the dependency on both S and A , and was previously achieved only by the **UCB-M-Q** algorithm at a price of a much higher storage burden.

Let us also briefly highlight the key ideas of our algorithm. As an optimistic variant of variance-reduced Q-learning, **Q-EarlySettled-Advantage** leverages the recently introduced reference-advantage decompositions for variance reduction [75]. As a distinguishing feature from prior algorithms, we employ an *early-stopped* reference update rule, with the assistance of two Q-learning sequences that incorporate upper and lower confidence bounds (LCBs), respectively. The design of our early-stopped variance reduction scheme, as well as its analysis framework, might be of independent interest to other settings that involve managing intricate exploration–exploitation trade-offs.

1.3 Related works

We now take a moment to discuss a small sample of other related works. We limit our discussions primarily to RL algorithms in the tabular setting with finite state and action spaces, which are the closest to our work. The readers interested in those model-free variants with function approximation are referred to [19, 22, 45] and the references therein.

Probably approximately correct (PAC) bounds for synchronous and asynchronous Q-learning. Q-learning is arguably among the most famous model-free algorithms developed in the RL literature [26, 57, 60, 65], which enjoys a space complexity $O(SAH)$. Non-asymptotic sample analysis and PAC bounds have seen extensive developments in the last several years, including but not limited to the works of [10, 12, 21, 37, 62] for the synchronous setting (the case with access to a generative model or a simulator), and the works of [10, 13, 21, 42, 50] for the asynchronous setting (where one observes a single Markovian trajectory induced by a behavior policy). Finite-time guarantees of other variants of Q-learning have also been developed; partial examples include speedy Q-learning [5], double Q-learning [68], variance-reduced Q-learning [42, 63], momentum Q-learning [67], pessimistic Q-learning [53] and Q-learning for linearly parameterized MDPs [64]. This line of works did not account for exploration, and hence the success of Q-learning in these settings heavily relies on the access to a simulator or a behavior policy with sufficient coverage over the state-action space.

Regret analysis for model-free RL with exploration. When it comes to online episodic RL (so that a simulator is unavailable), regret analysis is the prevailing analysis paradigm employed to capture the

trade-off between exploration and exploitation. A common theme is to augment the original model-free update rule (e.g. the Q-learning update rule) by an exploration bonus, which typically takes the form of, say, certain upper confidence bounds (UCBs) motivated by the bandit literature [4, 35]. In addition to the ones in Table 1 for episodic finite-horizon settings, sample-efficient model-free algorithms have been investigated for infinite-horizon MDPs as well [16, 27, 43, 70, 74, 76].

Variance reduction in RL. The seminal idea of variance reduction was originally proposed to accelerate finite-sum stochastic optimization, e.g. [24, 32, 46]. Thereafter, the variance reduction strategy has been imported to RL, which assists in improving the sample efficiency of RL algorithms in multiple contexts, including but not limited to policy evaluation [17, 34, 61, 69], RL with a generative model [54, 55, 63], asynchronous Q-learning [42] and offline RL [53, 71].

Model-based approach. Model-based RL is known to be minimax-optimal in the presence of a simulator [1, 6, 41], beating the state-of-the-art model-free algorithms by achieving optimality for the entire sample size range [41]. When it comes to online episodic RL, [7] was the first work that managed to achieve near-optimal regret (at least for large T); in fact, this was also the first result (for any algorithm) matching existing lower bounds for large T . The sample efficiency of the model-based approach has subsequently been established for other settings, including but not limited to discounted infinite-horizon MDPs [25], MDPs with bounded total reward [72, 74], offline RL [40, 52] and Markov games [39, 73].

Regret lower bound. Inspired by the classical lower bound argument developed for multi-armed bandits [3], the work [28] established a regret lower bound for MDPs with finite diameters (so that for an arbitrary pair of states, the expected time to transition between them is assumed to be finite as long as a suitable policy is used), which has been reproduced in the note [47] with the purpose of facilitating comparison with [9]. The way to construct hard MDPs in [28] has since been adapted by [29] to exhibit a lower bound on episodic MDPs (with a sketched proof provided therein). It was recently revisited in [15], which presented a detailed and rigorous proof argument with a different construction.

2. Problem formulation

In this section, we formally describe the problem setting. Here and throughout, we denote by $\Delta(\mathcal{S})$ the probability simplex over a set \mathcal{S} , and introduce the notation $[M] := \{1, \dots, M\}$ for any integer $M > 0$.

Basics of finite-horizon MDPs. Let $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \{P_h\}_{h=1}^H, \{r_h\}_{h=1}^H)$ represent a finite-horizon MDP, where $\mathcal{S} := \{1, \dots, S\}$ is the state space of size S , $\mathcal{A} := \{1, \dots, A\}$ is the action space of size A , H denotes the horizon length and $P_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ (resp. $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$) represents the probability transition kernel (resp. reward function) at the h -th time step, $1 \leq h \leq H$, respectively. More specifically, $P_h(\cdot | s, a) \in \Delta(\mathcal{S})$ stands for the transition probability vector from state s at time step h when action a is taken, while $r_h(s, a)$ indicates the immediate reward received at time step h for a state-action pair (s, a) (which is assumed to be deterministic and fall within the range $[0, 1]$). The MDP is said to be non-stationary when the P_h 's are not identical across $1 \leq h \leq H$. A policy of an agent is represented by $\pi = \{\pi_h\}_{h=1}^H$ with $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$ the action selection rule at time step h , so that $\pi_h(s)$ specifies which action to execute in state s at time step h . Throughout this paper, we concentrate on deterministic policies.

Value functions, Q-functions and Bellman equations. The value function $V_h^\pi(s)$ of a (deterministic) policy π at step h is defined as the expected cumulative rewards received between time steps h and H when executing this policy from an initial state s at time step h , namely,

$$V_h^\pi(s) := \mathbb{E}_{s_{t+1} \sim P_t(\cdot | s_t, \pi_t(s_t)), t \geq h} \left[\sum_{t=h}^H r_t(s_t, \pi_t(s_t)) \mid s_h = s \right], \quad (2.1)$$

where the expectation is taken over the randomness of the MDP trajectory $\{s_t \mid h \leq t \leq H\}$. The action-value function (a.k.a. the Q-function) $Q_h^\pi(s, a)$ of a policy π at step h can be defined analogously except that the action at step h is fixed to be a , that is,

$$Q_h^\pi(s, a) := r_h(s, a) + \mathbb{E}_{\substack{s_{h+1} \sim P_h(\cdot|s, a), \\ s_{t+1} \sim P_t(\cdot|s_t, \pi_t(s_t)), t > h}} \left[\sum_{t=h+1}^H r_t(s_t, \pi_t(s_t)) \mid s_h = s, a_h = a \right]. \quad (2.2)$$

In addition, we define $V_{H+1}^\pi(s) = Q_{H+1}^\pi(s, a) = 0$ for any policy π and any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$. By virtue of basic properties in dynamic programming [11], the value function and the Q-function satisfy the following Bellman equation:

$$Q_h^\pi(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot|s, a)} [V_{h+1}^\pi(s')]. \quad (2.3)$$

A policy $\pi^* = \{\pi_h^*\}_{h=1}^H$ is said to be an optimal policy if it maximizes the value function simultaneously for all states among all policies. The resulting optimal value function $V^* = \{V_h^*\}_{h=1}^H$ and optimal Q-functions $Q^* = \{Q_h^*\}_{h=1}^H$ satisfy

$$V_h^*(s) = V_h^{\pi^*}(s) = \max_{\pi} V_h^\pi(s) \quad \text{and} \quad Q_h^*(s, a) = Q_h^{\pi^*}(s, a) = \max_{\pi} Q_h^\pi(s, a) \quad (2.4)$$

for any $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$. It is well known that the optimal policy always exists [49], and satisfies the Bellman optimality equation:

$$\forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : \quad Q_h^*(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot|s, a)} [V_{h+1}^*(s')]. \quad (2.5)$$

Online episodic RL. This paper investigates the online episodic RL setting, where the agent is allowed to execute the MDP sequentially in a total number of K episodes each of length H . This amounts to collecting

$$T = KH \text{ samples}$$

in total. More specifically, in each episode $k = 1, \dots, K$, the agent is assigned an arbitrary initial state s_1^k (possibly by an adversary), and selects a policy $\pi^k = \{\pi_h^k\}_{h=1}^H$ learned based on the information collected up to the $(k-1)$ -th episode. The k -th episode is then executed following the policy π^k and the dynamic of the MDP \mathcal{M} , leading to a length- H sample trajectory.

Goal: regret minimization. In order to evaluate the quality of the learned policies $\{\pi^k\}_{1 \leq k \leq K}$, a frequently used performance metric is the cumulative regret defined as follows:

$$\text{Regret}(T) := \sum_{k=1}^K (V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k)). \quad (2.6)$$

In words, the regret reflects the sub-optimality gaps between the values of the optimal policy and those of the learned policies aggregated over K episodes. A natural objective is thus to design a sample-optimal algorithm, namely, an algorithm whose resulting regret scales optimally in the sample size T .

Accomplishing this goal requires carefully managing the trade-off between exploration and exploitation, which is particularly challenging in the sample-limited regime.

Notation. Before presenting our main results, we take a moment to introduce some convenient notation to be used throughout the remainder of this paper. For any vector $x \in \mathbb{R}^{SA}$ that constitutes certain quantities for all state-action pairs, we shall often use $x(s, a)$ to denote the entry associated with the state-action pair (s, a) , as long as it is clear from the context. We shall also let

$$P_{h,s,a} = P_h(\cdot | s, a) \in \mathbb{R}^{1 \times S} \quad (2.7)$$

abbreviate the transition probability vector given the (s, a) pair at time step h . Additionally, we denote by e_i the i -th standard basis vector, with the only non-zero element being in the i -th entry and equal to 1.

3. Algorithm and theoretical guarantees

In this section, we present the proposed algorithm called **Q-EarlySettled-Advantage**, as well as the accompanying theory confirming its sample and memory efficiency.

3.1 Review: *Q*-learning with UCB exploration and reference advantage

This subsection briefly reviews the *Q*-learning algorithm with UCB exploration proposed in [29], as well as a variant that further exploits the idea of variance reduction [75]. These two model-free algorithms inspire the algorithm design in this paper.

Q-learning with UCB exploration (UCB-Q or UCB-Q-Hoeffding). Recall that the classical *Q*-learning algorithm has been proposed as a stochastic approximation scheme [51] to solve the Bellman optimality equation (2.5), which consists of the following update rule [65, 66]:

$$Q_h(s, a) \leftarrow (1 - \eta)Q_h(s, a) + \eta \left\{ r_h(s, a) + \underbrace{\widehat{P}_{h,s,a} V_{h+1}}_{\text{stochastic estimate of } P_{h,s,a} V_{h+1}} \right\}. \quad (3.1)$$

Here, Q_h (resp. V_h) indicates the running estimate of Q_h^* (resp. V_h^*), η is the (possibly iteration-varying) learning rate or stepsize and $\widehat{P}_{h,s,a} V_{h+1}$ is a stochastic estimate of $P_{h,s,a} V_{h+1}$ (cf. (2.7)). For instance, if one has available a sample (s, a, s') transitioning from state s at step h to s' at step $h+1$ under action a , then a stochastic estimate of $P_{h,s,a} V_{h+1}$ can be taken as $V_{h+1}(s')$, which is unbiased in the sense that

$$\mathbb{E}[V_{h+1}(s')] = P_{h,s,a} V_{h+1}.$$

To further encourage exploration, the algorithm proposed in [29]—which shall be abbreviated as **UCB-Q** or **UCB-Q-Hoeffding** hereafter—augments the *Q*-learning update rule (3.1) in each episode via an additional exploration bonus:

$$Q_h^{\text{UCB}}(s, a) \leftarrow (1 - \eta)Q_h^{\text{UCB}}(s, a) + \eta \{ r_h(s, a) + \widehat{P}_{h,s,a} V_{h+1} + b_h \}. \quad (3.2)$$

The bonus term $b_h \geq 0$ is chosen to be a certain UCB for $(\widehat{P}_{h,s,a} - P_{h,s,a})V_{h+1}$, an exploration-efficient scheme that originated from the bandit literature [35, 36]. The algorithm then proceeds to the next episode by executing/sampling the MDP using a greedy policy w.r.t. the updated *Q*-estimate. These steps are repeated until the algorithm is terminated.

Q-learning with UCB exploration and reference advantage (UCB-Q-Advantage). The regret bounds derived for UCB-Q [29], however, fall short of being optimal, as they are at least a factor of \sqrt{H} away from the fundamental lower bound. In order to further shave this \sqrt{H} factor, one strategy is to leverage the idea of variance reduction to accelerate convergence [32, 42, 55, 63]. An instantiation of this idea for the regret setting is a variant of UCB-Q based on reference-advantage decomposition, which was put forward in [75] and shall be abbreviated as UCB-Q-Advantage throughout this paper.

To describe the key ideas of UCB-Q-Advantage, imagine that we are able to maintain a collection of reference values $V^R = \{V_h^R\}_{h=1}^H$, which form reasonable estimates of $V^* = \{V_h^*\}_{h=1}^H$ and become increasingly more accurate as the algorithm progresses.

At each time step h , the algorithm adopts the following update rule:

$$Q_h^R(s, a) \leftarrow (1 - \eta)Q_h^R(s, a) + \eta \left\{ r_h(s, a) + \underbrace{\widehat{P}_{h,s,a}(V_{h+1} - V_{h+1}^R)}_{\text{stochastic estimate of } P_{h,s,a}(V_{h+1} - V_{h+1}^R)} + \widehat{[P_h V_{h+1}^R]}(s, a) + b_h^R \right\}. \quad (3.3)$$

Two ingredients of this update rule are worth noting:

- Akin to the UCB-Q case, we can take $\widehat{P}_{h,s,a}(V_{h+1} - V_{h+1}^R)$ to be the stochastic estimate $V_{h+1}(s') - V_{h+1}^R(s')$ if we observe a sample transition (s, a, s') at time step h . If V_{h+1} is fairly close to the reference V_{h+1}^R , then this stochastic term can be less volatile than the stochastic term $\widehat{P}_{h,s,a}V_{h+1}$ in (3.2).
- Additionally, the term $\widehat{P_h V_{h+1}^R}$ indicates an estimate of the one-step look-ahead value $P_h V_{h+1}^R$, which shall be computed using a batch of samples.
- The variability of $\widehat{P_h V_{h+1}^R}$ can be well controlled through the use of batch data, at the price of an increased sample size.

Accordingly, the exploration bonus term b_h^R is taken to be a UCB for the above-mentioned two terms combined. Given that the uncertainty of (3.3) largely stems from these two terms (which can both be much smaller than the variability in (3.2)), the incorporation of the reference term helps accelerate convergence.

3.2 The proposed algorithm: Q-EarlySettled-Advantage

As alluded to previously, however, the sample size required for UCB-Q-Advantage to achieve optimal regret needs to exceed a large polynomial $S^6 A^4$ in the size of the state/action space. To overcome this sample complexity barrier, we come up with an improved variant called Q-EarlySettled-Advantage.

Motivation: early settlement of a reference value. An important insight obtained from previous algorithm designs is that in order to achieve low regret, it is desirable to maintain an estimate of Q -function that (i) provides an optimistic view (namely, an over-estimate) of the truth Q^* , and (ii) mitigates the bias $Q - Q^*$ as much as possible. With two additional optimistic Q-estimates in hand— Q_h^{UCB} based on UCB-Q, and a reference Q_h^R —it is natural to combine them as follows to further reduce the bias without violating the optimism principle:

$$Q_h(s_h, a_h) \leftarrow \min \left\{ Q_h^R(s_h, a_h), Q_h^{\text{UCB}}(s_h, a_h), Q_h(s_h, a_h) \right\}. \quad (3.4)$$

This is precisely what is conducted in UCB-Q-Advantage. However, while the estimate Q_h^R obtained with the aid of reference-advantage decomposition provides great promise, fully realizing its potential in the sample-limited regime relies on the ability to quickly *settle on* a desirable ‘reference’ during the initial stage of the algorithm. This leads us to a dilemma that requires careful thinking. On the one hand, the reference value V^R needs to be updated in a timely manner in order to better control the stochastic estimate of $P_{h,s,a}(V_{h+1} - V_{h+1}^R)$. On the other hand, updating V^R too frequently incurs an overly large sample size burden, as new samples need to be accumulated whenever V^R is updated.

Built upon the above insights, it is advisable to prevent frequent updating of the reference value V^R . In fact, it would be desirable to stop updating the reference value once a point of sufficient quality—denoted by V^R_{final} —has been obtained. Locking on a reasonable reference value early on means that (a) fewer samples will be wasted on estimating a drifting target $P_h V_{h+1}^R$, and (b) all ensuing samples can then be dedicated to estimating the key quantity of interest $P_h V^R_{\text{final},h+1}$.

REMARK 1. In [75], the algorithm UCB-Q-Advantage requires collecting $\tilde{O}(SAH^6)$ samples *for each state* before settling on the reference value, which inevitably contributes to the large burn-in cost.

Algorithm 1 Q-EarlySettled-Advantage

```

1: Parameters: some universal constant  $c_b > 0$  and probability of failure  $\delta \in (0, 1)$ ;
2: Initialize  $Q_h(s, a), Q_h^{\text{UCB}}(s, a), Q_h^R(s, a) \leftarrow H$ ;  $V_h(s), V_h^R(s) \leftarrow H$ ;  $Q_h^{\text{LCB}}(s, a) \leftarrow 0$ ;  $V_h^{\text{LCB}}(s) \leftarrow 0$ ;
    $N_h(s, a) \leftarrow 0$ ;  $\mu_h^{\text{ref}}(s, a), \sigma_h^{\text{ref}}(s, a), \mu_h^{\text{adv}}(s, a), \sigma_h^{\text{adv}}(s, a), \delta_h^R(s, a), B_h^R(s, a) \leftarrow 0$ ; and  $u_{\text{ref}}(s) = \text{True}$ 
   for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ .
3: for Episode  $k = 1$  to  $K$  do
4:   Set initial state  $s_1 \leftarrow s_1^k$ .
5:   for Step  $h = 1$  to  $H$  do
6:     Take action  $a_h = \pi_h^k(s_h) = \arg \max_a Q_h(s_h, a)$ , and draw  $s_{h+1} \sim P_h(\cdot | s_h, a_h)$ . // sampling
7:      $N_h(s_h, a_h) \leftarrow N_h(s_h, a_h) + 1$ ;  $n \leftarrow N_h(s_h, a_h)$ . // update the counter
8:      $\eta_n \leftarrow \frac{H+1}{H+n}$ . // update the learning rate
9:      $Q_h^{\text{UCB}}(s_h, a_h) \leftarrow \text{update-ucb-q}()$ . // run UCB-Q; see Algorithm 2
10:     $Q_h^{\text{LCB}}(s_h, a_h) \leftarrow \text{update-lcb-q}()$ . // run LCB-Q; see Algorithm 2
11:     $Q_h^R(s_h, a_h) \leftarrow \text{update-ucb-q-advantage}()$ . // estimate  $Q_h^R$ ; see Algorithm 2
        /* update Q-estimates using all estimates in hand, and update
        value estimates
        */
12:     $Q_h(s_h, a_h) \leftarrow \min \{Q_h^R(s_h, a_h), Q_h^{\text{UCB}}(s_h, a_h), Q_h(s_h, a_h)\}$ .
13:     $V_h(s_h) \leftarrow \max_a Q_h(s_h, a)$ .
14:     $V_h^{\text{LCB}}(s_h) \leftarrow \max \{\max_a Q_h^{\text{LCB}}(s_h, a), V_h^{\text{LCB}}(s_h)\}$ .
        /* update reference values
        */
15:    if  $V_h(s_h) - V_h^{\text{LCB}}(s_h) > 1$  then
16:       $V_h^R(s_h) \leftarrow V_h(s_h)$ ,  $u_{\text{ref}}(s_h) = \text{True}$ ,
17:    else if  $u_{\text{ref}}(s_h) = \text{True}$  then
18:       $V_h^R(s_h) \leftarrow V_h(s_h)$ ,  $u_{\text{ref}}(s_h) = \text{False}$ .
19:    end if
20:  end for
21: end for

```

The proposed Q-EarlySettled-Advantage algorithm. We now propose a new model-free algorithm that allows for early settlement of the reference value. A few key ingredients are as follows.

- *An auxiliary sequence based on LCB.* In addition to the two optimistic Q-estimates Q_h^R and Q_h^{UCB} described previously, we intend to maintain another *pessimistic* estimate $Q_h^{LCB} \leq Q_h^*$ using the subroutine `update-lcb-q`, based on LCBs. We will also maintain the corresponding value function V_h^{LCB} , which lower bounds V_h^* .
- *Termination rules for reference updates.* With $V_h^{LCB} \leq V_h^*$ in place, the updates of the references (lines 15-18 of Algorithm 1) are designed to terminate when

$$V_h(s_h) \leq V_h^{LCB}(s_h) + 1 \leq V_h^*(s_h) + 1. \quad (3.5)$$

Note that V_h^R keeps tracking the value of V_h before it stops being updated. In effect, when the additional condition in lines 15 is violated and thus (3.5) is satisfied, we claim that it is unnecessary to update the reference V_h^R afterwards, since it is of sufficient quality (being close enough to the optimal value V_h^*) and further drifting the reference does not appear beneficial. As we will make it rigorous shortly, this reference update rule is sufficient to ensure that $|V_h - V_h^R| \leq 2$ throughout the execution of the algorithm, which in turn suggests that the standard deviation of $\hat{P}_{h,s,a}(V_{h+1} - V_{h+1}^R)$ might be $O(H)$ times smaller than that of $\hat{P}_{h,s,a}V_{h+1}$ (i.e. the stochastic term used in (3.1) of UCB-Q). This is a key observation that helps shave the addition factor H in the regret bound of UCB-Q.

- *Update rules for Q_h^{UCB} and Q_h^R .* The two optimistic Q-estimates Q_h^{UCB} and Q_h^R are updated using the subroutine `update-ucb-q` (following the standard Q-learning with Hoeffding bonus [29]) and `update-ucb-q-advantage`, respectively. Note that Q_h^R continues to be updated even after V_h^R is no longer updated.

Q-learning with reference-advantage decomposition. The rest of this subsection is devoted to explaining the subroutine `update-ucb-q-advantage`, which produces a Q-estimate Q^R based on the reference-advantage decomposition similar to [75]. To facilitate the implementation, let us introduce the parameters associated with a reference value V^R , which include six different components, i.e.

$$[\mu_h^{\text{ref}}(s, a), \sigma_h^{\text{ref}}(s, a), \mu_h^{\text{adv}}(s, a), \sigma_h^{\text{adv}}(s, a), \delta_h^R(s, a), B_h^R(s, a)], \quad (3.6)$$

for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$. Here, $\mu_h^{\text{ref}}(s, a)$ and $\sigma_h^{\text{ref}}(s, a)$ estimate the running mean and the second moment of the reference $[P_h V_{h+1}^R](s, a)$; $\mu_h^{\text{adv}}(s, a)$ and $\sigma_h^{\text{adv}}(s, a)$ estimate the running (weighted) mean and the second moment of the advantage $[P_h(V_{h+1} - V_{h+1}^R)](s, a)$; $B_h^R(s, a)$ aggregates the empirical standard deviations of the reference and the advantage combined; and last but not least, $\delta_h^R(s, a)$ is the temporal difference between $B_h^R(s, a)$ and its previous value.

As alluded to previously, the Q-function estimation follows the strategy (3.3) at a high level. Upon observing a sample transition (s_h, a_h, s_{h+1}) , we compute the following estimates to update $Q^R(s_h, a_h)$.

- The term $\hat{P}_{h,s,a}(V_{h+1} - V_{h+1}^R)$ is set to be $V_{h+1}(s_{h+1}) - V_{h+1}^R(s_{h+1})$, which is an unbiased stochastic estimate of $P_{h,s,a}(V_{h+1} - V_{h+1}^R)$.

- The term $[P_h V_{h+1}^R](s, a)$ is estimated via μ^{ref}, R_h (cf. line 11). Given that this is estimated using all previous samples, we expect the variability of this term to be well-controlled as the sample size increases (especially after V^R is locked).
- The exploration bonus $b_h^R(s, a)$ is updated using $B_h^R(s_h, a_h)$ and $\delta_h^R(s_h, a_h)$ (cf. lines 7-8 of Algorithm 2), which is a confidence bound accounting for both the reference and the advantage. Let us also explain line 8 of Algorithm 2 a bit. If we augment the notation by letting $b_h^{R, n+1}(s, a)$ and $B_h^{R, n+1}(s, a)$ denote, respectively, $b_h^R(s, a)$ and $B_h^R(s, a)$ after (s, a) is visited for the n -th time, then this line is designed to ensure that

$$\eta_n b_h^{R, n+1}(s, a) + (1 - \eta_n) B_h^{R, n}(s, a) \approx B_h^{R, n+1}(s, a).$$

With the above updates implemented properly, Q_h^R provides the advantage-based update of the Q-function at time step h , according to the update rule (3.3).

Algorithm 2 Auxiliary functions

```

1: function update-ucb-q() :
2:    $Q_h^{\text{UCB}}(s_h, a_h) \leftarrow (1 - \eta_n) Q_h^{\text{UCB}}(s_h, a_h) + \eta_n \left( r_h(s_h, a_h) + V_{h+1}(s_{h+1}) + c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{n}} \right)$ .
3: function update-lcb-q() :
4:    $Q_h^{\text{LCB}}(s_h, a_h) \leftarrow (1 - \eta_n) Q_h^{\text{LCB}}(s_h, a_h) + \eta_n \left( r_h(s_h, a_h) + V_{h+1}^{\text{LCB}}(s_{h+1}) - c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{n}} \right)$ .
5: function update-ucb-q-advantage() :
   /* update the moment statistics of  $V_h^R$  */  

6:    $[\mu_h^{\text{ref}}, \sigma_h^{\text{ref}}, \mu_h^{\text{adv}}, \sigma_h^{\text{adv}}](s_h, a_h) \leftarrow \text{update-moments}()$ ;  

   /* update the accumulative bonus and bonus difference */  

7:    $[\delta_h^R, B_h^R](s_h, a_h) \leftarrow \text{update-bonus}()$ ;  

8:    $b_h^R \leftarrow B_h^R(s_h, a_h) + (1 - \eta_n) \frac{\delta_h^R(s_h, a_h)}{\eta_n} + c_b \frac{H^2 \log \frac{SAT}{\delta}}{n^{3/4}}$ ;  

   /* update the Q-estimate based on reference-advantage decomposition */  

9:    $Q_h^R(s_h, a_h) \leftarrow (1 - \eta_n) Q_h^R(s_h, a_h) + \eta_n (r_h(s_h, a_h) + V_{h+1}(s_{h+1}) - V_{h+1}^R(s_{h+1}) + \mu_h^{\text{ref}}(s_h, a_h) + b_h^R)$ ;  

  

10: function update-moments() :
11:    $\mu_h^{\text{ref}}(s_h, a_h) \leftarrow (1 - \frac{1}{n}) \mu_h^{\text{ref}}(s_h, a_h) + \frac{1}{n} V_{h+1}^R(s_{h+1})$ ; // mean of the reference  

12:    $\sigma_h^{\text{ref}}(s_h, a_h) \leftarrow (1 - \frac{1}{n}) \sigma_h^{\text{ref}}(s_h, a_h) + \frac{1}{n} (V_{h+1}^R(s_{h+1}))^2$ ; // 2nd moment of the reference  

13:    $\mu_h^{\text{adv}}(s_h, a_h) \leftarrow (1 - \eta_n) \mu_h^{\text{adv}}(s_h, a_h) + \eta_n (V_{h+1}(s_{h+1}) - V_{h+1}^R(s_{h+1}))$ ; // weighted average  

   of the advantage  

14:    $\sigma_h^{\text{adv}}(s_h, a_h) \leftarrow (1 - \eta_n) \sigma_h^{\text{adv}}(s_h, a_h) + \eta_n (V_{h+1}(s_{h+1}) - V_{h+1}^R(s_{h+1}))^2$ . // weighted 2nd  

   moment of the advantage  

15: function update-bonus() :
16:    $B_h^{\text{next}}(s_h, a_h) \leftarrow$   

17:    $c_b \sqrt{\frac{\log \frac{SAT}{\delta}}{n}} \left( \sqrt{\sigma_h^{\text{ref}}(s_h, a_h) - (\mu_h^{\text{ref}}(s_h, a_h))^2} + \sqrt{H} \sqrt{\sigma_h^{\text{adv}}(s_h, a_h) - (\mu_h^{\text{adv}}(s_h, a_h))^2} \right)$ ;  

18:    $\delta_h^R(s_h, a_h) \leftarrow B_h^{\text{next}}(s_h, a_h) - B_h^R(s_h, a_h)$ ;  

19:    $B_h^R(s_h, a_h) \leftarrow B_h^{\text{next}}(s_h, a_h)$ .

```

3.3 Main results

Encouragingly, the proposed **Q-EarlySettled-Advantage** algorithm manages to achieve near-optimal regret even in the sample-limited and memory-limited regime, as formalized by the following theorem.

THEOREM 3.1. Consider any $\delta \in (0, 1)$, and suppose that $c_b > 0$ is chosen to be a sufficiently large universal constant. Then there exists some absolute constant $C_0 > 0$ such that Algorithm 1 achieves

$$\text{Regret}(T) \leq C_0 \left(\sqrt{H^2 SAT \log^4 \frac{SAT}{\delta}} + H^6 SA \log^3 \frac{SAT}{\delta} \right) \quad (3.7)$$

with probability at least $1 - \delta$.

Theorem 3.1 delivers a non-asymptotic characterization of the performance of our algorithm **Q-EarlySettled-Advantage**. Several appealing features of the algorithm are noteworthy.

- *Regret optimality.* Our regret bound (3.7) simplifies to

$$\text{Regret}(T) \leq \tilde{O}(\sqrt{H^2 SAT}) \quad (3.8)$$

as long as the sample size T exceeds

$$T \geq SA \text{poly}(H). \quad (3.9)$$

This sublinear regret bound (3.8) is essentially optimal, as it coincides with the existing lower bound (1.1) modulo some logarithmic factor.

- *Sample complexity and substantially reduced burn-in cost.* As an interpretation of our theory (3.8), our algorithm attains ε average regret (i.e. $\frac{1}{K} \text{Regret}(T) \leq \varepsilon$) with a sample complexity

$$\tilde{O}\left(\frac{SAH^4}{\varepsilon^2}\right).$$

Crucially, the burn-in cost (3.9) is significantly lower than that of the state-of-the-art memory-efficient model-free algorithm [75] (whose optimality is guaranteed only in the range $T \geq S^6 A^4 \text{poly}(H)$).

- *Memory efficiency.* Our algorithm, which is model-free in nature, achieves a low space complexity $O(SAH)$. This is basically un-improvable for the tabular case, since even storing the optimal Q-values alone takes $O(SAH)$ units of space. In comparison, while [44] also accommodates the sample size range (3.9), the algorithm proposed therein incurs a space complexity of $O(S^2 AH)$ that is S times higher than ours.
- *Computational complexity.* An additional intriguing feature of our algorithm is its low computational complexity. The runtime of **Q-EarlySettled-Advantage** is no larger than $O(T)$, which is proportional to the time taken to read the samples. This matches the computational cost of the model-free algorithm UCB-Q proposed in [29], and is considerably lower than that of the UCB-M-Q algorithm in [44] (which has a computational cost of at least $O(ST)$).

4. Analysis

In this section, we outline the main steps needed to prove our main result in Theorem 3.1.

4.1 Preliminaries: basic properties about learning rates

Before continuing, let us first state some basic facts regarding the learning rates. Akin to [29], the proposed algorithm adopts the linearly rescaled learning rate

$$\eta_n = \frac{H+1}{H+n} \quad (4.1)$$

for the n -th visit of a state–action pair at any time step h . For notation convenience, we further introduce two sequences of related quantities defined for any integer $N \geq 0$ and $n \geq 1$:

$$\eta_n^N := \begin{cases} \eta_n \prod_{i=n+1}^N (1 - \eta_i), & \text{if } N > n, \\ \eta_n, & \text{if } N = n, \\ 0, & \text{if } N < n \end{cases} \quad \text{and} \quad \eta_0^N := \begin{cases} \prod_{i=1}^N (1 - \eta_i) = 0, & \text{if } N > 0, \\ 1, & \text{if } N = 0. \end{cases} \quad (4.2)$$

As can be easily verified, we have

$$\sum_{n=1}^N \eta_n^N = \begin{cases} 1, & \text{if } N > 0, \\ 0, & \text{if } N = 0. \end{cases} \quad (4.3)$$

The following properties play an important role in the analysis.

LEMMA 1. For any integer $N > 0$, the following properties hold:

$$\frac{1}{N^a} \leq \sum_{n=1}^N \frac{\eta_n^N}{n^a} \leq \frac{2}{N^a}, \quad \text{for all } \frac{1}{2} \leq a \leq 1, \quad (4.4a)$$

$$\max_{1 \leq n \leq N} \eta_n^N \leq \frac{2H}{N}, \quad \sum_{n=1}^N (\eta_n^N)^2 \leq \frac{2H}{N}, \quad \sum_{N=n}^{\infty} \eta_n^N \leq 1 + \frac{1}{H}. \quad (4.4b)$$

Proof. See Appendix B. □

4.2 Additional notation used in the proof

In order to enable a more concise description of the algorithm, we have suppressed the dependency of many quantities on the episode number k in Algorithms 1 and 2. This, however, becomes notationally inconvenient when presenting the proof. As a consequence, we shall adopt, throughout the analysis, a more complete set of notation, detailed below.

- (s_h^k, a_h^k) : the state–action pair encountered and chosen at time step h in the k -th episode.
- $k_h^n(s, a)$: the index of the episode in which (s, a) is visited for the n -th time at time step h ; for the sake of conciseness, we shall sometimes use the shorthand $k^n = k_h^n(s, a)$ whenever it is clear from the context.
- $k_h^n(s)$: the index of the episode in which state s is visited for the n -th time at time step h ; we might sometimes abuse the notation by abbreviating $k^n = k_h^n(s)$.

- $P_h^k \in \{0, 1\}^{1 \times |\mathcal{A}|}$: the empirical transition at time step h in the k -th episode, namely,

$$P_h^k(s) = \mathbb{1}(s = s_{h+1}^k). \quad (4.5)$$

In addition, for several parameters of interest in Algorithm 1, we introduce the following set of augmented notation.

- $N_h^k(s, a)$ denotes $N_h(s, a)$ by the end of the k -th episode; for the sake of conciseness, we shall often abbreviate $N^k = N_h^k(s, a)$ or $N^k = N_h^k(s_h^k, a_h^k)$ (depending on which result we are proving).
- $Q_h^k(s, a)$, $V_h^k(s)$ and $Q_h^{\text{UCB}, k}(s, a)$ denote, respectively, $Q_h(s, a)$, $V_h(s)$ and $Q_h^{\text{UCB}}(s, a)$ at the beginning of the k -th episode.
- $Q_h^{\text{LCB}, k}(s, a)$ and $V_h^{\text{LCB}, k}(s)$ denote, respectively, $Q_h^{\text{LCB}}(s, a)$ and $V_h^{\text{LCB}}(s)$ at the beginning of the k -th episode.
- $Q_h^R, k(s, a)$, $V_h^R, k(s)$ and $u_{\text{ref}}^k(s)$ denote, respectively, $Q_h^R(s, a)$, $V_h^R(s)$ and $u_{\text{ref}}(s)$ at the beginning of the k -th episode.
- $[\mu^{\text{ref}}, k_h, \sigma^{\text{ref}}, k_h, \mu^{\text{adv}}, k_h, \sigma^{\text{adv}}, k_h, \delta^R, k_h, B^R, k_h]$ denotes $[\mu_h^{\text{ref}}, \sigma_h^{\text{ref}}, \mu_h^{\text{adv}}, \sigma_h^{\text{adv}}, \delta_h^R, B_h^R]$ at the beginning of the k -th episode.

Further, for any matrix $P = [P_{i,j}]_{1 \leq i \leq m, 1 \leq j \leq n}$, we define $\|P\|_1 := \max_{1 \leq i \leq m} \sum_{j=1}^n |P_{i,j}|$. For any vector $V = [V_i]_{1 \leq i \leq n}$, we define its ℓ_∞ norm as $\|V\|_\infty := \max_{1 \leq i \leq n} |V_i|$. We often overload scalar functions and expressions to take vector-valued arguments, with the understanding that they are applied in an entrywise manner. For example, for a vector $x = [x_i]_{1 \leq i \leq n}$, we denote $x^2 = [x_i^2]_{1 \leq i \leq n}$. For any two vectors $x = [x_i]_{1 \leq i \leq n}$ and $y = [y_i]_{1 \leq i \leq n}$, the notation $x \leq y$ (resp. $x \geq y$) means $x_i \leq y_i$ (resp. $x_i \geq y_i$) for all $1 \leq i \leq n$. For any given vector $V \in \mathbb{R}^S$, we define the variance parameter w.r.t. $P_{h,s,a}$ (cf. (2.7)) as follows:

$$\text{Var}_{h,s,a}(V) := \mathbb{E}_{s' \sim P_{h,s,a}} \left[(V(s') - P_{h,s,a}V)^2 \right] = P_{h,s,a}(V^2) - (P_{h,s,a}V)^2. \quad (4.6)$$

Finally, let $\mathcal{X} := (S, A, H, T, \frac{1}{\delta})$. The notation $f(\mathcal{X}) \lesssim g(\mathcal{X})$ (resp. $f(\mathcal{X}) \gtrsim g(\mathcal{X})$) means that there exists a universal constant $C_0 > 0$ such that $f(\mathcal{X}) \leq C_0 g(\mathcal{X})$ (resp. $f(\mathcal{X}) \geq C_0 g(\mathcal{X})$); the notation $f(\mathcal{X}) \asymp g(\mathcal{X})$ means that $f(\mathcal{X}) \lesssim g(\mathcal{X})$ and $f(\mathcal{X}) \gtrsim g(\mathcal{X})$ hold simultaneously.

4.3 Key properties of Q -estimates and auxiliary sequences

In this subsection, we introduce several key properties of our Q -estimates and value estimates, which play a crucial role in the proof of Theorem 3.1. The proofs for this subsection are deferred to Appendix C.s

Properties of the Q -estimate Q_h^k : monotonicity and optimism. We first make an important observation regarding the monotonicity of the value estimates Q_h^k and V_h^k . To begin with, it is straightforward to see that the update rule in Algorithm 3 (cf. line 12) ensures the following monotonicity property:

$$Q_h^{k+1}(s, a) \leq Q_h^k(s, a) \quad \text{for all } (s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H], \quad (4.7a)$$

Algorithm 3 Q-EarlySettled-Advantage (rewrite of Algorithm 1)

1: **Parameters:** some universal constant $c_b > 0$ and probability of failure $\delta \in (0, 1)$;

2: **Initialize** $Q_h^1(s, a), Q_h^{\text{UCB}, 1}(s, a), Q_h^{\text{R}, 1}(s, a) \leftarrow H; Q_h^{\text{LCB}, 1}(s, a) \leftarrow 0; N_h^0(s, a) \leftarrow 0; V_h^1(s), V_h^{\text{R}, 1}(s) \leftarrow H; \mu_h^{\text{ref}}(s, a), \sigma_h^{\text{ref}}(s, a), \mu_h^{\text{adv}}(s, a), \sigma_h^{\text{adv}}(s, a), \delta_h^{\text{R}}(s, a), B_h^{\text{R}}(s, a) \leftarrow 0$; and $u_{\text{ref}}^1(s) = \text{True}$, for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$.

3: **for** Episode $k = 1$ **to** K **do**

4: Set initial state $s_1 \leftarrow s_1^k$.

5: **for** Step $h = 1$ **to** H **do**

6: Take action $a_h^k = \pi_h^k(s_h) = \arg \max_a Q_h^k(s_h^k, a)$, and draw $s_{h+1}^k \sim P_h(\cdot | s_h^k, a_h^k)$. // sampling

7: $N_h^k(s_h^k, a_h^k) \leftarrow N_h^{k-1}(s_h^k, a_h^k) + 1; n \leftarrow N_h^k(s_h^k, a_h^k)$. // update the counter

8: $\eta_n \leftarrow \frac{H+1}{H+n}$. // update the learning rate

9: $Q_h^{\text{UCB}, k+1}(s_h^k, a_h^k) \leftarrow \text{update-ucb-q}()$. // run UCB-Q;

10: $Q_h^{\text{LCB}, k+1}(s_h^k, a_h^k) \leftarrow \text{update-lcb-q}()$. // run LCB-Q;

11: $Q_h^{\text{R}, k+1}(s_h^k, a_h^k) \leftarrow \text{update-ucb-q-advantage}()$. // estimate Q_h^{R} ;

12: /* update Q-estimates using all estimates in hand, and update value estimates */

13: $Q_h^{k+1}(s_h^k, a_h^k) \leftarrow \min \{Q_h^{\text{R}, k+1}(s_h^k, a_h^k), Q_h^{\text{UCB}, k+1}(s_h^k, a_h^k), Q_h^{\text{LCB}, k+1}(s_h^k, a_h^k)\}$;

14: $V_h^{k+1}(s_h^k) \leftarrow \max_a Q_h^{k+1}(s_h^k, a)$.

15: $V_h^{\text{LCB}, k+1}(s_h^k) \leftarrow \max \left\{ \max_a Q_h^{\text{LCB}, k+1}(s_h^k, a), V_h^{\text{R}, k+1}(s_h^k) \right\}$.

16: /* update reference values */

17: **if** $V_h^{k+1}(s_h^k) - V_h^{\text{LCB}, k+1}(s_h^k) > 1$ **then**

18: $V_h^{\text{R}, k+1}(s_h^k) \leftarrow V_h^{k+1}(s_h^k), u_{\text{ref}}^{k+1}(s_h^k) = \text{True}$;

19: **else if** $u_{\text{ref}}^k(s_h^k) = \text{True}$ **then**

20: $V_h^{\text{R}, k+1}(s_h^k) \leftarrow V_h^{k+1}(s_h^k), u_{\text{ref}}^{k+1}(s_h^k) = \text{False}$.

21: **end if**

22: **end for**

23: **end for**

which combined with line 13 of Algorithm 3 leads to monotonicity of $V_h(s)$ as follows:

$$V_h^{k+1}(s) = Q_h^{k+1}(s, \pi_h^{k+1}(s)) \leq Q_h^k(s, \pi_h^{k+1}(s)) \leq V_h^k(s). \quad (4.7b)$$

Moreover, by virtue of the update rule in line 12 of Algorithm 3, we can immediately obtain (via induction) the following useful property:

$$Q_h^{\text{R}, k}(s, a) \geq Q_h^k(s, a) \quad \text{for all } (k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}. \quad (4.8)$$

In addition, Q_h^k and V_h^k form an ‘optimistic view’ of Q_h^* and V_h^* , respectively, as asserted by the following lemma.

LEMMA 2. Consider any $\delta \in (0, 1)$. Suppose that $c_b > 0$ is some sufficiently large constant. Then with probability at least $1 - \delta$,

$$Q_h^k(s, a) \geq Q_h^*(s, a) \quad \text{and} \quad V_h^k(s) \geq V_h^*(s) \quad (4.9)$$

hold simultaneously for all $(s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]$.

Lemma 2 implies that Q_h^k (resp. V_h^k) is a pointwise upper bound on Q_h^* (resp. V_h^*). Taking this result together with the non-increasing property (4.7), we see that Q_h^k (resp. V_h^k) becomes an increasingly tighter estimate of Q_h^* (resp. V_h^*) as the number of episodes k increases. This important fact forms the basis of the subsequent proof, allowing us to replace V_h^* with V_h^k when upper bounding the regret. Combining Lemma 2 with (4.8), we can straightforwardly see that with probability at least $1 - \delta$:

$$Q_h^R, k(s, a) \geq Q_h^*(s, a) \quad \text{for all } (k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}. \quad (4.10)$$

Properties of the Q-estimate $Q_h^{\text{LCB}, k}$: pessimism and proximity. In parallel, we formalize the fact that $Q_h^{\text{LCB}, k}$ and $V_h^{\text{LCB}, k}$ provide a ‘pessimistic view’ of Q_h^* and V_h^* , respectively. Furthermore, it becomes increasingly more likely for $Q_h^{\text{LCB}, k}$ and Q_h^k to stay close to each other as k increases, which indicates that the confidence interval that contains the optimal value Q_h^* becomes shorter and shorter. These properties are summarized in the following lemma.

LEMMA 3. Consider any $\delta \in (0, 1)$, and suppose that $c_b > 0$ is some sufficiently large constant. Then with probability at least $1 - \delta$,

$$Q_h^{\text{LCB}, k}(s, a) \leq Q_h^*(s, a) \quad \text{and} \quad V_h^{\text{LCB}, k}(s) \leq V_h^*(s) \quad (4.11)$$

hold for all $(s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]$, and

$$\sum_{h=1}^H \sum_{k=1}^K \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB}, k}(s_h^k, a_h^k) > \varepsilon \right) \lesssim \frac{H^6 S A \log \frac{SAT}{\delta}}{\varepsilon^2} \quad (4.12)$$

holds for all $\varepsilon \in (0, H]$.

Interestingly, the upper bound (4.12) only scales logarithmically in the number K of episodes, thus implying the closeness of $Q_h^{\text{LCB}, k}$ and Q_h^k for a large fraction of episodes. Note that it is straightforward to ensure the monotonicity property of $V_h^{\text{LCB}, k}$ from the update rule in Algorithm 3 (cf. line 14):

$$V_h^{\text{LCB}, k+1}(s) \geq V_h^{\text{LCB}, k}(s) \quad \text{for all } (s, k, h) \in \mathcal{S} \times [K] \times [H], \quad (4.13)$$

which in conjunction with (4.11), implies that $V_h^{\text{LCB}, k}(s)$ gets closer to $V_h^*(s)$ as the number of episodes k increases. Together with the monotonicity of V_h^k (cf. (4.7b)), an important consequence is that the reference value V_h^R will stop being updated shortly after the following condition is met for the first time

(according to lines 15–18 of Algorithm 1)

$$V_h^k(s) \leq V_h^{\text{LCB},k}(s) + 1 \leq V_h^*(s) + 1 \quad \text{for all } s \in \mathcal{S}. \quad (4.14)$$

Properties of the reference V_h^R, k . The above fact ensures that V_h^R, k will not be updated too many times. In fact, its value stays reasonably close to V_h^k even after being locked to a fixed value, which ensures its fidelity as a reference signal. Moreover, the aggregate difference between V^R, k_h and the final reference V^R, k_h over the entire trajectory can be bounded in a reasonably tight fashion (owing to (4.12)), as formalized in the next lemma. These properties play a key role in reducing the burn-in cost of the proposed algorithm.

LEMMA 4. Consider any $\delta \in (0, 1)$. Suppose that $c_b > 0$ is some sufficiently large constant. Then with probability exceeding $1 - \delta$, one has

$$|V_h^k(s) - V^R, k_h(s)| \leq 2 \quad (4.15)$$

for all $(k, h, s) \in [K] \times [H] \times \mathcal{S}$, and

$$\begin{aligned} & \sum_{h=1}^H \sum_{k=1}^K \left(V^R, k_h(s_h^k) - V^R, k_h(s_h^k) \right) \\ & \leq H^2 S + \sum_{h=1}^H \sum_{k=1}^K \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) \right) \mathbb{1}(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > 1) \quad (4.16) \end{aligned}$$

$$\lesssim H^6 S A \log \frac{SAT}{\delta}. \quad (4.17)$$

In words, Lemma 4 guarantees that (i) our value function estimate and the reference value are always sufficiently close (cf. (4.15)), and (ii) the aggregate difference between V_h^R, k and the final reference value V^R, k_h is nearly independent of the sample size T (except for some logarithmic scaling).

4.4 Main steps of the proof

We are now ready to embark on the regret analysis for Q-EarlySettled-Advantage, which consists of multiple steps as follows.

Step 1: regret decomposition. Lemma 2 allows one to upper bound the regret as follows:

$$\text{Regret}(T) := \sum_{k=1}^K (V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k)) \leq \sum_{k=1}^K (V_1^k(s_1^k) - V_1^{\pi^k}(s_1^k)). \quad (4.18)$$

To continue, it boils down to controlling $V_1^k(s_1^k) - V_1^{\pi^k}(s_1^k)$. Toward this end, we intend to examine $V_h^k(s_h^k) - V_h^{\pi^k}(s_h^k)$ across all time steps $1 \leq h \leq H$, which admits the following decomposition:

$$\begin{aligned}
V_h^k(s_h^k) - V_h^{\pi^k}(s_h^k) &\stackrel{(i)}{=} Q_h^k(s_h^k, a_h^k) - Q_h^{\pi^k}(s_h^k, a_h^k) \\
&= Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) + Q_h^*(s_h^k, a_h^k) - Q_h^{\pi^k}(s_h^k, a_h^k) \\
&\stackrel{(ii)}{=} Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) + P_{h, s_h^k, a_h^k}(V_{h+1}^* - V_{h+1}^{\pi^k}) \\
&\stackrel{(iii)}{=} Q_h^k(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) + (P_{h, s_h^k, a_h^k} - P_h^k)(V_{h+1}^* - V_{h+1}^{\pi^k}) + V_{h+1}^*(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k) \\
&\leq Q^R, k_h(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) + (P_{h, s_h^k, a_h^k} - P_h^k)(V_{h+1}^* - V_{h+1}^{\pi^k}) + V_{h+1}^*(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k).
\end{aligned} \tag{4.19}$$

Here, (i) holds since π_h^k is a greedy policy w.r.t. Q_h^k and $\pi_h^k(s_h^k) = a_h^k$, (ii) comes from the Bellman equations

$$Q_h^{\pi^k}(s, a) - Q_h^*(s, a) = (r_h(s, a) + P_{h, s, a} V_{h+1}^{\pi^k}) - (r_h(s, a) + P_{h, s, a} V_{h+1}^*) = P_{h, s, a}(V_{h+1}^{\pi^k} - V_{h+1}^*),$$

(iii) follows from $P_h^k(V_{h+1}^* - V_{h+1}^{\pi^k}) = V_{h+1}^*(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k)$ (see the notation (4.5)), whereas the last inequality comes from (4.8). Summing (4.19) over $1 \leq k \leq K$ and using Lemma 2, we obtain

$$\begin{aligned}
\sum_{k=1}^K (V_h^*(s_h^k) - V_h^{\pi^k}(s_h^k)) &\leq \sum_{k=1}^K (V_h^k(s_h^k) - V_h^{\pi^k}(s_h^k)) \\
&\leq \sum_{k=1}^K (Q^R, k_h(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) + \sum_{k=1}^K (P_{h, s_h^k, a_h^k} - P_h^k)(V_{h+1}^* - V_{h+1}^{\pi^k}) \\
&\quad + \sum_{k=1}^K (V_{h+1}^*(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k)).
\end{aligned} \tag{4.20}$$

This allows us to establish a connection between $\sum_k (V_h^*(s_h^k) - V_h^{\pi^k}(s_h^k))$ for step h and $\sum_k (V_{h+1}^*(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k))$ for step $h+1$.

Step 2: managing regret by recursion. The regret can be further manipulated by leveraging the update rule of Q^R, k_h as well as recursing over the time steps $h = 1, 2, \dots, H$ with the terminal condition $V_{H+1}^k = V_{H+1}^{\pi^k} = 0$. This leads to a key decomposition as summarized in the lemma below, whose proof is provided in Appendix D.

LEMMA 5. Fix $\delta \in (0, 1)$. Suppose that $c_b > 0$ is some sufficiently large constant. Then with probability at least $1 - \delta$, one has

$$\sum_{k=1}^K (V_1^k(s_1^k) - V_1^{\pi^k}(s_1^k)) \leq \mathcal{R}_1 + \mathcal{R}_2 + \mathcal{R}_3, \quad (4.21)$$

where

$$\mathcal{R}_1 := \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^{h-1} \left(HSA + 8c_b H^2 (SA)^{3/4} K^{1/4} \log \frac{SAT}{\delta} + \sum_{k=1}^K (P_{h,s_h^k, a_h^k} - P_h^k) (V_{h+1}^* - V_{h+1}^{\pi^k}) \right), \quad (4.22a)$$

$$\mathcal{R}_2 := \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^{h-1} \sum_{k=1}^K B^R, k_h(s_h^k, a_h^k), \quad (4.22b)$$

$$\begin{aligned} \mathcal{R}_3 := & \sum_{h=1}^H \sum_{k=1}^K \lambda_h^k \left((P_h^k - P_{h,s_h^k, a_h^k}) (V_{h+1}^* - V^R, k_{h+1}) \right. \\ & \left. + \frac{\sum_{i=1}^{N_h^k(s_h^k, a_h^k)} (V_{h+1}^{R, k_h^i(s_h^k, a_h^k)} (s_{h+1}^{k_i(s_h^k, a_h^k)}) - P_{h,s_h^k, a_h^k} V^R, k_{h+1})}{N_h^k(s_h^k, a_h^k)} \right), \end{aligned} \quad (4.22c)$$

with

$$\lambda_h^k := \left(1 + \frac{1}{H}\right)^{h-1} \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \eta_{N_h^k(s_h^k, a_h^k)}^n.$$

This lemma attempts to upper bound the target quantity $\sum_{k=1}^K (V_1^k(s_1^k) - V_1^{\pi^k}(s_1^k))$ via three terms (see (4.21)). Informally, these terms reflect (i) the influence of the initialization as well as the finite-sample uncertainty of $P_h^k (V_{h+1}^* - V_{h+1}^{\pi^k})$, (ii) the influence of the size of the bonus terms and (iii) the discrepancy term when the running value iterates are replaced by the reference values. As we shall see in the analysis, the key to obtaining these terms lies in properly expanding the component $\sum_{k=1}^K (Q^R, k_h(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k))$ in (4.20), as well as applying induction across all $h = 1, \dots, H$.

Step 3: controlling the terms in (4.22) separately. As it turns out, each of the terms in (4.22) can be well controlled. We provide the bounds for these terms in the following lemma.

LEMMA 6. Consider any $\delta \in (0, 1)$. With probability at least $1 - \delta$, we have the following upper bounds:

$$\begin{aligned}\mathcal{R}_1 &\leq C_r \left\{ \sqrt{H^2 SAT \log \frac{SAT}{\delta}} + H^{4.5} SA \log^2 \frac{SAT}{\delta} \right\}, \\ \mathcal{R}_2 &\leq C_r \left\{ \sqrt{H^2 SAT \log \frac{SAT}{\delta}} + H^4 SA \log^2 \frac{SAT}{\delta} \right\}, \\ \mathcal{R}_3 &\leq C_r \left\{ \sqrt{H^2 SAT \log^4 \frac{SAT}{\delta}} + H^6 SA \log^3 \frac{SAT}{\delta} \right\}\end{aligned}$$

for some universal constant $C_r > 0$.

In order to derive the above bounds, the main strategy is to apply the Bernstein-type concentration inequalities carefully, and to upper bound the sum of variance in a careful manner. The proofs are deferred to Appendix E.

Step 4: putting all this together. We now have everything in place to establish our main result. Taking the preceding bounds in Lemma 6 together with (4.22), we see that with probability exceeding $1 - \delta$, one has

$$\text{Regret}(T) \leq \mathcal{R}_1 + \mathcal{R}_2 + \mathcal{R}_3 \lesssim \sqrt{H^2 SAT \log^4 \frac{SAT}{\delta}} + H^6 SA \log^3 \frac{SAT}{\delta}$$

as claimed.

5. Discussion

In this paper, we have proposed a novel model-free RL algorithm—tailored to online episodic settings—that attains near-optimal regret $\tilde{O}(\sqrt{H^2 SAT})$ and near-minimal memory complexity $O(SAH)$ at once. Remarkably, the near-optimality of the algorithm comes into effect as soon as the sample size rises above $O(SA \text{poly}(H))$, which has significantly improved upon the sample size requirements (or burn-in cost) for any prior regret-optimal model-free algorithm (based on the definition of the model-free algorithm in [29]). Given that online data collection could be expensive, time-consuming or high-stakes in a variety of contemporary applications (e.g. clinical trials, autonomous driving, online advertisement), reducing burn-in sample sizes compromising sample optimality is crucial in enabling sample-efficient solutions in these sample-constrained applications.

The results in this paper naturally suggest a number of possible extensions and directions for future investigation. We close the paper by listing a few of them.

- While the proposed algorithm provably enables minimal burn-in cost in terms of the dependency on S and A , our current theory falls short of delivering optimal horizon dependency of the burn-in cost. More specifically, even though our burn-in cost improves upon the state-of-the-art theory for sample-optimal model-free algorithms by a factor of at least $S^5 A^3 H^{18}$ (see [75]), the way we cope with the dependency on H remains inadequate. This calls for more refined analysis tools to optimize the horizon dependency.
- This paper focuses primarily on MDPs with non-stationary probability transition kernels. Another important scenario is concerned with MDPs with stationary transition kernels (i.e. the case where

P_h is identical across different h). It is worth noting that the algorithm developed herein is incapable of attaining optimal regret for the stationary case (i.e. the resulting regret might be off by a factor of \sqrt{H}). While our analysis already contains multiple key ingredients that are useful for analyzing the stationary case, how to complete the picture is non-trivial, which we leave for future work.

- Admittedly, even though we are now able to settle the sample size dependency on the state-action space, the size of SA might remain prohibitively large in many modern RL applications. As a result, parsimonious function representation/approximation of the underlying MDP is needed in order to further reduce the sample complexity. Prominent examples of this kind include linearly parameterized or realizable MDPs [18, 31, 38]. We hope that the method and analysis framework developed herein might inspire further development of sample-efficient algorithms that can effectively accommodate low-dimensional function approximation.

5. Data availability

No new data were generated or analyzed in support of this research.

Acknowledgment

L. Shi and Y. Chi are supported in part by the grants Office of Naval Research N00014-19-1-2404, National Science Foundation CCF-2106778, CCF-2007911 and DMS-2134080. L. Shi is also gratefully supported by the Leo Finzi Memorial Fellowship, Wei Shen and Xuehong Zhang Presidential Fellowship, and Liang Ji-Dian Graduate Fellowship at Carnegie Mellon University. Y. Chen is supported in part by the the Alfred P. Sloan Research Fellowship, the Air Force Office of Scientific Research grant FA9550-22-1-0198, the ONR grant N00014-22-1-2354, and the NSF grants CCF-2221009, CCF-1907661, IIS-2218713 and IIS-2218773, and the Google Research Scholar Award. Part of this work was done while Y. Chen and G. Li were visiting the Simons Institute for the Theory of Computing.

REFERENCES

1. AGARWAL, A., KAKADE, S. & YANG, L. F. (2020) Model-based reinforcement learning with a generative model is minimax optimal. *Conference on Learning Theory*, **125**, 67–83.
2. AGRAWAL, S. & JIA, R. (2017) Posterior sampling for reinforcement learning: worst-case regret bounds. arXiv preprint arXiv:1705.07041.
3. AUER, P., CESA-BIANCHI, N., FREUND, Y. & SCHAPIRE, R. E. (2002) The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, **32**, 48–77.
4. AUER, P. & ORTNER, R. (2010) UCB revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Period. Math. Hungar.*, **61**, 55–65.
5. AZAR, M. G., KAPPEN, H. J., GHAVAMZADEH, M. & MUNOS, R. (2011) Speedy Q-learning. *Advances in neural information processing systems*, pp. 2411–2419.
6. AZAR, M. G., MUNOS, R. & KAPPEN, H. J. (2013) Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, **91**, 325–349.
7. AZAR, M. G., OSBAND, I. & MUNOS, R. (2017) Minimax regret bounds for reinforcement learning. *Proceedings of the 34th International Conference on Machine Learning*, vol. **70**. JMLR. org, pp. 263–272.
8. BAI, Y., XIE, T., JIANG, N. & WANG, Y.-X. (2019) Provably efficient q -learning with low switching cost. *Advances in Neural Information Processing Systems*, pp. 8002–8011.
9. BARTLETT, P. & TEWARI, A. (2009) Regal: a regularization based algorithm for reinforcement learning in weakly communicating MDPs. *Uncertainty in Artificial Intelligence: Proceedings of the 25th Conference*. AUAI Press, pp. 35–42.

10. BECK, C. L. & SRIKANT, R. (2012) Error bounds for constant step-size Q-learning. *Systems Control Lett.*, **61**, 1203–1208.
11. BERTSEKAS, D. P. (2017) *Dynamic programming and optimal control (4th edition)*. Athena Scientific.
12. CHEN, Z., MAGULURI, S. T., SHAKKOTTAI, S. & SHANMUGAM, K. (2020) Finite-sample analysis of stochastic approximation using smooth convex envelopes. *arXiv preprint arXiv:2002.00874*.
13. CHEN, Z., MAGULURI, S. T., SHAKKOTTAI, S. & SHANMUGAM, K. (2021) A Lyapunov theory for finite-sample guarantees of asynchronous Q-learning and TD-learning variants. *arXiv preprint arXiv:2102.01567*.
14. DANN, C., LATTIMORE, T. & BRUNSKILL, E. (2017) Unifying PAC and regret: Uniform PAC bounds for episodic reinforcement learning. *arXiv preprint arXiv:1703.07710*.
15. DOMINGUES, O. D., MÉNARD, P., KAUFMANN, E. & VALKO, M. (2021) Episodic reinforcement learning in finite MDPs: Minimax lower bounds revisited. *Algorithmic Learning Theory*. PMLR, pp. 578–598.
16. DONG, K., WANG, Y., CHEN, X. & WANG, L. (2019) Q-learning with UCB exploration is sample efficient for infinite-horizon MDP. *arXiv preprint arXiv:1901.09311*.
17. DU, S. S., CHEN, J., LI, L., XIAO, L. & ZHOU, D. (s2017) Stochastic variance reduction methods for policy evaluation. *Proceedings of the 34th International Conference on Machine Learning*, vol. **70**. JMLR. org, pp. 1049–1058.
18. DU, S. S., KAKADE, S. M., WANG, R. & YANG, L. F. (2020) Is a good representation sufficient for sample efficient reinforcement learning? *International Conference on Learning Representations*.
19. DU, S. S., LUO, Y., WANG, R. & ZHANG, H. (2019) Provably efficient Q-learning with function approximation via distribution shift error checking oracle. *Advances in Neural Information Processing Systems*, pp. 8058–8068.
20. EFRONI, Y., MERLIS, N., GHAVAMZADEH, M. & MANNER, S. (2019) *Tight regret bounds for model-based reinforcement learning with greedy policies*. *arXiv preprint arXiv:1905.11527*.
21. EVEN-DAR, E. & MANSOUR, Y. (2003) Learning rates for Q-learning. *Journal of machine learning Research*, **5**, 1–25.
22. FAN, J., WANG, Z., XIE, Y. & YANG, Z. (2020) A theoretical analysis of deep Q-learning. *Learning for Dynamics and Control*. PMLR, 486–489.
23. FREEDMAN, D. A. (1975) On tail probabilities for martingales. *the Annals of Probability*, pp. 100–118.
24. GOWER, R. M., SCHMIDT, M., BACH, F. & RICHTÁRIK, P. (2020) Variance-reduced methods for machine learning. *Proceedings of the IEEE*, **108**, 1968–1983.
25. HE, J., ZHOU, D. & GU, Q. (2020) *Nearly minimax optimal reinforcement learning for discounted MDPs*. *arXiv preprint arXiv:2010.00587*.
26. JAAKKOLA, T., JORDAN, M. I. & SINGH, S. P. (1994) Convergence of stochastic iterative dynamic programming algorithms. *Advances in neural information processing systems*, pp. 703–710.
27. JAFARNIA-JAHROMI, M., WEI, C.-Y., JAIN, R. & LUO, H. (2020) *A model-free learning algorithm for infinite-horizon average-reward MDPs with near-optimal regret*. *arXiv preprint arXiv:2006.04354*.
28. JAKSCH, T., ORTNER, R. & AUER, P. (2010) Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, **11**, 1563–1600.
29. JIN, C., ALLEN-ZHU, Z., BUBECK, S. & JORDAN, M. I. (2018a) Is Q-learning provably efficient? *Advances in Neural Information Processing Systems*, pp. 4863–4873.
30. JIN, C., ALLEN-ZHU, Z., BUBECK, S. & JORDAN, M. I. (2018b) *Is Q-learning provably efficient?* *arXiv preprint arXiv:1807.03765*.
31. JIN, C., YANG, Z., WANG, Z. & JORDAN, M. I. (2020) Provably efficient reinforcement learning with linear function approximation. *Conference on Learning Theory*. PMLR, pp. 2137–2143.
32. JOHNSON, R. & ZHANG, T. (2013) Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, pp. 315–323.
33. KAKADE, S. (2003) *On the sample complexity of reinforcement learning* PhD thesis,. University of London.
34. KHAMARU, K., PANANJADY, A., RUAN, F., WAINWRIGHT, M. J. & JORDAN, M. I. (2021) Is temporal difference learning optimal? an instance-dependent analysis. *SIAM Journal on Mathematics of Data Science*, SIAM, **3**(4), 1013–1040.

35. LAI, T. L. & ROBBINS, H. (1985) Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, **6**, 4–22.
36. LATTIMORE, T. & SZEPESVÁRI, C. (2020) *Bandit algorithms*. Cambridge University Press.
37. LI, G., CAI, C., CHEN, Y., GU, Y., WEI, Y. & CHI, Y. (2021a) *Is Q-learning minimax optimal? a tight sample complexity analysis*. arXiv preprint arXiv:2102.06548.
38. LI, G., CHEN, Y., CHI, Y., GU, Y. & WEI, Y. (2021b) Sample-efficient reinforcement learning is feasible for linearly realizable MDPs with limited revisiting. *Advances in Neural Information Processing Systems*, **34**, 16671–16685.
39. LI, G., CHI, Y., WEI, Y. & CHEN, Y. (2022a) *Minimax-optimal multi-agent RL in markov games with a generative model*. arXiv preprint arXiv:2208.10458.
40. LI, G., SHI, L., CHEN, Y., CHI, Y. & WEI, Y. (2022b) *Settling the sample complexity of model-based offline reinforcement learning*. arXiv preprint arXiv:2204.05275.
41. LI, G., WEI, Y., CHI, Y., GU, Y. & CHEN, Y. (2020a) Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Advances in Neural Information Processing Systems*, vol. **33**.
42. LI, G., WEI, Y., CHI, Y., GU, Y. & CHEN, Y. (2020b) Sample complexity of asynchronous Q-learning: Sharper analysis and variance reduction. *Advances in Neural Information Processing Systems (NeurIPS)*.
43. LIU, S. & SU, H. (2020) γ -regret for non-episodic reinforcement learning. arXiv:2002.05138.
44. MÉNARD, P., DOMINGUES, O. D., SHANG, X. & VALKO, M. (2021) UCB Momentum Q-learning: Correcting the bias without forgetting. *International Conference on Machine Learning*, PMLR, 7609–7618.
45. MURPHY, S. (2005) A generalization error for Q-learning. *Journal of Machine Learning Research*, **6**, 1073–1097.
46. NGUYEN, L. M., LIU, J., SCHEINBERG, K. & TAKÁČ, M. (2017) SARAH: A novel method for machine learning problems using stochastic recursive gradient. *International Conference on Machine Learning*. PMLR, pp. 2613–2621.
47. OSBAND, I. & VAN ROY, B. (2016) *On lower bounds for regret in reinforcement learning*. arXiv preprint arXiv:1608.02732.
48. PACCHIANO, A., BALL, P., PARKER-HOLDER, J., CHOROMANSKI, K. & ROBERTS, S. (2020) *On optimism in model-based reinforcement learning*. arXiv preprint arXiv:2006.11911.
49. PUTERMAN, M. L. (2014) *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
50. QU, G. & WIERMAN, A. (2020) Finite-time analysis of asynchronous stochastic approximation and Q-learning. *Conference on Learning Theory*, **125**, 3185–3205.
51. ROBBINS, H. & MONRO, S. (1951) A stochastic approximation method. *The annals of mathematical statistics*, **22**(3), 400–407.
52. SHI, L. & CHI, Y. (2022) *Distributionally robust model-based offline reinforcement learning with near-optimal sample complexity*. arXiv preprint arXiv:2208.05767.
53. SHI, L., LI, G., WEI, Y., CHEN, Y. & CHI, Y. (2022) Pessimistic Q-learning for offline reinforcement learning: Towards optimal sample complexity. *Proceedings of the 39th International Conference on Machine Learning*, PMLR, **162**, 19967–20025.
54. SIDFORD, A., WANG, M., WU, X., YANG, L. & YE, Y. (2018a) Near-optimal time and sample complexities for solving markov decision processes with a generative model. *Advances in Neural Information Processing Systems*, pp. 5186–5196.
55. SIDFORD, A., WANG, M., WU, X. & YE, Y. (2018b) Variance reduced value iteration and faster algorithms for solving Markov decision processes. *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, pp. 770–787.
56. STREHL, A. L., LI, L., WIEWIORA, E., LANGFORD, J. & LITTMAN, M. L. (2006) PAC model-free reinforcement learning. *Proceedings of the 23rd international conference on Machine learning*, pp. 881–888.
57. SZEPESVÁRI, C. (1997) The asymptotic convergence-rate of Q-learning. *NIPS*, vol. **10**. Citeseer, pp. 1064–1070.

58. TAO, T. (2012). *Topics in Random Matrix Theory*. Graduate Studies in Mathematics. American Mathematical Society, Providence, Rhode Island.
59. TROPP, J. (2011) Freedman's inequality for matrix martingales. *Electron. Comm. Probab.*, **16**, 262–270.
60. TSITSIKLIS, J. N. (1994) Asynchronous stochastic approximation and Q-learning. *Machine learning*, **16**, 185–202.
61. WAI, H.-T., HONG, M., YANG, Z., WANG, Z. & TANG, K. (2019) Variance reduced policy evaluation with smooth function approximation. *Advances in Neural Information Processing Systems*, **32**, 5784–5795.
62. WAINWRIGHT, M. J. (2019a) *Stochastic approximation with cone-contractive operators: Sharp ℓ_∞ -bounds for Q-learning*. arXiv preprint arXiv:1905.06265.
63. WAINWRIGHT, M. J. (2019b) *Variance-reduced Q-learning is minimax optimal*. arXiv preprint arXiv:1906.04697.
64. WANG, B., YAN, Y. & FAN, J. (2021) Sample-efficient reinforcement learning for linearly-parameterized mdps with a generative model. *Advances in Neural Information Processing Systems*, **34**, 23009–23022.
65. WATKINS, C. J. & DAYAN, P. (1992) Q-learning. *Machine learning*, vol. **8**, pp. 279–292.
66. WATKINS, C. J. C. H. (1989) Learning from delayed rewards. *PhD thesis, King's College, University of Cambridge*.
67. WENG, B., XIONG, H., ZHAO, L., LIANG, Y. & ZHANG, W. (2020) *Momentum Q-learning with finite-sample convergence guarantee*. arXiv preprint arXiv:2007.15418.
68. XIONG, H., ZHAO, L., LIANG, Y. & ZHANG, W. (2020) Finite-time analysis for double Q-learning. *Advances in Neural Information Processing Systems*, **33**.
69. XU, T., WANG, Z., ZHOU, Y. & LIANG, Y. (2019) Reanalysis of variance reduced temporal difference learning. *International Conference on Learning Representations*.
70. YANG, K., YANG, L. & DU, S. (2021) Q-learning with logarithmic regret. *International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 1576–1584.
71. YIN, M., BAI, Y. & WANG, Y.-X. (2021) Near-optimal offline reinforcement learning via double variance reduction. *Advances in neural information processing systems*, **34**, 7677–7688.
72. ZANETTE, A. & BRUNSKILL, E. (2019) Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. *International Conference on Machine Learning*. PMLR, pp. 7304–7312.
73. ZHANG, K., KAKADE, S., BASAR, T. & YANG, L. (2020a) Model-based multi-agent RL in zero-sum Markov games with near-optimal sample complexity. *Advances in Neural Information Processing Systems*, **33**.
74. ZHANG, Z., JI, X. & DU, S. S. (2021) Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. *Conference on Learning Theory*, PMLR, 4528–4531.
75. ZHANG, Z., ZHOU, Y. & JI, X. (2020) Almost optimal model-free reinforcement learning via reference-advantage decomposition. *Advances in Neural Information Processing Systems*, **33**, 15198–15207.
76. ZHANG, Z., ZHOU, Y. & JI, X. (2021) Model-free reinforcement learning: from clipped pseudo-regret to sample complexity. *International Conference on Machine Learning*, PMLR, 12653–12662.

A. Freedman's inequality

A.1 A user-friendly version of Freedman's inequality

Due to the Markovian structure of the problem, our analysis relies heavily on the celebrated Freedman's inequality [23, 59], which extends the Bernstein's inequality to accommodate martingales. For ease of reference, we state below a user-friendly version of Freedman's inequality as provided in (37, Section C).

THEOREM A.1. Freedman's inequality. Consider a filtration $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots$, and let \mathbb{E}_k stand for the expectation conditioned on \mathcal{F}_k . Suppose that $Y_n = \sum_{k=1}^n X_k \in \mathbb{R}$, where $\{X_k\}$ is a real-valued scalar sequence obeying

$$|X_k| \leq R \quad \text{and} \quad \mathbb{E}_{k-1}[X_k] = 0 \quad \text{for all } k \geq 1$$

for some quantity $R < \infty$. We also define

$$W_n := \sum_{k=1}^n \mathbb{E}_{k-1}[X_k^2].$$

In addition, suppose that $W_n \leq \sigma^2$ holds deterministically for some given quantity $\sigma^2 < \infty$. Then for any positive integer $m \geq 1$, with probability at least $1 - \delta$ one has

$$|Y_n| \leq \sqrt{8 \max \left\{ W_n, \frac{\sigma^2}{2^m} \right\} \log \frac{2m}{\delta}} + \frac{4}{3} R \log \frac{2m}{\delta}. \quad (\text{A.1})$$

A.2 Application of Freedman's inequality

We now develop several immediate consequences of Freedman's inequality, which lend themselves well to our context. Before proceeding, we recall that $N_h^i(s, a)$ denotes the number of times that the state-action pair (s, a) has been visited at step h by the end of the i -th episode, and $k_h^n(s, a)$ stands for the episode index when (s, a) is visited at step h for the n -th time (see Section 4.2).

Our first result is concerned with a martingale concentration bound as follows:

LEMMA 7. Let $\{W_h^i \in \mathbb{R}^S \mid 1 \leq i \leq K, 1 \leq h \leq H+1\}$ and $\{u_h^i(s, a, N) \in \mathbb{R} \mid 1 \leq i \leq K, 1 \leq h \leq H+1\}$ be the collections of vectors and scalars, respectively, and suppose that they obey the following properties:

- W_h^i is fully determined by the samples collected up to the end of the $(h-1)$ -th step of the i -th episode;
- $\|W_h^i\|_\infty \leq C_w$;
- $u_h^i(s, a, N)$ is fully determined by the samples collected up to the end of the $(h-1)$ -th step of the i -th episode, and a given positive integer $N \in [K]$;
- $0 \leq u_h^i(s, a, N) \leq C_u$;
- $0 \leq \sum_{n=1}^{N_h^i(s, a)} u_h^{k_h^n(s, a)}(s, a, N) \leq 2$.

In addition, consider the following sequence:

$$X_i(s, a, h, N) := u_h^i(s, a, N) (P_h^i - P_{h, s, a}) W_{h+1}^i \mathbb{1}\{(s_h^i, a_h^i) = (s, a)\}, \quad 1 \leq i \leq K, \quad (\text{A.2})$$

with P_h^i defined in (4.5). Consider any $\delta \in (0, 1)$. Then with probability at least $1 - \delta$,

$$\left| \sum_{i=1}^k X_i(s, a, h, N) \right| \lesssim \sqrt{C_u \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k(s,a)} u_h^{k_h^n(s,a)}(s, a, N) \text{Var}_{h,s,a}(W_{h+1}^{k_h^n(s,a)})} + \left(C_u C_w + \sqrt{\frac{C_u}{N}} C_w \right) \log^2 \frac{SAT}{\delta} \quad (\text{A.3})$$

holds simultaneously for all $(k, h, s, a, N) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A} \times [K]$.

Proof. For the sake of notational convenience, we shall abbreviate $X_i(s, a, h, N)$ as X_i throughout the proof of this lemma, as long as it is clear from the context. The plan is to apply Freedman's inequality (cf. Theorem A.1) to control the term $\sum_{i=1}^k X_i$ of interest.

Consider any given $(k, h, s, a, N) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A} \times [K]$. It can be easily verified that

$$\mathbb{E}_{i-1}[X_i] = 0,$$

where \mathbb{E}_{i-1} denotes the expectation conditioned on everything happening up to the end of the $(h-1)$ -th step of the i -th episode. Additionally, we make note of the following crude bound:

$$\begin{aligned} |X_i| &\leq u_h^i(s, a, N) \left| (P_h^i - P_{h,s,a}) W_{h+1}^i \right| \\ &\leq u_h^i(s, a, N) \left(\|P_h^i\|_1 + \|P_{h,s,a}\|_1 \right) \|W_{h+1}^i\|_\infty \leq 2C_w C_u, \end{aligned} \quad (\text{A.4})$$

which results from the assumptions $\|W_{h+1}^i\|_\infty \leq C_w$, $0 \leq u_h^i(s, a, N) \leq C_u$ as well as the basic facts $\|P_h^i\|_1 = \|P_{h,s,a}\|_1 = 1$. To continue, recalling the definition of the variance parameter in (4.6), we obtain

$$\begin{aligned} \sum_{i=1}^k \mathbb{E}_{i-1} [|X_i|^2] &= \sum_{i=1}^k (u_h^i(s, a, N))^2 \mathbb{1}\{(s_h^i, a_h^i) = (s, a)\} \mathbb{E}_{i-1} [| (P_h^i - P_{h,s,a}) W_{h+1}^i |^2] \\ &= \sum_{n=1}^{N_h^k(s,a)} (u_h^{k_h^n(s,a)}(s, a, N))^2 \text{Var}_{h,s,a}(W_{h+1}^{k_h^n(s,a)}) \\ &\leq C_u \left(\sum_{n=1}^{N_h^k(s,a)} u_h^{k_h^n(s,a)}(s, a, N) \right) \|W_{h+1}^{k_h^n(s,a)}\|_\infty^2 \\ &\leq 2C_u C_w^2, \end{aligned} \quad (\text{A.5})$$

where the inequalities hold true due to the assumptions $\|W_h^i\|_\infty \leq C_w$, $0 \leq u_h^i(s, a, N) \leq C_u$, and $0 \leq \sum_{n=1}^{N_h^k(s,a)} u_h^{k_h^n(s,a)}(s, a, N) \leq 1$.

With (A.4) and (A.5) in place, we can invoke Theorem A.1 (with $m = \lceil \log_2 N \rceil$) and take the union bound over all $(k, h, s, a, N) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A} \times [K]$ to show that, with probability at least $1 - \delta$,

$$\begin{aligned} \left| \sum_{i=1}^k X_i \right| &\lesssim \sqrt{\max \left\{ C_u \sum_{n=1}^{N_h^k(s,a)} u_h^{k_h^n(s,a)}(s, a, N) \text{Var}_{h,s,a}(W_{h+1}^{k_h^n(s,a)}), \frac{C_u C_w^2}{N} \right\} \log \frac{SAT^2 \log N}{\delta}} \\ &\quad + C_u C_w \log \frac{SAT^2 \log N_h^k}{\delta} \\ &\lesssim \sqrt{C_u \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k(s,a)} u_h^{k_h^n(s,a)}(s, a, N) \text{Var}_{h,s,a}(W_{h+1}^{k_h^n(s,a)})} + \left(C_u C_w + \sqrt{\frac{C_u}{N}} C_w \right) \log^2 \frac{SAT}{\delta} \end{aligned}$$

holds simultaneously for all $(k, h, s, a, N) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A} \times [K]$. \square

The next result is concerned with martingale concentration bounds for another type of sequences of interest.

LEMMA 8. Let $\{N(s, a, h) \in [K] \mid (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]\}$ be a collection of positive integers, and let $\{c_h : 0 \leq c_h \leq e, h \in [H]\}$ be a collection of fixed and bounded universal constants. Moreover, let $\{W_h^i \in \mathbb{R}^S \mid 1 \leq i \leq K, 1 \leq h \leq H + 1\}$ and $\{u_h^i(s_h^i, a_h^i) \in \mathbb{R} \mid 1 \leq i \leq K, 1 \leq h \leq H + 1\}$ represent, respectively, the collections of random vectors and scalars, which obey the following properties.

- W_h^i is fully determined by the samples collected up to the end of the $(h - 1)$ -th step of the i -th episode;
- $\|W_h^i\|_\infty \leq C_w$ and $W_h^i \geq 0$;
- $u_h^i(s_h^i, a_h^i)$ is fully determined by the integer $N(s_h^i, a_h^i, h)$ and all samples collected up to the end of the $(h - 1)$ -th step of the i -th episode;
- $0 \leq u_h^i(s_h^i, a_h^i) \leq C_u$.

Consider any $\delta \in (0, 1)$, and introduce the following sequences:

$$X_{i,h} := u_h^i(s_h^i, a_h^i)(P_h^i - P_{h,s_h^i, a_h^i})W_{h+1}^i, \quad 1 \leq i \leq K, 1 \leq h \leq H + 1, \quad (\text{A.6})$$

$$Y_{i,h} := c_h(P_h^i - P_{h,s_h^i, a_h^i})W_{h+1}^i, \quad 1 \leq i \leq K, 1 \leq h \leq H + 1. \quad (\text{A.7})$$

Then with probability at least $1 - \delta$,

$$\begin{aligned} \left| \sum_{h=1}^H \sum_{i=1}^K X_{i,h} \right| &\lesssim \sqrt{C_u^2 \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[\left| (P_h^i - P_{h,s_h^i, a_h^i}) W_{h+1}^i \right|^2 \right] \log \frac{T^{HSA}}{\delta}} + C_u C_w \log \frac{T^{HSA}}{\delta} \\ &\lesssim \sqrt{C_u^2 C_w \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[P_h^i W_{h+1}^i \right] \log \frac{T^{HSA}}{\delta}} + C_u C_w \log \frac{T^{HSA}}{\delta} \\ \left| \sum_{h=1}^H \sum_{i=1}^K Y_{i,h} \right| &\lesssim \sqrt{T C_w^2 \log \frac{1}{\delta}} + C_w \log \frac{1}{\delta} \end{aligned}$$

holds simultaneously for all possible collections $\{N(s, a, h) \in [K] \mid (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]\}$.

Proof. This lemma can be proved by Freedman's inequality (cf. Theorem A.1).

- We start by controlling the first term of interest $\sum_{h=1}^H \sum_{i=1}^K X_{i,h}$. As can be easily seen, $a_h^i = \arg \max Q_h^i(s_h^i, a)$ is fully determined by what happens before step h of the i -th episode. Consider any given $\{N(s, a, h) \in [K] \mid (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]\}$. It is readily seen that

$$\mathbb{E}_{i,h-1} [X_i] = \mathbb{E}_{i,h-1} \left[u_h^i(s_h^i, a_h^i) (P_h^i - P_{h,s_h^i, a_h^i}) W_{h+1}^i \right] = 0,$$

where $\mathbb{E}_{i,h-1}$ denotes the expectation conditioned on everything happening before step h of the i -th episode. In addition, we make note of the following crude bound:

$$\begin{aligned} |X_{i,h}| &\leq u_h^i(s_h^i, a_h^i) \left| (P_h^i - P_{h,s_h^i, a_h^i}) W_{h+1}^i \right| \\ &\leq u_h^i(s_h^i, a_h^i) \left(\|P_h^i\|_1 + \|P_{h,s_h^i, a_h^i}\|_1 \right) \|W_{h+1}^i\|_\infty \leq 2C_w C_u, \end{aligned} \tag{A.8}$$

which arises from the assumptions $\|W_{h+1}^i\|_\infty \leq C_w$, $0 \leq u_h^i(s, a, N) \leq C_u$ together with the basic facts $\|P_h^i\|_1 = \|P_{h,s_h^i, a_h^i}\|_1 = 1$. Additionally, we can calculate that

$$\begin{aligned} \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} [|X_{i,h}|^2] &= \sum_{h=1}^H \sum_{i=1}^K (u_h^i(s_h^i, a_h^i))^2 \mathbb{E}_{i,h-1} \left[\left| (P_h^i - P_{h,s_h^i, a_h^i}) W_{h+1}^i \right|^2 \right] \\ &\stackrel{(i)}{\leq} C_u^2 \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[\left| (P_h^i - P_{h,s_h^i, a_h^i}) W_{h+1}^i \right|^2 \right] \end{aligned} \tag{A.9}$$

$$\begin{aligned}
&\leq C_u^2 \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[|P_h^i W_{h+1}^i|^2 \right] \\
&\stackrel{(ii)}{=} C_u^2 \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[P_h^i (W_{h+1}^i)^2 \right] \\
&\stackrel{(iii)}{\leq} C_u^2 \sum_{h=1}^H \sum_{i=1}^K \|W_{h+1}^i\|_\infty \mathbb{E}_{i,h-1} \left[P_h^i W_{h+1}^i \right] \\
&\stackrel{(iv)}{\leq} C_u^2 C_w \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[P_h^i W_{h+1}^i \right] \\
&\stackrel{(A.10)}{=} C_u^2 C_w \sum_{h=1}^H \sum_{i=1}^K \|W_{h+1}^i\|_\infty \stackrel{(v)}{\leq} HK C_u^2 C_w^2 = T C_u^2 C_w^2.
\end{aligned}$$

$$\leq C_u^2 C_w \sum_{h=1}^H \sum_{i=1}^K \|W_{h+1}^i\|_\infty \stackrel{(v)}{\leq} HK C_u^2 C_w^2 = T C_u^2 C_w^2. \quad (A.11)$$

Here, (i) holds true due to the assumption $0 \leq u_h^i(s_h^i, a_h^i) \leq C_u$, (ii) is valid since P_h^i only has one non-zero entry (cf. (4.5)), (iii) relies on the assumptions that W_h^i is non-negative, whereas (iv) and (v) follow since $\|W_h^i\|_\infty \leq C_w$,

- With (A.8), (A.10) and (A.11) in mind, we can invoke Theorem A.1 (with $m = \lceil \log_2 T \rceil$) and take the union bound over all possible collections $\{N(s, a, h) \in [K] \mid (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]\}$ —which has at most K^{HSA} possibilities—to show that, with probability at least $1 - \delta$,

$$\begin{aligned}
\left| \sum_{h=1}^H \sum_{i=1}^K X_{i,h} \right| &\lesssim \sqrt{\max \left\{ C_u^2 \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[|(P_h^i - P_{h,s_h^i, a_h^i}) W_{h+1}^i|^2 \right], \frac{T C_u^2 C_w^2}{2^m} \right\} \log \frac{K^{HSA} \log T}{\delta}} \\
&\quad + C_u C_w \log \frac{K^{HSA} \log T}{\delta} \\
&\lesssim \sqrt{C_u^2 \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[|(P_h^i - P_{h,s_h^i, a_h^i}) W_{h+1}^i|^2 \right] \log \frac{T^{HSA}}{\delta}} + C_u C_w \log \frac{T^{HSA}}{\delta} \\
&\lesssim \sqrt{C_u^2 C_w \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[P_h^i W_{h+1}^i \right] \log \frac{T^{HSA}}{\delta}} + C_u C_w \log \frac{T^{HSA}}{\delta}
\end{aligned}$$

holds simultaneously for all $\{N(s, a, h) \in [K] \mid (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]\}$.

- Then we turn to control the second term $\left| \sum_{h=1}^H \sum_{i=1}^K Y_{i,h} \right|$ of interest. Similar to $\left| \sum_{h=1}^H \sum_{i=1}^K X_{i,h} \right|$, we have

$$\begin{aligned}
|Y_{i,h}| &\leq 2eC_w, \\
\sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[|Y_{i,h}|^2 \right] &\leq e^2 T C_w^2.
\end{aligned}$$

Invoke Theorem A.1 (with $m = 1$) to arrive at

$$\left| \sum_{h=1}^H \sum_{i=1}^K Y_{i,h} \right| \lesssim \sqrt{TC_w^2 \log \frac{1}{\delta}} + C_w \log \frac{1}{\delta} \quad (\text{A.12})$$

with probability at least $1 - \delta$. \square

B. Proof of Lemma 1

First of all, the properties in (4.4b) follow directly from (29, Lemma 4.1). Therefore, it suffices to establish the property in (4.4a), which forms the remainder of this subsection.

When $N = 1$, the statement holds trivially since

$$\sum_{n=1}^N \frac{\eta_n^N}{n^a} = \eta_1^1 = 1 \in [1, 2].$$

Now suppose that $N \geq 2$. Using the basic relation $\eta_n^N = (1 - \eta_N)\eta_n^{N-1}$ for all $n = 1, \dots, N-1$, we observe the following identity:

$$\sum_{n=1}^N \frac{\eta_n^N}{n^a} = \frac{\eta_N}{N^a} + (1 - \eta_N) \sum_{n=1}^{N-1} \frac{\eta_n^{N-1}}{n^a}. \quad (\text{B.1})$$

We now prove the property in (4.4a) by induction. Suppose for the moment that the property holds for $N-1$, namely,

$$\frac{1}{(N-1)^a} \leq \sum_{n=1}^{N-1} \frac{\eta_n^{N-1}}{n^a} \leq \frac{2}{(N-1)^a}. \quad (\text{B.2})$$

Then it is readily seen from (B.1) that

$$\sum_{n=1}^N \frac{\eta_n^N}{n^a} = \frac{\eta_N}{N^a} + (1 - \eta_N) \sum_{n=1}^{N-1} \frac{\eta_n^{N-1}}{n^a} \geq \frac{\eta_N}{N^a} + \frac{1 - \eta_N}{(N-1)^a} \geq \frac{\eta_N}{N^a} + \frac{1 - \eta_N}{N^a} = \frac{1}{N^a}, \quad (\text{B.3})$$

where the first inequality comes from (B.2). Similarly, one can upper bound

$$\begin{aligned} \sum_{n=1}^N \frac{\eta_n^N}{n^a} &= \frac{\eta_N}{N^a} + (1 - \eta_N) \sum_{n=1}^{N-1} \frac{\eta_n^{N-1}}{n^a} \stackrel{\text{(i)}}{\leq} \frac{\eta_N}{N^a} + \frac{2(1 - \eta_N)}{(N-1)^a} \stackrel{\text{(ii)}}{=} \frac{H+1}{N^a(H+N)} + \frac{2(N-1)^{1-a}}{H+N} \\ &\stackrel{\text{(iii)}}{\leq} \frac{H+1}{N^a(H+N)} + \frac{2N^{1-a}}{H+N} = \frac{1}{N^a} \left(\frac{H+1}{H+N} + \frac{2N}{H+N} \right) \stackrel{\text{(iv)}}{\leq} \frac{2}{N^a}, \end{aligned}$$

where (i) arises from (B.2), (ii) follows from the choice $\eta_N = \frac{H+1}{H+N}$, (iii) holds since $a \leq 1$ and (iv) follows since $H \geq 1$. Consequently, we can immediately establish the advertised property (4.4a) by induction.

C. Proof of key lemmas in Section 4.3

C.1 Proof of Lemma 2

To begin with, suppose that we can prove

$$Q_h^k(s, a) \geq Q_h^*(s, a) \quad \text{for all } (k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}. \quad (\text{C.1})$$

Then this property would immediately lead to the claim w.r.t. V_h^k , namely,

$$V_h^k(s) \geq Q_h^k(s, \pi_h^*(s)) \geq Q_h^*(s, \pi_h^*(s)) = V_h^*(s) \quad \text{for all } (k, h, s) \in [K] \times [H] \times \mathcal{S}. \quad (\text{C.2})$$

As a result, it suffices to focus on justifying the claim (C.1), which we shall accomplish by induction.

- *Base case.* Given that the initialization obeys $Q_h^1(s, a) = H \geq Q_h^*(s, a)$ for all $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$, the claim (C.1) holds trivially when $k = 1$.
- *Induction.* Suppose that the claim (C.1) holds all the way up to the k -th episode, and we wish to establish it for the $(k + 1)$ -th episode as well. To complete the induction argument, it suffices to justify

$$\min \left\{ Q_h^{\text{UCB}, k+1}(s, a), Q_h^{\text{R}, k+1}(s, a) \right\} \geq Q_h^*(s, a)$$

according to line 12 of Algorithm 3. Recognizing that $Q_h^{\text{UCB}, k+1}$ is computed via the standard UCB-Q update rule (see line 2 of Algorithm 2), we can readily invoke the argument in (29, Lemma 4.3) to show that with probability at least $1 - \delta$,

$$Q_h^{\text{UCB}, k+1}(s, a) \geq Q_h^*(s, a)$$

holds simultaneously for all $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$. Therefore, it is sufficient to prove that

$$Q_h^{\text{R}, k+1}(s, a) \geq Q_h^*(s, a). \quad (\text{C.3})$$

The remainder of the proof is thus devoted to justifying (C.3), assuming that the claim (C.1) holds all the way up to k .

Since Q_h^{R} , $k(s_h^k, a_h^k)$ is updated in the k -th episode while other entries of Q_h^{R} , k remain fixed, it suffices to verify

$$Q_h^{\text{R}, k+1}(s_h^k, a_h^k) \geq Q_h^*(s_h^k, a_h^k).$$

We remind the readers of two important short-hand notation that shall be used when it is clear from the context:

- $N_h^k = N_h^k(s_h^k, a_h^k)$ denotes the number of times that the state-action pair (s_h^k, a_h^k) has been visited at step h by the end of the k -th episode;
- $k^n = k_h^n(s_h^k, a_h^k)$ denotes the index of the episode in which the state-action pair (s_h^k, a_h^k) is visited for the n -th time at step h .

Step 1: decomposing $Q_h^{R,k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)$.

To begin with, the above definition of N_h^k and k^n allows us to write

$$Q_h^{R,k+1}(s_h^k, a_h^k) = Q_h^{R,k^N_h+1}(s_h^k, a_h^k), \quad (\text{C.4})$$

since $k^N_h = k^{N_h^k(s_h^k, a_h^k)} = k$. According to the update rule (i.e. line 11 in Algorithm 3 and line 9 in Algorithm 2), we obtain

$$\begin{aligned} Q_h^{R,k+1}(s_h^k, a_h^k) &= Q_h^{R,k^N_h+1}(s_h^k, a_h^k) = (1 - \eta_{N_h^k})Q_h^{R,k^N_h}(s_h^k, a_h^k) \\ &\quad + \eta_{N_h^k} \left\{ r_h(s_h^k, a_h^k) + V_{h+1}^{k^N_h}(s_{h+1}^{k^N_h}) - V_{h+1}^{R,k^N_h}(s_{h+1}^{k^N_h}) + \mu_h^{\text{ref}, k^N_h+1}(s_h^k, a_h^k) + b_h^{R,k^N_h+1} \right\} \\ &= (1 - \eta_{N_h^k})Q_h^{R,k^{N_h^k-1}+1}(s_h^k, a_h^k) \\ &\quad + \eta_{N_h^k} \left\{ r_h(s_h^k, a_h^k) + V_{h+1}^{k^N_h}(s_{h+1}^{k^N_h}) - V_{h+1}^{R,k^N_h}(s_{h+1}^{k^N_h}) + \mu_h^{\text{ref}, k^N_h+1}(s_h^k, a_h^k) + b_h^{R,k^N_h+1} \right\}, \end{aligned}$$

where the last identity again follows from our argument for justifying (C.4). Applying this relation recursively and invoking the definitions of η_0^N and η_n^N in (4.2), we are left with

$$\begin{aligned} Q_h^{R,k+1}(s_h^k, a_h^k) &= \eta_0^{N_h^k} Q_h^{R,1}(s_h^k, a_h^k) \\ &\quad + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left\{ r_h(s_h^k, a_h^k) + V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \mu_h^{\text{ref}, k^n+1}(s_h^k, a_h^k) + b_h^{R,k^n+1} \right\}. \end{aligned} \quad (\text{C.5})$$

Additionally, the basic relation $\eta_0^{N_h^k} + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} = 1$ (see (4.2) and (4.3)) tells us that

$$Q_h^*(s_h^k, a_h^k) = \eta_0^{N_h^k} Q_h^*(s_h^k, a_h^k) + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} Q_h^*(s_h^k, a_h^k), \quad (\text{C.6})$$

which combined with (C.5) leads to

$$\begin{aligned} Q_h^{R,k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) &= \eta_0^{N_h^k} (Q_h^{R,1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)) \\ &\quad + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left\{ r_h(s_h^k, a_h^k) + V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \mu_h^{\text{ref}, k^n+1}(s_h^k, a_h^k) + b_h^{R,k^n+1} - Q_h^*(s_h^k, a_h^k) \right\}. \end{aligned} \quad (\text{C.7})$$

To continue, invoking the Bellman optimality equation

$$Q_h^*(s_h^k, a_h^k) = r_h(s_h^k, a_h^k) + P_{h,s_h^k,a_h^k} V_{h+1}^* \quad (\text{C.8})$$

and using the construction of μ_h^{ref} in line 11 of Algorithm 2 (which is the running mean of V_{h+1}^R), we reach

$$\begin{aligned} r_h(s_h^k, a_h^k) + V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R, k^n}(s_{h+1}^{k^n}) + \mu_h^{\text{ref}, k^n+1}(s_h^k, a_h^k) + b_h^{R, k^n+1} - Q_h^*(s_h^k, a_h^k) \\ = V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R, k^n}(s_{h+1}^{k^n}) + \frac{\sum_{i=1}^n V_{h+1}^{R, k^i}(s_{h+1}^{k^i})}{n} - P_{h, s_h^k, a_h^k} V_{h+1}^* + b_h^{R, k^n+1} \end{aligned} \quad (\text{C.9})$$

$$\begin{aligned} &= P_{h, s_h^k, a_h^k} \left\{ V_{h+1}^{k^n} - V_{h+1}^{R, k^n} \right\} + \frac{\sum_{i=1}^n P_{h, s_h^k, a_h^k} (V_{h+1}^{R, k^i})}{n} - P_{h, s_h^k, a_h^k} V_{h+1}^* + b_h^{R, k^n+1} + \xi_h^{k^n}, \\ &= P_{h, s_h^k, a_h^k} \left\{ V_{h+1}^{k^n} - V_{h+1}^* + \frac{\sum_{i=1}^n (V_{h+1}^{R, k^i} - V_{h+1}^{R, k^n})}{n} \right\} + b_h^{R, k^n+1} + \xi_h^{k^n}. \end{aligned} \quad (\text{C.10})$$

Here, we have introduced the following quantity:

$$\xi_h^{k^n} := (P_h^{k^n} - P_{h, s_h^k, a_h^k}) (V_{h+1}^{k^n} - V_{h+1}^{R, k^n}) + \frac{1}{n} \sum_{i=1}^n (P_h^{k^i} - P_{h, s_h^k, a_h^k}) V_{h+1}^{R, k^i}, \quad (\text{C.11})$$

with the notation P_h^k defined in (4.5). Putting (C.10) and (C.7) together leads to the following decomposition:

$$\begin{aligned} Q_h^{R, k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) &= \eta_0^{N_h^k} \left(Q_h^{R, 1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \right) \\ &\quad + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left\{ P_{h, s_h^k, a_h^k} \left(V_{h+1}^{k^n} - V_{h+1}^* + \frac{\sum_{i=1}^n (V_{h+1}^{R, k^i} - V_{h+1}^{R, k^n})}{n} \right) + b_h^{R, k^n+1} + \xi_h^{k^n} \right\}. \end{aligned} \quad (\text{C.12})$$

Step 2: two key quantities for lower bounding $Q_h^{R, k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)$.

In order to develop a lower bound on $Q_h^{R, k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)$ based on the decomposition (C.12), we make note of several simple facts as follows:

- (i) The initialization satisfies $Q_h^{R, 1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \geq 0$.
- (ii) For any $1 \leq k^n \leq k$, one has

$$V_{h+1}^{k^n} \geq V_{h+1}^*, \quad (\text{C.13})$$

owing to the induction hypotheses (C.1) and (C.2) that hold up to k .

- (iii) For all $0 \leq i \leq n$ and any $s \in \mathcal{S}$, one has

$$V_{h+1}^{R, k^i}(s) - V_{h+1}^{R, k^n}(s) \geq 0, \quad (\text{C.14})$$

which holds since the reference value $V_h^R(s)$ is monotonically non-increasing in view of the monotonicity of $V_h(s)$ in (4.7b) and the update rule in line 16 of Algorithm 3.

The above three facts taken collectively with (C.12) allow one to drop several terms and yield

$$Q_h^{R,k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \geq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (b_h^{R,k^n+1} + \xi_h^{k^n}). \quad (\text{C.15})$$

In the sequel, we aim to establish $Q_h^{R,k+1}(s_h^k, a_h^k) \geq Q_h^*(s_h^k, a_h^k)$ based on this inequality (C.15).

As it turns out, if one could show that

$$\left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \xi_h^{k^n} \right| \leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} b_h^{R,k^n+1}, \quad (\text{C.16})$$

then taking this together with (C.15) and the triangle inequality would immediately lead to the desired result

$$Q_h^{R,k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \geq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} b_h^{R,k^n+1} - \left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \xi_h^{k^n} \right| \geq 0. \quad (\text{C.17})$$

As a result, the remaining steps come down to justifying the claim (C.16). In order to do so, we need to control the following two quantities (in view of (C.11)):

$$I_1 := \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (P_h^{k^n} - P_{h,s_h^k, a_h^k}) (V_{h+1}^{k^n} - V_{h+1}^{R,k^n}), \quad (\text{C.18a})$$

$$I_2 := \sum_{n=1}^{N_h^k} \frac{1}{n} \eta_n^{N_h^k} \sum_{i=1}^n (P_h^{k^i} - P_{h,s_h^k, a_h^k}) V_{h+1}^{R,k^i} \quad (\text{C.18b})$$

separately, which constitutes the next two steps. As will be seen momentarily, these two terms can be controlled in a similar fashion using Freedman's inequality.

Step 3: controlling I_1 . In the following text, we intend to invoke Lemma 7 to control the term I_1 defined in (C.18a). To begin with, consider any $(N, h) \in [K] \times [H]$, and introduce

$$W_{h+1}^i := V_{h+1}^i - V_{h+1}^{R,i} \quad \text{and} \quad u_h^i(s, a, N) := \eta_{N_h^i(s, a)}^N \geq 0. \quad (\text{C.19})$$

Accordingly, we can derive and define

$$\|W_{h+1}^i\|_\infty \leq \|V_{h+1}^{R,i}\|_\infty + \|V_{h+1}^i\|_\infty \leq 2H =: C_w, \quad (\text{C.20})$$

and

$$\max_{N, h, s, a \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}} \eta_{N_h^i(s, a)}^N \leq \frac{2H}{N} =: C_u, \quad (\text{C.21})$$

where the last inequality follows since (according to Lemma 1 and the definition in (4.2))

$$\begin{aligned} \eta_{N_h^i(s, a)}^N &\leq \frac{2H}{N}, & \text{if } 1 \leq N_h^i(s, a) \leq N; \\ \eta_{N_h^i(s, a)}^N &= 0, & \text{if } N_h^i(s, a) > N. \end{aligned}$$

Moreover, observed from (4.3), we have

$$0 \leq \sum_{n=1}^N u_h^{k^n(s,a)}(s,a,N) = \sum_{n=1}^N \eta_n^N \leq 1 \quad (\text{C.22})$$

holds for all $(N, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$. Therefore, choosing $(N, s, a) = (N_h^k, s_h^k, a_h^k)$ and applying Lemma 7 with the quantities (C.19) implies that, with probability at least $1 - \delta$,

$$\begin{aligned} |I_1| &= \left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (P_h^{k^n} - P_{h,s_h^k, a_h^k}) (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right| = \left| \sum_{i=1}^k X_i(s_h^k, a_h^k, h, N_h^k) \right| \\ &\lesssim \sqrt{C_u \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} u_h^{k^n}(s_h^k, a_h^k, N_h^k) \text{Var}_{h,s_h^k, a_h^k}(W_{h+1}^{k^n})} + \left(C_u C_w + \sqrt{\frac{C_u}{N}} C_w \right) \log^2 \frac{SAT}{\delta} \\ &\asymp \sqrt{\frac{H}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \text{Var}_{h,s_h^k, a_h^k}(V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n})} + \frac{H^2 \log^2 \frac{SAT}{\delta}}{N_h^k} \quad (\text{C.23}) \end{aligned}$$

$$\lesssim \sqrt{\frac{H}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sigma_h^{\text{adv}, k^N_h + 1}(s_h^k, a_h^k) - (\mu_h^{\text{adv}, k^N_h + 1}(s_h^k, a_h^k))^2} + \frac{H^2 \log^2 \frac{SAT}{\delta}}{(N_h^k)^{3/4}}, \quad (\text{C.24})$$

where the proof of the last inequality (C.24) needs additional explanation and is postponed to Appendix C.1.1 to streamline the presentation.

Step 4: controlling I_2 . Next, we turn attention to the quantity I_2 defined in (C18b). Rearranging terms in the definition (C18b), we are left with

$$I_2 = \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \frac{\sum_{i=1}^n (P_h^{k^i} - P_{h,s_h^k, a_h^k}) V_{h+1}^{\text{R}, k^i}}{n} = \sum_{i=1}^{N_h^k} \left(\sum_{n=i}^{N_h^k} \frac{\eta_n^{N_h^k}}{n} \right) (P_h^{k^i} - P_{h,s_h^k, a_h^k}) V_{h+1}^{\text{R}, k^i},$$

which can again be controlled by invoking Lemma 7. To do so, we abuse the notation by taking

$$W_{h+1}^i := V_{h+1}^{\text{R}, i} \quad \text{and} \quad u_h^i(s, a, N) := \sum_{n=N_h^i(s, a)}^N \frac{\eta_n^N}{n} \geq 0. \quad (\text{C.25})$$

These quantities satisfy

$$\|W_{h+1}^i\|_\infty \leq \|V_{h+1}^{\text{R}, i}\|_\infty \leq H =: C_w \quad (\text{C.26})$$

and, according to Lemma 1,

$$\max_{N, h, s, a \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}} \sum_{n=N_h^i(s, a)}^N \frac{\eta_n^N}{n} \leq \sum_{n=1}^N \frac{\eta_n^N}{n} \leq \frac{2}{N} =: C_u. \quad (\text{C.27})$$

Then it is readily seen from (C.27) that

$$0 \leq \sum_{n=1}^N u_h^{k^n(s,a)}(s, a, N) \leq \sum_{n=1}^N \frac{2}{N} \leq 2 \quad (\text{C.28})$$

holds for all $(N, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$.

With the above relations in mind, Taking $(N, s, a) = (N_h^k, s_h^k, a_h^k)$ and applying Lemma 7 w.r.t. the quantities (C.25) reveals that

$$|I_2| = \left| \sum_{i=1}^{N_h^k} \sum_{n=i}^{N_h^k} \frac{\eta_n^{N_h^k}}{n} (P_h^{k^i} - P_{h, s_h^k, a_h^k}) V_{h+1}^{\text{R}, k^i} \right| = \left| \sum_{i=1}^k X_i(s_h^k, a_h^k, h, N_h^k) \right| \quad (\text{C.29})$$

$$\begin{aligned} & \lesssim \sqrt{C_u \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} u_h^{k^n}(s_h^k, a_h^k, N_h^k) \text{Var}_{h, s_h^k, a_h^k}(W_{h+1}^{k^n})} + \left(C_u C_w + \sqrt{\frac{C_u}{N}} C_w \right) \log^2 \frac{SAT}{\delta} \\ & \lesssim \sqrt{\frac{1}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} \text{Var}_{h, s_h^k, a_h^k}(V_{h+1}^{\text{R}, k^n})} + \frac{H}{N_h^k} \log^2 \frac{SAT}{\delta} \end{aligned} \quad (\text{C.30})$$

$$\lesssim \sqrt{\frac{1}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sigma_h^{\text{ref}, k^{N_h^k+1}}(s_h^k, a_h^k) - (\mu_h^{\text{ref}, k^{N_h^k+1}}(s_h^k, a_h^k))^2} + \frac{H}{(N_h^k)^{3/4}} \log^2 \frac{SAT}{\delta} \quad (\text{C.31})$$

with probability exceeding $1 - \delta$, where the proof of the last inequality (C.31) is deferred to Appendix C.1.2 in order to streamline presentation.

Step 5: combining the above bounds. Summing up the results in (C.24) and (C.31), we arrive at an upper bound on $|\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \xi_h^{k^n}|$ as follows:

$$\begin{aligned} & \left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \xi_h^{k^n} \right| \leq |I_1| + |I_2| \\ & \lesssim \sqrt{\frac{H}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sigma_h^{\text{adv}, k^{N_h^k+1}}(s_h^k, a_h^k) - (\mu_h^{\text{adv}, k^{N_h^k+1}}(s_h^k, a_h^k))^2} \\ & \quad + \sqrt{\frac{1}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sigma_h^{\text{ref}, k^{N_h^k+1}}(s_h^k, a_h^k) - (\mu_h^{\text{ref}, k^{N_h^k+1}}(s_h^k, a_h^k))^2} + \frac{H^2 \log^2 \frac{SAT}{\delta}}{(N_h^k)^{3/4}} \\ & \leq B_h^{\text{R}, k^{N_h^k+1}}(s_h^k, a_h^k) + c_b \frac{H^2 \log^2 \frac{SAT}{\delta}}{(N_h^k)^{3/4}} \end{aligned} \quad (\text{C.32})$$

for some sufficiently large constant $c_b > 0$, where the last line follows from the definition of $B_h^{\text{R}, k^{N_h^k+1}}(s_h^k, a_h^k)$ in line 17 of Algorithm 2.

In order to establish the desired bound (C.16), we still need to control the sum $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} b_h^{\text{R},k^n+1}$. Toward this end, the definition of b_h^{R,k^n+1} (resp. δ_h^{R}) in line 8 (resp. line 18) of Algorithm 2 yields

$$b_h^{\text{R},k^n+1} = \left(1 - \frac{1}{\eta_n}\right) B_h^{\text{R},k^n}(s_h^k, a_h^k) + \frac{1}{\eta_n} B_h^{\text{R},k^n+1}(s_h^k, a_h^k) + \frac{c_b}{n^{3/4}} H^2 \log^2 \frac{SAT}{\delta}. \quad (\text{C.33})$$

This taken collectively with the definition (4.2) of η_n^N allows us to expand

$$\begin{aligned} & \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} b_h^{\text{R},k^n+1} \\ &= \sum_{n=1}^{N_h^k} \eta_n \prod_{i=n+1}^{N_h^k} (1 - \eta_i) \left(\left(1 - \frac{1}{\eta_n}\right) B_h^{\text{R},k^n}(s_h^k, a_h^k) + \frac{1}{\eta_n} B_h^{\text{R},k^n+1}(s_h^k, a_h^k) \right) + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \log^2 \frac{SAT}{\delta} \\ &= \sum_{n=1}^{N_h^k} \prod_{i=n+1}^{N_h^k} (1 - \eta_i) \left(- (1 - \eta_n) B_h^{\text{R},k^n}(s_h^k, a_h^k) + B_h^{\text{R},k^n+1}(s_h^k, a_h^k) \right) + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \log^2 \frac{SAT}{\delta} \\ &= \sum_{n=1}^{N_h^k} \left(\prod_{i=n+1}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n+1}(s_h^k, a_h^k) - \prod_{i=n}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n}(s_h^k, a_h^k) \right) + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \log^2 \frac{SAT}{\delta} \\ &\stackrel{(i)}{=} \sum_{n=1}^{N_h^k} \prod_{i=n+1}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n+1}(s_h^k, a_h^k) - \sum_{n=2}^{N_h^k} \prod_{i=n}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n}(s_h^k, a_h^k) + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \log^2 \frac{SAT}{\delta} \\ &\stackrel{(ii)}{=} \sum_{n=1}^{N_h^k} \prod_{i=n+1}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n+1}(s_h^k, a_h^k) - \sum_{n=1}^{N_h^k-1} \prod_{i=n+1}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n+1}(s_h^k, a_h^k) + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \log^2 \frac{SAT}{\delta} \\ &= B_h^{\text{R},k^{N_h^k+1}}(s_h^k, a_h^k) + c_b \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} H^2 \log^2 \frac{SAT}{\delta}. \end{aligned} \quad (\text{C.34})$$

Here, (i) is valid due to the fact that $B_h^{\text{R},k^1}(s_h^k, a_h^k) = 0$; (ii) follows from the fact that

$$\begin{aligned} \sum_{n=2}^{N_h^k} \prod_{i=n}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n}(s_h^k, a_h^k) &= \sum_{n=1}^{N_h^k-1} \prod_{i=n+1}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^{n+1}}(s_h^k, a_h^k) \\ &= \sum_{n=1}^{N_h^k-1} \prod_{i=n+1}^{N_h^k} (1 - \eta_i) B_h^{\text{R},k^n+1}(s_h^k, a_h^k), \end{aligned}$$

where the first relation can be seen by replacing n with $n+1$, and the last relation holds true since the state-action pair (s_h^k, a_h^k) has not been visited at step h between the (k^n+1) -th episode and the $(k^{n+1}-1)$ -th episode. Combining the above identity (C.34) with the following property (see Lemma 1)

$$\frac{1}{(N_h^k)^{3/4}} \leq \sum_{n=1}^{N_h^k} \frac{\eta_n^{N_h^k}}{n^{3/4}} \leq \frac{2}{(N_h^k)^{3/4}},$$

we can immediately demonstrate that

$$B_h^{\text{R},k^{N_h^k}+1}(s_h^k, a_h^k) + c_b \frac{H^2 \log^2 \frac{SAT}{\delta}}{(N_h^k)^{3/4}} \leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} b_h^{\text{R},k^n+1} \leq B_h^{\text{R},k^{N_h^k}+1}(s_h^k, a_h^k) + 2c_b \frac{H^2 \log^2 \frac{SAT}{\delta}}{(N_h^k)^{3/4}}. \quad (\text{C.35})$$

Taking (C.32) and (C.35) collectively demonstrates that

$$\left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \xi_h^{k^n} \right| \leq B_h^{\text{R},k^{N_h^k}+1}(s_h^k, a_h^k) + c_b \frac{H^2 \log^2 \frac{SAT}{\delta}}{(N_h^k)^{3/4}} \leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} b_h^{\text{R},k^n+1} \quad (\text{C.36})$$

as claimed in (C.16). We have thus concluded the proof of Lemma 2 based on the argument in Step 2.

C.1.1 Proof of the inequality (C24). In order to establish the inequality (C.24), it suffices to look at the following term:

$$I_3 := \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \text{Var}_{h,s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}) - \sigma_h^{\text{adv},k^{N_h^k}+1}(s_h^k, a_h^k) + (\mu_h^{\text{adv},k^{N_h^k}+1}(s_h^k, a_h^k))^2, \quad (\text{C.37})$$

which forms the main content of this subsection.

First of all, the update rules of $\mu_h^{\text{adv},k^{n+1}}$ and $\sigma_h^{\text{adv},k^{n+1}}$ in lines 13-14 of Algorithm 2 tell us that

$$\begin{aligned} \mu_h^{\text{adv},k^{n+1}}(s_h^k, a_h^k) &= \mu_h^{\text{adv},k^n+1}(s_h^k, a_h^k) = (1 - \eta_n) \mu_h^{\text{adv},k^n}(s_h^k, a_h^k) + \eta_n (V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n})), \\ \sigma_h^{\text{adv},k^{n+1}}(s_h^k, a_h^k) &= \sigma_h^{\text{adv},k^n+1}(s_h^k, a_h^k) = (1 - \eta_n) \sigma_h^{\text{adv},k^n}(s_h^k, a_h^k) + \eta_n (V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}))^2. \end{aligned}$$

Applying this relation recursively and invoking the definitions of η_n^N (resp. P_h^k) in (4.2) (resp. (4.5)) give

$$\mu_h^{\text{adv},k^{N_h^k}+1}(s_h^k, a_h^k) \stackrel{\text{(i)}}{=} \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n})) = \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_h^{k^n} (V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}), \quad (\text{C.38a})$$

$$\sigma_h^{\text{adv},k^{N_h^k}+1}(s_h^k, a_h^k) \stackrel{\text{(ii)}}{=} \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}))^2 = \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_h^{k^n} (V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n})^2. \quad (\text{C.38b})$$

Recognizing that $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} = 1$ (see (4.3)), we can immediately apply Jensen's inequality to the expressions (i) and (ii) to yield

$$\sigma_h^{\text{adv},k^{N_h^k}+1}(s_h^k, a_h^k) \geq \left(\mu_h^{\text{adv},k^{N_h^k}+1}(s_h^k, a_h^k) \right)^2. \quad (\text{C.39})$$

Further, in view of the definition (4.6), we have

$$\text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}) = P_{h,s_h^k,a_h^k}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n})^2 - \left(P_{h,s_h^k,a_h^k}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}) \right)^2,$$

which allows one to decompose and bound I_3 as follows:

$$\begin{aligned} I_3 &= \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_{h,s_h^k,a_h^k}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n})^2 - \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_h^{k^n}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n})^2 \\ &\quad + \left(\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_h^{k^n}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}) \right)^2 - \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(P_{h,s_h^k,a_h^k}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}) \right)^2 \\ &\leq \underbrace{\left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (P_h^{k^n} - P_{h,s_h^k,a_h^k})(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n})^2 \right|}_{=:I_{3,1}} \\ &\quad + \underbrace{\left(\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_h^{k^n}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}) \right)^2 - \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(P_{h,s_h^k,a_h^k}(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n}) \right)^2}_{=:I_{3,2}}. \end{aligned} \quad (\text{C.40})$$

It then boils down to controlling the above two terms in (C.40) separately.

Step 1: bounding $I_{3,1}$. To upper bound the term $I_{3,1}$ in (C.40), we resort to Lemma 7 by setting

$$W_{h+1}^i := (V_{h+1}^i - V_{h+1}^{\text{R},i})^2 \quad \text{and} \quad u_h^i(s, a, N) := \eta_{N_h^i(s, a)}^N. \quad (\text{C.41})$$

It is easily seen that

$$\|W_{h+1}^i\|_\infty \leq \left(\|V_{h+1}^{\text{R},i}\|_\infty + \|V_{h+1}^i\|_\infty \right)^2 \leq 4H^2 =: C_w, \quad (\text{C.42})$$

and it follows from (C.21) that

$$\max_{N, h, s, a \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}} \eta_{N_h^i(s, a)}^N \leq \frac{2H}{N} =: C_u. \quad (\text{C.43})$$

Armed with the properties (C.42) and (C.43) and recalling (C.28), we can invoke Lemma 7 w.r.t. (C.41) and set $(N, s, a) = (N_h^k, s_h^k, a_h^k)$ to yield

$$\begin{aligned} I_{3,1} &= \left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (P_h^{k^n} - P_{h,s_h^k,a_h^k})(V_{h+1}^{k^n} - V_{h+1}^{\text{R},k^n})^2 \right| = \left| \sum_{i=1}^k X_i(s_h^k, a_h^k, h, N_h^k) \right| \\ &\lesssim \sqrt{C_u \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} u_h^{k^n}(s_h^k, a_h^k, N_h^k) \text{Var}_{h,s_h^k,a_h^k}(W_{h+1}^{k^n}) + \left(C_u C_w + \sqrt{\frac{C_u}{N}} C_w \right) \log^2 \frac{SAT}{\delta}} \end{aligned}$$

$$\begin{aligned}
&\lesssim \sqrt{\frac{H}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \text{Var}_{h, s_h^k, a_h^k} ((V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n})^2)} + \frac{H^3 \log^2 \frac{SAT}{\delta}}{N_h^k} \\
&\lesssim \sqrt{\frac{H^5}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^3}{N_h^k} \log^2 \frac{SAT}{\delta}
\end{aligned} \tag{C.44}$$

with probability at least $1 - \delta$. Here, the last inequality results from the fact $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \leq 1$ (see (4.3)) and the following trivial result:

$$\text{Var}_{h, s_h^k, a_h^k} ((V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n})^2) \leq \| (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n})^4 \|_{\infty} \leq 16H^4. \tag{C.45}$$

Step 2: bounding $I_{3,2}$. Jensen's inequality tells us that

$$\begin{aligned}
\left(\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_{h, s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right)^2 &= \left(\sum_{n=1}^{N_h^k} (\eta_n^{N_h^k})^{1/2} \cdot (\eta_n^{N_h^k})^{1/2} P_{h, s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right)^2 \\
&\leq \left\{ \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \right\} \left\{ \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(P_{h, s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right)^2 \right\} \\
&\leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(P_{h, s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right)^2,
\end{aligned}$$

where the last line arises from (4.3). Substitution into $I_{3,2}$ (cf. (C.40)) gives

$$\begin{aligned}
I_{3,2} &\leq \left(\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_h^{k^n} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right)^2 - \left(\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} P_{h, s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right)^2 \\
&= \left\{ \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (P_h^{k^n} - P_{h, s_h^k, a_h^k}) (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right\} \left\{ \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (P_h^{k^n} + P_{h, s_h^k, a_h^k}) (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right\}. \tag{C.46}
\end{aligned}$$

In what follows, we would like to use this relation to show that

$$I_{3,2} \leq C_{32} \left\{ \sqrt{\frac{H^5}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^3}{N_h^k} \log^2 \frac{SAT}{\delta} \right\} \tag{C.47}$$

for some universal constant $C_{32} > 0$.

If $I_{3,2} \leq 0$, then (C.47) holds true trivially. Consequently, it is sufficient to study the case where $I_{3,2} > 0$. To this end, we first note that the term in the first pair of curly brackets of (C.46) is exactly I_1

(see (C18a)), which can be bounded by recalling (C.23):

$$\begin{aligned}
|I_1| &\lesssim \sqrt{\frac{H}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \text{Var}_{h,s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n})} + \frac{H^2 \log^2 \frac{SAT}{\delta}}{N_h^k} \\
&\lesssim \sqrt{\frac{H^3}{N_h^k} \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} + \frac{H^2 \log^2 \frac{SAT}{\delta}}{N_h^k}} \\
&\lesssim \sqrt{\frac{H^3}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^2}{N_h^k} \log^2 \frac{SAT}{\delta}, \tag{C.48}
\end{aligned}$$

with probability at least $1 - \delta$. Here, the second inequality arises from the following property:

$$\text{Var}_{h,s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \leq \| (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n})^2 \|_\infty \leq 4H^2, \tag{C.49}$$

whereas the last inequality (C.48) holds as a result of the fact $\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \leq 1$ (see (4.3)).

Moreover, the term in the second pair of curly brackets of (C.46) can be bounded straightforwardly as follows:

$$\begin{aligned}
&\left| \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (P_h^{k^n} + P_{h,s_h^k, a_h^k}) (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right| \\
&\leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} (\|P_h^{k^n}\|_1 + \|P_{h,s_h^k, a_h^k}\|_1) \|V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}\|_\infty \leq 2H, \tag{C.50}
\end{aligned}$$

where we have used the property (4.3), as well as the elementary facts $\|V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}\|_\infty \leq H$ and $\|P_h^{k^n}\|_1 = \|P_{h,s_h^k, a_h^k}\|_1 = 1$. Substituting the above two results (C.48) and (C.50) back into (C.46), we arrive at the bound (C.47) as long as $I_{3,2} > 0$. Putting all cases together, we have established the claim (C.47).

Step 3: putting all this together. To finish up, plugging the bounds (C.44) and (C.47) into (C.40), we can conclude that

$$I_3 \leq I_{3,1} + I_{3,2} \leq C_3 \left\{ \sqrt{\frac{H^5}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^3}{N_h^k} \log^2 \frac{SAT}{\delta} \right\}$$

for some constant $C_3 > 0$. This together with the definition (C.37) of I_3 results in

$$\begin{aligned}
&\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \text{Var}_{h,s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \\
&\leq \left\{ \sigma_h^{\text{adv}, k^{N_h^k}+1} (s_h^k, a_h^k) - (\mu_h^{\text{adv}, k^{N_h^k}+1} (s_h^k, a_h^k))^2 \right\} + C_3 \left(\sqrt{\frac{H^5}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^3}{N_h^k} \log^2 \frac{SAT}{\delta} \right),
\end{aligned}$$

which combined with the elementary inequality $\sqrt{u+v} \leq \sqrt{u} + \sqrt{v}$ for any $u, v \geq 0$ and (C.39) yields

$$\begin{aligned} & \left\{ \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \text{Var}_{h, s_h^k, a_h^k} (V_{h+1}^{k^n} - V_{h+1}^{\text{R}, k^n}) \right\}^{1/2} \\ & \lesssim \left\{ \sigma_h^{\text{adv}, k^{N_h^k}+1}(s_h^k, a_h^k) - (\mu_h^{\text{adv}, k^{N_h^k}+1}(s_h^k, a_h^k))^2 \right\}^{1/2} + \frac{H^{5/4}}{(N_h^k)^{1/4}} \log^{1/2} \frac{\text{SAT}}{\delta} + \frac{H^{3/2}}{(N_h^k)^{1/2}} \log \frac{\text{SAT}}{\delta}. \end{aligned}$$

Substitution into (C.23) establishes the desired result (C.24).

C.1.2 Proof of the inequality (C31) In order to prove the inequality (C.31), it suffices to look at the following term:

$$I_4 := \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} \text{Var}_{h, s_h^k, a_h^k} (V_{h+1}^{\text{R}, k^n}) - \left(\sigma_h^{\text{ref}, k^{N_h^k}+1}(s_h^k, a_h^k) - (\mu_h^{\text{ref}, k^{N_h^k}+1}(s_h^k, a_h^k))^2 \right). \quad (\text{C.51})$$

In view of the update rules of $\mu_h^{\text{ref}, k^{n+1}}$ and $\sigma_h^{\text{ref}, k^{n+1}}$ in lines 11–12 of Algorithm 2, we have

$$\begin{aligned} \mu_h^{\text{ref}, k^{n+1}}(s_h^k, a_h^k) &= \mu_h^{\text{ref}, k^n+1}(s_h^k, a_h^k) = \left(1 - \frac{1}{n}\right) \mu_h^{\text{ref}, k^n}(s_h^k, a_h^k) + \frac{1}{n} V_{h+1}^{\text{R}, k^n}(s_{h+1}^n), \\ \sigma_h^{\text{ref}, k^{n+1}}(s_h^k, a_h^k) &= \sigma_h^{\text{ref}, k^n+1}(s_h^k, a_h^k) = \left(1 - \frac{1}{n}\right) \sigma_h^{\text{ref}, k^n}(s_h^k, a_h^k) + \frac{1}{n} (V_{h+1}^{\text{R}, k^n}(s_{h+1}^n))^2, \end{aligned}$$

Through simple recursion, these identities together with the definition (4.5) of P_h^k lead to

$$\mu_h^{\text{ref}, k^{N_h^k}+1}(s_h^k, a_h^k) \stackrel{\text{(i)}}{=} \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^{\text{R}, k^n}(s_{h+1}^n) = \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_h^{k^n} V_{h+1}^{\text{R}, k^n}, \quad (\text{C.52a})$$

$$\sigma_h^{\text{ref}, k^{N_h^k}+1}(s_h^k, a_h^k) \stackrel{\text{(ii)}}{=} \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (V_{h+1}^{\text{R}, k^n}(s_{h+1}^n))^2 = \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_h^{k^n} (V_{h+1}^{\text{R}, k^n})^2, \quad (\text{C.52b})$$

The expressions (i) and (ii) combined with Jensen's inequality give

$$\sigma_h^{\text{ref}, k^{N_h^k}+1}(s_h^k, a_h^k) \geq \left(\mu_h^{\text{ref}, k^{N_h^k}+1}(s_h^k, a_h^k) \right)^2. \quad (\text{C.53})$$

Taking these together with the definition

$$\text{Var}_{h, s_h^k, a_h^k}(V_{h+1}^{\text{R}, k^n}) = P_{h, s_h^k, a_h^k} (V_{h+1}^{\text{R}, k^n})^2 - (P_{h, s_h^k, a_h^k} V_{h+1}^{\text{R}, k^n})^2,$$

we obtain

$$\begin{aligned}
I_4 &= \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} \left(P_{h,s_h^k, a_h^k} (V_{h+1}^{\text{R}, k^n})^2 - (P_{h,s_h^k, a_h^k} V_{h+1}^{\text{R}, k^n})^2 \right) - \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_h^{k^n} (V_{h+1}^{\text{R}, k^n})^2 + \left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_h^{k^n} V_{h+1}^{\text{R}, k^n} \right)^2 \\
&= \underbrace{\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_{h,s_h^k, a_h^k} - P_h^{k^n}) (V_{h+1}^{\text{R}, k^n})^2}_{=: I_{4,1}} + \underbrace{\left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_h^{k^n} V_{h+1}^{\text{R}, k^n} \right)^2 - \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_{h,s_h^k, a_h^k} V_{h+1}^{\text{R}, k^n})^2}_{=: I_{4,2}}. \quad (\text{C.54})
\end{aligned}$$

In what follows, we shall bound the terms $I_{4,1}$ and $I_{4,2}$ in (C.54) separately.

Step 1: bounding $I_{4,1}$. The first term $I_{4,1}$ in (C.54) can be bounded by means of Lemma 7 in an almost identical fashion as $I_{3,1}$ in (C.44). Specifically, let us set

$$W_{h+1}^i := (V_{h+1}^{\text{R}, i})^2 \quad \text{and} \quad u_h^i(s, a, N) := \frac{1}{N},$$

which clearly obey

$$|u_h^i(s, a, N)| = \frac{1}{N} =: C_u \quad \text{and} \quad \|W_{h+1}^i\|_\infty \leq H^2 =: C_w.$$

It is easily verified that

$$\sum_{n=1}^N u_h^{k^n(s,a)}(s, a, N) = \sum_{n=1}^N \frac{1}{N} = 1$$

holds for all $(N, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$. Hence we can take $(N, s, a) = (N_h^k, s_h^k, a_h^k)$ and apply Lemma 7 to yield

$$\begin{aligned}
|I_{4,1}| &= \left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_h^{k^n} - P_{h,s_h^k, a_h^k}) (V_{h+1}^{\text{R}, k^n})^2 \right| = \left| \sum_{i=1}^k X_i(s_h^k, a_h^k, h, N_h^k) \right| \\
&\lesssim \sqrt{C_u \log^2 \frac{SAT}{\delta}} \sqrt{\sum_{n=1}^{N_h^k} u_h^{k^n}(s_h^k, a_h^k, N_h^k) \text{Var}_{h,s_h^k, a_h^k}(W_{h+1}^{k^n})} + \left(C_u C_w + \sqrt{\frac{C_u}{N}} C_w \right) \log^2 \frac{SAT}{\delta} \\
&\vee \sqrt{\frac{H^4 \log^2 \frac{SAT}{\delta}}{N_h^k} + \frac{H^2 \log^2 \frac{SAT}{\delta}}{N_h^k}}
\end{aligned} \quad (\text{C.55})$$

with probability at least $1 - \delta$, where the last inequality results from the fact that

$$\text{Var}_{h,s_h^k, a_h^k}(W_{h+1}^{k^n}) \leq \|W_{h+1}^{k^n}\|_\infty^2 \leq C_w^2 = H^4.$$

Step 2: bounding $I_{4,2}$. We now turn to the other term $I_{4,2}$ defined in (C.54). Toward this, we first make the observation that

$$\left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_{h,s_h^k, a_h^k} V_{h+1}^{R,k^n} \right)^2 \leq \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} \left(P_{h,s_h^k, a_h^k} V_{h+1}^{R,k^n} \right)^2, \quad (\text{C.56})$$

which follows from Jensen's inequality. Equipped with this relation, we can upper bound $I_{4,2}$ as follows:

$$\begin{aligned} I_{4,2} &\leq \left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_h^{k^n} V_{h+1}^{R,k^n} \right)^2 - \left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} P_{h,s_h^k, a_h^k} V_{h+1}^{R,k^n} \right)^2 \\ &= \left\{ \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_h^{k^n} - P_{h,s_h^k, a_h^k}) V_{h+1}^{R,k^n} \right\} \left\{ \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_h^{k^n} + P_{h,s_h^k, a_h^k}) V_{h+1}^{R,k^n} \right\}. \end{aligned} \quad (\text{C.57})$$

In the following text, we would like to apply this relation to prove

$$I_{4,2} \leq C_{42} \left(\sqrt{\frac{H^4}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^2}{N_h^k} \log^2 \frac{SAT}{\delta} \right) \quad (\text{C.58})$$

for some constant $C_{42} > 0$.

When $I_{4,2} \leq 0$, the claim (C.58) holds trivially. As a result, we shall focus on the case where $I_{4,2} > 0$. Let us begin with the term in the first pair of curly brackets of (C.57). Toward this, let us abuse the notation and set

$$W_{h+1}^i := V_{h+1}^{R,i} \quad \text{and} \quad u_h^i(s, a, N) := \frac{1}{N},$$

which satisfy

$$|u_h^i(s, a, N)| = \frac{1}{N} =: C_u \quad \text{and} \quad \|W_{h+1}^i\|_\infty \leq H =: C_w.$$

Akin to our argument for bounding $I_{4,1}$, invoking Lemma 7 and setting $(N, s, a) = (N_h^k, s_h^k, a_h^k)$ imply that

$$\left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_h^{k^n} - P_{h,s_h^k, a_h^k}) V_{h+1}^{R,k^n} \right| \lesssim \sqrt{\frac{H^2 \log^2 \frac{SAT}{\delta}}{N_h^k} + \frac{H \log^2 \frac{SAT}{\delta}}{N_h^k}}$$

with probability at least $1 - \delta$. In addition, the term in the second pair of curly brackets of (C.57) can be bounded straightforwardly by

$$\left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_h^{k^n} + P_{h,s_h^k, a_h^k}) V_{h+1}^{R,k^n} \right| \leq \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (\|P_h^{k^n}\|_1 + \|P_{h,s_h^k, a_h^k}\|_1) \|V_{h+1}^{R,k^n}\|_\infty \leq 2H,$$

where we have used $\|V_{h+1}^{R,k^n}\|_\infty \leq H$ and $\|P_h^{k^n}\|_1 = \|P_{h,s_h^k, a_h^k}\|_1 = 1$. Substituting the preceding facts into (C.57) validates the bound (C.58) as long as $I_{4,2} > 0$. We have thus finished the proof of the claim (C.58).

Step 3: putting all pieces together. Combining the results (C.55) and (C.58) with (C.54) yields

$$I_4 \leq |I_{4,1}| + I_{4,2} \leq C_4 \left\{ \sqrt{\frac{H^4}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^2}{N_h^k} \log^2 \frac{SAT}{\delta} \right\}$$

for some constant $C_4 > 0$. This bound taken together with the definition (C.51) of I_4 gives

$$\begin{aligned} \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} \text{Var}_{h,s_h^k, a_h^k}(V_{h+1}^{R,k^n}) &\leq \left\{ \sigma_h^{\text{ref}, k^N_h + 1}(s_h^k, a_h^k) - (\mu_h^{\text{ref}, k^N_h + 1}(s_h^k, a_h^k))^2 \right\} \\ &\quad + C_4 \left\{ \sqrt{\frac{H^4}{N_h^k} \log^2 \frac{SAT}{\delta}} + \frac{H^2}{N_h^k} \log^2 \frac{SAT}{\delta} \right\}. \end{aligned}$$

Invoke the elementary inequality $\sqrt{u+v} \leq \sqrt{u} + \sqrt{v}$ for any $u, v \geq 0$ and use the property (C.53) to obtain

$$\begin{aligned} &\left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} \text{Var}_{h,s_h^k, a_h^k}(V_{h+1}^{R,k^n}) \right)^{1/2} \\ &\lesssim \left\{ \sigma_h^{\text{ref}, k^N_h + 1}(s_h^k, a_h^k) - (\mu_h^{\text{ref}, k^N_h + 1}(s_h^k, a_h^k))^2 \right\}^{1/2} + \frac{H}{(N_h^k)^{1/4}} \log^{1/2} \frac{SAT}{\delta} + \frac{H}{(N_h^k)^{1/2}} \log \frac{SAT}{\delta}. \end{aligned}$$

Substitution into (C.30) directly establishes the desired result (C.31).

C.2 Proof of Lemma 3

C.2.1 Proof of the inequalities (4.11) Suppose that we can verify the following inequality:

$$Q_h^{\text{LCB},k}(s, a) \leq Q_h^*(s, a) \quad \text{for all } (s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H], \quad (\text{C.59})$$

which in turn yields

$$\max_a Q_h^{\text{LCB},k}(s, a) \leq \max_a Q_h^*(s, a) = V_h^*(s) \quad \text{for all } (k, h, s) \in [K] \times [H] \times \mathcal{S}. \quad (\text{C.60})$$

In addition, the construction of $V_h^{\text{LCB},k}$ (see line 14 of Algorithm 3) allows us to show that

$$V_h^{\text{LCB},k+1}(s) \leq \max \left\{ \max_{j:j \leq k+1} \max_a Q_h^{\text{LCB},j}(s, a), \max_{j:j \leq k} V_h^{\text{LCB},j}(s) \right\}.$$

This taken together with the initialization $V_h^{\text{LCB},1} = 0$ and a simple induction argument yields

$$V_h^{\text{LCB},k}(s) \leq V_h^*(s) \quad \text{for all } (k, h, s) \in [K] \times [H] \times \mathcal{S}. \quad (\text{C.61})$$

As a consequence, everything comes down to proving the claim (C.59), which we shall accomplish by induction.

Base case. Given our initialization, we have

$$Q_h^{\text{LCB},1}(s, a) - Q_h^*(s, a) = 0 - Q_h^*(s, a) \leq 0,$$

and hence the claim (C.59) holds trivially when $k = 1$.

Induction step. Suppose now that the claim (C.59) holds all the way up to k for all (s, a, h) , and we would like to validate it for the $(k + 1)$ -th episode as well. Toward this end, recall that the state-action pair (s_h^k, a_h^k) is visited in the k -th episode at time step h ; this means that $Q_h^{\text{LCB}}(s_h^k, a_h^k)$ is updated once we collect samples in the k -th episode, with all other entries Q_h^{LCB} frozen. It thus suffices to verify that

$$Q_h^{\text{LCB},k+1}(s_h^k, a_h^k) \leq Q_h^*(s_h^k, a_h^k).$$

In what follows, we shall adopt the short-hand notation (see also Section 4.2)

$$N_h^k = N_h^k(s_h^k, a_h^k) \quad \text{and} \quad k^n = k_h^n(s_h^k, a_h^k)$$

which will be used throughout this subsection as long as it is clear from the context.

The update rule of $Q_h^{\text{LCB},k}$ (cf. line 4 of Algorithm 2) and the Bellman optimality equation in (C.8) tell us the following identities:

$$\begin{aligned} Q_h^{\text{LCB},k+1}(s_h^k, a_h^k) &= Q_h^{\text{LCB},k^{N_h^k}+1}(s_h^k, a_h^k) \\ &= (1 - \eta_{N_h^k})Q_h^{\text{LCB},k^{N_h^k}}(s_h^k, a_h^k) + \eta_{N_h^k}\left(r_h(s_h^k, a_h^k) + V_{h+1}^{\text{LCB},k^{N_h^k}}(s_{h+1}^{k^{N_h^k}}) - b_h^{k^{N_h^k}}\right), \\ Q_h^*(s_h^k, a_h^k) &= (1 - \eta_{N_h^k})Q_h^*(s_h^k, a_h^k) + \eta_{N_h^k}Q_h^*(s_h^k, a_h^k) \\ &= (1 - \eta_{N_h^k})Q_h^*(s_h^k, a_h^k) + \eta_{N_h^k}\left(r(s_h^k, a_h^k) + P_{h,s_h^k,a_h^k}V_{h+1}^*\right), \end{aligned}$$

which taken collectively lead to the following identity

$$\begin{aligned} Q_h^{\text{LCB},k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) &= Q_h^{\text{LCB},k^{N_h^k}+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \\ &= (1 - \eta_{N_h^k})\left(Q_h^{\text{LCB},k^{N_h^k}}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)\right) + \eta_{N_h^k}\left(V_{h+1}^{\text{LCB},k^{N_h^k}}(s_{h+1}^{k^{N_h^k}}) - P_{h,s_h^k,a_h^k}V_{h+1}^* - b_h^{k^{N_h^k}}\right) \\ &= (1 - \eta_{N_h^k})\left(Q_h^{\text{LCB},k^{N_h^k-1}+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)\right) + \eta_{N_h^k}\left(V_{h+1}^{\text{LCB},k^{N_h^k}}(s_{h+1}^{k^{N_h^k}}) - P_{h,s_h^k,a_h^k}V_{h+1}^* - b_h^{k^{N_h^k}}\right). \end{aligned}$$

Recall the definitions of η_0^N and η_n^N in (4.2). Applying the above relation recursively and using the decomposition of $Q_h^*(s_h^k, a_h^k)$ in (C.6) result in

$$\begin{aligned} Q_h^{\text{LCB},k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) &= \eta_0^{N_h^k}\left(Q_h^{\text{LCB},1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k)\right) + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(V_{h+1}^{\text{LCB},k^n}(s_{h+1}^{k^n}) - P_{h,s_h^k,a_h^k}V_{h+1}^* - b_h^{k^n}\right) \\ &\leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(V_{h+1}^{\text{LCB},k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) + (P_h^{k^n} - P_{h,s_h^k,a_h^k})V_{h+1}^* - b_h^{k^n}\right), \end{aligned} \tag{C.62}$$

where the inequality follows from the initialization $Q_h^{\text{LCB},1}(s_h^k, a_h^k) = 0 \leq Q_h^*(s_h^k, a_h^k)$ and the definition of P_h^k in (4.5). To continue, we invoke a result established in (29, proof of Lemma 4.3), which guarantees

that with probability at least $1 - \delta$,

$$\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(P_h^{k^n} - P_{h, s_h^k, a_h^k} \right) V_{h+1}^* \lesssim \sqrt{\frac{H^3 \log(\frac{SAT}{\delta})}{N_h^k}} \leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} b_h^{k^n},$$

provided that c_b is some sufficiently large constant. Substituting the above relation into (C.62) implies that

$$Q_h^{\text{LCB}, k+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \leq \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(V_{h+1}^{\text{LCB}, k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) \right) \leq 0, \quad (\text{C.63})$$

where the last inequality follows from the induction hypothesis

$$V_{h+1}^{\text{LCB}, j}(s) \leq V_{h+1}^*(s) \quad \text{for all } s \in \mathcal{S} \text{ and } j \leq k.$$

The proof is thus completed by induction.

C.2.2 Proof of the inequality (4.12) The proof of (4.12) essentially follows the same arguments of (70, Lemma 4.2) (see also (30, Lemma C.7)), an algebraic result leveraging certain relations w.r.t. the Q-value estimates. Accounting for the difference between our algorithm and the one in [70], we paraphrase (70, Lemma 4.2) into the following form that is convenient for our purpose.

LEMMA 9. paraphrased from Lemma 4.2 in [70] Assume that there exists a constant $c_b > 0$ such that for all $(s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]$, it holds that

$$\begin{aligned} 0 &\leq Q_h^{k+1}(s, a) - Q_h^{\text{LCB}, k+1}(s, a) \\ &\leq \eta_0^{N_h^k(s, a)} H + \sum_{n=1}^{N_h^k(s, a)} \eta_n^{N_h^k(s, a)} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\text{LCB}, k^n}(s_{h+1}^{k^n}) \right) + 4c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k(s, a)}}. \end{aligned} \quad (\text{C.64})$$

Consider any $\varepsilon \in (0, H]$. Then for all $\beta = 1, \dots, \lceil \log_2 \frac{H}{\varepsilon} \rceil$, one has

$$\left| \sum_{h=1}^H \sum_{k=1}^K \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB}, k}(s_h^k, a_h^k) \in [2^{\beta-1} \varepsilon, 2^\beta \varepsilon) \right) \right| \lesssim \frac{H^6 S A \log \frac{SAT}{\delta}}{4^\beta \varepsilon^2}. \quad (\text{C.65})$$

We first show how to justify (4.12) if the inequality (C.65) holds. As can be seen, the fact (C.65) immediately leads to

$$\sum_{h=1}^H \sum_{k=1}^K \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB}, k}(s_h^k, a_h^k) > \varepsilon \right) \lesssim \sum_{\beta=1}^{\lceil \log_2 \frac{H}{\varepsilon} \rceil} \frac{H^6 S A \log \frac{SAT}{\delta}}{4^\beta \varepsilon^2} \leq \frac{H^6 S A \log \frac{SAT}{\delta}}{2 \varepsilon^2} \quad (\text{C.66})$$

as desired.

We now return to justify the claim (C.65), toward which it suffices to demonstrate that (C.64) holds. Lemma 2 and Lemma 3 directly verify the left-hand side of (C.64) since

$$Q_h^k(s, a) \geq Q_h^*(s, a) \geq Q_h^{\text{LCB}, k}(s, a) \quad \text{for all } (s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]. \quad (\text{C.67})$$

The remainder of the proof is thus devoted to justifying the upper bound on $Q_h^{k+1}(s, a) - Q_h^{\text{LCB}, k+1}(s, a)$ in (C.64). In view of the update rule in line 12 of Algorithm 3, we have the following basic fact:

$$Q_h^{k+1}(s, a) \leq Q_h^{\text{UCB}, k+1}(s, a).$$

This enables us to obtain

$$Q_h^{k+1}(s, a) - Q_h^{\text{LCB}, k+1}(s, a) \leq Q_h^{\text{UCB}, k+1}(s, a) - Q_h^{\text{LCB}, k+1}(s, a) = Q_h^{\text{UCB}, k^N_h+1}(s, a) - Q_h^{\text{LCB}, k^N_h+1}(s, a), \quad (\text{C.68})$$

where we abbreviate

$$N_h^k = N_h^k(s, a)$$

throughout this subsection as long as it is clear from the context. Using the update rules of $Q_h^{\text{UCB}, k}$ and $Q_h^{\text{LCB}, k}$ in line 2 and line 4 of Algorithm 2, we reach

$$\begin{aligned} & Q_h^{\text{UCB}, k^N_h+1}(s, a) - Q_h^{\text{LCB}, k^N_h+1}(s, a) \\ &= (1 - \eta_{N_h^k})Q_h^{\text{UCB}, k^N_h}(s, a) + \eta_{N_h^k} \left(r_h(s, a) + V_{h+1}^{k^N_h}(s_{h+1}^{k^N_h}) + c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k}} \right) \\ &\quad - (1 - \eta_{N_h^k})Q_h^{\text{LCB}, k^N_h}(s, a) - \eta_{N_h^k} \left(r_h(s, a) + V_{h+1}^{\text{LCB}, k^N_h}(s_{h+1}^{k^N_h}) - c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k}} \right) \\ &= (1 - \eta_{N_h^k}) \left(Q_h^{\text{UCB}, k^N_h}(s, a) - Q_h^{\text{LCB}, k^N_h}(s, a) \right) \\ &\quad + \eta_{N_h^k} \left(V_{h+1}^{k^N_h}(s_{h+1}^{k^N_h}) - V_{h+1}^{\text{LCB}, k^N_h}(s_{h+1}^{k^N_h}) + 2c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k}} \right) \\ &= (1 - \eta_{N_h^k}) \left(Q_h^{\text{UCB}, k^{N_h-1}+1}(s, a) - Q_h^{\text{LCB}, k^N_h}(s, a) \right) \\ &\quad + \eta_{N_h^k} \left(V_{h+1}^{k^N_h}(s_{h+1}^{k^N_h}) - V_{h+1}^{\text{LCB}, k^N_h}(s_{h+1}^{k^N_h}) + 2c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k}} \right). \end{aligned}$$

Applying this relation recursively leads to the desired result

$$\begin{aligned}
& Q_h^{\text{UCB},k^h+1}(s, a) - Q_h^{\text{LCB},k^h+1}(s, a) \\
&= \eta_0^{N_h^k} \left(Q_h^{\text{UCB},1}(s, a) - Q_h^{\text{LCB},1}(s, a) \right) + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\text{LCB},k^n}(s_{h+1}^{k^n}) + 2c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{n}} \right) \\
&\leq \eta_0^{N_h^k} H + \sum_{n=1}^{N_h^k} \eta_n^{N_h^k} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\text{LCB},k^n}(s_{h+1}^{k^n}) \right) + 4c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k}}.
\end{aligned}$$

Here, the last line is valid due to the property $0 \leq Q_h^{\text{LCB},1}(s, a) \leq Q_h^{\text{UCB},1}(s, a) \leq H$ and the following fact:

$$\sum_{n=1}^{N_h^k} \eta_n^{N_h^k} c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k}} \leq 2c_b \sqrt{\frac{H^3 \log \frac{SAT}{\delta}}{N_h^k}},$$

which is an immediate consequence of the elementary property $\sum_{n=1}^N \frac{\eta_n^N}{\sqrt{n}} \leq \frac{2}{\sqrt{N}}$ (see Lemma 1). This combined with (C.68) establishes the condition (C.64), thus concluding the proof of the inequality (4.12).

C.3 Proof of Lemma 4

C.3.1 Proof of the inequality (4.15). Consider any state s that has been visited at least once during the K episodes. Throughout this proof, we shall adopt the shorthand notation

$$k^i = k_h^i(s),$$

which denotes the index of the episode in which state s is visited for the i -th time at step h . Given that $V_h(s)$ and $V_h^R(s)$ are only updated during the episodes with indices coming from $\{i \mid 1 \leq k^i \leq K\}$, it suffices to show that for any s and the corresponding $1 \leq k^i \leq K$, the claim (4.15) holds in the sense that

$$|V_h^{k^i+1}(s) - V_h^{R,k^i+1}(s)| \leq 2. \quad (\text{C.69})$$

Toward this end, we look at three scenarios separately.

Case 1. Suppose that k^i obeys

$$V_h^{k^i+1}(s) - V_h^{\text{LCB},k^i+1}(s) > 1 \quad (\text{C.70})$$

or

$$V_h^{k^i+1}(s) - V_h^{\text{LCB},k^i+1}(s) \leq 1 \quad \text{and} \quad u_{\text{ref}}^{k^i}(s) = \text{True} \quad (\text{C.71})$$

The above conditions correspond to the ones in line 15 and line 17 of Algorithm 3 (meaning that V_h^R is updated during the k^i -th episode), thus resulting in

$$V_h^{k^i+1}(s) = V_h^{R,k^i+1}(s).$$

This clearly satisfies (C.69).

Case 2. Suppose that k^{i_0} is the first time such that (C.70) and (C.71) are violated, namely,

$$i_0 := \min \left\{ j \mid V_h^{k^j+1}(s) - V_h^{\text{LCB}, k^j+1}(s) \leq 1 \text{ and } u_{\text{ref}}^{k^j}(s) = \text{False} \right\}. \quad (\text{C.72})$$

We make three observations.

- The definition (C.72) taken together with the update rules (lines 15–18 of Algorithm 3) reveals that V_h^R has been updated in the k^{i_0-1} -th episode, thus indicating that

$$V_h^{R, k^{i_0}}(s) = V_h^{R, k^{i_0-1}+1}(s) = V_h^{k^{i_0-1}+1}(s) = V_h^{k^{i_0}}(s). \quad (\text{C.73})$$

- Additionally, note that under the definition (C.72), $V_h^R(s)$ is not updated during the k^{i_0} -th episode, namely,

$$V_h^{R, k^{i_0}+1}(s) = V_h^{R, k^{i_0}}(s). \quad (\text{C.74})$$

- The definition of k^{i_0} indicates that either (C.70) or (C.71) is satisfied in the previous episode $k^i = k^{i_0-1}$ in which s was visited. If (C.70) is satisfied, then lines 15–16 in Algorithm 3 tell us that

$$\text{True} = u_{\text{ref}}^{k^{i_0-1}+1}(s) = u_{\text{ref}}^{k^{i_0}}(s), \quad (\text{C.75})$$

which, however, contradicts the assumption $u_{\text{ref}}^{k^{i_0}}(s) = \text{False}$ in (C.72). Therefore, in the k^{i_0-1} -th episode, (C.71) is satisfied, thus leading to

$$V_h^{k^{i_0}}(s) - V_h^{\text{LCB}, k^{i_0}}(s) = V_h^{k^{i_0-1}+1}(s) - V_h^{\text{LCB}, k^{i_0-1}+1}(s) \leq 1. \quad (\text{C.76})$$

We see from (C.73), (C.74) and (C.76) that

$$V_h^{R, k^{i_0}+1}(s) - V_h^{k^{i_0}+1}(s) = V_h^{R, k^{i_0}}(s) - V_h^{k^{i_0}+1}(s) = V_h^{k^{i_0}}(s) - V_h^{k^{i_0}+1}(s) \quad (\text{C.77})$$

$$\stackrel{\text{(i)}}{\leq} V_h^{k^{i_0}}(s) - V_h^{\text{LCB}, k^{i_0}}(s) \stackrel{\text{(ii)}}{\leq} 1, \quad (\text{C.78})$$

where (i) holds since $V_h^{k^{i_0}+1}(s) \geq V_h^*(s) \geq V_h^{\text{LCB}, k^{i_0}}(s)$, and (ii) follows from (C.76). In addition, we make note of the fact that

$$V_h^{R, k^{i_0}+1}(s) - V_h^{k^{i_0}+1}(s) = V_h^{k^{i_0}}(s) - V_h^{k^{i_0}+1}(s) \geq 0, \quad (\text{C.79})$$

which follows from (C.77) and the monotonicity of $V_h^k(s)$ in k . With the above results in place, we arrive at the advertised bound (C.69) when $i = i_0$.

Case 3. Consider any $i > i_0$. It is easily verified that

$$V_h^{k^i+1}(s) - V_h^{\text{LCB}, k^i+1}(s) \leq 1 \quad \text{and} \quad u_{\text{ref}}^{k^i}(s) = \text{False}. \quad (\text{C.80})$$

It then follows that

$$\begin{aligned} V_h^{R,k+1}(s) &\stackrel{(i)}{\leq} V_h^{R,k^i+1}(s) \stackrel{(ii)}{\leq} V_h^{k^i+1}(s) + 1 \stackrel{(iii)}{\leq} V_h^{LCB,k^i+1}(s) + 2 \\ &\stackrel{(iv)}{\leq} V_h^*(s) + 2 \stackrel{(v)}{\leq} V_h^{k^i+1}(s) + 2. \end{aligned} \quad (C.81)$$

Here, (i) holds due to the monotonicity of V_h^R and V_h^k (see line 14 of Algorithm 3), (ii) is a consequence of (C.78), (iii) comes from the definition (C.72) of i_0 , (iv) arises since V_h^{LCB} is a lower bound on V_h^* (see Lemma 3) and (v) is valid since $V_h^{k^i+1}(s) \geq V_h^*(s)$ (see Lemma 2). In addition, in view of the monotonicity of V_h^k (see line 14 of Algorithm 3) and the update rule in line 16 of Algorithm 3, we know that

$$V_h^{R,k+1}(s) \geq V_h^{k^i+1}(s).$$

The preceding two bounds taken collectively demonstrate that

$$0 \leq V_h^{R,k+1}(s) - V_h^{k^i+1}(s) \leq 2,$$

thus justifying (C.69) for this case.

Therefore, we have established (C.69)—and hence (4.15)—for all cases.

C.3.2 Proof of the inequality (4.16)

Suppose that

$$V^R, k_h(s_h^k) - V^R, k_h(s_h^k) \neq 0 \quad (C.82)$$

holds for some $k < K$. Then there are two possible scenarios to look at:

(a) *Case 1: the condition in line 15 and line 17 of Algorithm 3 are violated at step h of the k -th episode.* This means that we have

$$V_h^{k+1}(s_h^k) - V_h^{LCB,k+1}(s_h^k) \leq 1 \quad \text{and} \quad u_{\text{ref}}^k(s_h^k) = \text{False} \quad (C.83)$$

in this case. Then for any $k' > k$, one necessarily has

$$\begin{cases} V_h^{k'}(s_h^k) - V_h^{LCB,k'}(s_h^k) \leq V_h^{k+1}(s_h^k) - V_h^{LCB,k+1}(s_h^k) \leq 1, \\ u_{\text{ref}}^{k'}(s_h^k) = u_{\text{ref}}^k(s_h^k) = \text{False}, \end{cases} \quad (C.84)$$

where the first property uses the monotonicity of V_h^k and $V_h^{LCB,k}$ (see (4.7b) and line 14 of Algorithm 3). In turn, Condition (C.84) implies that V_h^R will no longer be updated after the k -th episode (see line 15 of Algorithm 3), thus indicating that

$$V^R, k_h(s_h^k) = V_h^{R,k+1}(s_h^k) = \dots = V^R, k_h(s_h^k). \quad (C.85)$$

This, however, contradicts the assumption (C.82).

(b) *Case 2: the condition in either line 15 or line 17 of Algorithm 3 is satisfied at step h of the k -th episode.* If this occurs, then the update rule in line 15 of Algorithm 3 implies that

$$V_h^{k+1}(s_h^k) - V_h^{LCB,k+1}(s_h^k) > 1, \quad (C.86)$$

or

$$V_h^{k+1}(s_h^k) - V_h^{LCB,k+1}(s_h^k) \leq 1 \quad \text{and} \quad u_{\text{ref}}^k(s_h^k) = \text{True}. \quad (C.87)$$

To summarize, the above argument demonstrates that (C.82) can only occur if either (C.86) or (C.87) holds.

With the above observation in place, we can proceed with the following decomposition:

$$\begin{aligned}
 \sum_{h=1}^H \sum_{k=1}^K (V_h^{R,k}(s_h^k) - V_h^{R,K}(s_h^k)) &= \sum_{h=1}^H \sum_{k=1}^K (V_h^{R,k}(s_h^k) - V_h^{R,K}(s_h^k)) \mathbb{1}(V_h^{R,k}(s_h^k) - V_h^{R,K}(s_h^k) \neq 0) \\
 &\leq \sum_{h=1}^H \sum_{k=1}^K (V_h^{R,k}(s_h^k) - V_h^{R,K}(s_h^k)) \mathbb{1}(V_h^{k+1}(s_h^k) - V_h^{\text{LCB},k+1}(s_h^k) \leq 1 \text{ and } u_{\text{ref}}^k(s_h^k) = \text{True}) \\
 &\quad + \underbrace{\sum_{h=1}^H \sum_{k=1}^K (V_h^k(s_h^k) - V_h^{\text{LCB},k}(s_h^k)) \mathbb{1}(V_h^k(s_h^k) - V_h^{\text{LCB},k}(s_h^k) > 1)}_{=: \omega}. \tag{C.88}
 \end{aligned}$$

Regarding the first term in (C.88), it is readily seen that for all $s \in \mathcal{S}$,

$$\sum_{k=1}^K \mathbb{1}(V_h^{k+1}(s) - V_h^{\text{LCB},k+1}(s) \leq 1 \text{ and } u_{\text{ref}}^k(s) = \text{True}) \leq 1, \tag{C.89}$$

which arises since, for each $s \in \mathcal{S}$, the above condition is satisfied in at most one episode, owing to the monotonicity property of V_h, V_h^{LCB} and the update rule for u_{ref} in (17). As a result, one has

$$\begin{aligned}
 \sum_{h=1}^H \sum_{k=1}^K (V_h^{R,k}(s_h^k) - V_h^{R,K}(s_h^k)) \mathbb{1}(V_h^{k+1}(s_h^k) - V_h^{\text{LCB},k+1}(s_h^k) \leq 1 \text{ and } u_{\text{ref}}^k(s_h^k) = \text{True}) \\
 &\leq H \sum_{h=1}^H \sum_{k=1}^K \mathbb{1}(V_h^{k+1}(s_h^k) - V_h^{\text{LCB},k+1}(s_h^k) \leq 1 \text{ and } u_{\text{ref}}^k(s_h^k) = \text{True}) \\
 &= H \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{k=1}^K \mathbb{1}(V_h^{k+1}(s) - V_h^{\text{LCB},k+1}(s) \leq 1 \text{ and } u_{\text{ref}}^k(s) = \text{True}) \\
 &\leq H \sum_{h=1}^H \sum_{s \in \mathcal{S}} 1 = H^2 S,
 \end{aligned}$$

where the first inequality holds since $\|V_h^R, k - V_h^R, k\|_\infty \leq H$. Substitution into (C.88) yields

$$\sum_{h=1}^H \sum_{k=1}^K (V_h^{R,k}(s_h^k) - V_h^{R,K}(s_h^k)) \leq H^2 S + \omega. \tag{C.90}$$

To complete the proof, it boils down to bounding the term ω defined in (C.88). To begin with, note that

$$V_h^R, k_h(s_h^k) \geq V_h^*(s_h^k) \geq V_h^{\text{LCB},k}(s_h^k),$$

where we use the optimism of $V_h^R, k(s_h^k)$ stated in Lemma 2 (cf. (4.9)) and the pessimism of V_h^{LCB} in Lemma 3 (see (4.11)). As a result, we can obtain

$$\begin{aligned} \omega &= \sum_{h=1}^H \sum_{k=1}^K \left(V_h^k(s_h^k) - V_h^{\text{LCB},k}(s_h^k) \right) \mathbb{1} \left(V_h^k(s_h^k) - V_h^{\text{LCB},k}(s_h^k) > 1 \right) \\ &\leq \sum_{h=1}^H \sum_{k=1}^K \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) \right) \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > 1 \right), \end{aligned} \quad (\text{C.91})$$

where the second line arises from the properties $V_h^k(s_h^k) = Q_h^k(s_h^k, a_h^k)$ (given that $a_h^k = \arg \max_a Q_h^k(s_h^k, a)$) as well as the following fact (see line 14 of Algorithm 3)

$$V_h^{\text{LCB},k}(s_h^k) \geq \max_a Q_h^{\text{LCB},k}(s_h^k, a) \geq Q_h^{\text{LCB},k}(s_h^k, a_h^k).$$

Further, let us make note of the following elementary identity:

$$Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) = \int_0^\infty \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > t \right) dt.$$

This allows us to obtain

$$\begin{aligned} \omega &\leq \sum_{h=1}^H \sum_{k=1}^K \left\{ \int_0^\infty \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > t \right) dt \right\} \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > 1 \right) \\ &= \int_1^H \sum_{h=1}^H \sum_{k=1}^K \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > 1 \right) dt \\ &\lesssim \int_1^H \frac{H^6 SA \log \frac{SAT}{\delta}}{t^2} dt \lesssim H^6 SA \log \frac{SAT}{\delta}, \end{aligned} \quad (\text{C.92})$$

where the last line follows from the property (4.12) in Lemma 3. Combining the above bounds (C.91) and (C.92) with (C.90) establishes

$$\begin{aligned} &\sum_{h=1}^H \sum_{k=1}^K \left(V^R, k_h(s_h^k) - V^R, k_h(s_h^k) \right) \\ &\leq H^2 S + \sum_{h=1}^H \sum_{k=1}^K \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) \right) \mathbb{1} \left(Q_h^k(s_h^k, a_h^k) - Q_h^{\text{LCB},k}(s_h^k, a_h^k) > 1 \right) \\ &\leq H^6 SA \log \frac{SAT}{\delta} \end{aligned}$$

as claimed.

D. Proof of Lemma 5

For notational simplicity, we shall adopt the short-hand notation

$$k^n = k_h^n(s_h^k, a_h^k)$$

throughout this section. A starting point for proving this lemma is the upper bound already derived in (4.20), and we intend to further bound the first term on the right-hand side of (4.20). Recalling the expression of $Q_h^{R,k+1}(s_h^k, a_h^k)$ in (C.7) and (C.9), we can derive

$$\begin{aligned}
Q_h^{R,k}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) &= Q_h^{R,k^{N_h^{k-1}(s_h^k, a_h^k)}+1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \\
&= \eta_0^{N_h^{k-1}(s_h^k, a_h^k)} \left(Q_h^{R,1}(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \right) + \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} b_h^{R,k^n+1} \\
&\quad + \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \frac{1}{n} \sum_{i=1}^n V_{h+1}^{R,k^i}(s_{h+1}^{k^i}) - P_{h,s_h^k, a_h^k} V_{h+1}^* \right) \\
&\leq \eta_0^{N_h^{k-1}(s_h^k, a_h^k)} H + B_h^{R,k}(s_h^k, a_h^k) + \frac{2c_b H^2}{(N_h^{k-1}(s_h^k, a_h^k))^{3/4}} \log \frac{SAT}{\delta} \\
&\quad + \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \frac{1}{n} \sum_{i=1}^n V_{h+1}^{R,k^i}(s_{h+1}^{k^i}) - P_{h,s_h^k, a_h^k} V_{h+1}^* \right),
\end{aligned} \tag{D.1}$$

where the last line follows from (C.35) with $B_h^{R,k^{N_h^{k-1}}+1} = B^R, k_h$ and the initialization $Q_h^{R,1}(s_h^k, a_h^k) = H$. Summing over all $1 \leq k \leq K$ gives

$$\begin{aligned}
&\sum_{k=1}^K \left(Q^R, k_h(s_h^k, a_h^k) - Q_h^*(s_h^k, a_h^k) \right) \\
&\leq \sum_{k=1}^K \left(H \eta_0^{N_h^{k-1}(s_h^k, a_h^k)} + B^R, k_h(s_h^k, a_h^k) + \frac{2c_b H^2}{(N_h^{k-1}(s_h^k, a_h^k))^{3/4}} \log \frac{SAT}{\delta} \right) \\
&\quad + \sum_{k=1}^K \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \frac{\sum_{i=1}^n V_{h+1}^{R,k^i}(s_{h+1}^{k^i})}{n} - P_{h,s_h^k, a_h^k} V_{h+1}^* \right) \\
&\leq \sum_{k=1}^K \left(H \eta_0^{N_h^{k-1}(s_h^k, a_h^k)} + B^R, k_h(s_h^k, a_h^k) + \frac{2c_b H^2}{(N_h^{k-1}(s_h^k, a_h^k))^{3/4}} \log \frac{SAT}{\delta} \right) \\
&\quad + \sum_{k=1}^K \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} \left(V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) \right) \\
&\quad + \sum_{k=1}^K \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} \left(V_{h+1}^*(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) + \frac{1}{n} \sum_{i=1}^n V_{h+1}^{R,k^i}(s_{h+1}^{k^i}) - P_{h,s_h^k, a_h^k} V_{h+1}^* \right).
\end{aligned} \tag{D.2}$$

Next, we control each term in (D.2) separately.

- Regarding the first term of (D.2), we make two observations. To begin with,

$$\sum_{k=1}^K \eta_0^{N_h^{k-1}(s_h^k, a_h^k)} \leq \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{n=0}^{N_h^{K-1}(s,a)} \eta_0^n \leq SA, \quad (\text{D.3})$$

where the last inequality follows since $\eta_0^n = 0$ for all $n > 0$ (see (4.2)). Next, it is also observed that

$$\begin{aligned} \sum_{k=1}^K \frac{1}{(N_h^{k-1}(s_h^k, a_h^k))^{3/4}} &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{n=1}^{N_h^{K-1}(s,a)} \frac{1}{n^{3/4}} \\ &\leq \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} 4(N_h^{K-1}(s,a))^{1/4} \leq 4(SA)^{3/4} K^{1/4}, \end{aligned} \quad (\text{D.4})$$

where the last inequality comes from Holder's inequality

$$\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} (N_h^{K-1}(s,a))^{1/4} \leq \left[\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} 1 \right]^{3/4} \left[\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} N_h^{K-1}(s,a) \right]^{1/4} \leq (SA)^{3/4} K^{1/4}.$$

Combining the above bounds yields

$$\begin{aligned} \sum_{k=1}^K &\left(H \eta_0^{N_h^{k-1}(s_h^k, a_h^k)} + B_h^{\mathbb{R},k}(s_h^k, a_h^k) + \frac{2c_b H^2}{(N_h^{k-1}(s_h^k, a_h^k))^{3/4}} \log \frac{SAT}{\delta} \right) \\ &\leq HSA + \sum_{k=1}^K B_h^{\mathbb{R},k}(s_h^k, a_h^k) + 8c_b (SA)^{3/4} K^{1/4} H^2 \log \frac{SAT}{\delta}. \end{aligned} \quad (\text{D.5})$$

- We now turn to the second term of (D.2). A little algebra gives

$$\begin{aligned} &\sum_{k=1}^K \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} (V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\star}(s_{h+1}^{k^n})) \\ &= \sum_{l=1}^K \sum_{N=N_h^l(s_h^l, a_h^l)}^{N_h^{k-1}(s_h^l, a_h^l)} \eta_N^{N_h^l(s_h^l, a_h^l)} (V_{h+1}^l(s_{h+1}^l) - V_{h+1}^{\star}(s_{h+1}^l)) \\ &\leq \left(1 + \frac{1}{H} \right) \sum_{l=1}^K (V_{h+1}^l(s_{h+1}^l) - V_{h+1}^{\star}(s_{h+1}^l)) \\ &= \left(1 + \frac{1}{H} \right) \left[\sum_{k=1}^K (V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\star k}(s_{h+1}^k)) - \sum_{k=1}^K (V_{h+1}^{\star k}(s_{h+1}^k) - V_{h+1}^{\star k}(s_{h+1}^k)) \right]. \end{aligned} \quad (\text{D.6})$$

Here, the second line replaces k^n (resp. n) with l (resp. $N_h^l(s_h^l, a_h^l)$), the third line is due to the property $\sum_{N=n}^{\infty} \eta_N^N \leq 1 + 1/H$ (see Lemma 1), while the last relation replaces l with k again.

- When it comes to the last term of (D.2), we can derive

$$\begin{aligned}
& \sum_{k=1}^K \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} \left(V_{h+1}^*(s_{h+1}^{k^n}) - V_{h+1}^{R, k^n}(s_{h+1}^{k^n}) + \frac{1}{n} \sum_{i=1}^n V_{h+1}^{R, k^i}(s_{h+1}^{k^i}) - P_{h, s_h^k, a_h^k} V_{h+1}^* \right) \\
&= \sum_{k=1}^K \sum_{n=1}^{N_h^{k-1}(s_h^k, a_h^k)} \eta_n^{N_h^{k-1}(s_h^k, a_h^k)} \left((P_h^{k^n} - P_{h, s_h^k, a_h^k}) (V_{h+1}^* - V_{h+1}^{R, k^n}) + \frac{1}{n} \sum_{i=1}^n (V_{h+1}^{R, k^i}(s_{h+1}^{k^i}) - P_{h, s_h^k, a_h^k} V_{h+1}^{R, k^n}) \right) \\
&= \sum_{k=1}^K \sum_{N=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \eta_{N_h^k(s_h^k, a_h^k)}^N \left((P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^* - V^R, k_{h+1}) \right. \\
&\quad \left. + \frac{\sum_{i=1}^{N_h^k(s_h^k, a_h^k)} (V_{h+1}^{R, k^i}(s_{h+1}^{k^i}) - P_{h, s_h^k, a_h^k} V^R, k_{h+1})}{N_h^k(s_h^k, a_h^k)} \right).
\end{aligned}$$

Here, the first equality holds since $V_{h+1}^*(s_{h+1}^{k^n}) - V_{h+1}^{R, k^n}(s_{h+1}^{k^n}) = P_h^{k^n} (V_{h+1}^* - V_{h+1}^{R, k^n})$ (in view of the definition of P_h^k in (4.5)), the second equality can be seen via simple rearrangement of the terms, while in the last line we replace k^n (resp. n) with k (resp. $N_h^k(s_h^k, a_h^k)$).

Taking the above bounds together with (D.2) and (4.20), we can rearrange terms to reach

$$\begin{aligned}
& \sum_{k=1}^K (V_h^k(s_h^k) - V_h^{\pi^k}(s_h^k)) \\
& \leq \left(1 + \frac{1}{H} \right) \sum_{k=1}^K (V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k)) + \sum_{k=1}^K B^R, k_h(s_h^k, a_h^k) \\
& \quad + HSA + 8c_b H^2 (SA)^{3/4} K^{1/4} \log \frac{SAT}{\delta} + \sum_{k=1}^K (P_{h, s_h^k, a_h^k} - P_h^k) (V_{h+1}^* - V_{h+1}^{\pi^k}) \\
& \quad + \sum_{k=1}^K \sum_{N=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \eta_{N_h^k(s_h^k, a_h^k)}^N \left[(P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^* - V^R, k_{h+1}) \right. \\
& \quad \left. + \frac{\sum_{i=1}^{N_h^k(s_h^k, a_h^k)} (V_{h+1}^{R, k^i}(s_{h+1}^{k^i}) - P_{h, s_h^k, a_h^k} V^R, k_{h+1})}{N_h^k(s_h^k, a_h^k)} \right], \tag{D.7}
\end{aligned}$$

where we have dropped the term $-\frac{1}{H} \sum_k (V_{h+1}^*(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k))$ owing to the fact that $V_{h+1}^* \geq V_{h+1}^{\pi^k}$.

Thus far, we have established a crucial connection between $\sum_{k=1}^K (V_h^k(s_h^k) - V_h^{\pi^k}(s_h^k))$ at step h and $\sum_{k=1}^K (V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k))$ at step $h+1$. Clearly, the term $V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k)$ can be further bounded in the same manner. As a result, by recursively applying the above relation (D.7) over the time steps $h = 1, 2, \dots, H$ and using the terminal condition $V_{H+1}^k = V_{H+1}^{\pi^k} = 0$, we can immediately arrive at the advertised bound in Lemma 5.

E. Proof of Lemma 6

E.1 Bounding the term \mathcal{R}_1

First of all, let us look at the first two terms of \mathcal{R}_1 in (4.22a). Recognizing the following elementary inequality:

$$\left(1 + \frac{1}{H}\right)^{h-1} \leq \left(1 + \frac{1}{H}\right)^H \leq e \quad \text{for all } h = 1, 2, \dots, H+1, \quad (\text{E.1})$$

we are allowed to upper bound the first two terms in (4.22a) as follows:

$$\begin{aligned} \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^{h-1} \left\{ HSA + 8c_b H^2 (SA)^{3/4} K^{1/4} \log \frac{SAT}{\delta} \right\} &\lesssim H^2 SA + H^3 (SA)^{3/4} K^{1/4} \log \frac{SAT}{\delta} \\ &\lesssim H^{4.5} SA \log^2 \frac{SAT}{\delta} + \sqrt{H^3 SAK} = H^{4.5} SA \log^2 \frac{SAT}{\delta} + \sqrt{H^2 SAT}, \end{aligned} \quad (\text{E.2})$$

where the last inequality can be shown using the AM–GM inequality as follows:

$$H^3 (SA)^{3/4} K^{1/4} \log \frac{SAT}{\delta} = \left(H^{9/4} \sqrt{SA} \log \frac{SAT}{\delta} \right) \cdot (H^3 SAK)^{1/4} \leq H^{4.5} SA \log^2 \frac{SAT}{\delta} + \sqrt{H^3 SAK}.$$

We are now left with the last term of \mathcal{R}_1 in (4.22a). Toward this, we resort to Lemma 8 by setting

$$W_{h+1}^i := V_{h+1}^* - V_{h+1}^{\pi^k} \quad \text{and} \quad c_h := \left(1 + \frac{1}{H}\right)^{h-1}.$$

In view of (E.1) and the property $H \geq V^*(s) \geq V^\pi(s) \geq 0$, we see that

$$0 \leq c_h \leq e, \quad W_{h+1}^i \geq 0, \quad \text{and} \quad \|W_{h+1}^i\|_\infty \leq H =: C_w.$$

Therefore, applying Lemma 8 yields

$$\begin{aligned} \left| \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^{h-1} \sum_{k=1}^K (P_{h,s_h^k, a_h^k} - P_h^k) (V_{h+1}^* - V_{h+1}^{\pi^k}) \right| &= \left| \sum_{h=1}^H \sum_{k=1}^K Y_{k,h} \right| \\ &\lesssim \sqrt{TC_w^2 \log \frac{1}{\delta}} + C_w \log \frac{1}{\delta} = \sqrt{H^2 T \log \frac{1}{\delta}} + H \log \frac{1}{\delta} \end{aligned} \quad (\text{E.3})$$

with probability exceeding $1 - \delta$.

Combining (E.2) and (E.3) with the definition (4.22a) of \mathcal{R}_1 immediately leads to the claimed bound.

E.2 Bounding the term \mathcal{R}_2

In view of the definition of $B^R, k_h(s_h^k, a_h^k)$ in line ?? of Algorithm 2, we can decompose \mathcal{R}_2 (cf. (4.22b)) as follows:

$$\begin{aligned} \mathcal{R}_2 &= \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^{h-1} c_b \sqrt{H \log \frac{SAT}{\delta}} \sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{adv},k}(s_h^k, a_h^k) - (\mu_h^{\text{adv},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \\ &\quad + \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^{h-1} c_b \sqrt{\log \frac{SAT}{\delta}} \sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \\ &\lesssim \sqrt{H \log \frac{SAT}{\delta}} \sum_{h=1}^H \sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{adv},k}(s_h^k, a_h^k) - (\mu_h^{\text{adv},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \\ &\quad + \sqrt{\log \frac{SAT}{\delta}} \sum_{h=1}^H \sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}}, \end{aligned} \quad (\text{E.4})$$

where the last relation holds due to (E.1). In what follows, we intend to bound these two terms separately.

Step 1: upper bounding the first term in (E.4). Toward this, we make the observation that

$$\begin{aligned} \sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{adv},k}(s_h^k, a_h^k) - (\mu_h^{\text{adv},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} &\leq \sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{adv},k}(s_h^k, a_h^k)}{N_h^k(s_h^k, a_h^k)}} \\ &= \sum_{k=1}^K \sqrt{\frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} \eta_n^k(s_h^k, a_h^k) (V_{h+1}^{k^n}(s_{h+1}^{k^n}) - V_{h+1}^{R,k^n}(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)}}, \end{aligned} \quad (\text{E.5})$$

where the second line follows from the update rule of σ_h^{adv} , k in (C38). Combining the relation $|V_{h+1}^k(s_h^k) - V_{h+1}^R(s_h^k)| \leq 2$ (cf. (4.15)) and the property $\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} \eta_n^k(s_h^k, a_h^k) \leq 1$ (cf. (4.3)) with (E.5) yields

$$\sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{adv},k}(s_h^k, a_h^k) - (\mu_h^{\text{adv},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \leq \sum_{k=1}^K \sqrt{\frac{4}{N_h^k(s_h^k, a_h^k)}} \leq 2\sqrt{SAK}. \quad (\text{E.6})$$

Here, the last inequality holds due to the following fact:

$$\begin{aligned} \sum_{k=1}^K \sqrt{\frac{1}{N_h^k(s_h^k, a_h^k)}} &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{n=1}^{N_h^K(s,a)} \sqrt{\frac{1}{n}} \leq 2 \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sqrt{N_h^K(s,a)} \\ &\leq 2 \sqrt{\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} 1} \cdot \sqrt{\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} N_h^K(s,a)} = 2\sqrt{SAK}, \end{aligned} \quad (\text{E.7})$$

where the last line arises from Cauchy–Schwarz and the basic fact that $\sum_{(s,a)} N_h^K(s,a) = K$.

Step 2: upper bounding the second term in (E.4). Recalling the update rules of $\mu_h^{\text{ref},k}$ and $\sigma_h^{\text{ref},k}$ in (C52), we have

$$\begin{aligned} &\sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \\ &= \sum_{k=1}^K \sqrt{\frac{1}{N_h^k(s_h^k, a_h^k)}} \underbrace{\sqrt{\frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} (V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)}} - \left(\frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2}_{=:J_h^k}. \end{aligned} \quad (\text{E.8})$$

Additionally, the quantity J_h^k defined in (E.8) obeys

$$\begin{aligned} (J_h^k)^2 &\leq \frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} (V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}))^2 - (V_{h+1}^{\star}(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)} \\ &\quad + \frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} (V_{h+1}^{\star}(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)} - \left(\frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} V_{h+1}^{\star}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2 \\ &\leq \underbrace{\frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} 2H(V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\star}(s_{h+1}^{k^n}))}{N_h^k(s_h^k, a_h^k)}}_{=:J_1} \\ &\quad + \underbrace{\frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} (V_{h+1}^{\star}(s_{h+1}^{k^n}))^2}{N_h^k(s_h^k, a_h^k)} - \left(\frac{\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} V_{h+1}^{\star}(s_{h+1}^{k^n})}{N_h^k(s_h^k, a_h^k)} \right)^2}_{=:J_2}, \end{aligned} \quad (\text{E.9})$$

which arises from the fact that $H \geq V_{h+1}^{\text{R},k^n} \geq V_{h+1}^{\star} \geq 0$ for all $k^n \leq K$ and hence

$$\begin{aligned} (V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}))^2 - (V_{h+1}^{\star}(s_{h+1}^{k^n}))^2 &= (V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}) + V_{h+1}^{\star}(s_{h+1}^{k^n}))(V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\star}(s_{h+1}^{k^n})) \\ &\leq 2H(V_{h+1}^{\text{R},k^n}(s_{h+1}^{k^n}) - V_{h+1}^{\star}(s_{h+1}^{k^n})). \end{aligned}$$

With (E.9) in mind, we shall proceed to bound each term in (E.9) separately.

- The first term J_1 can be straightforwardly bounded as follows:

$$\begin{aligned}
J_1 &= \frac{2H}{N_h^k(s_h^k, a_h^k)} \left(\sum_{n=1}^{N_h^k(s_h^k, a_h^k)} \left(V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) \right) \right. \\
&\quad \left. \mathbb{1}\left(V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) \leq 3 \right) + \Phi_h^k(s_h^k, a_h^k) \right) \\
&\leq 6H + \frac{2H}{N_h^k(s_h^k, a_h^k)} \Phi_h^k(s_h^k, a_h^k), \tag{E.10}
\end{aligned}$$

where $\Phi_h^k(s_h^k, a_h^k)$ is defined as

$$\Phi_h^k(s_h^k, a_h^k) := \sum_{n=1}^{N_h^k(s_h^k, a_h^k)} \left(V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) \right) \mathbb{1}\left(V_{h+1}^{R,k^n}(s_{h+1}^{k^n}) - V_{h+1}^*(s_{h+1}^{k^n}) > 3 \right). \tag{E.11}$$

- When it comes to the second term J_2 , we claim that

$$J_2 \lesssim \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*) + H^2 \sqrt{\frac{\log \frac{SAT}{\delta}}{N_h^k(s_h^k, a_h^k)}}, \tag{E.12}$$

which will be justified in Appendix E.2.1.

Plugging (E.10) and (E.12) into (E.9) and (E.8) allows one to demonstrate that

$$\begin{aligned}
&\sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \\
&\lesssim \sum_{k=1}^K \sqrt{\frac{1}{N_h^k(s_h^k, a_h^k)}} \sqrt{H + \frac{H\Phi_h^k(s_h^k, a_h^k)}{N_h^k(s_h^k, a_h^k)} + \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*) + H^2 \sqrt{\frac{\log \frac{SAT}{\delta}}{N_h^k(s_h^k, a_h^k)}}} \\
&\leq \sum_{k=1}^K \left(\sqrt{\frac{H}{N_h^k(s_h^k, a_h^k)}} + \frac{\sqrt{H\Phi_h^k(s_h^k, a_h^k)}}{N_h^k(s_h^k, a_h^k)} + \sqrt{\frac{\text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*)}{N_h^k(s_h^k, a_h^k)} + \frac{H \log^{1/4} \frac{SAT}{\delta}}{(N_h^k(s_h^k, a_h^k))^{3/4}}} \right) \\
&\lesssim \sqrt{HS\bar{A}K} + \sum_{k=1}^K \frac{\sqrt{H\Phi_h^k(s_h^k, a_h^k)}}{N_h^k(s_h^k, a_h^k)} + \sum_{k=1}^K \sqrt{\frac{\text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*)}{N_h^k(s_h^k, a_h^k)} + H(SA)^{3/4} \left(K \log \frac{SAT}{\delta} \right)^{1/4}}, \tag{E.13}
\end{aligned}$$

where the last line follows from (E.7) and (D.4).

Step 3: putting together the preceding results.

Finally, the above results in (E.6) and (E.13) taken collectively with (E.4) lead to

$$\begin{aligned}
\mathcal{R}_2 &\lesssim \sqrt{H^3 SAK \log \frac{SAT}{\delta}} + \sum_{h=1}^H \sqrt{\log \frac{SAT}{\delta}} \sum_{k=1}^K \sqrt{\frac{\sigma_h^{\text{ref},k}(s_h^k, a_h^k) - (\mu_h^{\text{ref},k}(s_h^k, a_h^k))^2}{N_h^k(s_h^k, a_h^k)}} \\
&\lesssim \sqrt{H^3 SAK \log \frac{SAT}{\delta}} + H^2 (SA)^{3/4} K^{1/4} \log^{5/4} \frac{SAT}{\delta} + \sqrt{\log \frac{SAT}{\delta}} \sum_{h=1}^H \sum_{k=1}^K \sqrt{\frac{\text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*)}{N_h^k(s_h^k, a_h^k)}} \\
&\quad + \sqrt{H \log \frac{SAT}{\delta}} \sum_{h=1}^H \sum_{k=1}^K \frac{\sqrt{\Phi_h^k(s_h^k, a_h^k)}}{N_h^k(s_h^k, a_h^k)} \\
&\stackrel{(i)}{\lesssim} \sqrt{H^3 SAK \log \frac{SAT}{\delta}} + H^2 (SA)^{3/4} K^{1/4} \log^{5/4} \frac{SAT}{\delta} + H^4 SA \log^2 \frac{SAT}{\delta} \\
&\stackrel{(ii)}{\lesssim} \sqrt{H^3 SAK \log \frac{SAT}{\delta}} + H^4 SA \log^2 \frac{SAT}{\delta} = \sqrt{H^2 SAT \log \frac{SAT}{\delta}} + H^4 SA \log^2 \frac{SAT}{\delta}.
\end{aligned}$$

Here, (i) holds due to the following two claimed inequalities:

$$\sum_{h=1}^H \sum_{k=1}^K \sqrt{\frac{\text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*)}{N_h^k(s_h^k, a_h^k)}} \lesssim \sqrt{H^2 SAT \log \frac{SAT}{\delta}} + H^4 SA \log \frac{SAT}{\delta}, \quad (\text{E.14})$$

$$\sum_{h=1}^H \sum_{k=1}^K \frac{\sqrt{\Phi_h^k(s_h^k, a_h^k)}}{N_h^k(s_h^k, a_h^k)} \lesssim H^{7/2} SA \log^{3/2} \frac{SAT}{\delta}, \quad (\text{E.15})$$

whose proofs are postponed to Appendix E.2.2 and Appendix E.2.3, respectively. Additionally, the inequality (ii) above is valid since

$$\begin{aligned}
H^2 (SA)^{3/4} K^{1/4} \log^{5/4} \frac{SAT}{\delta} &= \left(H^{5/4} (SA)^{1/2} \log \frac{SAT}{\delta} \right) \cdot \left(H^3 SAK \log \frac{SAT}{\delta} \right)^{1/4} \\
&\lesssim H^{2.5} SA \log^2 \frac{SAT}{\delta} + \sqrt{H^3 SAK \log \frac{SAT}{\delta}} = H^{2.5} SA \log^2 \frac{SAT}{\delta} + \sqrt{H^2 SAT \log \frac{SAT}{\delta}}
\end{aligned}$$

due to the Cauchy–Schwarz inequality. This concludes the proof of the advertised upper bound on \mathcal{R}_2 .

E.2.1 *Proof of the inequality (E12).* Akin to the proof of I_4^1 in (C.55), let

$$W_{h+1}^i := (V_{h+1}^*)^2 \quad \text{and} \quad u_h^i(s, a, N) := \frac{1}{N}.$$

By observing and setting

$$C_u := \frac{1}{N}, \quad \|W_{h+1}^i\|_\infty \leq H^2 =: C_w,$$

we can apply Lemma 7 to yield

$$\left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (V_{h+1}^*(s_{h+1}^{k^n}))^2 - P_{h,s_h^k, a_h^k} (V_{h+1}^*)^2 \right| = \left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_h^{k^n} - P_{h,s_h^k, a_h^k}) (V_{h+1}^*)^2 \right| \lesssim H^2 \sqrt{\frac{\log^2 \frac{SAT}{\delta}}{N_h^k}}$$

with probability at least $1 - \delta$. Similarly, by applying the trivial bound $\|V_{h+1}^*\|_\infty \leq H$ and Lemma 7, we can obtain

$$\left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n}) - P_{h,s_h^k, a_h^k} V_{h+1}^* \right| = \left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} (P_h^{k^n} - P_{h,s_h^k, a_h^k}) V_{h+1}^* \right| \lesssim H \sqrt{\frac{\log \frac{SAT}{\delta}}{N_h^k}}$$

with probability at least $1 - \delta$.

Recalling from (4.6) the definition

$$\text{Var}_{h,s_h^k, a_h^k} (V_{h+1}^*) = P_{h,s_h^k, a_h^k} (V_{h+1}^*)^2 - (P_{h,s_h^k, a_h^k} V_{h+1}^*)^2,$$

we can use the preceding two bounds and the triangle inequality to show that

$$\begin{aligned} & \left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n})^2 - \left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n}) \right)^2 - \text{Var}_{h,s_h^k, a_h^k} (V_{h+1}^*) \right| \\ & \leq \left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n})^2 - P_{h,s_h^k, a_h^k} (V_{h+1}^*)^2 \right| + \left| \left(\frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n}) \right)^2 - (P_{h,s_h^k, a_h^k} V_{h+1}^*)^2 \right| \\ & \lesssim H^2 \sqrt{\frac{\log \frac{SAT}{\delta}}{N_h^k}} + \left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n}) - P_{h,s_h^k, a_h^k} V_{h+1}^* \right| \cdot \left| \frac{1}{N_h^k} \sum_{n=1}^{N_h^k} V_{h+1}^*(s_{h+1}^{k^n}) + P_{h,s_h^k, a_h^k} V_{h+1}^* \right| \\ & \lesssim H^2 \sqrt{\frac{\log \frac{SAT}{\delta}}{N_h^k}} \end{aligned}$$

with probability at least $1 - \delta$, where the last line also uses the fact that $\|V_{h+1}^*\|_\infty \leq H$.

E.2.2 Proof of the inequality (E14) To begin with, we make the observation that

$$\sum_{k=1}^K \sqrt{\frac{\text{Var}_{h,s_h^k, a_h^k} (V_{h+1}^*)}{N_h^k (s_h^k, a_h^k)}} = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{n=1}^{N_h^K (s,a)} \sqrt{\frac{\text{Var}_{h,s,a} (V_{h+1}^*)}{n}} \leq 2 \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sqrt{N_h^K (s,a) \text{Var}_{h,s,a} (V_{h+1}^*)},$$

which relies on the fact that $\sum_{n=1}^N 1/\sqrt{n} \leq 2\sqrt{N}$. It then follows that

$$\begin{aligned}
 \sum_{h=1}^H \sum_{k=1}^K \sqrt{\frac{\text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*)}{N_h^k(s_h^k, a_h^k)}} &\leq 2 \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sqrt{N_h^K(s, a) \text{Var}_{h,s,a}(V_{h+1}^*)} \\
 &\leq 2 \sqrt{\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} 1} \cdot \sqrt{\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} N_h^K(s, a) \text{Var}_{h,s,a}(V_{h+1}^*)} \\
 &= 2\sqrt{HSA} \sqrt{\sum_{h=1}^H \sum_{k=1}^K \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*)}, \tag{E.16}
 \end{aligned}$$

where the second inequality invokes the Cauchy–Schwarz inequality.

The rest of the proof is then dedicated to bounding (E.16). Toward this end, we first decompose

$$\begin{aligned}
 \sum_{h=1}^H \sum_{k=1}^K \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*) &\leq \sum_{h=1}^H \sum_{k=1}^K \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^{\pi^k}) + \sum_{h=1}^H \sum_{k=1}^K \left| \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*) - \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^{\pi^k}) \right| \\
 &\stackrel{(ii)}{\lesssim} HT + H^3 \log \frac{SAT}{\delta} + \sum_{h=1}^H \sum_{k=1}^K \left| \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*) - \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^{\pi^k}) \right|, \tag{E.17}
 \end{aligned}$$

where (ii) follows directly from (30, Lemma C.5). The second term on the right-hand side of (E.17) can be bounded as follows:

$$\begin{aligned}
 &\sum_{h=1}^H \sum_{k=1}^K \left| \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*) - \text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^{\pi^k}) \right| \\
 &= \sum_{h=1}^H \sum_{k=1}^K \left| P_{h,s_h^k,a_h^k}(V_{h+1}^*)^2 - (P_{h,s_h^k,a_h^k} V_{h+1}^*)^2 - P_{h,s_h^k,a_h^k}(V_{h+1}^{\pi^k})^2 + (P_{h,s_h^k,a_h^k} V_{h+1}^{\pi^k})^2 \right| \\
 &\leq \sum_{h=1}^H \sum_{k=1}^K \left\{ \left| P_{h,s_h^k,a_h^k} ((V_{h+1}^* - V_{h+1}^{\pi^k})(V_{h+1}^* + V_{h+1}^{\pi^k})) \right| + \left| (P_{h,s_h^k,a_h^k} V_{h+1}^*)^2 - (P_{h,s_h^k,a_h^k} V_{h+1}^{\pi^k})^2 \right| \right\} \\
 &\stackrel{(i)}{\leq} 4H \sum_{h=1}^H \sum_{k=1}^K P_{h,s_h^k,a_h^k} (V_{h+1}^* - V_{h+1}^{\pi^k}) \\
 &= 4H \sum_{h=1}^H \sum_{k=1}^K \left\{ V_{h+1}^*(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k) + (P_{h,s_h^k,a_h^k} - P_h^k) (V_{h+1}^* - V_{h+1}^{\pi^k}) \right\}
 \end{aligned}$$

$$\begin{aligned}
&\stackrel{(ii)}{\leq} 4H \sum_{h=1}^H \sum_{k=1}^K (\phi_{h+1}^k + \delta_{h+1}^k) \stackrel{(iii)}{\lesssim} H^2 \sqrt{T \log \frac{SAT}{\delta}} + H^4 \sqrt{SAT \log \frac{SAT}{\delta}} + H^4 SA \\
&\asymp H^4 \sqrt{SAT \log \frac{SAT}{\delta}} + H^4 SA,
\end{aligned} \tag{E.18}$$

where we define

$$\delta_{h+1}^k := V_{h+1}^{\text{UCB},k}(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k), \quad \phi_{h+1}^k := (P_{h,s_h^k,a_h^k} - P_h^k)(V_{h+1}^* - V_{h+1}^{\pi^k}). \tag{E.19}$$

We shall take a moment to explain how we derive (E.18). The inequality (i) holds by observing that $V_{h+1}^* - V_{h+1}^{\pi^k} \geq 0$ and

$$\begin{aligned}
\left| P_{h,s_h^k,a_h^k} \left((V_{h+1}^* - V_{h+1}^{\pi^k})(V_{h+1}^* + V_{h+1}^{\pi^k}) \right) \right| &\leq P_{h,s_h^k,a_h^k} (V_{h+1}^* - V_{h+1}^{\pi^k}) (\|V_{h+1}^*\|_\infty + \|V_{h+1}^{\pi^k}\|_\infty) \\
&\leq 2HP_{h,s_h^k,a_h^k} (V_{h+1}^* - V_{h+1}^{\pi^k}), \\
\left| (P_{h,s_h^k,a_h^k} V_{h+1}^*)^2 - (P_{h,s_h^k,a_h^k} V_{h+1}^{\pi^k})^2 \right| &\leq \left| P_{h,s_h^k,a_h^k} (V_{h+1}^* - V_{h+1}^{\pi^k}) \right| \cdot \left| P_{h,s_h^k,a_h^k} (V_{h+1}^* + V_{h+1}^{\pi^k}) \right| \\
&\leq 2HP_{h,s_h^k,a_h^k} (V_{h+1}^* - V_{h+1}^{\pi^k});
\end{aligned}$$

(ii) is valid since $V_{h+1}^{\text{UCB}} \geq V_{h+1}^*$; and (iii) results from the following two bounds:

$$\sum_{h=1}^H \sum_{k=1}^K \delta_{h+1}^k \lesssim H^3 \sqrt{SAT \log \frac{SAT}{\delta}} + H^3 SA, \tag{E.20a}$$

$$\sum_{h=1}^H \sum_{k=1}^K \phi_{h+1}^k \lesssim H \sqrt{T \log \frac{SAT}{\delta}}, \tag{E.20b}$$

which come, respectively, from (30, Eqn. (C.13)) and the argument for (30, Eqn. (C.12)).⁴

As a consequence, substituting (E.17) and (E.18) into (E.16), we reach

$$\begin{aligned}
\sum_{h=1}^H \sum_{k=1}^K \sqrt{\frac{\text{Var}_{h,s_h^k,a_h^k}(V_{h+1}^*)}{N_h^k(s_h^k, a_h^k)}} &\lesssim \sqrt{HSA} \sqrt{HT + H^4 \sqrt{SAT \log \frac{SAT}{\delta}} + H^4 SA} \\
&\lesssim \sqrt{H^2 SAT} + H^{5/2} (SA)^{3/4} \left(T \log \frac{SAT}{\delta} \right)^{1/4} + H^{2.5} SA \\
&= \sqrt{H^2 SAT} + \left(H^2 SAT \log \frac{SAT}{\delta} \right)^{1/4} (H^4 SA)^{1/2} + H^{2.5} SA \\
&\lesssim \sqrt{H^2 SAT \log \frac{SAT}{\delta}} + H^4 SA \log \frac{SAT}{\delta},
\end{aligned}$$

⁴ Note that the notation δ_h^k used in ([30], Section C.2) and the one in the proof of ([30], Theorem 1) are different; here, we need to adopt the notation used in the proof of ([30], Theorem 1).

where we have applied the basic inequality $2ab \leq a^2 + b^2$ for any $a, b \geq 0$.

E.2.3 Proof of the inequality (E15)

First, it is observed that

$$\begin{aligned} \sum_{k=1}^K \frac{\sqrt{\Phi_h^k(s_h^k, a_h^k)}}{N_h^k(s_h^k, a_h^k)} &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{n=1}^{N_h^K(s,a)} \frac{\sqrt{\Phi_h^{k^n}(s,a)}(s,a)}{n} \\ &\leq \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sqrt{\Phi_h^{N_h^K(s,a)}(s,a) \log T} \leq \sqrt{SA \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \Phi_h^{N_h^K(s,a)}(s,a) \log T}. \end{aligned} \quad (\text{E.21})$$

Here, the first inequality holds by the monotonicity property of $\Phi_h^k(s_h, a_h)$ with respect to k (see its definition in (E.11)) due to the same property of V_{h+1}^R, k , while the second inequality comes from Cauchy–Schwarz.

To continue, note that

$$\begin{aligned} &\sum_{h=1}^H \sqrt{\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \Phi_h^{N_h^K(s,a)}(s,a)} \\ &= \sum_{h=1}^H \sqrt{\sum_{k=1}^K \left(V_{h+1}^{R,k}(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k) \right) \mathbb{1}\left(V_{h+1}^{R,k}(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k) > 3 \right)} \\ &\leq \sum_{h=1}^H \sqrt{\sum_{k=1}^K \left(V_{h+1}^k(s_{h+1}^k) + 2 - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) \right) \mathbb{1}\left(V_{h+1}^k(s_{h+1}^k) + 2 - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) > 3 \right)} \\ &= \sum_{h=1}^H \sqrt{\sum_{k=1}^K \left(V_{h+1}^k(s_{h+1}^k) + 2 - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) \right) \mathbb{1}\left(V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) > 1 \right)} \\ &\leq \sum_{h=1}^H \sqrt{\sum_{k=1}^K 3 \left(V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) \right) \mathbb{1}\left(V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) > 1 \right)} \\ &\leq \sqrt{H} \sqrt{\sum_{h=1}^H \sum_{k=1}^K 3 \left(V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) \right) \mathbb{1}\left(V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) > 1 \right)}, \end{aligned} \quad (\text{E.22})$$

where the first inequality follows from Lemma 4 (cf. (4.15)) and Lemma 3 (so that $V_{h+1}^R, k(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k) \leq V_{h+1}^k(s_{h+1}^k) + 2 - V_{h+1}^{\text{LCB},k}(s_{h+1}^k)$), the penultimate inequality holds since $1 \leq V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k)$ when $\mathbb{1}\left(V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) > 1 \right) \neq 0$, and the last inequality is a consequence of the Cauchy–Schwarz inequality.

Combining the above relation with (C.91) and applying the triangle inequality, we can demonstrate that

$$\begin{aligned} & \sum_{h=1}^H \sqrt{\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \Phi_h^{N_h^K(s,a)}(s,a)} \\ & \lesssim \sqrt{H} \sqrt{\sum_{h=1}^H \sum_{k=1}^K \left(Q_{h+1}^k(s_{h+1}^k, a_{h+1}^k) - Q_{h+1}^{\text{LCB},k}(s_{h+1}^k, a_{h+1}^k) \right) \mathbb{1} \left(Q_{h+1}^k(s_{h+1}^k, a_{h+1}^k) - Q_{h+1}^{\text{LCB},k}(s_{h+1}^k, a_{h+1}^k) > 1 \right)} \\ & \lesssim \sqrt{H^7 S A \log \frac{SAT}{\delta}}, \end{aligned}$$

where the second inequality follows directly from (4.16), and the first inequality is valid since

$$V_{h+1}^k(s_{h+1}^k) - V_{h+1}^{\text{LCB},k}(s_{h+1}^k) \leq Q_{h+1}^k(s_{h+1}^k, a_{h+1}^k) - Q_{h+1}^{\text{LCB},k}(s_{h+1}^k, a_{h+1}^k).$$

Substitution into (E.21) gives

$$\sum_{h=1}^H \sum_{k=1}^K \frac{\sqrt{\Phi_h^k(s_h^k, a_h^k)}}{N_h^k(s_h^k, a_h^k)} \lesssim \left(\sqrt{S A \log T} \right) \cdot \sqrt{H^7 S A \log \frac{SAT}{\delta}} \asymp H^{7/2} S A \log^{3/2} \frac{SAT}{\delta},$$

thus concluding the proof.

E.3 Bounding the term \mathcal{R}_3

For notational convenience, we shall use the short-hand notation

$$k^i := k_h^i(s_h^k, a_h^k)$$

whenever it is clear from the context. This allows us to decompose the expression of \mathcal{R}_3 in (4.22c) as follows:

$$\begin{aligned} \mathcal{R}_3 &:= \underbrace{\sum_{h=1}^H \sum_{k=1}^K \lambda_h^k (P_h^k - P_{h,s_h^k, a_h^k}) (V_{h+1}^{\star} - V^{\text{R}}, k_{h+1})}_{=: \mathcal{R}_3^1} \\ &\quad + \underbrace{\sum_{h=1}^H \sum_{k=1}^K \lambda_h^k \frac{\sum_{i \leq N_h^k(s_h^k, a_h^k)} (V_{h+1}^{\text{R}, k^i}(s_{h+1}^{k^i}) - P_{h,s_h^k, a_h^k} V^{\text{R}}, k_{h+1})}{N_h^k(s_h^k, a_h^k)}}_{=: \mathcal{R}_3^2} \end{aligned}$$

with

$$\lambda_h^k := \left(1 + \frac{1}{H}\right)^{h-1} \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \eta_{N_h^k(s_h^k, a_h^k)}^n \leq \left(1 + \frac{1}{H}\right)^h \leq \left(1 + \frac{1}{H}\right)^H \leq e. \quad (\text{E.23})$$

Here, the first inequality in (E.23) follows from the property $\sum_{N=n}^{\infty} \eta_n^N \leq 1 + 1/H$ in Lemma 1, while the last inequality in (E.23) results from (E.1). In the sequel, we shall control each of these two terms separately.

Step 1: upper bounding \mathcal{R}_3^1 . We plan to control this term by means of Lemma 8. For notational simplicity, let us define

$$N(s, a, h) := N_h^{K-1}(s, a)$$

and set

$$W_{h+1}^i := V^R, k_{h+1} - V_{h+1}^* \quad \text{and} \quad u_h^i(s_h^i, a_h^i) := \lambda_h^i = \left(1 + \frac{1}{H}\right)^{h-1} \sum_{n=N_h^i(s_h^i, a_h^i)}^{N(s_h^i, a_h^i, h)} \eta_n^i.$$

Given the fact that $V^R, k_{h+1}(s), V_{h+1}^*(s) \in [0, H]$ and the condition (E.23), it is readily seen that

$$|u_h^i(s_h^i, a_h^i)| \leq e =: C_u \quad \text{and} \quad \|W_{h+1}^i\|_\infty \leq H =: C_w.$$

Apply Lemma 8 to yield

$$\begin{aligned} \left| \sum_{h=1}^H \sum_{k=1}^K \lambda_h^k (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^* - V_{h+1}^{R,k}) \right| &= \left| \sum_{h=1}^H \sum_{k=1}^K X_{k,h} \right| \\ &\lesssim \sqrt{C_u^2 C_w HSA \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} [P_h^i W_{h+1}^i] \log \frac{K}{\delta} + C_u C_w HSA \log \frac{K}{\delta}} \\ &\lesssim \sqrt{H^2 SA \sum_{h=1}^H \sum_{k=1}^K \mathbb{E}_{i,h-1} [P_h^k (V_{h+1}^{R,k} - V_{h+1}^*)] \log \frac{T}{\delta} + H^2 SA \log \frac{T}{\delta}} \\ &\asymp \sqrt{H^2 SA \left\{ \sum_{h=1}^H \sum_{k=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^{R,k} - V_{h+1}^*) \right\} \log \frac{T}{\delta} + H^2 SA \log \frac{T}{\delta}} \end{aligned} \quad (\text{E.24})$$

with probability at least $1 - \delta/2$.

It then comes down to controlling the sum $\sum_{h=1}^H \sum_{k=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^R, k - V_{h+1}^*)$. Toward this end, we first single out the following useful fact:

$$\begin{aligned} \sum_{h=1}^H \sum_{k=1}^K P_h^k (V^R, k_{h+1} - V_{h+1}^*) &\stackrel{(i)}{\leq} \sum_{h=1}^H \sum_{k=1}^K P_h^k (V_{h+1}^k + 2 - V_{h+1}^*) \\ &\leq 2HK + \sum_{h=1}^H \sum_{k=1}^K (V_{h+1}^k(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k)) \stackrel{(ii)}{\lesssim} \sqrt{H^7 SAK \log \frac{SAT}{\delta}} + H^3 SA + HK \end{aligned} \quad (\text{E.25})$$

with probability at least $1 - \delta/4$, where (i) holds according to (4.15), and (ii) is valid since

$$\begin{aligned} \sum_{h=1}^H \sum_{k=1}^K (V_{h+1}^k(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k)) &\leq \sum_{h=1}^H \sum_{k=1}^K (V_{h+1}^{\text{UCB},k}(s_{h+1}^k) - V_{h+1}^{\pi^k}(s_{h+1}^k)) \\ &\lesssim \sqrt{H^7 SAK \log \frac{SAT}{\delta}} + H^3 SA, \end{aligned}$$

where the first inequality follows since $V_{h+1}^{\text{UCB},k} \geq V_{h+1}^k$ and $V_{h+1}^* \geq V_{h+1}^{\pi^k}$, and the second inequality comes from (E20a). Additionally, invoking Freedman's inequality (see Lemma 8) with $c_h = 1$ and $\tilde{W}_h^i = V_{h+1}^R - V_{h+1}^*$ (so that $0 \leq \tilde{W}_h^i(s) \leq H$) directly leads to

$$\left| \sum_{h=1}^H \sum_{k=1}^K (P_h^k - P_{h,s_h^k, a_h^k})(V_{h+1}^{\text{R},k} - V_{h+1}^*) \right| \lesssim \sqrt{TH^2 \log \frac{1}{\delta}} + H \log \frac{1}{\delta} \asymp \sqrt{H^3 K \log \frac{1}{\delta}}$$

with probability at least $1 - \delta/4$, which taken collectively with (E.25) reveals that

$$\begin{aligned} \sum_{h=1}^H \sum_{k=1}^K P_{s_h^k, a_h^k, h} (V_{h+1}^{\text{R},k} - V_{h+1}^*) &\leq \sum_{h=1}^H \sum_{k=1}^K P_h^k (V_{h+1}^{\text{R},k} - V_{h+1}^*) + \left| \sum_{h=1}^H \sum_{k=1}^K (P_h^k - P_{s_h^k, a_h^k, h})(V_{h+1}^{\text{R},k} - V_{h+1}^*) \right| \\ &\lesssim \sqrt{H^7 SAK \log \frac{SAT}{\delta}} + H^3 SA + HK \end{aligned} \quad (\text{E.26})$$

with probability at least $1 - \delta/2$. Substitution into (E.24) then gives

$$\begin{aligned} &\left| \sum_{h=1}^H \sum_{k=1}^K \lambda_h^k (P_h^k - P_{h,s_h^k, a_h^k})(V_{h+1}^* - V_{h+1}^{\text{R},k}) \right| \\ &\lesssim \sqrt{H^2 SA \sum_{h=1}^H \sum_{k=1}^K P_{h,s_h^k, a_h^k} (V_{h+1}^{\text{R},k} - V_{h+1}^*) \log \frac{T}{\delta} + H^2 SA \log \frac{T}{\delta}} \\ &\lesssim \sqrt{H^2 SA \left(\sqrt{H^7 SAK \log \frac{SAT}{\delta}} + H^3 SA + HK \right) \log \frac{T}{\delta} + H^2 SA \log \frac{T}{\delta}} \\ &\asymp \sqrt{H^2 SA \left(H^6 SA \log \frac{SAT}{\delta} + H^3 SA + HK \right) \log \frac{T}{\delta} + H^2 SA \log \frac{T}{\delta}} \\ &\lesssim \sqrt{H^3 SAK \log \frac{SAT}{\delta} + H^4 SA \log \frac{SAT}{\delta}} \\ &= \sqrt{H^2 SAT \log \frac{SAT}{\delta} + H^4 SA \log \frac{SAT}{\delta}} \end{aligned} \quad (\text{E.27})$$

with probability exceeding $1 - \delta$, where the third line holds since (due to Cauchy–Schwarz)

$$\sqrt{H^7 SAK \log \frac{SAT}{\delta}} = \sqrt{H^6 SA \log \frac{SAT}{\delta}} \sqrt{HK} \lesssim H^6 SA \log \frac{SAT}{\delta} + HK.$$

Step 2: upper bounding \mathcal{R}_3^2 . We start by making the following observation:

$$\begin{aligned}
\mathcal{R}_3^2 &\leq \sum_{h=1}^H \sum_{k=1}^K \frac{\lambda_h^k}{N_h^k(s_h^k, a_h^k)} \sum_{i \leq N_h^k(s_h^k, a_h^k)} (V_{h+1}^{\mathbb{R}, k^i}(s_{h+1}^{k^i}) - P_{h, s_h^k, a_h^k} V^{\mathbb{R}}, k_{h+1}) \\
&= \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (V^{\mathbb{R}}, k_{h+1}(s_{h+1}^k) - V^{\mathbb{R}}, k_{h+1}(s_{h+1}^k) + (P_h^k - P_{h, s_h^k, a_h^k}) V^{\mathbb{R}}, k_{h+1}) \\
&\leq (e \log T) \sum_{h=1}^H \sum_{k=1}^K (V^{\mathbb{R}}, k_{h+1}(s_{h+1}^k) - V^{\mathbb{R}}, k_{h+1}(s_{h+1}^k)) + \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) V_{h+1}^* \\
&\quad + \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V^{\mathbb{R}}, k_{h+1} - V_{h+1}^*), \tag{E.28}
\end{aligned}$$

where the first inequality comes from the monotonicity property $V^{\mathbb{R}}, k_{h+1} \geq V_{h+1}^{\mathbb{R}, k+1} \geq \dots \geq V^{\mathbb{R}}, k_{h+1}$, and the last line follows from the facts that $\sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{1}{n} \leq \log T$ and $\lambda_h^k \leq e$ (cf. (E.23)). In what follows, we shall control the three terms in (E.28) separately.

- The first term in (E.28) can be controlled by Lemma 4 (cf. (4.16)) as follows:

$$\sum_{h=1}^H \sum_{k=1}^K (V^{\mathbb{R}}, k_{h+1}(s_{h+1}^k) - V^{\mathbb{R}}, k_{h+1}(s_{h+1}^k)) \lesssim H^6 S A \log \frac{S A T}{\delta} \tag{E.29}$$

with probability at least $1 - \delta/3$.

- To control the second term in (E.28), we abuse the notation by setting

$$N(s, a, h) := N_h^{K-1}(s, a)$$

and

$$W_{h+1}^i := V_{h+1}^*, \quad \text{and} \quad u_h^i(s_h^i, a_h^i) := \sum_{n=N_h^i(s_h^i, a_h^i)}^{N(s_h^i, a_h^i, h)} \frac{\lambda_h^i}{n},$$

which clearly satisfy

$$|u_h^i(s_h^i, a_h^i)| \leq e \sum_{n=N_h^i(s_h^i, a_h^i)}^{N(s_h^i, a_h^i, h)} \frac{1}{n} \leq e \log T =: C_u \quad \text{and} \quad \|W_{h+1}^i\|_\infty \leq H =: C_w.$$

Here, we have used the properties $\sum_{n=N_h^i(s_h^i, a_h^i)}^{N_h^{K-1}(s_h^i, a_h^i)} \frac{1}{n} \leq \log T$ and $\lambda_h^k \leq e$ (cf. (E.23)). With these in place, applying Lemma 8 reveals that

$$\begin{aligned}
& \left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) V_{h+1}^* \right| = \left| \sum_{h=1}^H \sum_{k=1}^K X_{k,h} \right| \\
& \stackrel{(i)}{\lesssim} \sqrt{C_u^2 HSA \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} \left[|(P_h^i - P_{h, s_h^i, a_h^i}) W_{h+1}^i|^2 \right] \log \frac{T}{\delta} + C_u C_w HSA \log \frac{T}{\delta}} \\
& \stackrel{(ii)}{\lesssim} \sqrt{\sum_{h=1}^H \sum_{k=1}^K \text{Var}_{h, s_h^k, a_h^k}(V_{h+1}^*) \cdot HSA \log^3 \frac{T}{\delta} + H^2 SA \log^2 \frac{T}{\delta}} \\
& \stackrel{(iii)}{\lesssim} \sqrt{HSA(HT + H^4 \sqrt{SAT}) \log^4 \frac{SAT}{\delta} + H^2 SA \log^2 \frac{T}{\delta}} \\
& \lesssim \sqrt{HSA(HT + H^7 SA) \log^4 \frac{SAT}{\delta} + H^2 SA \log^2 \frac{T}{\delta}} \\
& \stackrel{(E.30)}{\lesssim} \sqrt{H^2 SAT \log^4 \frac{SAT}{\delta} + H^4 SA \log^2 \frac{SAT}{\delta}}
\end{aligned}$$

with probability at least $1 - \delta/3$. Here, (i) comes from the definition in (4.6), (ii) holds due to (E.17) and (E.18) and (iii) is valid since

$$HT + H^4 \sqrt{SAT} = HT + \sqrt{H^7 SA} \cdot \sqrt{HT} \lesssim HT + H^7 SA$$

due to the Cauchy–Schwarz inequality.

- Turning attention the third term of (E.28), we need to properly cope with the dependency between P_h^k and V_{h+1}^R , k . Toward this, we shall resort to the standard epsilon-net argument (see, e.g. [58]), which will be presented in Appendix E.3.1. The final bound reads like

$$\left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^{R,K} - V_{h+1}^*) \right| \lesssim H^4 SA \log^2 \frac{SAT}{\delta} + \sqrt{H^3 SAK \log^3 \frac{SAT}{\delta}}. \quad (E.31)$$

- Combining (E.29), (E.30) and (E.31) with (E.28), we can use the union bound to demonstrate that

$$\mathcal{R}_3^2 \leq C_{3,2} \left\{ H^6 SA \log^3 \frac{SAT}{\delta} + \sqrt{H^2 SAT \log^4 \frac{SAT}{\delta}} \right\} \quad (E.32)$$

with probability at least $1 - \delta$, where $C_{3,2} > 0$ is some constant.

Step 3: final bound of \mathcal{R}_3 . Putting the above results (E.27) and (E.32) together, we immediately arrive at

$$\mathcal{R}_3 \leq |\mathcal{R}_3^1| + \mathcal{R}_3^2 \leq C_{r,3} \left\{ H^6 SA \log^3 \frac{SAT}{\delta} + \sqrt{H^2 SAT \log^4 \frac{SAT}{\delta}} \right\} \quad (\text{E.33})$$

with probability at least $1 - 2\delta$, where $C_{r,3} > 0$ is some constant. This immediately concludes the proof.

E.3.1 Proof of (E31). **Step 1: concentration bounds for a fixed group of vectors.** Consider a fixed group of vectors $\{V_{h+1}^d \in \mathbb{R}^S \mid 1 \leq h \leq H\}$ obeying the following properties:

$$V_{h+1}^* \leq V_{h+1}^d \leq H \quad \text{for } 1 \leq h \leq H. \quad (\text{E.34})$$

We intend to control the following sum:

$$\sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^d - V_{h+1}^*).$$

To do so, we shall resort to Lemma 8. For the moment, let us take $N(s, a, h) := N_h^{K-1}(s, a)$ and

$$W_{h+1}^i := V_{h+1}^d - V_{h+1}^*, \quad u_h^i(s_h^i, a_h^i) := \sum_{n=N_h^i(s_h^i, a_h^i)}^{N(s_h^i, a_h^i, h)} \frac{\lambda_h^i}{n}.$$

It is easily seen that

$$|u_h^i(s_h^i, a_h^i)| \leq e \sum_{n=N_h^i(s_h^i, a_h^i)}^{N(s_h^i, a_h^i, h)} \frac{1}{n} \leq e \log T =: C_u \quad \text{and} \quad \|W_{h+1}^i\|_\infty \leq H =: C_w,$$

which hold due to the facts $\sum_{n=N_h^i(s_h^i, a_h^i)}^{N_h^{K-1}(s_h^i, a_h^i)} \frac{1}{n} \leq \log T$ and $\lambda_h^i \leq e$ (cf. (E.23)) as well as the property that $V_{h+1}^d(s), V_{h+1}^*(s) \in [0, H]$. Thus, invoking Lemma 8 yields

$$\begin{aligned} & \left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^d - V_{h+1}^*) \right| = \left| \sum_{h=1}^H \sum_{k=1}^K X_{k,h} \right| \\ & \stackrel{\mathcal{D}}{\sim} \sqrt{C_u^2 C_w \sum_{h=1}^H \sum_{i=1}^K \mathbb{E}_{i,h-1} [P_h^i W_{h+1}^i] \log \frac{K^{HSA}}{\delta_0} + C_u C_w \log \frac{K^{HSA}}{\delta_0}} \\ & \stackrel{\mathcal{D}}{\sim} \sqrt{H \sum_{h=1}^H \sum_{i=1}^K P_{h, s_h^i, a_h^i} (V_{h+1}^d - V_{h+1}^*) (\log^2 T) \log \frac{K^{HSA}}{\delta_0} + H (\log T) \log \frac{K^{HSA}}{\delta_0}} \end{aligned} \quad (\text{E.35})$$

with probability at least $1 - \delta_0$, where the choice of δ_0 will be revealed momentarily.

Step 2: constructing and controlling an epsilon net. Our argument in Step 1 is only applicable to a fixed group of vectors. The next step is then to construct an epsilon net that allows one to cover the set

of interest. Specifically, let us construct an epsilon net $\mathcal{N}_{h+1,\alpha}$ (the value of α will be specified shortly) for each $h \in [H]$ such that:

a) for any $V_{h+1} \in [0, H]^S$, one can find a point $V_{h+1}^{\text{net}} \in \mathcal{N}_{h+1,\alpha}$ obeying

$$0 \leq V_{h+1}(s) - V_{h+1}^{\text{net}}(s) \leq \alpha \quad \text{for all } s \in \mathcal{S};$$

b) its cardinality obeys

$$|\mathcal{N}_{h+1,\alpha}| \leq \left(\frac{H}{\alpha}\right)^S. \quad (\text{E.36})$$

Clearly, this also means that

$$|\mathcal{N}_{2,\alpha} \times \mathcal{N}_{3,\alpha} \times \cdots \times \mathcal{N}_{H+1,\alpha}| \leq \left(\frac{H}{\alpha}\right)^{SH}.$$

Set $\delta_0 = \frac{1}{6}\delta/\left(\frac{H}{\alpha}\right)^{SH}$. Taking (E.35) together the union bound implies that: with probability at least $1 - \delta_0\left(\frac{H}{\alpha}\right)^{SH} = 1 - \delta/6$, one has

$$\begin{aligned} & \left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^{\text{net}} - V_{h+1}^{\star}) \right| \\ & \lesssim \sqrt{H \sum_{h=1}^H \sum_{i=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^{\text{net}} - V_{h+1}^{\star}) (\log^2 T) \log \frac{K^{HSA}}{\delta_0}} + H(\log T) \log \frac{K^{HSA}}{\delta_0} \\ & \lesssim \sqrt{H^2 SA \sum_{h=1}^H \sum_{i=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^{\text{net}} - V_{h+1}^{\star}) (\log^2 T) \log \frac{SAT}{\delta\alpha}} + H^2 SA \log^2 \frac{SAT}{\delta\alpha} \end{aligned} \quad (\text{E.37})$$

simultaneously for all $\{V_{h+1}^{\text{net}} \mid 1 \leq h \leq H\}$ obeying $V_{h+1}^{\text{d}} \in \mathcal{N}_{h+1,\alpha}$ ($h \in [H]$).

Step 3: obtaining uniform bounds. We are now positioned to establish a uniform bound over the entire set of interest. Consider an arbitrary group of vectors $\{V_{h+1}^{\text{u}} \in \mathbb{R}^S \mid 1 \leq h \leq H\}$ obeying (E.34). By construction, one can find a group of points $\{V_{h+1}^{\text{net}} \in \mathcal{N}_{h+1,\alpha} \mid h \in [H]\}$ such that

$$0 \leq V_{h+1}^{\text{u}}(s) - V_{h+1}^{\text{net}}(s) \leq \alpha \quad \text{for all } (h, s) \in \mathcal{S} \times [H]. \quad (\text{E.38})$$

It is readily seen that

$$\begin{aligned}
& \left| \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^u - V_{h+1}^{\text{net}}) \right| \\
& \leq \left| \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} \left(\|P_h^k\|_1 + \|P_{h, s_h^k, a_h^k}\|_1 \right) \|V_{h+1}^u - V_{h+1}^{\text{net}}\|_\infty \right| \\
& \leq 2eK\alpha \log T,
\end{aligned} \tag{E.39}$$

where the last inequality follows from $\sum_{n=N_h^i(s_h^i, a_h^i)}^{N_h^{K-1}(s_h^i, a_h^i)} \frac{1}{n} \leq \log T$ and $\lambda_h^k \leq e$ (cf. (E.23)). Consequently, by taking $\alpha = 1/(SAT)$, we can deduce that

$$\begin{aligned}
& \left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^u - V_{h+1}^{\star}) \right| \\
& \leq \left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^{\text{net}} - V_{h+1}^{\star}) \right| \\
& \quad + \sum_{h=1}^H \left| \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^u - V_{h+1}^{\text{net}}) \right| \\
& \leq \left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^{\text{net}} - V_{h+1}^{\star}) \right| + HK\alpha \log T \\
& \lesssim \sqrt{H^2 SA \sum_{h=1}^H \sum_{i=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^{\text{net}} - V_{h+1}^{\star}) (\log^2 T) \log \frac{SAT}{\delta\alpha}} + H^2 SA \log^2 \frac{SAT}{\delta\alpha} + HK\alpha \log T \\
& \asymp \sqrt{H^2 SA \sum_{h=1}^H \sum_{i=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^u - V_{h+1}^{\star}) (\log^2 T) \log \frac{SAT}{\delta} + H^2 SA \log^2 \frac{SAT}{\delta}},
\end{aligned} \tag{E.40}$$

where the last line holds due to the condition (E.38) and our choice of α . To summarize, with probability exceeding $1 - \delta/6$, the property (E.40) holds simultaneously for all $\{V_{h+1}^u \in \mathbb{R}^S \mid 1 \leq h \leq H\}$ obeying (E.34).

Step 4: controlling the original term of interest. With the above union bound in hand, we are ready to control the original term of interest

$$\sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V^R, k_{h+1} - V_{h+1}^*). \quad (\text{E.41})$$

To begin with, it can be easily verified using (4.10) that

$$V_{h+1}^* \leq V^R, k_{h+1} \leq H \quad \text{for all } 1 \leq h \leq H. \quad (\text{E.42})$$

Moreover, we make the observation that

$$\begin{aligned} \sum_{h=1}^H \sum_{k=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^{R, k} - V_{h+1}^*) &\stackrel{\text{(i)}}{\leq} \sum_{h=1}^H \sum_{k=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^{R, k} - V_{h+1}^*) \\ &\stackrel{\text{(ii)}}{\leq} \sqrt{H^7 SAK \log \frac{SAT}{\delta}} + H^3 SA + HK \end{aligned} \quad (\text{E.43})$$

with probability exceeding $1 - \delta/6$, where (i) holds because V_{h+1}^R is monotonically non-increasing (in view of the monotonicity of $V_h(s)$ in (4.7b) and the update rule in line 16 of Algorithm 3), and (ii) follows from (E.26). Substitution into (E.40) yields

$$\begin{aligned} &\left| \sum_{h=1}^H \sum_{k=1}^K \sum_{n=N_h^k(s_h^k, a_h^k)}^{N_h^{K-1}(s_h^k, a_h^k)} \frac{\lambda_h^k}{n} (P_h^k - P_{h, s_h^k, a_h^k}) (V_{h+1}^{R, k} - V_{h+1}^*) \right| \\ &\lesssim \sqrt{H^2 SA \sum_{h=1}^H \sum_{i=1}^K P_{h, s_h^k, a_h^k} (V_{h+1}^{R, k} - V_{h+1}^*) (\log^2 T) \log \frac{SAT}{\delta} + H^2 SA \log^2 \frac{SAT}{\delta}} \\ &\lesssim \sqrt{H^2 SA \left\{ \sqrt{H^7 SAK \log \frac{SAT}{\delta}} + H^3 SA + HK \right\} (\log^2 T) \log \frac{SAT}{\delta} + H^2 SA \log^2 \frac{SAT}{\delta}} \\ &\lesssim \sqrt{H^2 SA \left\{ H^6 SA \log \frac{SAT}{\delta} + H^3 SA + HK \right\} \log^3 \frac{SAT}{\delta} + H^2 SA \log^2 \frac{SAT}{\delta}} \\ &\lesssim H^4 SA \log^2 \frac{SAT}{\delta} + \sqrt{H^3 SAK \log^3 \frac{SAT}{\delta}}, \end{aligned} \quad (\text{E.44})$$

where the penultimate line holds since

$$\sqrt{H^7 SAK \log \frac{SAT}{\delta}} = \sqrt{H^6 SA \log \frac{SAT}{\delta}} \sqrt{HK} \lesssim H^6 SA \log \frac{SAT}{\delta} + HK.$$