Feasible Q-Learning for Average Reward Reinforcement Learning

Ying Jin Stanford University Jose Blanchet
Stanford University

Ramki Gummadi Google

Zhengyuan Zhou² New York University

Abstract

Average reward reinforcement learning (RL) provides a suitable framework for capturing the objective (i.e. long-run average reward) for continuing tasks, where there is often no natural way to identify a discount factor. However, existing average reward RL algorithms with sample complexity guarantees are not feasible, as they take as input the (unknown) mixing time of the Markov decision process (MDP). In this paper, we make initial progress towards addressing this open problem. We design a feasible average-reward Qlearning framework that requires no knowledge of any problem parameter as input. Our framework is based on discounted Qlearning, while we dynamically adapt the discount factor (and hence the effective horizon) to progressively approximate the average reward. In the synchronous setting, we solve three tasks: (i) learn a policy that is ϵ -close to optimal, (ii) estimate optimal average reward with ϵ -accuracy, and (iii) estimate the bias function (similar to Q-function in discounted case) with ϵ -accuracy. We show that with carefully designed adaptation schemes, (i) can be achieved with $\widetilde{O}(\frac{SAt_{\text{mix}}^8}{\epsilon^8})$ samples, (ii) with $\widetilde{O}(\frac{SAt_{\text{mix}}^5}{\epsilon^5})$ samples, and (iii) with $\widetilde{O}(\frac{SAB}{s^9})$ samples, where $t_{\rm mix}$ is the mixing time, and B>0 is an MDP-dependent constant. To our knowledge, we provide the first finite-sample guarantees that are polynomial in $S, A, t_{\text{mix}}, \epsilon$ for a feasible variant of Qlearning. That said, the sample complexity bounds have tremendous room for improve-

Proceedings of the 27th International Conference on Artificial Intelligence and Statistics (AISTATS) 2024, Valencia, Spain. PMLR: Volume 238. Copyright 2024 by the author(s).

ment, which we leave for the community's best minds. Preliminary simulations verify that our framework is effective without prior knowledge of parameters as input.

1 Introduction

Reinforcement learning (RL) has achieved remarkable success in simulated environments such as beating world human champions in playing poker, chess and Go (Brown and Sandholm, 2018; Schrittwieser et al., 2020; Silver et al., 2018; Mnih et al., 2015; Schaeffer et al., 1992; Campbell et al., 2002). Such empirical successes—and the excitements generated therefrom—have motivated a remarkably fruitful line of research on the sample complexity of various RL algorithms, which characterizes how many samples (from a generative model such as a simulator) are needed to obtain an ϵ -optimal policy.

Regarding sample complexity, an important setting is the discounted infinite-horizon RL, where the goal is to maximize the (infinite) sum of all future γ -discounted rewards for a given factor 0 < γ < 1. Broadly speaking, there are (at least) two main classes of methods to tackle the problem: model-based methods and model-free methods. Model-free algorithms learn to select actions without model estimation. Compared with model-based ones (Azar et al., 2013; Sidford et al., 2018b,a; Wang, 2020; Agarwal et al., 2020), they are often more computationally efficient, have less storage overhead, and are easy to generalize to RL with function approximation (Sutton and Barto, 2018; François-Lavet et al., 2018). In particular, Q-learning (Watkins and Dayan, 1992), as the prototypical model-free algorithm, has been studied³ extensively in discounted infinite-horizon RL (Kearns and Singh, 1999; Even-Dar et al., 2003; Beck and Srikant, 2012; Wainwright, 2019a; Chen et al., 2020; Li et al., 2021).

 $^{^2{\}rm This}$ work is generously supported by National Science Foundation grants CCF-2312205 and CCF-2312204. Correspond to ying 531@stanford.edu and zz26@nyu.edu.

³Earlier works on discounted *Q*-learning focused on characterizing its asymptotic convergence without any finite-sample guarantees; see Jaakkola et al. (1994); Tsitsiklis (1994); Szepesvári (1998); Borkar and Meyn (2000).

Such remarkable research efforts notwithstanding, average-reward infinite-horizon RL remains far from being resolved: As recognized by the RL community (Sutton and Barto, 2018; Wan et al., 2021; Mahadevan, 1996; Dewanto et al., 2020), average reward RL, which provides a more natural objective for many continuing tasks, is a much more challenging problem and has remained largely under-explored. The sparsity of theoretical understanding also limits the implementation. Indeed, discount factors are often used in practical instantiations of RL algorithms, even when the final objective of interest is clearly the average undiscounted reward in the long term. This is partly due to the fact that algorithms to directly optimize average reward are much more challenging to characterize than the discounted ones. As a result, popular RL algorithms make an ad-hoc choice of a fixed discount factor that are only partially understood (Tang et al., 2021). However, many RL applications do not have a natural discount factor that can be endogenously identified.

Classical results (Blackwell, 1962; Mahadevan, 1996) connect the discounted and average reward formulations via fundamental relations between optimal value functions asymptotically, yet there is only limited finite-sample guarantees related to algorithmic transfer between the two frameworks, given the ubiquitous practice of implementing the former to solve the latter. Recent works leverage similar relations to devise algorithms that optimizes the γ -discounted reward to approximate optimal average reward for a carefully designed, fixed γ (Jin and Sidford, 2020, 2021). However, these (model-based) algorithms that have finitesample guarantees all require knowing the mixing time t_{mix} and are hence *infeasible*. In parallel, finite-sample guarantees for model-free methods (sometimes called R-learning, see Section 1.1) for average reward are even more scarce; the only one we know of is very recent (Zhang et al., 2021), and the algorithm therein also relies on unknown parameters as input. As such, and in light of the merits of model-free RL algorithms mentioned before, we are naturally led to the open question:

Can we design a feasible model-free averagereward RL algorithm with finite-sample guarantees for the sample complexity?

In this paper, we provide new theoretical results that highlight an important algorithmic role for the discounted formulation as a subroutine in solving several learning tasks in the average reward scenario. In contrast to the typical use of a fixed discount factor, we highlight the importance of a carefully selected schedule for progressively finer discount factors to obtain a solution to the average reward objective. Our high-level idea is in concordance with Hordijk and Tijms (1975) which provides an asymptotic analysis for such paradigm when planning with knowledge of the MDP; in contrast, we provide finite-sample analysis for a learning problem, which requires judiciously chosen schemes of both discount factors and learning rates to tackle random data.

1.1 Related Work

We include more detailed discussion of related work in this part, extending Section 1.2 in the main text.

Model-based average-reward RL. (Kearns and Singh, 2002) is an early work on average-reward RL that proposes a model-based learning algorithm and establishes a sample complexity bound of $\widetilde{O}(\frac{S^5 \operatorname{poly}(A) t_{mix}^5}{\epsilon^6})$, where $\operatorname{poly}(A)$ is an unspecified polynomial of A. From an algorithmic standpoint, an important issue of the algorithm is that it requires the knowledge of both $t_{\rm mix}$ and the optimal average reward⁴. More recently, several different model-based average reward RL algorithms have been proposed Wang (2017); Zhang and Xie (2023); Jin and Sidford (2020, 2021), with the latter two achieving state-of-the-art sample-complexity bounds of $\widetilde{O}(\frac{SAt_{mix}^2}{\epsilon^2})$ and $\widetilde{O}(\frac{SAt_{mix}}{\epsilon^3})$, respectively. All three algorithms rely on knowing t_{mix} or similar quantities. Jin and Sidford (2021) further provides a lower bound of $\widetilde{O}(\frac{SAt_{mix}}{\epsilon^2})$ for average-reward RL when t_{mix} is known. Consequently, while having the pleasing optimal dependence on SA, a minimax optimal algorithm for average-reward RL is not yet known even under known t_{mix} , an impractical assumption to begin with.

Model-free average-reward RL. More recently, driven by the merits of model-free algorithms, Zhang et al. (2021) provides the very first finite-sample analysis of a average reward Q-learning variant, yielding a sample complexity bound of $\widetilde{O}(\frac{SAJ^3}{(1-\delta)^5\epsilon^2})$, where J and δ are two unknown MDP-dependent constants that may arbitrarily depend on t_{mix} . This pioneering bound itself is valuable in light of the difficulty in characterizing finite-sample guarantees of Q-learning algorithms for average reward RL. However, the issue is that Zhang et al. (2021) assumes J and δ are known and given, rendering the algorithm infeasible.

Feasible average-reward RL with asymptotic guarantees. The only feasible average-reward Q-learning (sometimes also called R-learning) variants

⁴How to learn this quantity efficiently also remains under-explored. We provide an answer on this as well.

that we know of are Wan et al. (2021); Abounadi et al. (2001), which attempt to directly learn the Q-values (called bias function in average-reward RL) associated with the Bellman equation for average reward MDP (discussed in more detail in Section 2). However, for both algorithms, only asymptotic consistency is established. It is largely unclear whether finite-sample sample complexity bounds can be established for Rlearning algorithms, since the underlying update is not a contraction. Remotely related is Hordijk and Tijms (1975), which shares a similar spirit as our methods for gradually adjusting discount factors to approximate average reward; however, it requires the exact knowledge of the MDP and hence solves a planning – rather than learning - problem (see the next subsection for more details), and only provide asymptotic analysis.

Other metrics for average-reward RL. Finally, we mention in passing that the literature has also studied other metrics such as regret Jaksch et al. (2010); Jin et al. (2018); Dong et al. (2019); Fruit et al. (2020); Dong et al. (2021); Wei et al. (2020), which is not the focus here. While there are online-to-batch tricks to turn regrets to sample complexity (e.g., via a randomly sample from history of policies), they can not tell which policy is ϵ -close to optimal; this deviates from our goal of designing a practically feasible algorithm and analyzing its sample complexity. Instead, we focus on algorithms with a deterministic output given data.

1.2 Our Contributions and Related Work

First, we design a feasible average-reward Q-learning algorithmic framework that requires no knowledge of any problem parameter. In contrast to the aforementioned Q-learning variants for average-reward RL Wan et al. (2021); Abounadi et al. (2001); Zhang et al. (2021) that all aim to directly learn the bias function (which determines the optimal policy; see Section 2) in the average-reward Bellman equation, our algorithmic framework uses discounted Q-learning but dynamically adjusts the discount factor towards 1 – and hence gradually enlarging the effective horizon – to yield progressively finer approximations of the average-reward setting. This idea is quite simple and intuitive; however, the challenge lies in designing the specific horizon adaptation scheme (involving the simultaneous adjustment of learning rate and discount factor; see Algorithm 1) such that the finite-sample analysis goes through. Note that unlike in the discounted setting with a fixed γ , we now have a "moving target" problem as the discount factor is constantly shifting and weaker contraction, thereby making the analysis much more difficult. On this note, an early work back to 1970s Hordijk and Tijms (1975) shares a similar idea of adjusting discount factors; however, it only provides asymptotic guarantees (and does so without any specific adaptation scheme); further, it requires the exact knowledge of the MDP and hence solves a planning problem.

Second, we offer two concrete instantiations of our framework under (distinct) judiciously chosen adaptation schemes. Using the first scheme, stopping our algorithm at any iteration T yields a policy whose average reward differs from the optimal by at most $\widetilde{O}(t_{\rm mix}/T^{1/8})$. Under the synchronous setting we consider (see Section 2), it translates to a sample complexity bound of $\widetilde{O}(SAt_{mix}^8/\epsilon^8)$. Furthermore, if the goal is to simply estimate the optimal average reward, then we can do so at a faster rate by our second scheme: stopping at any iteration T yields an estimation accuracy of $\widetilde{O}(t_{\rm mix}/T^{1/5})$, which translates to a sample complexity of $\widetilde{O}(SAt_{\rm mix}^5/\epsilon^5)$. These results also indicate that our algorithm is any time: we do not need to know T beforehand when running the algorithm⁵. To the best of our knowledge, these are the first finitesample guarantees for a feasible Q-learning algorithm. That said, we believe these bounds can be improved; and we view our results as an open invitation for further progress in feasible average-reward RL.

Additionally, although not the main focus here, we also extend our algorithmic framework to estimate the bias function in the average-reward MDP Bellman equation (Section 5). As mentioned before, existing average-reward Q-learning algorithms Wan et al. (2021); Abounadi et al. (2001); Zhang and Ji (2019); Zhang et al. (2021) learn a policy through estimating the bias function. Viewed through this lens, we provide an alternative way for estimating this quantity that has a sample complexity of $\widetilde{O}(\frac{SAB}{\epsilon^9})$. Similar to Zhang et al. (2021) (which is the only work that has a finite-sample guarantee for estimating the bias function), our bound has an unpleasing dependence on an unknown MDP-dependent constant B. However, different from Zhang et al. (2021), our algorithm itself is feasible (and does not need to know B or any other problem parameter). Note that the two bounds – ours and the one in Zhang et al. (2021) – are incomparable as all of those problem-dependent constants are unknown.

Finally, through preliminary simulations, we verify that our algorithm learns the near-optimal policy with satisfactory convergence, and performs well across MDPs with various values of mixing time without the need of any prior knowledge of the mixing time.

 $^{^{5}}$ Algorithms without this property cannot provably adapt to additional samples beyond an initially chosen T.

2 Problem Setup

We consider an infinite-horizon tabular MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r)$, with finite state space $\mathcal{S} = \{1, \ldots, S\}$, finite action space $\mathcal{A} = \{1, \ldots, A\}$, transition probability $P \colon \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$ (i.e., $P(s' \mid s, a)$ is the probability of transiting to s' from a state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$), and reward function $r \colon \mathcal{S} \times \mathcal{A} \to [0, 1]$ (i.e., r(s, a) is the immediate reward at state $s \in \mathcal{S}$ if action $a \in \mathcal{A}$ is taken). We let $\pi \colon \mathcal{S} \to \Delta(\mathcal{A})$ denote a policy, i.e., $\pi(a \mid s)$ is the probability of taking a at state s. When π is a deterministic policy, $\pi(s)$ denotes the action chosen at state s.

Learning objective. Given a policy π , we define its long-term average reward as

$$V^{\pi}(s) = \liminf_{T \to \infty} \mathbb{E}_{\pi} \left[\frac{1}{T} \sum_{k=1}^{T} r(s_k, a_k) \, \middle| \, s_1 = s \right]$$

for all $s \in \mathcal{S}$. Here \mathbb{E}_{π} denotes the expectation over the trajectory $\{(s_k, a_k)\}_{k \geq 0}$ of the MDP under policy π . Under sufficient generality⁷, the standard MDP theory (Puterman, 2014) shows that for any policy π , there exists a constant $J^{\pi} \in [0, 1]$ such that $V^{\pi}(s) = J^{\pi}$ for all $s \in \mathcal{S}$. The long-term average reward of π initialized at any state-action pair also equals J^{π} , that is, $J^{\pi} \equiv \liminf_{T \to \infty} \mathbb{E}_{\pi} \left[\frac{1}{T} \sum_{k=1}^{T} r(s_k, a_k) \mid s_1 = s, a_1 = a \right]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. We use $J^{\pi} \in [0, 1]$ to denote the average reward of policy π .

The optimal policy π^* for average reward attains $J^{\pi^*} = \max_{\pi} J^{\pi}$, and we denote $J^* = J^{\pi^*}$. Furthermore, there exists a function $q^* \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ so that the following Bellman equation holds for all $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$J^* + q^*(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot \mid s, a)} [v^*(s')], \qquad (1)$$

where $v^*(s') = \max_{a' \in \mathcal{A}} q^*(s', a')$ for all $s' \in \mathcal{S}$. The optimal policy π^* is greedy w.r.t. q^* , i.e., $\pi^*(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}}_{a \in \mathcal{A}} q^*(s, a)$ for all $s \in \mathcal{S}$. The solution (q^*, v^*) to (1) is unique up to a constant; one solution is

$$v^*(s) = \mathbb{E}_{\pi^*} \left[\sum_{k=1}^{\infty} \left(r(s_k, a_k) - J^* \right) \, \middle| \, s_1 = s \right], \tag{2}$$

$$q^*(s, a) = \mathbb{E}_{\pi^*} \left[\sum_{k=1}^{\infty} (r(s_k, a_k) - J^*) \mid s_1 = s, a_1 = a \right],$$

where v^* and q^* are called the value bias function and q-value bias function respectively (Puterman, 2014).

In this work, we will consider three tasks: (1) learning the optimal reward J^* up to ϵ -accuracy, (2) learning a

policy whose average reward is ϵ -close to J^* , and (3) learning the bias functions q^* and v^* up to constants.

2.1 Sampling Scheme: Synchronous Setting with a Generative Model

Throughout this paper, we work under the synchronous scenario with a generative model (or simulator) (Even-Dar et al., 2003). That is, we consider algorithms that proceed in multiple iterations. In each iteration t, we receive an independent sample $s' \sim P(\cdot \mid s, a)$ for all state-action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$. Despite being relatively simple, it serves as an idealistic sampling protocol that has received much attention for various RL algorithms throughout the RL literature (Kearns et al., 2002; Kakade, 2003; Even-Dar et al., 2003; Azar et al., 2013; Beck and Srikant, 2012; Sidford et al., 2018a,b: Wainwright, 2019a,b: Yang and Wang, 2019; Zanette et al., 2019; Agarwal et al., 2020; Li et al., 2021). We focus on this setting as a starting point for studying sample complexity of feasible average-reward RL algorithms.

2.2 Mixing time and other related notions

Definition 2.1. The mixing time t_{mix} of an MDP is

$$\max_{\pi} \min \left\{ t \colon \max_{q \in \Delta(\mathcal{S})} d_{\text{TV}} \left((P^{\pi})^t(q), \nu^{\pi} \right) \le 1/4 \right\}, \quad (3)$$

where $d_{\text{TV}}(\mu, \nu) := \frac{1}{2} \sum_{s \in \mathcal{S}} |\mu(s) - \nu(s)|$ is the total variation distance, and $(P^{\pi})^t(q)$ is the distribution of s_t induced by policy π with initial distribution $s_0 \sim q$. We suppose for any policy π there exists a stationary distribution $\nu^{\pi} \in \Delta(\mathcal{S})$; otherwise $t_{\text{mix}} = \infty$.

The notion of mixing time (3) is widely adopted in the literature of average reward MDPs (Wei et al., 2020; Jin and Sidford, 2021), and our assumption is standard in the literature Wang (2017); Jin and Sidford (2020). While we adopt such notion to be consistent with the literature, our analysis framework also works under other regularity conditions.

Remark 2.2. Our results remain true if the mixing time (3) is replaced by the reward averaging time similar to De Farias and Van Roy (2006); Dong et al. (2021); the latter is also closely related to the notion of averaging time considered in Kearns and Singh (2002). To be specific, the reward averaging time is defined as $\tau = \max_{\pi} \sup_{T \geq 1} \left\{ T \cdot |J^{\pi} - V^{\pi}(s,T)| \right\}$, where $V^{\pi}(s,T) = \mathbb{E}_{\pi} \left[\frac{1}{T} \sum_{t=1}^{T} r(s_{t},a_{t}) \, \middle| \, s_{1} = s \right]$ is the T-step average reward value function of π . Our results carry over to this setting with t_{mix} replaced by $\tau < \infty$ in the upper bound.

Remark 2.3. Our framework also applies to weakly communicating MDPs (Wei et al., 2020). Define the

⁷Puterman (2014) shows that if the state space S is finite or countable, then the limit (instead of liminf) $V^{\pi}(s)$ exists; if the chain induced by π is irreducible or has a single recurrent class, then $V^{\pi}(s)$ is a constant function. We operate in this scenario.

span of the optimal γ -discounted value function as $\operatorname{sp}(\gamma) = (\max_s V_\gamma^*(s) - \min_s V_\gamma^*(s))/(1-\gamma)$, where V_γ^* is the optimal discounted value function with a scaling; the exact definition of V_γ^* is deferred to (4). It is argued in Wei et al. (2020) that the span is bounded by the diameter of the MDP (Lattimore and Szepesvári, 2020) for weakly-communicating MDPs. Our bound applies to MDPs with bounded spans with t_{mix} replaced by $\operatorname{sup}_{\gamma \in (0,1)} \operatorname{sp}(\gamma)$.

3 Algorithm: Dynamic Horizon Q-Learning

3.1 Recap: Discounted Q-Learning

We start with a brief recap on the celebrated Q-learning algorithm for discounted value functions. Given any discount factor $\gamma \in (0,1)$ and any policy π , with rewards $r_k = r(s_k, a_k)$, we denote the rescaled γ -discounted value function and Q-function of π as

$$V_{\gamma}^{\pi}(s) = (1 - \gamma) \mathbb{E}_{\pi} \Big[\sum_{k=1}^{\infty} \gamma^{t-1} r_k \, \big| \, s_1 = s \Big],$$
$$Q_{\gamma}^{\pi}(s, a) = (1 - \gamma) \mathbb{E}_{\pi} \Big[\sum_{k=1}^{\infty} \gamma^{k-1} r_k \, \big| \, s_1 = s, a_1 = a \Big]$$

for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. The optimal discounted value function V_{γ}^* and Q-function Q_{γ}^* are

$$V_{\gamma}^*(s) = \max_{\pi} V_{\gamma}^{\pi}(s), \quad Q_{\gamma}^*(s,a) = \max_{\pi} Q_{\gamma}^{\pi}(s,a) \quad (4)$$

for $(s, a) \in \mathcal{S} \times \mathcal{A}$. For preparation, we also denote the unscaled γ -discounted value and Q-functions of π as

$$v_{\gamma}^{\pi} = V_{\gamma}^{\pi}/(1-\gamma), \quad q_{\gamma}^{\pi} = Q_{\gamma}^{\pi}/(1-\gamma),$$

and the optimal ones as

$$v_{\gamma}^* = V_{\gamma}^*/(1-\gamma), \quad q_{\gamma}^* = Q_{\gamma}^*/(1-\gamma).$$
 (5)

The optimal policy for γ -discounted reward is denoted as π_{γ}^* . It is well-known that π_{γ}^* is a deterministic policy with $\pi_{\gamma}^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q_{\gamma}^*(s, a) = \operatorname{argmax}_{a \in \mathcal{A}} q_{\gamma}^*(s, a)$.

In synchronous setting, the Q-learning algorithm for γ -discounted rewards maintains an estimate $Q_t \colon \mathcal{S} \times \mathcal{A} \to [0,1]$ for the optimal Q-function Q_{γ}^* . In each iteration t, it updates all entries of the estimate at once, according to $Q_t = (1 - \eta_t)Q_{t-1} + \eta_t \mathcal{T}_t(Q_{t-1})$. Here $\eta_t \in (0,1]$ is the learning rate or the step size, and \mathcal{T}_t is the empirical Bellman operator depending on the samples collected in the t-th iteration (with proper scaling): $\mathcal{T}_t(Q)(s,a) = (1-\gamma)r(s,a) + \gamma \max_{a' \in \mathcal{A}} Q(s',a')$, where $s' \sim P(\cdot | s,a)$ is the independent sample collected for (s,a) from the generative model.

3.2 Dynamic Horizon Q-Learning Framework

In the existing works of Wei et al. (2020); Jin and Sidford (2021), algorithms for the discounted setting are applied with a properly chosen discounted factor γ that depends on the *known* mixing time and a prespecified sample size. Rather distinct from them, we avoid the knowledge of the mixing time and a prespecified sample size by employing a series of dynamic discount factors.

To be specific, given a sequence of discount factors $\{\gamma_t\}_{t\geq 1}$, we maintain an estimate $Q_t: \mathcal{S}\times\mathcal{A}\to\mathbb{R}$ in the t-th iteration. In each iteration, the algorithm updates all entries of the Q-function estimate via

$$Q_t(s, a) = (1 - \eta_t)Q_{t-1}(s, a) + \eta_t \left[(1 - \gamma_t)r(s, a) + \gamma_t \max_{a' \in \mathcal{A}} Q_{t-1}(s_t(s, a), a') \right]$$
(6)

for all state-action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$. Here $s_t(s, a)$ is the independent sample from the generative model, γ_t is the discount factor and $\eta_t \in (0, 1]$ is the learning rate in the t-th iteration. Correspondingly, we define the estimate of value function in the t-th iteration as $V_t(s) := \max_{a \in \mathcal{A}} Q_t(s, a)$ for all $s \in \mathcal{S}$, and the associated greedy policy as

$$\pi_t(s) := \operatorname*{argmax}_{a \in A} Q_t(s, a), \tag{7}$$

so that $V_t(s) = Q_t(s, \pi_t(s))$ for all $s \in \mathcal{S}$. The complete algorithm is summarized as follows.

Algorithm 1 Dynamic Horizon Q-Learning

 $\{\eta_t\}_{t\geq 1}, \{\gamma_t\}_{t\geq 1}$. Initialization: $Q_0 \equiv 0$. for $t=1,2,\ldots$ do

- 3: Generate $s_t(s, a) \sim P(\cdot \mid s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.
- 4: Set $Q_t(s, a) = (1 \eta_t)Q_{t-1}(s, a) + \eta_t[(1 \gamma_t)r(s, a) + \gamma_t V_{t-1}(s_t(s, a))], \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$
- 5: Set $V_t(s) = \max_{a \in \mathcal{A}} Q_t(s, a)$ for all $s \in \mathcal{S}$.
- 6: end for

3.3 Theory for estimating optimal reward

We provide theoretical guarantee of Algorithm 1 for learning the optimal reward J^* as follows. The proof of Theorem 3.1 is sketched in Section 4.1 and 4.2, while detailed in Appendix A.1.

Theorem 3.1. In Algorithm 1, we set $\eta_t = (1 + \frac{c_1 t^{3/5}}{(\log t)^3})^{-1}$, $\gamma_t = 1 - t^{-1/5}$, $t \geq 2$, for some constant $c_1 > 0$ and set $\eta_1 = \eta_2$, $\gamma_1 = \gamma_2$. Let $\varepsilon \in (0,1)$ and $\delta \in (0,1)$. Suppose T is sufficiently large such that $T/\log T \geq 300$, $T^{1/5}(\log T)^2 \geq 4c_2$, $(\log T)^2 \geq 12(10 + c_1)$, $c_2 T^{1/5} \geq 24(10 + c_1)(\log T)^{2/5}$

and $T^{2/5} \geq \frac{64(\log T)^4}{9c_1}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}$ for some constant $c_2 > 0$. Then with probability at least $1 - \delta$, after T iterations, Algorithm 1 achieves

$$\left| V_T(s) - J^* \right| \le \frac{c \cdot t_{mix} (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}}$$

simultaneously for all $s \in \mathcal{S}$, where $J^* \in [0,1]$ is the optimal average reward defined in Section 2, and c > 0 is an absolute constant that only depends on c_1, c_2 .

We take a moment to discuss the general idea of our learning framework. At a high level, we approximate the ultimate targets (i.e., the optimal average reward J^*) with a sequence of proxies $\{Q_{\gamma_t}^*\}$ that our estimates $\{Q_t\}$ eventually converge to. In particular,

$$\left| Q_t(\cdot, \cdot) - J^* \right| \le \left| Q_{\gamma_t}^*(\cdot, \cdot) - J^* \right| + \left| Q_t(\cdot, \cdot) - Q_{\gamma_t}^*(\cdot, \cdot) \right|$$
(8)

for all inputs in $S \times A$. The first term in (8) is the approximation error by discounted Q-functions that is controlled as follows (proof in Appendix A.2).

Lemma 3.2.
$$|V_{\gamma}^{*}(s) - J^{*}| \leq 3(1 - \gamma)t_{mix}$$
 and $|Q_{\gamma}^{*}(s, a) - J^{*}| \leq 3(1 - \gamma)t_{mix}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

The second term of (8) is the estimation error of $\{Q_t\}$ to the dynamic targets $\{Q_{\gamma_t}^*\}$. We provide the sketch of a recursive analysis for this term in Section 4.1.

3.4 Theory for Learning ϵ -Optimal Policy

Algorithm 1 can be adapted to learn a policy whose average reward is ϵ -close to optimal. The idea is still to take the greedy policy from our estimate Q_t . However, due to the subtlety between Q-function approximation and the value of the greedy policy, we need another adaptation scheme of $\{\gamma_t\}$ and $\{\eta_t\}$. A partial proof sketch of Theorem 3.3 is in Section 4.1; the detailed proof is in Appendix B.

Theorem 3.3. In Algorithm 1, we set $\eta_t = \left(1 + \frac{c_1 t^{5/8}}{(\log t)^2}\right)^{-1}$, $\gamma_t = 1 - t^{-1/8}$, $t \geq 2$ for some constant $c_1 > 0$ and set $\eta_1 = \eta_2$, $\gamma_1 = \gamma_2$. Let π_t be the greedy policy with respect to Q_t from Algorithm 1 for all $t \geq 1$. Let $\varepsilon \in (0,1)$ and $\delta \in (0,1)$. Suppose T is sufficiently large such that $(\log T)^2 \geq 11(2+c_1)/2$, $T/\log T \geq 100$, $c_2 T^{1/4} \geq 4(2+c_1)(\log T)^{11/8}$ and $T^{5/8} \geq \frac{64(\log T)^3}{9c_1}\log \frac{|S||A|T}{\delta}$ for some constant $c_2 > 0$. Then with probability at least $1 - \delta$, after T iterations, π_T in Algorithm 1 satisfies

$$\left| J^{\pi_T} - J^* \right| \le \frac{c \cdot t_{mix} (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/8}}.$$

4 Analysis Framework

We provide the analysis framework for Theorems 3.1 and 3.3. In Section 4.1, we provide a general decomposition of the estimation error; we then give the proof

sketch for Theorem 3.1 in Section 4.2. Theorem 3.3 follows similar ideas, with details in Appendix B.

Notations. Bold letters denote vectors and matrices. For any matrix M, we denote $||M||_1 =$ $\max_{i} \sum_{j} |M_{ij}|$ and $||\boldsymbol{M}||_{\infty} = \max_{i,j} |M_{ij}|$. For vectors $\boldsymbol{a} = [a_i], \boldsymbol{b} = [b_i] \in \mathbb{R}^n, \, \boldsymbol{a} \leq \boldsymbol{b} \text{ (resp. } \boldsymbol{a} \geq \boldsymbol{b} \text{) means } a_i \leq \boldsymbol{b}$ b_i (resp. $a_i \geq b_i$) for all i. We let $\boldsymbol{a} \circ \boldsymbol{b} = [a_i b_i] \in \mathbb{R}^n$. For vector $\mathbf{a} = [a_i]$, we denote $|\mathbf{a}| = [|a_i|]$. For a group of vectors $\{a_i: i \in \mathcal{I}\}$, we denote $\max_{i \in \mathcal{I}} a_i$ as the vector of entrywise maximum, that is, $[\max_{i \in \mathcal{I}} a_i]_i =$ $\max_{i \in \mathcal{I}} [a_i]_j$, $\forall j$. We use vector $r \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ to represent reward functions, so that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, the (s,a)-th entry of r is given by r(s,a). We represent value and Q-functions in vectors: for example, the s-th entry of $V^{\pi} \in \mathbb{R}^{|\mathcal{S}|}$ is given by $V^{\pi}(s)$; we define $Q_t, Q^{\pi}, Q^{\pi}_{\gamma}, Q^{*}_{\gamma}, \mathbf{q}_t, \mathbf{q}^{\pi}, \mathbf{q}^{*}, \mathbf{q}^{\pi}_{\alpha}, \mathbf{q}^{*}_{\alpha} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, and $V_t, V^{\pi}, V_{\gamma}^{\pi}, V_{\gamma}^{*}, \mathbf{v}_t, \mathbf{v}^{\pi}, \mathbf{v}^{*}, \mathbf{v}_{\alpha}^{\pi}, \mathbf{v}_{\alpha}^{*} \in \mathbb{R}^{|\mathcal{S}|}$, analogously. We use a matrix $P \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|\times|\mathcal{S}|}$ to represent the probability transition kernel P, whose (s, a)-row $P_{s,a}$ represents the vector $P(\cdot | s, a)$. For any vector $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$, we define $Var_{\mathbf{P}}(\mathbf{V}) = \mathbf{P}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}\mathbf{V}) \circ (\mathbf{P}\mathbf{V}) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|},$ that is, the (s, a)-entry of $Var_{\mathbf{P}}(\mathbf{V})$ is Var(V(s')) for $s' \sim P(\cdot | s, a)$. We also define the square probability transition matrix $P^{\pi} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ (resp. $P_{\pi} \in$ $\mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$) induced by a deterministic policy π over the state-action pairs (resp. states) as $P^{\pi} := P\Pi^{\pi}$, $P_{\pi} := \Pi^{\pi} P$, where $\Pi^{\pi} \in \{0,1\}^{|\mathcal{S}| \times |\mathcal{S}| |\mathcal{A}|}$ is the projection matrix associated with π , whose s-th row consists of $|\mathcal{S}|$ blocks each of length $|\mathcal{A}|$, among which the s-th block is $e_{\pi(s)}^{\top}$, and e_i is the *i*-th standard basis vector. Given samples $s_t(s,a) \sim P(\cdot | s,a)$ collected in the t-th iteration, we define the empirical transition matrix $P_t \in \{0,1\}^{|S||A| \times |S|}$ by $P_t((s,a),s') = \mathbb{1}\{s' =$ $s_t(s,a)$.

4.1 Framework of Analysis for Algorithm 1

In this section, we provide a sketch of analysis for Algorithm 1. To begin with, we denote $\Delta_t = Q_t - Q_{\gamma_t}^*$, which is the estimation error of Q_t for the rescaled optimal Q-function with γ_t in the t-th iteration. Our updating rule (6) in the t-th iteration admits the representation:

$$\mathbf{Q}_t = (1 - \eta_t)\mathbf{Q}_{t-1} + \eta_t [(1 - \gamma_t)\mathbf{r} + \gamma_t \mathbf{P}_t \mathbf{V}_{t-1}].$$

Employing the Bellman equation $Q_{\gamma_t}^* = (1 - \gamma_t)r + \gamma_t P V_{\gamma_t}^*$ (note the rescaling of the value and Q-functions compared to conventional notations),

$$\begin{split} & \boldsymbol{\Delta}_{t} = \boldsymbol{Q}_{t} - \boldsymbol{Q}_{\gamma_{t}}^{*} \\ &= (1 - \eta_{t}) \boldsymbol{Q}_{t-1} + \eta_{t} \big[(1 - \gamma_{t}) \boldsymbol{r} + \gamma_{t} \boldsymbol{P}_{t} \boldsymbol{V}_{t-1} \big] - \boldsymbol{Q}_{\gamma_{t}}^{*} \\ &= (1 - \eta_{t}) \big[\boldsymbol{Q}_{t-1} - \boldsymbol{Q}_{\gamma_{t}}^{*} \big] \\ &+ \eta_{t} \big[(1 - \gamma_{t}) \boldsymbol{r} + \gamma_{t} \boldsymbol{P}_{t} \boldsymbol{V}_{t-1} - \boldsymbol{Q}_{\gamma_{t}}^{*} \big] \end{split}$$

$$= (1 - \eta_t) [\mathbf{Q}_{t-1} - \mathbf{Q}_{\gamma_t}^*]$$

$$+ \eta_t [(1 - \gamma_t)\mathbf{r} + \gamma_t \mathbf{P}_t \mathbf{V}_{t-1} - (1 - \gamma_t)\mathbf{r} - \gamma_t \mathbf{P} \mathbf{V}_{\gamma_t}^*]$$

$$= (1 - \eta_t) \Delta_{t-1} + (1 - \eta_t) [\mathbf{Q}_{\gamma_{t-1}}^* - \mathbf{Q}_{\gamma_t}^*]$$

$$+ \eta_t \gamma_t [\mathbf{P}_t \mathbf{V}_{t-1} - \mathbf{P} \mathbf{V}_{\gamma_t}^*].$$
(9)

The decomposition (9) is similar to what appears in the standard analysis of discounted Q-learning updates, such as Li et al. (2021); Wainwright (2019b); however, we note a few key technical challenges due to a dynamic discount factor and a moving target $Q_{\gamma_{\ell}}^{*}$, which require considerably more efforts and techniques to address.

- (1) First, due to the dynamic discount factor $\gamma_t \to 1$, the contraction of (9) is much weaker than that with a fixed discount factor, so that the last term $\eta_t \gamma_t \left[\mathbf{P}_t \mathbf{V}_{t-1} \mathbf{P} \mathbf{V}_{\gamma_t}^* \right]$ is much more difficult to control. As a result, γ_t cannot converge to 1 too fast.
- (2) Due to the moving target $Q_{\gamma_t}^*$, there is a bias $(1 \eta_t)[Q_{\gamma_{t-1}}^* Q_{\gamma_t}^*]$ in (9). It is large if γ_t converges to 1 too quickly or if η_t is too small. A similar bias occurs in $\eta_t \gamma_t [P_t V_{t-1} P V_{\gamma_t}^*]$. In light of (1) and (2) and the bias in Lemma 3.2 (large if γ_t converges to 1 too slowly), we need a careful choice of γ_t to balance different sources of bias.
- (3) Similar to the discounted setting, the learning rate η_t needs to balance the bias from earlier updates and variance from recent updates. Moreover, as η_t is coupled with γ_t in the third term, we have to admit a less aggressive learning rate due to additionally balancing the bias&variance in random update with η_t and the bias from γ_t .

To summarize, in sharp distinction from the discounted Q-learning, we need to address the bias due to $\gamma_t \to 1$ in conjunction with the statistical error in random learning update. As a result, our dynamic Q-learning algorithm typically yield slower convergence compared with the discounted counterpart.

Our analysis partially builds on the recent progress in sharp analysis for discounted Q-learning Li et al. (2021). However, we additionally tackle the bias due to $\gamma_t \neq \gamma_{t-1}$, which leads to distinct analysis. The last term in the bracket in (9) can be further decomposed as

$$egin{aligned} P_t V_{t-1} - P V_{\gamma_t}^* = & (P_t - P) V_{t-1} + P (V_{t-1} - V_{\gamma_{t-1}}^*) \ & + P (V_{\gamma_{t-1}}^* - V_{\gamma_t}^*), \end{aligned}$$

where

$$egin{aligned} m{P}(m{V}_{t-1} - m{V}_{\gamma_{t-1}}^*) &= m{P}^{\pi_{t-1}}m{Q}_{t-1} - m{P}^{\pi_{\gamma_{t-1}}^*}m{Q}_{\gamma_{t-1}}^* \ &\leq m{P}^{\pi_{t-1}}m{Q}_{t-1} - m{P}^{\pi_{t-1}}m{Q}_{\gamma_{t-1}}^* &= m{P}^{\pi_{t-1}}m{\Delta}_{t-1}, \end{aligned}$$

$$egin{aligned} m{P}(m{V}_{t-1} - m{V}_{\gamma_{t-1}}^*) &= m{P}^{\pi_{t-1}}m{Q}_{t-1} - m{P}^{\pi_{\gamma_{t-1}}^*}m{Q}_{\gamma_{t-1}}^* \ &\geq m{P}^{\pi_{\gamma_{t-1}}^*}m{Q}_{t-1} - m{P}^{\pi_{\gamma_{t-1}}^*}m{Q}_{\gamma_{t-1}}^* &= m{P}^{\pi_{\gamma_{t-1}}^*}m{\Delta}_{t-1}. \end{aligned}$$

Above two inequalities use the fact that $P_{\gamma t-1}^*$ is greedy w.r.t. $Q_{\gamma_{t-1}}^*$, while π_{t-1} is greedy w.r.t. Q_{t-1} . Plugging them back into (9) leads to

$$\Delta_{t} \leq (1 - \eta_{t}) \Delta_{t}
+ d_{t} + \eta_{t} \gamma_{t} \left[\mathbf{P}^{\pi_{t-1}} \Delta_{t-1} + (\mathbf{P}_{t} - \mathbf{P}) \mathbf{V}_{t-1} \right]; \quad (11a)
\Delta_{t} \geq (1 - \eta_{t}) \Delta_{t}
+ d_{t} + \eta_{t} \gamma_{t} \left[\mathbf{P}^{\pi_{\gamma_{t-1}}^{*}} \Delta_{t-1} + (\mathbf{P}_{t} - \mathbf{P}) \mathbf{V}_{t-1} \right], \quad (11b)$$

where we define the switching error in the t-th iteration

$$d_t = (1 - \eta_t)[Q_{\gamma_{t-1}}^* - Q_{\gamma_t}^*] + \eta_t \gamma_t P(V_{\gamma_{t-1}}^* - V_{\gamma_t}^*).$$

Applying (11a) and (11b) recursively, we arrive at

$$\Delta_{t} \leq \sum_{i=1}^{t} \eta_{i}^{(t)} \gamma_{i} \left[(\boldsymbol{P}_{i} - \boldsymbol{P}) \boldsymbol{V}_{i-1} + \boldsymbol{P}^{\pi_{t-1}} \Delta_{i-1} \right]
+ \eta_{0}^{(t)} \Delta_{0} + \sum_{i=1}^{t} \eta_{i}^{(t)} \boldsymbol{d}_{i} / \eta_{i},
\Delta_{t} \geq \sum_{i=1}^{t} \eta_{i}^{(t)} \gamma_{i} \left[(\boldsymbol{P}_{i} - \boldsymbol{P}) \boldsymbol{V}_{i-1} + \boldsymbol{P}^{\pi_{\gamma_{i-1}}^{*}} \Delta_{i-1} \right]
+ \eta_{0}^{(t)} \Delta_{0} + \sum_{i=1}^{t} \eta_{i}^{(t)} \boldsymbol{d}_{i} / \eta_{i},$$
(12)

where we define $\eta_t^{(t)} = \eta_t$, $\eta_0^{(t)} = \prod_{j=1}^t (1 - \eta_j)$ and $\eta_i^{(t)} = \eta_i \cdot \prod_{j=i+1}^t (1 - \eta_j)$ for $i \ge 1$.

We now proceed to bound (12). Let $\beta \in (0,1)$ be a constant whose value will be specified later. Writing $t_{\beta} = |(1-\beta)t|$, the upper bound in (12) is

$$\Delta_t \le \zeta_t + \xi_t + \sum_{i=1+t_\beta}^t \eta_i^{(t)} \gamma_i P^{\pi_{i-1}} \Delta_{i-1} + \delta_t, \quad (13)$$

where

$$egin{aligned} oldsymbol{\zeta}_t &:= \eta_0^{(t)} oldsymbol{\Delta}_0 + \sum_{i=1}^{t_eta} \eta_i^{(t)} \gamma_i ig[(oldsymbol{P}_i - oldsymbol{P}) oldsymbol{V}_{i-1} + oldsymbol{P}^{\pi_{t-1}} oldsymbol{\Delta}_{i-1} ig], \ oldsymbol{\xi}_t &:= \sum_{i=1+t_eta}^t \eta_i^{(t)} \gamma_i (oldsymbol{P}_i - oldsymbol{P}) oldsymbol{V}_{i-1}, \ oldsymbol{\delta}_t &:= \sum_{i=1}^t \eta_i^{(t)} oldsymbol{d}_i / \eta_i. \end{aligned}$$

The convergence rates for these quantities then depend on the choice of $\{\eta_t\}$ and $\{\gamma_t\}$, which we analyze in a case-by-case fashion for all our theoretical results.

4.2 Sketch of Analysis for Theorem 3.1

As an example of our analysis, we now bound the terms δ_t , ζ_t and ξ_t in the decomposition (13), which leads to

a recursive bound on Δ_t for Theorem 3.1. To begin with, for the constant $c_2 > 0$ in Theorem 3.1, we set

$$\beta = \frac{c_2}{T^{1/5}(\log T)^2}.$$

The lemmas throughout this subsection will be under the same conditions as Theorem 3.1.

Bounding switching error. First, the dynamic discount factors leads to the switching error δ_t , since the estimation targets $Q_{\gamma_t}^*$ keeps moving.

Lemma 4.1. Let T satisfy the conditions of Theorem 3.1. Then $\|\boldsymbol{\delta}_t\|_{\infty} \leq 2(\frac{\log T}{T})^{2/5}$ for all t obeying $T/\log T \leq t \leq T$.

Bounding ζ_t . The second term ζ_t is the cumulative estimation error up to $t_{\beta} = \lfloor (1 - \beta)t \rfloor$. It can be bounded via appropriate contraction with appropriate choices of t_{β} . We bound its ℓ_{∞} -norm as follows. The proof of Lemma 4.2 is in Appendix A.5.

Lemma 4.2. Let T satisfy the conditions of Theorem 3.1. Then $\|\zeta_t\|_{\infty} \leq \frac{2}{T}$ for all $T/\log T \leq t \leq T$.

Bounding ξ_t . Finally, ξ_t is a random error term, for which we derive a high-probability bound, adapting the sharp analysis strategies introduced in Li et al. (2021). The proof of Lemma 4.3 is in Appendix A.6.

Lemma 4.3. For any fixed t obeying $T/\log T \le t \le T$, it holds with probability at least $1 - \delta$ that

$$|\boldsymbol{\xi}_{t}| \leq 5\sqrt{\frac{(\log T)^{\frac{19}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{c_{1}T^{3/5}}}$$

$$\times \sqrt{\left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-1/5}\mathbf{1}\right)}.$$
(14)

A recursive bound. With the above three bounds in place, we put them together and obtain a recursive bound on Δ_t , whose detailed proof is in Appendix A.3. **Proposition 4.4.** Under the conditions in Theorem 3.1, with probability at least $1 - \delta$,

$$\Delta_{t} \leq \sqrt{\frac{c_{3}(\log T)^{4} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5} (T \log T)^{1/5}}}
\times \sqrt{\left(1 + 9t_{mix}(\log T)^{1/5} + T^{1/5} \max_{\lfloor t/2 \rfloor \leq i < t} \|\Delta_{i}\|_{\infty}\right)} \cdot \mathbf{1}$$

 \bigvee $\lfloor t/2 \rfloor \leq i < t$ "

holds simultaneously for all $\frac{3T}{2 \cdot 1 \cdot T} < t < T$, where

holds simultaneously for all $\frac{3T}{2 \log T} \leq t \leq T$, where $c_3 > 0$ is a constant that only depends on c_1, c_2 .

The techniques of Proposition 4.4 are related to Li et al. (2021), but controlling multiple sources of bias and variance relies on quite different ideas.

A similar lower bound is in Proposition 4.5, whose proof is essentially the same as that of Proposition 4.4 hence omitted here.

Proposition 4.5. Suppose T satisfies the conditions in Theorem 3.1. Then with probability at least $1 - \delta$,

$$\Delta_{t} \geq -\sqrt{\frac{c_{3}(\log T)^{4}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}}(T\log T)^{-1/5}}$$

$$\times \sqrt{\left(1 + 9t_{mix}(\log T)^{1/5} + T^{1/5}\max_{\lfloor t/2 \rfloor \leq i < t} \|\Delta_{i}\|_{\infty}\right)} \mathbf{1}$$

holds simultaneously for all t obeying $\frac{3T}{2 \log T} \leq t \leq T$, where $c_3 > 0$ is a constant that only depends on c_1, c_2 .

Finally, solving the recursive bounds proves Theorem 3.1; see Appendix A.1 for details.

5 Extension to Learning Bias Function

We now propose a variant of our framework that learns (up to a constant) $q^*(s, a)$ in (1) hence the optimal policy for average reward. In particular, our approach overcomes the non-contraction of the empirical Bellman updates that are often considered in the literature (Abounadi et al., 2001; Zhang et al., 2021), enabling finite-sample analysis of convergence to solutions to (1).

Our new algorithm maintains an estimate $q_t : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ for all $t \geq 1$; in each iteration t, it updates all entries of the estimate for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ via

$$q_t(s, a) = (1 - \theta_t)q_{t-1}(s, a) + \theta_t [r(s, a) + \alpha_t \max_{a' \in \mathcal{A}} q_{t-1}(s_t(s, a), a')].$$
 (17)

Here $s_t(s,a)$ is the independent sample from the generative model, and α_t and θ_t are the discount factor and learning rate in the t-th iteration, respectively. In addition, we define the value function $v_t \colon \mathcal{S} \to \mathbb{R}$ in the t-th iteration by $v_t(s) := \max_{a \in \mathcal{A}} q_t(s,a)$ for all $s \in \mathcal{S}$. See Algorithm 2 for a formal statement.

Algorithm 2 Dynamic Horizon q-Learning

 $\{\theta_t\}_{t\geq 1}, \{\alpha_t\}_{t\geq 0}$. Initialization: $q_0\equiv 0$. for $t=1,2,\ldots$ do

- **3**: Generate $s_t(s, a) \sim P(\cdot \mid s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.
- 4: Update $q_t(s, a) = (1 \theta_t)q_{t-1}(s, a) + \theta_t [r(s, a) + \alpha_t v_{t-1}(s_t(s, a))]$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.
- 5: Set $v_t = \max_{a \in \mathcal{A}} q_t(s, a)$ for all $s \in \mathcal{S}$.
- 6: end for

The distinction between this procedure and Algorithm 1 is that we use r(s, a) instead of $(1-\alpha_t)r(s, a)$ in the update. As a result, our estimates approximate another set of dynamic targets: the unscaled discounted

value functions $\{q_{\alpha_t}^*\}$ and $\{v_{\alpha_t}^*\}$ (c.f. (5)), which further approximate solutions to (1).

Lemma 5.1. There exists a constant B>0 which only depends on the underlying MDP such that $\left|v_{\alpha}^{*}(s)-\frac{J^{*}}{1-\alpha}-v^{*}(s)\right|\leq (1-\alpha)B$ and $\left|q_{\alpha}^{*}(s,a)-\frac{J^{*}}{1-\alpha}-q^{*}(s,a)\right|\leq (1-\alpha)B$ for all state-action pairs $(s,a)\in\mathcal{S}\times\mathcal{A}$, where J^{*} is the optimal average reward, q_{α}^{*} and v_{α}^{*} are the unscaled optimal α -discounted Q- (value) functions, and q^{*} and v^{*} are the bias functions defined in (2).

The proof of Lemma 5.1 is in Appendix D.2. Thus,

$$|q_t(s, a) - q^*(s, a) - J^*/(1 - \alpha_t)|$$

 $\leq (1 - \alpha)B + |q_t(s, a) - q^*_{\alpha_t}(s, a)|,$

where, similar to the previous case, it remains to control the estimation error to the dynamic targets.

Theorem 5.2. In Algorithm 2, we set $\theta_t = (1 + \frac{c_1 t^{2/3}}{(\log t)^2})^{-1}$, $\alpha_t = 1 - t^{-1/9}$, $t \geq 2$, for some constant $c_1 > 0$, and set $\theta_1 = \theta_2$, $\alpha_1 = \alpha_2$. Let $\varepsilon \in (0,1)$ and $\delta \in (0,1)$. Suppose T is sufficiently large such that $T/\log T \geq \max\{e^{2+c_1}, 100, 64(2+c_1)^3\}$, $c_3(\log T)^5 \geq 5$, $c_1 T^3 \geq 3(\log T)^5$, $2c_3 T \geq (1+3/c_1)^2(\log T)^2$ and $c_2 \log T \geq 5(2+c_1)(\log T)^{1/3} + \log\log T$ for some constant $c_2 > 0$ and $c_3 = 16/c_1 + 144/c_1^2$. Then with probability at least $1 - \delta$, after T iterations, Algorithm 2 achieves

$$\left| q_T(s,a) - q^*(s,a) - \frac{J^*}{1 - \alpha_T} \right| \le \frac{B + c(\log T)^4 (\log \frac{T|S||A|}{2\delta})^2}{T^{1/9}}$$
and
$$\left| v_T(s) - v^*(s) - \frac{J^*}{1 - \alpha_T} \right| \le \frac{B + c(\log T)^4 (\log \frac{T|S||A|}{2\delta})^2}{T^{1/9}}$$

simultaneously for all $(s, a) \in S \times A$, where c > 0 is an absolute constant that only depends on c_1, c_2 , B is the constant given in Lemma 5.1, q^* and v^* are the bias functions defined in (2).

See Appendix C.1 for a sketch of analysis, and Appendix C.2 for a detailed proof.

With the bias function estimator q_T , a natural idea for the task of policy learning is to take the greedy policy with respect to q_T :

$$\pi_T(s) = \operatorname*{argmax}_{a \in \mathcal{A}} \in q_T(s, a).$$

Its suboptimality from the true optimal reward turns out to be of the same scale as the estimation error bound of q_T . We thus have its near-optimality property. The proof of Corollary 5.3 is in Appendix C.3.

Corollary 5.3. Under the same conditions of Theorem 5.2, the greedy policy π_T with respect to q_T obeys

$$J^* - J^{\pi_T} \leq \frac{2B + 2c(\log T)^4(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{1/9}}.$$

6 Empirical Validation

We support our theoretical results with preliminary simulations. We design MDPs with $|\mathcal{S}|=10$ and $|\mathcal{A}|=8$ according to the construction in Jin and Sidford (2021), which is the hardest instance (in the sense of information-theoretic lower bound) for learning a policy. We vary the mixing time $t_{\text{mix}} \in \{O(10), O(100), O(1000)\}$, which is obtained by tuning the parameter $\gamma \in \{0.1, 0.01, 0.001\}$ in Jin and Sidford (2021). On each MDP instance, we run our algorithm for Theorem 3.3 in 200 independent experiments for $T=10^7$. We keep the same scheduling of $\gamma_t=1-t^{-1/8}$ and $\eta_t=(1+\frac{t^{5/8}}{(\log t)^2})^{-1}$ for all instances, without adjusting algorithm inputs using any prior knowledge.

We compute the frequency where the algorithm finds the optimal action for each state at $t \in \{10, 10^2, \dots, 10^6, 10^7\}$ in Figure 1 (for the first 4 states for visualization). The results suggest that our algorithm performs stably well, finding the optimal policy within a reasonable time of training (which might be better than what is predicted by our theory). As promised, it does not need any prior information as input, but works well for MDPs with different $t_{\rm mix}$. Interestingly, the recovery frequency turns out to be relatively stable across different values of $t_{\rm mix}$.

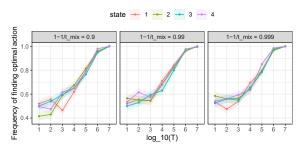


Figure 1: Frequency of recovering optimal action.

Conclusions

In this paper, we provide a feasible framework for average-reward RL that, distinct from existing works, does not require any problem-dependent parameters as input. Our results highlight the algorithmic role of sequentially adjusted discounted factors, along with a carefully selected adaptation scheme, in achieving several average reward learning objectives. We provide finite-sample guarantees for three popular learning tasks. We envision this work as initial progress towards algorithmic design and theoretical understanding of feasible model-free average-reward RL, and an invitation for subsequent efforts in these aspects.

References

- Abounadi, J., Bertsekas, D., and Borkar, V. S. (2001). Learning algorithms for markov decision processes with average cost. *SIAM Journal on Control and Optimization*, 40(3):681–698.
- Agarwal, A., Kakade, S., and Yang, L. F. (2020). Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR.
- Azar, M. G., Munos, R., and Kappen, H. J. (2013). Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349.
- Beck, C. L. and Srikant, R. (2012). Error bounds for constant step-size q-learning. Systems & control letters, 61(12):1203–1208.
- Blackwell, D. (1962). Discrete dynamic programming. The Annals of Mathematical Statistics, pages 719–726.
- Borkar, V. S. and Meyn, S. P. (2000). The ode method for convergence of stochastic approximation and reinforcement learning. *SIAM Journal on Control and Optimization*, 38(2):447–469.
- Brown, N. and Sandholm, T. (2018). Superhuman ai for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424.
- Campbell, M., Hoane Jr, A. J., and Hsu, F.-h. (2002). Deep blue. *Artificial intelligence*, 134(1-2):57–83.
- Chen, Z., Theja Maguluri, S., Shakkottai, S., and Shanmugam, K. (2020). Finite-sample analysis of stochastic approximation using smooth convex envelopes. *arXiv e-prints*, pages arXiv–2002.
- De Farias, D. P. and Van Roy, B. (2006). A cost-shaping linear program for average-cost approximate dynamic programming with performance guarantees. *Mathematics of Operations Research*, 31(3):597–620.
- Dewanto, V., Dunn, G., Eshragh, A., Gallagher, M., and Roosta, F. (2020). Average-reward model-free reinforcement learning: a systematic review and literature mapping. arXiv preprint arXiv:2010.08920.
- Dong, S., Van Roy, B., and Zhou, Z. (2019). Provably efficient reinforcement learning with aggregated states. arXiv preprint arXiv:1912.06366.
- Dong, S., Van Roy, B., and Zhou, Z. (2021). Simple agent, complex environment: Efficient reinforcement learning with agent state. arXiv preprint arXiv:2102.05261.

- Even-Dar, E., Kakade, S. M., and Mansour, Y. (2009). Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736.
- Even-Dar, E., Mansour, Y., and Bartlett, P. (2003). Learning rates for q-learning. *Journal of machine learning Research*, 5(1).
- Feinberg, E. A. and Shwartz, A. (2012). *Handbook of Markov decision processes: methods and applications*, volume 40. Springer Science & Business Media.
- François-Lavet, V., Henderson, P., Islam, R., Bellemare, M. G., and Pineau, J. (2018). An introduction to deep reinforcement learning. arXiv preprint arXiv:1811.12560.
- Freedman, D. A. (1975). On tail probabilities for martingales. the Annals of Probability, pages 100–118.
- Fruit, R., Pirotta, M., and Lazaric, A. (2020). Improved analysis of ucrl2 with empirical bernstein inequality. arXiv preprint arXiv:2007.05456.
- Hordijk, A. and Tijms, H. (1975). A modified form of the iterative method of dynamic programming. *The Annals of Statistics*, pages 203–208.
- Jaakkola, T., Jordan, M. I., and Singh, S. P. (1994). On the convergence of stochastic iterative dynamic programming algorithms. *Neural computation*, 6(6):1185–1201.
- Jaksch, T., Ortner, R., and Auer, P. (2010). Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(51):1563–1600.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is q-learning provably efficient? arXiv preprint arXiv:1807.03765.
- Jin, Y. and Sidford, A. (2020). Efficiently solving mdps with stochastic mirror descent. In *International Conference on Machine Learning*, pages 4890–4900. PMLR.
- Jin, Y. and Sidford, A. (2021). Towards tight bounds on the sample complexity of average-reward mdps. arXiv preprint arXiv:2106.07046.
- Kakade, S. and Langford, J. (2002). Approximately optimal approximate reinforcement learning. In *IN PROC. 19TH INTERNATIONAL CONFERENCE ON MACHINE LEARNING*, pages 267–274.
- Kakade, S. M. (2003). On the sample complexity of reinforcement learning. University of London, University College London (United Kingdom).

- Kearns, M., Mansour, Y., and Ng, A. Y. (2002). A sparse sampling algorithm for near-optimal planning in large markov decision processes. *Machine learning*, 49(2):193–208.
- Kearns, M. and Singh, S. (1999). Finite-sample convergence rates for q-learning and indirect algorithms. Advances in neural information processing systems, pages 996–1002.
- Kearns, M. and Singh, S. (2002). Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2):209–232.
- Lattimore, T. and Szepesvári, C. (2020). Bandit algorithms. Cambridge University Press.
- Li, G., Cai, C., Chen, Y., Gu, Y., Wei, Y., and Chi, Y. (2021). Is q-learning minimax optimal? a tight sample complexity analysis. arXiv preprint arXiv:2102.06548.
- Mahadevan, S. (1996). Average reward reinforcement learning: Foundations, algorithms, and empirical results. *Machine learning*, 22(1):159–195.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533.
- Puterman, M. L. (2014). Markov decision processes: discrete stochastic dynamic programming. John Wiley & Sons.
- Schaeffer, J., Culberson, J., Treloar, N., Knight, B., Lu, P., and Szafron, D. (1992). A world championship caliber checkers program. *Artificial Intelligence*, 53(2-3):273–289.
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., et al. (2020). Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609.
- Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. (2018a). Near-optimal time and sample complexities for solving markov decision processes with a generative model. *Advances in Neural Information Processing Systems*, 31.
- Sidford, A., Wang, M., Wu, X., and Ye, Y. (2018b). Variance reduced value iteration and faster algorithms for solving markov decision processes. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 770–787. SIAM.

- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., et al. (2018). A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144.
- Sutton, R. S. and Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.
- Szepesvári, C. (1998). The asymptotic convergencerate of q-learning. In *Proceedings of the 1997 confer*ence on Advances in neural information processing systems 10, pages 1064–1070.
- Tang, Y., Rowland, M., Munos, R., and Valko, M. (2021). Taylor expansions of discount factors. In *Proceedings of the 29th International Conference on Machine Learning (ICML-21)*.
- Tsitsiklis, J. N. (1994). Asynchronous stochastic approximation and q-learning. *Machine learning*, 16(3):185–202.
- Wainwright, M. J. (2019a). Stochastic approximation with cone-contractive operators: Sharp ℓ_{∞} -bounds for *q*-learning. arXiv preprint arXiv:1905.06265.
- Wainwright, M. J. (2019b). Variance-reduced q-learning is minimax optimal. arXiv preprint arXiv:1906.04697.
- Wan, Y., Naik, A., and Sutton, R. S. (2021). Learning and planning in average-reward markov decision processes. In *International Conference on Machine Learning*, pages 10653–10662. PMLR.
- Wang, M. (2017). Primal-dual π learning: Sample complexity and sublinear run time for ergodic markov decision problems. $arXiv\ preprint\ arXiv:1710.06100$.
- Wang, M. (2020). Randomized linear programming solves the markov decision problem in nearly linear (sometimes sublinear) time. *Mathematics of Operations Research*, 45(2):517–546.
- Watkins, C. J. and Dayan, P. (1992). Q-learning. *Machine learning*, 8(3-4):279–292.
- Wei, C.-Y., Jahromi, M. J., Luo, H., Sharma, H., and Jain, R. (2020). Model-free reinforcement learning in infinite-horizon average-reward markov decision processes. In *International conference on machine learning*, pages 10170–10180. PMLR.
- Yang, L. and Wang, M. (2019). Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004. PMLR.

Zanette, A., Kochenderfer, M. J., and Brunskill, E. (2019). Almost horizon-free structure-aware best policy identification with a generative model. *Advances in Neural Information Processing Systems*, 32.

Zhang, S., Zhang, Z., and Maguluri, S. T. (2021). Finite sample analysis of average-reward td learning and q-learning. Advances in Neural Information Processing Systems, 34.

Zhang, Z. and Ji, X. (2019). Regret minimization for reinforcement learning by evaluating the optimal bias function. Advances in Neural Information Processing Systems, 32.

Zhang, Z. and Xie, Q. (2023). Sharper model-free reinforcement learning for average-reward markov decision processes. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 5476–5477. PMLR.

Checklist

- 1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes]
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes]
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Not Applicable]
- 2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. [Yes]
 - (b) Complete proofs of all theoretical results. [Yes]
 - (c) Clear explanations of any assumptions. [Yes]
- 3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Not Applicable] (Explanation: we will do so when it is possible to de-anonymize the paper.)
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes]
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes]

- (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Not Applicable]
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
 - (a) Citations of the creator If your work uses existing assets. [Not Applicable]
 - (b) The license information of the assets, if applicable. [Not Applicable]
 - (c) New assets either in the supplemental material or as a URL, if applicable. [Not Applicable]
 - (d) Information about consent from data providers/curators. [Not Applicable]
 - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable]
- 5. If you used crowdsourcing or conducted research with human subjects, check if you include:
 - (a) The full text of instructions given to participants and screenshots. [Not Applicable]
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable]
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable]

A Technical Proofs Regarding Theorem 3.1

In this section, we provide detailed proofs for results regarding Theorem 3.1. Appendix A.1 provides the detailed proof for Theorem 3.1, and the remaining subsections provide proofs for intermediate results in the analysis sketch of Section 4.

Throughout the section, the learning rates and discount factors are

$$\eta_t = \frac{1}{1 + \frac{c_1 t^{3/5}}{(\log t)^3}}, \quad \gamma_t = 1 - t^{-1/5},$$

and we choose the proportion

$$\beta = \frac{c_2}{T^{1/5}(\log T)^2}$$

for the constant $c_2 > 0$ in Theorem 3.1.

A.1 Proof of Theorem 3.1

Proof of Theorem 3.1. We solve the recursive bounds in Propositions 4.4 and 4.5 to obtain the final high-probability bound. For any positive integer k, we define

$$u_k = \max \left\{ \| \boldsymbol{\Delta}_i \|_{\infty} \colon 2^k \frac{3T}{2 \log T} \le i \le T \right\}.$$

A naive upper bound is $u_k \leq 1$ for all k. Furthermore, by (15) and the definition of u_k , with probability at least $1 - \delta$, it holds simultaneously for all $k \geq 1$ that

$$u_{k+1} \le \sqrt{\frac{c_3(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}} \left(\frac{(\log T)^{-1/5} + 9t_{\text{mix}}}{T^{1/5}} + \frac{u_k}{(\log T)^{1/5}}\right)} \mathbf{1}$$

for some absolute constant $c_3 > 0$. If there exists some k such that

$$u_k \le c_4 \left((\log T)^{-1/5} + 9t_{\text{mix}} \right) \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}}$$

for some sufficiently large constant $c_4 > 0$, then

$$u_{k+1} \le \sqrt{\frac{c_3(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}} \frac{(\log T)^{-1/5} + 9t_{\text{mix}}}{T^{1/5}} \left(1 + c_4(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}\right)}{\frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}}}$$

as well. By induction, the above bound holds for all $j \geq k$, hence

$$\|\Delta_T\|_{\infty} \le c_4 ((\log T)^{-1/5} + 9t_{\text{mix}}) \frac{(\log T)^4 \log \frac{|S||S|T}{\delta}}{T^{1/5}}.$$

On the other hand, suppose $u_j > c_4((\log T)^{-1/5} + 9t_{\text{mix}}) \frac{(\log T)^4 \log \frac{|S||A|T}{\delta}}{T^{1/5}}$ for any $1 \le j \le k$. Then we know

$$u_{j+1} \le \sqrt{\frac{2c_3(\log T)^4 \log \frac{|\mathcal{A}||\mathcal{S}|T}{\delta}}{T^{1/5}} u_j}, \quad \text{for all } 1 \le j \le k,$$

which implies

$$\log u_{j+1} \le \frac{1}{2} \log u_j + \frac{1}{2} \log \theta, \quad \text{where} \quad \theta = 2c_3 \cdot T^{-1/5} (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}$$

for some constant $c_6 > 0$. Recursively applying this relation for $1 \le j < k$ yields

$$2^k \log u_k \le \log u_0 + \sum_{j=0}^{k-1} 2^j \log \theta \le (2^k - 1) \log \theta,$$

where we used the fact that $u_0 \leq 1$. Hence

$$u_k \le \theta^{1-1/2^k} = \left(2c_3 \cdot T^{-1/5} (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}\right)^{1-1/2^k}$$

$$= \left(2c_3 \cdot (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}\right)^{1-1/2^k} \cdot T^{-1/5} \cdot T^{\frac{2^{-k}}{5}}$$

$$\le 2c_3 \cdot \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}} \cdot \exp\left(\log T \cdot 2^{-k}/5\right).$$

Now we can set $2^k \ge \log T/5$, so that the above upper bound translates to

$$\|\mathbf{\Delta}_T\|_{\infty} \le u_k \le 6c_3 \cdot \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}}.$$

Combining the two cases above, we arrive at

$$\|\Delta_T\|_{\infty} \le \left[c_4\left((\log T)^{-1/5} + 9t_{\min}\right) + 6c_3\right] \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}}$$

with probability at least $1 - \delta$ for T obeying the conditions of Proposition 4.4. Furthermore, since $t_{\text{mix}} \ge 1$, the above bound translates to

$$\|\mathbf{\Delta}_T\|_{\infty} \le ct_{\text{mix}} \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/5}}$$

for some constant c > 0 that only depends on c_1, c_2 , which completes the proof of Theorem 3.1.

A.2 Proof of Lemma 3.2

Proof of Lemma 3.2. For any discounted factor $\gamma \in (0,1)$, the optimal policy π_{γ}^* with respect to γ -discounted reward attains the optimal Q-function, i.e., $\mathbf{Q}_{\gamma}^* = \mathbf{Q}_{\gamma}^{\pi_{\gamma}^*}$. Similarly, the optimal policy π^* with respect to average reward satisfies $J^* = J^{\pi^*}$. Applying Lemma D.1 to π_{γ}^* and γ leads to

$$Q_{\gamma}^{*}(s,a) = Q_{\gamma}^{\pi_{\gamma}^{*}}(s,a) \le J^{\pi_{\gamma}^{*}} + 3(1-\gamma)t_{\text{mix}} \le J^{*} + 3(1-\gamma)t_{\text{mix}}$$

for any state-action pair $(s,a) \in \mathcal{S} \times \mathcal{A}$. On the other hand, applying Lemma D.1 to π^* and γ leads to

$$Q_{\gamma}^{*}(s, a) \ge Q_{\gamma}^{\pi^{*}}(s, a) \ge J^{\pi^{*}} - 3(1 - \gamma)t_{\text{mix}} \ge J^{*} - 3(1 - \gamma)t_{\text{mix}}$$

for any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$. The last steps of the two equations follow from the bounded reward averaging time for π_{γ}^* and π^* . Combining the two inequalities, we complete the proof of Lemma 3.2.

A.3 Proof of Proposition 4.4

Proof of Proposition 4.4. Combining Lemma 4.1, 4.2 and 4.3, for sufficiently larget T obeying the conditions of Theorem 3.1 and any fixed t within $T/\log T \le t \le T$, it holds with probability at least $1-\delta$ that

$$\Delta_{t} \leq \frac{2(\log T)^{2/5}}{T^{2/5}} + \frac{2}{T} + 5\sqrt{\frac{(\log T)^{\frac{19}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{c_{1}T^{3/5}}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-1/5}\mathbf{1}\right)$$
$$+ \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \eta_{i}^{(t)} \gamma_{i} \boldsymbol{P}^{\pi_{t-1}} \Delta_{i-1}$$

$$\leq \left(1 + \frac{5}{\sqrt{c_1}}\right) \sqrt{\frac{(\log T)^{\frac{19}{5}} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{3/5}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-1/5} \mathbf{1}\right)} + \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \eta_i^{(t)} \gamma_i \boldsymbol{P}^{\pi_{i-1}} \boldsymbol{\Delta}_{i-1}, \tag{18}$$

where the second inequality follows from the fact that

$$\frac{2(\log T)^{2/5}}{T^{2/5}} + \frac{2}{T} \le \frac{4(\log T)^{2/5}}{T^{2/5}} \le \sqrt{\frac{(\log T)^{\frac{19}{5}} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{4/5}}}$$

for T so large that $T/\log T \geq 300$. Once $(1-\beta) \geq 3/4$, for any fixed t obeying $3T/(2\log T) \leq t \leq T$, we apply (18) with a union bound over $\{k \colon \frac{2t}{3} \leq k \leq t\}$ —as a result, with probability at least $1-\delta$, one has

$$\Delta_k \le \sqrt{\varphi_t} + \sum_{i=1+\lfloor (1-\beta)k\rfloor}^k \eta_i^{(k)} \gamma_i \mathbf{P}^{\pi_{i-1}} \Delta_{i-1}, \quad \text{for all } \frac{2t}{3} \le k \le t,$$
(19)

where we define

$$\varphi_{t} = \left(1 + \frac{5}{\sqrt{c_{1}}}\right)^{2} \cdot \frac{\left(\log T\right)^{\frac{19}{5}} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{3/5}} \left(\max_{\lfloor t/2 \rfloor < i \le t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-1/5}\mathbf{1}\right)
\geq 2\left(1 + \frac{5}{\sqrt{c_{1}}}\right)^{2} \cdot \frac{\left(\log T\right)^{\frac{19}{5}} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta/T}}{T^{3/5}} \left(\max_{\lfloor (1-\beta)k \rfloor < i \le k} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-1/5}\mathbf{1}\right).$$
(20)

Before proceeding to recursively bound Δ_t , we define $\{\alpha_i^{(t)}\}$ as

$$\alpha_i^{(t)} := \frac{\eta_{i+1}^{(t)}}{\sum_{j=\lfloor (1-\beta)t\rfloor}^{t-1} \eta_{j+1}^{(t)}},\tag{21}$$

which, following Li et al. (2021), satisfies

$$\alpha_i^{(t)} \ge \eta_{i+1}^{(t)} \quad \text{and} \quad \sum_{j=\lfloor (1-\beta)t \rfloor}^{t-1} \alpha_{j+1}^{(t)} = 1$$
 (22)

for all t. When T is sufficiently large so that $1 - \beta \ge 2/3$, We decompose (19) as

$$\boldsymbol{\Delta}_{k} \leq \sqrt{\boldsymbol{\varphi}_{t}} + \sum_{i=1+|(1-\beta)k|}^{k} \eta_{i}^{(k)} \gamma_{i} \boldsymbol{P}^{\pi_{i-1}} \boldsymbol{\Delta}_{i-1} = \sum_{i_{1}=|(1-\beta)k|}^{k-1} \left(\alpha_{i_{1}}^{(k)} \sqrt{\boldsymbol{\varphi}_{t}} + \eta_{i_{1}+1}^{(k)} \gamma_{i_{1}+1} \boldsymbol{P}^{\pi_{i_{1}}} \boldsymbol{\Delta}_{i_{1}} \right).$$

We recursively apply the above relation in a manner similar to Equation 68 of Li et al. (2021), yet with a sequence of dynamic discount factor $\{\gamma_t\}$:

$$\begin{split} & \boldsymbol{\Delta}_{t} \leq \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \left(\alpha_{i_{1}}^{(t)}\sqrt{\varphi_{t}} + \eta_{i_{1}+1}^{(t)}\gamma_{i_{1}+1}\boldsymbol{P}^{\pi_{i_{1}}}\boldsymbol{\Delta}_{i_{1}}\right) \\ & \leq \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \left[\alpha_{i_{1}}^{(t)}\sqrt{\varphi_{t}} + \eta_{i_{1}+1}^{(t)}\gamma_{i_{1}+1}\boldsymbol{P}^{\pi_{i_{1}}}\sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \left(\alpha_{i_{2}}^{(i_{1})}\sqrt{\varphi_{t}} + \eta_{i_{2}+1}^{(i_{1})}\gamma_{i_{2}+1}\boldsymbol{P}^{\pi_{i_{2}}}\boldsymbol{\Delta}_{i_{2}}\right)\right] \\ & \leq \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \alpha_{i_{1}}^{(t)}\sqrt{\varphi_{t}} + \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \alpha_{i_{1}}^{(t)}\alpha_{i_{2}}^{(i_{1})}\left(\gamma_{i_{1}+1}\boldsymbol{P}^{\pi_{i_{1}}}\right)\sqrt{\varphi_{t}} \\ & + \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \eta_{i_{1}+1}^{(t)}\eta_{i_{2}+1}^{(i_{1})} \prod_{k=1}^{2} \left(\gamma_{i_{k}+1}\boldsymbol{P}^{\pi_{i_{k}}}\right)\boldsymbol{\Delta}_{i_{2}}\right) \end{split}$$

$$\leq \sum_{i_{1}=\lfloor (1-\beta)t \rfloor}^{t-1} \sum_{i_{2}=\lfloor (1-\beta)i_{1} \rfloor}^{i_{1}-1} \alpha_{i_{1}}^{(t)} \alpha_{i_{2}}^{(i_{1})} \{ \boldsymbol{I} + \gamma_{i_{1}+1} \boldsymbol{P}^{\pi_{i_{1}}} \} \sqrt{\varphi_{t}}
+ \sum_{i_{1}=\lfloor (1-\beta)t \rfloor}^{t-1} \sum_{i_{2}=\lfloor (1-\beta)i_{1} \rfloor}^{i_{1}-1} \eta_{i_{1}+1}^{(t)} \eta_{i_{2}+1}^{(i_{1})} \prod_{k=1}^{2} (\gamma_{i_{k}+1} \boldsymbol{P}^{\pi_{i_{k}}}) \Delta_{i_{2}},$$
(23)

where the third line uses $\eta_{i_1+1}^{(t)} \leq \alpha_i^{(t)}$ in (22), and the last line uses $\sum_{i_2=\lfloor (1-\beta)i_1\rfloor}^{i_1-1} \alpha_{i_2}^{(i_1)} = 1$ in (22). We shall further recursively apply the above relation. To begin with, we set

$$H:=T^{1/5}(\log T)^2 \quad \text{and} \quad \alpha_{\{i_k\}_{k=1}^H}:=\alpha_{i_1}^{(t)}\alpha_{i_2}^{(i_1)}\dots\alpha_{i_H}^{(i_{H-1})}\geq 0$$

for any $t > i_1 > i_2 > \cdots > i_H$, which (according to (22)) satisfies

$$\alpha_{\{i_k\}_{k=1}^H} \ge \eta_{i_1+1}^{(t)} \eta_{i_2+1}^{(i_1)} \cdots \eta_{i_H+1}^{(i_{H-1})}.$$

We define the index set

$$\mathcal{I}_t = \{(i_1, \dots, i_H) : \lfloor (1 - \beta)t \rfloor \le i_1 \le t - 1, \ \lfloor (1 - \beta)i_{j-1} \rfloor \le i_j \le i_{j-1} - 1, \ \forall 1 \le j < H \},\$$

which satisfies

$$\sum_{(i_1, \dots, i_H) \in \mathcal{I}_t} \alpha_{\{i_k\}_{k=1}^H} = 1.$$

By definition of β , we have

$$(1 - \beta)^H \ge \exp(-\beta H) \ge \frac{2}{3}$$

for sufficiently small $c_2 \leq \log(3/2)$, which implies

$$i_1 > i_2 > \dots > i_H \ge 2t/3$$
, for all $(i_1, \dots, i_H) \in \mathcal{I}_t$.

Recursively invoking (23), we obtain

$$\Delta_{t} \leq \sum_{(i_{1},\dots,i_{H})\in\mathcal{I}_{t}} \alpha_{\{i_{k}\}_{k=1}^{H}} \left\{ \left(\boldsymbol{I} + \sum_{h=1}^{H-1} \prod_{k=1}^{h} (\gamma_{i_{k}+1} \boldsymbol{P}^{\pi_{i_{k}}}) \right) \sqrt{\varphi_{t}} + \prod_{h=1}^{H} (\gamma_{i_{h}+1} \boldsymbol{P}^{\pi_{i_{h}}}) |\Delta_{i_{H}}| \right\}$$

$$\leq \max_{(i_{1},\dots,i_{H})\in\mathcal{I}_{t}} \left\{ \underbrace{\left(\boldsymbol{I} + \sum_{h=1}^{H-1} \prod_{k=1}^{h} (\gamma_{i_{k}+1} \boldsymbol{P}^{\pi_{i_{k}}}) \right) \sqrt{\varphi_{t}}}_{=:\beta_{1}} + \underbrace{\prod_{h=1}^{H} (\gamma_{i_{h}+1} \boldsymbol{P}^{\pi_{i_{h}}}) |\Delta_{i_{H}}|}_{=:\beta_{2}} \right\}, \tag{24}$$

where the second line uses $\sum_{i_2=\lfloor (1-\beta)i_1\rfloor}^{i_1-1} \alpha_{i_2}^{(i_1)} = 1$ in (22). In the following, we are to bound β_1 and β_2 in (24) separately. The easier term is β_2 : noting that by definition of the discount factor $\{\gamma_t\}$, we have $\gamma_j \leq \gamma_T$ for all $2t/3 \leq j \leq T$, indicating

$$|\beta_{2}| \leq \gamma_{t}^{H} \prod_{h=1}^{H} \mathbf{P}^{\pi_{i_{h}}} |\mathbf{\Delta}_{H}| \leq \gamma_{T}^{H} \left\| \prod_{h=1}^{H} \mathbf{P}^{\pi_{i_{h}}} \right\|_{1} \|\mathbf{\Delta}_{H}\|_{\infty} \stackrel{\text{(i)}}{\leq} \gamma_{T}^{H} \stackrel{\text{(ii)}}{\leq} \frac{1}{T}, \tag{25}$$

where (i) follows from the bounded magnitude $\|\Delta_H\|_{\infty} \leq 1$ and the fact that $\prod_{h=1}^H P^{\pi_{i_h}}$ is a probability transition matrix; (ii) follows from

$$\gamma_T^H = (1 - T^{-1/5})^{T^{1/5}(\log T)^2} \le \exp(-(\log T)^2) \le \frac{1}{T}.$$

Moving on to β_1 , its entrywise square can be upper bounded as

$$|oldsymbol{eta}_1|^2 = igg|\sum_{h=0}^{H-1}\prod_{k=1}^h(\gamma_{i_k+1}oldsymbol{P}^{\pi_{i_k}})\sqrt{oldsymbol{arphi}_t}igg|^2 \leq igg|\sum_{h=0}^{H-1}\gamma_T^h\Big(\prod_{k=1}^holdsymbol{P}^{\pi_{i_k}}\Big)\sqrt{oldsymbol{arphi}_t}igg|^2$$

$$\stackrel{\text{(i)}}{\leq} \left| \sum_{h=0}^{H-1} \gamma_{T}^{h/2} \gamma_{T}^{h/2} \sqrt{\prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \boldsymbol{\varphi}_{t}} \right|^{2} \\
\stackrel{\text{(ii)}}{\leq} \sum_{h=0}^{H-1} \gamma_{T}^{h} \cdot \sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \boldsymbol{\varphi}_{t} \\
\stackrel{\text{(iii)}}{\leq} \frac{1}{1 - \gamma_{T}} \sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \frac{2(1 + \frac{5}{\sqrt{c_{1}}})^{2} (\log T)^{\frac{19}{5}} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta/T}}{T^{3/5}} \left(\max_{\lfloor (1-\beta)k \rfloor < i \leq k} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i-1}) + T^{-1/5} \mathbf{1} \right) \\
\stackrel{\text{(iv)}}{\leq} \frac{4(1 + \frac{5}{\sqrt{c_{1}}})^{2} (\log T)^{\frac{19}{5}} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{2/5}} \left(\sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i}) + \mathbf{1} \right), \tag{26}$$

where (i) follows from Jensen's inequality noting that $\prod_{k=1}^{h} P^{\pi_{i_k}}$ is a probability transition matrix, (ii) follows from the Cauchy-Schwarz inequality, (iii) follows from the definition of φ_t in (20), and (iv) follows from the fact that $\prod_{k=1}^{h} P^{\pi_{i_k}} \mathbf{1} = \mathbf{1}$ since $\prod_{k=1}^{h} P^{\pi_{i_k}}$ is a probability transition matrix. We employ the following lemma to bound the first term of (26), whose proof is in Appendix D.4.

Lemma A.1. Suppose $T \geq 160$, then it holds for all $T/\log T \leq t \leq T$ that

$$\sum_{h=0}^{H-1} \gamma_T^h \prod_{k=1}^h \mathbf{P}^{\pi_{i_k}} \max_{\lfloor t/2 \rfloor \le i < t} \text{Var}_{\mathbf{P}}(\mathbf{V}_i) \le \left(7 + 72t_{mix} (\log T)^{1/5} + 8T^{1/5} \max_{\lfloor t/2 \rfloor \le i < t} \|\mathbf{\Delta}_i\|_{\infty}\right) \mathbf{1}. \tag{27}$$

Invoking Lemma A.1, we obtain an upper bound

$$|\beta_1|^2 \le \frac{4(1+\frac{5}{\sqrt{c_1}})^2(\log T)^{\frac{19}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{2/5}} \left(8+72t_{\text{mix}}(\log T)^{1/5}+8T^{1/5}\max_{|t/2|\le i< t}\|\boldsymbol{\Delta}_i\|_{\infty}\right) \mathbf{1}. \tag{28}$$

Recalling (24), by the upper bound (25) on $|\beta_2|$ and (28) on $|\beta_1|$, for any fixed t such that $3T/(2\log T) \le t \le T$,

$$\Delta_{t} \leq \frac{1}{T} + \sqrt{\frac{4(1 + \frac{5}{\sqrt{c_{1}}})^{2}(\log T)^{\frac{19}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{2/5}} \left(8 + 72t_{\text{mix}}(\log T)^{1/5} + 8T^{1/5}\max_{\lfloor t/2\rfloor \leq i < t} \|\Delta_{i}\|_{\infty}\right) \mathbf{1}}$$

$$\leq 4\left(1 + \frac{5}{\sqrt{c_{1}}}\right) \sqrt{\frac{(\log T)^{4}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{2/5}}} \left(\frac{8}{(\log T)^{1/5}} + 72t_{\text{mix}} + \frac{8T^{1/5}}{(\log T)^{1/5}}\max_{\lfloor t/2\rfloor \leq i < t} \|\Delta_{i}\|_{\infty}\right) \mathbf{1}} \tag{29}$$

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

We now further take a union bound for (29) over $\{t: 3T/(2\log T) \le t \le T\}$, which leads to the simultaneous $(1-\delta)$ high-probability bound

$$\boldsymbol{\Delta}_t \leq 16 \left(1 + \frac{5}{\sqrt{c_1}}\right) \sqrt{\frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{2/5}} \left((\log T)^{-1/5} + 9t_{\text{mix}} + (T/\log T)^{1/5} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\Delta}_i\|_{\infty}\right)} \mathbf{1}$$

for all t such that $3T/(2\log T) \le t \le T$, as $\log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta/T} \le 2\log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}$. This completes the proof of Proposition 4.4.

A.4 Proof of Lemma 4.1

Proof of Lemma 4.1. Employing Lemma D.4, each term d_t can be bounded as

$$\|\boldsymbol{d}_{t}\|_{\infty} \leq (1 - \eta_{t}) \|\boldsymbol{Q}_{\gamma_{t-1}}^{*} - \boldsymbol{Q}_{\gamma_{t}}^{*}\|_{\infty} + \eta_{t} \gamma_{t} \|\boldsymbol{P}(\boldsymbol{V}_{\gamma_{t-1}}^{*} - \boldsymbol{V}_{\gamma_{t}}^{*})\|_{\infty}$$

$$\leq (1 - \eta_{t}) \|\boldsymbol{Q}_{\gamma_{t-1}}^{*} - \boldsymbol{Q}_{\gamma_{t}}^{*}\|_{\infty} + \eta_{t} \gamma_{t} \|\boldsymbol{P}\|_{1} \|\boldsymbol{V}_{\gamma_{t-1}}^{*} - \boldsymbol{V}_{\gamma_{t}}^{*}\|_{\infty}$$

$$\leq (1 - \eta_{t}) \|\boldsymbol{Q}_{\gamma_{t-1}}^{*} - \boldsymbol{Q}_{\gamma_{t}}^{*}\|_{\infty} + \eta_{t} \gamma_{t} \|\boldsymbol{V}_{\gamma_{t-1}}^{*} - \boldsymbol{V}_{\gamma_{t}}^{*}\|_{\infty} \leq \frac{\gamma_{t-1} - \gamma_{t}}{1 - \gamma_{t-1}},$$
(30)

_

where the second line uses $\|\mathbf{A}\mathbf{B}\|_{\infty} \leq \|\mathbf{A}\|_1 \|\mathbf{B}\|_{\infty}$ for matrices \mathbf{A}, \mathbf{B} , the third line follows from $\|\mathbf{P}\|_1 = 1$ since each row of \mathbf{P} is a probability vector. For any t > 1, we let $\alpha = t^{-2/5} \in (0, 1)$. By the triangular inequality, the switching error can be bounded as

$$\left\| \sum_{i=1}^{t} \eta_{i}^{(t)} \boldsymbol{d}_{i} / \eta_{i} \right\|_{\infty} \leq \sum_{i=1}^{t} \prod_{j=i}^{t} (1 - \eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty}$$

$$= \sum_{i=1}^{\lfloor (1 - \alpha)t \rfloor} \prod_{j=i}^{t} (1 - \eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty} + \sum_{1 + \lfloor (1 - \alpha)t \rfloor}^{t} \prod_{j=i}^{t} (1 - \eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty}.$$

Firstly, when t obeys $\lfloor (1-\alpha)t \rfloor = \lfloor t-t^{2/5} \rfloor \geq 150$ (which is satisfied by $t \geq 160$), for any i such that $1 \leq i \leq \lfloor (1-\alpha)t \rfloor$, we have

$$\prod_{j=i}^{t} (1 - \eta_j) \leq \prod_{j=\lfloor (1-\alpha)t \rfloor}^{t} (1 - \eta_j) \leq \exp\left(-\sum_{j=\lfloor (1-\alpha)t \rfloor}^{t} \eta_j\right)
\stackrel{\text{(i)}}{\leq} \exp\left(-\sum_{j=\lfloor (1-\alpha)t \rfloor}^{t} \frac{(\log j)^3}{(10 + c_1)j^{3/5}}\right) \stackrel{\text{(ii)}}{\leq} \exp\left(-\frac{\alpha t (\log t)^3}{(10 + c_1)t^{3/5}}\right),$$

where (i) follows from (99), and (ii) follows from (100) for $j \ge \lfloor (1-\alpha)t \rfloor \ge 150$. Together with the fact that $\|\boldsymbol{d}_i\|_{\infty} \le 1$, the above bound implies

$$\sum_{i=1}^{\lfloor (1-\alpha)t\rfloor} \prod_{j=i}^{t} (1-\eta_j) \|\boldsymbol{d}_i\|_{\infty} \le t \cdot \exp\left(-\frac{\alpha t^{2/5} (\log t)^3}{10+c_1}\right) = \exp\left(\log t - (\log t)^3/(10+c_1)\right). \tag{31}$$

On the other hand, (30) implies

$$\|\boldsymbol{d}_i\|_{\infty} \le \frac{(i-1)^{-1/5} - i^{-1/5}}{(i-1)^{-1/5}} \le 1 - (1-1/i)^{1/5} \le \frac{2}{5i} \le \frac{2}{5(1-\alpha)t},$$

for all $i \ge 1 + \lfloor (1 - \alpha)t \rfloor$. The second inequality above is implied by

$$1 - (1 - 1/i)^{1/5} = 1 - \exp\left(\frac{\log(1 - 1/i)}{5}\right) \le 1 - \exp\left(-2/(5i)\right) \le 2/(5i),$$

where for $i \geq 2$, we use the fact that $\log(1-x) \geq -2x$ for $x \leq \log 2/2$ and $e^x \geq 1+x$ for $x \in \mathbb{R}$. Therefore, noting that $\alpha \leq 0.2$ for $t \geq 150$, we have

$$\sum_{1+\lfloor (1-\alpha)t\rfloor}^{t} \prod_{j=i}^{t} (1-\eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty} \leq \sum_{1+\lfloor (1-\alpha)t\rfloor}^{t} \|\boldsymbol{d}_{i}\|_{\infty} \leq \frac{2\alpha t}{5(1-\alpha)t} \leq \alpha/2 = \frac{1}{2t^{2/5}}.$$
 (32)

Combining (31) and (32), once $T/\log T \ge 160$, we have

$$\left\| \sum_{i=1}^{t} \eta_{i}^{(t)} \mathbf{d}_{i} / \eta_{i} \right\|_{\infty} \leq \exp\left(\log t - (\log t)^{3} / (10 + c_{1})\right) + \frac{1}{2t^{2/5}}$$

$$\leq \exp\left(\log T - \frac{(\log T - \log \log T)^{3}}{10 + c_{1}}\right) + \frac{(\log T)^{2/5}}{2T^{2/5}}$$

$$\leq \exp\left(\log T - \frac{(\log T)^{3}}{8(10 + c_{1})}\right) + \frac{(\log T)^{2/5}}{2T^{2/5}} \leq \frac{2(\log T)^{2/5}}{T^{2/5}}$$

for all $T/\log T \le t \le T$, where the third inequality follows from $\log T \ge 2\log\log T$, and the last inequality follows from $(\log T)^3/(80+8c_1) \ge \frac{7}{5}\log T$. Therefore, we complete the proof of Lemma 4.1.

A.5 Proof of Lemma 4.2

Proof of Lemma 4.2. The ℓ_{∞} -norm of ζ_t can be bounded as

$$\begin{split} \|\boldsymbol{\zeta}_{t}\|_{\infty} &\leq \|\eta_{0}^{(t)}\boldsymbol{\Delta}_{0}\|_{\infty} + \sum_{i=1}^{\lfloor (1-\beta)t\rfloor} \eta_{i}^{(t)} \gamma_{i} \|(\boldsymbol{P}_{i}-\boldsymbol{P})\boldsymbol{V}_{i-1} + \boldsymbol{P}^{\pi_{t-1}}\boldsymbol{\Delta}_{i-1}\|_{\infty} \\ &\stackrel{\text{(i)}}{\leq} \eta_{0}^{(t)} + \sum_{i=1}^{\lfloor (1-\beta)t\rfloor} \eta_{i}^{(t)} \gamma_{i} [\|\boldsymbol{P}_{i}-\boldsymbol{P}\|_{1} \|\boldsymbol{V}_{i-1}\|_{\infty} + \|\boldsymbol{P}^{\pi_{t-1}}\|_{1} \|\boldsymbol{\Delta}_{i-1}\|_{\infty}] \\ &\stackrel{\text{(ii)}}{\leq} \eta_{0}^{(t)} + 3 \sum_{i=1}^{\lfloor (1-\beta)t\rfloor} \eta_{i}^{(t)} \gamma_{i}. \end{split}$$

Here (i) uses $\|\boldsymbol{A}\boldsymbol{B}\|_{\infty} \leq \|\boldsymbol{A}\|_{1}\|\boldsymbol{B}\|_{\infty}$ for matrices $\boldsymbol{A}, \boldsymbol{B}$, and (ii) follows from the fact that $\|\boldsymbol{P}_{i} - \boldsymbol{P}\|_{1} \leq \|\boldsymbol{P}_{i}\|_{1} + \|\boldsymbol{P}\|_{1} \leq 2$ and $\|\boldsymbol{P}^{\pi_{t-1}}\|_{1} \leq 1$ (since they are all probability matrices), and the bounded magnitudes $\|\boldsymbol{V}_{i-1}\|_{\infty} \leq 1$ and $\|\boldsymbol{\Delta}_{i-1}\|_{\infty} \leq 1$.

By definition of $\{\eta_t\}$, as long as $\beta < 1/2$ and $t \ge 300$,

$$\eta_0^{(t)} = \prod_{j=1}^t (1 - \eta_j) \le \exp\left(-\sum_{j=1}^t \eta_j\right)
\le \exp\left(-\sum_{j=150}^t \frac{(\log j)^3}{(10 + c_1)j^{3/5}}\right) \le \exp\left(-\frac{(t - 150)(\log t)^3}{(10 + c_1)t^{3/5}}\right),$$
(33)

where the second inequality follows from (99), and the last inequality follows from the monotonicity (100).

Similarly, when $\lfloor (1-\beta)t \rfloor \geq 150$ (satisfied as long as $\beta < 1/2$ and $t \geq 300$), for any $1 \leq i \leq \lfloor (1-\beta)t \rfloor$, it holds that

$$\eta_i^{(t)} \le \prod_{j=1+\lfloor (1-\beta)t\rfloor}^t (1-\eta_j) \le \exp\left(-\sum_{j=1+\lfloor (1-\beta)t\rfloor}^t \eta_j\right) \\
\le \exp\left(-\sum_{j=1+\lfloor (1-\beta)t\rfloor}^t \frac{(\log j)^3}{(10+c_1)j^{3/5}}\right) \le \exp\left(-\frac{\beta t(\log t)^3}{(10+c_1)t^{3/5}}\right),$$
(34)

where the last inequality follows from (100). Combining (33) and (34), once $\beta < 1/2$, we have

$$\|\zeta_t\|_{\infty} \le \exp\left(-\frac{(t-150)(\log t)^3}{(10+c_1)t^{3/5}}\right) + 3\exp\left(\log t - \frac{\beta t(\log t)^3}{(10+c_1)t^{3/5}}\right)$$

$$\le \exp\left(-\frac{t(\log t)^3}{2(10+c_1)t^{3/5}}\right) + 3\exp\left(\log t - \frac{\beta t(\log t)^3}{(10+c_1)t^{3/5}}\right)$$

for all $t \ge 300$. Note that when $T/\log T \ge 300$, for all t such that $T/\log T \le t \le T$,

$$t^{2/5}(\log t)^3 \ge (T/\log T)^{2/5}(\log T - \log\log T)^3 \ge (T/\log T)^{2/5} \cdot (\log T)^3/8 \ge T^{2/5}(\log T)^2/8,$$

where the second inequality follows from $\log T \geq 2 \log \log T$. Meanwhile, for $T/\log T \leq t \leq T$,

$$\log t - \frac{\beta t (\log t)^3}{(10+c_1)t^{3/5}} \le \log T - \frac{c_2}{T^{1/5}(\log T)^2} \frac{T^{2/5}}{(\log T)^{2/5}} \frac{(\log T)^3}{8(10+c_1)} \le \log T - \frac{c_2 T^{1/5}(\log T)^{3/5}}{8(10+c_1)}.$$

Therefore, for any $T/\log T \le t \le \log T$,

$$\|\zeta_t\|_{\infty} \le \exp\left(-\frac{T^{2/5}(\log T)^2}{16(10+c_1)}\right) + \exp\left(\log(3T) - \frac{c_2T^{1/5}(\log T)^{3/5}}{8(10+c_1)}\right) \le \frac{2}{T}$$

if T obeys (satisfied by the conditions of Theorem 3.1)

$$T^{2/5} \log T \ge 16(10 + c_1)$$
 and $c_2 T^{1/5} \ge 24(10 + c_1)(\log T)^{2/5}$,

which completes the proof of Lemma 4.2.

A.6 Proof of Lemma 4.3

Proof of Lemma 4.3. We begin with some basic notations for convenience: write

$$\boldsymbol{\xi}_{t} = \sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \boldsymbol{z}_{i}, \quad \text{where} \quad \boldsymbol{z}_{i} = \eta_{i}^{(t)} \gamma_{i} (\boldsymbol{P}_{i} - \boldsymbol{P}) \boldsymbol{V}_{i-1}. \tag{35}$$

Then (entries of) $\{z_i\}$ are martingale differences in the sense that

$$\mathbb{E}[\boldsymbol{z}_i \mid \boldsymbol{V}_{i-1}, \dots, \boldsymbol{V}_0] = \boldsymbol{0} \quad \forall i > |(1-\beta)t|.$$

We shall obtain high-probability bound on $\|\zeta_t\|_{\infty}$ via Freedman's inequality in Lemma D.6, for which we compute some basic quantities that would be of use.

Firstly, when $\lfloor (1-\beta)t \rfloor \geq 150$ and $1-\beta \geq 1/2$.

$$\max_{\lfloor (1-\beta)t\rfloor < i \le t} \|\boldsymbol{z}_i\|_{\infty} \le \max_{\lfloor (1-\beta)t\rfloor < i \le t} \eta_i^{(t)} \gamma_i \|\boldsymbol{P}_i - \boldsymbol{P}\|_1 \|\boldsymbol{V}_{i-1}\|_{\infty}$$

$$\le \max_{\lfloor (1-\beta)t\rfloor < i \le t} 2\eta_i$$

$$\le \frac{4(\log t)^3}{c_1 t^{3/5}} =: R,$$
(36)

where the second inequality follows from $\|P_i - P\|_1 \le \|P_i\|_1 + \|P\|_1 \le 2$ and the bounded magnitude $\|V_{i-1}\|_{\infty} \le 1$, and the third inequality follows from (103).

We denote the entrywise conditional variance of z_i as $\operatorname{Var}(z_i | V_{i-1}, \dots, V_0) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, whose j-th element is the variance of $[z_i]_j$ conditional on V_{i-1}, \dots, V_0 . By the definition of z_i in (35),

$$W_{t} = \sum_{i=1+\lfloor(1-\beta)t\rfloor}^{t} \operatorname{Var}(\boldsymbol{z}_{i} \mid \boldsymbol{V}_{i-1}, \dots, \boldsymbol{V}_{0})$$

$$= \sum_{i=1+\lfloor(1-\beta)t\rfloor}^{t} \gamma_{i}^{2} (\eta_{i}^{(t)})^{2} \operatorname{Var}((\boldsymbol{P}_{i} - \boldsymbol{P})\boldsymbol{V}_{i-1} \mid \boldsymbol{V}_{i-1}) = \sum_{i=1+\lfloor(1-\beta)t\rfloor}^{t} \gamma_{i}^{2} (\eta_{i}^{(t)})^{2} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}).$$

It can be (entrywisely) upper bounded as

$$\begin{aligned} \boldsymbol{W}_{t} &\leq \left(\max_{\lfloor (1-\beta)t \rfloor < i \leq t} \eta_{i}^{(t)} \right) \sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \gamma_{i}^{2} \eta_{i}^{(t)} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) \\ &\leq \left(\max_{\lfloor (1-\beta)t \rfloor < i \leq t} \eta_{i}^{(t)} \right) \left(\sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \eta_{i}^{(t)} \right) \left(\max_{\lfloor (1-\beta)t \rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) \right) \\ &\leq \frac{2(\log t)^{3}}{c_{1} t^{3/5}} \max_{\lfloor (1-\beta)t \rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}), \end{aligned}$$

where the last inequality follows from (103) and the fact that $\sum_{i=1+\lfloor (1-\beta)t\rfloor}^t \eta_i^{(t)} \leq 1$ due to (101). In addition, since $\|V_{i-1}\|_{\infty} \leq 1$, a deterministic upper bound of W_t is given by

$$\|\mathbf{W}_t\|_{\infty} \le \frac{2(\log t)^3}{c_1 t^{3/5}} =: \sigma^2.$$
 (37)

To apply Freedman's inequality, we further choose the (smallest) positive integer K such that for some constant $c_2 \geq 1$,

$$\frac{\sigma^2}{2^K} \le \frac{2(\log t)^3}{c_1 t^{4/5}},$$

which leads to

$$\frac{\log t}{5\log 2} \le K \le \frac{\log t}{5\log 2} + 1.$$

In view of (36) and (37) together with a union bound over all |S||A| entries, Freedman's inequality (Lemma D.6) implies that for any fixed t obeying $|(1-\beta)t| \ge 150$ and $\beta < 1/2$,

$$|\boldsymbol{\xi}_{t}| \leq \frac{8(\log t)^{3}}{3c_{1}t^{3/5}} \log \frac{2|\mathcal{S}||\mathcal{A}| \log t}{\delta \log 2} \mathbf{1} + \sqrt{8\left(\boldsymbol{W}_{t} + \frac{2(\log t)^{3}}{c_{1}t^{4/5}} \mathbf{1}\right) \log \frac{2|\mathcal{S}||\mathcal{A}| \log t}{\delta \log 2}}$$

$$\leq \frac{8(\log t)^{3} \log \frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{3c_{1}t^{3/5}} \mathbf{1} + \sqrt{\frac{16(\log t)^{3} \log \frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{c_{1}t^{3/5}}} \left(\max_{\lfloor (1-\beta)t\rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + t^{-1/5} \mathbf{1}\right)}$$
(38)

holds with probability at least $1 - \delta$, where the second line follows from the fact that $\log t/t \le \log 2/2$ for $t \ge 4$. Finally, for any fixed t obeying $T/\log T \le t \le \log T$, we have

$$\frac{(\log t)^3\log\frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{t^{3/5}}\leq \frac{(\log T)^{\frac{18}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{3/5}},\quad \frac{(\log t)^3\log\frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{t^{4/5}}\leq \frac{(\log T)^{\frac{19}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{4/5}}.$$

The bound in (38) thus translates to

$$|\boldsymbol{\xi}_t| \le 5\sqrt{\frac{(\log T)^{\frac{19}{5}}\log\frac{|A||\mathcal{S}|T}{\delta}}{c_1 T^{3/5}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \le t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-1/5}\mathbf{1}\right)},$$
 (39)

as long as t additionally satisfies

$$T^{2/5} \ge \frac{64(\log T)^4}{9c_1} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}.$$

To see (39), note that under the conditions of Theorem 3.1, one has $T^{2/5} \ge \frac{64(\log T)^4}{9c_1} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}$, hence

$$\begin{split} \frac{8(\log T)^{\frac{18}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{3c_{1}T^{3/5}}\mathbf{1} &\leq \sqrt{\frac{(\log T)^{\frac{19}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{c_{1}T^{4/5}}}\mathbf{1} \\ &\leq \sqrt{\frac{(\log T)^{\frac{19}{5}}\log\frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{c_{1}T^{3/5}}\Big(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \mathrm{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-1/5}\mathbf{1}\Big)}, \end{split}$$

which completes the proof of Lemma 4.3.

B Proof of Theorem 3.3

In this section, we provide the detailed proof of Theorem 3.3. It follows the same idea as that of Theorem 3.1, yet with different analysis on the convergence rates.

Recall that in Theorem 3.3, we set $\gamma_t = 1 - t^{-1/8}$ and

$$\eta_t = \frac{1}{1 + \frac{c_1 t^{5/8}}{(\log t)^2}}, \quad t \ge 2$$

for some constant $c_1 > 0$ in Algorithm 1. In the following, we are to set

$$\beta = \frac{c_2}{T^{1/8}(\log T)^2}$$

for the constant $c_2 > 0$ in Theorem 3.3. In this section, we still use the same notations as in the decomposition of Section 4.1.

B.1 Bound on the switching error

Following exactly the same arguments as in the proof of Lemma 4.1, the switching error can be bounded as

$$\left\| \sum_{i=1}^{t} \eta_{i}^{(t)} \boldsymbol{d}_{i} / \eta_{i} \right\|_{\infty} \leq \sum_{i=1}^{t} \prod_{j=i}^{t} (1 - \eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty}$$

$$= \sum_{i=1}^{\lfloor (1 - \alpha)t \rfloor} \prod_{j=i}^{t} (1 - \eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty} + \sum_{1 + \lfloor (1 - \alpha)t \rfloor}^{t} \prod_{j=i}^{t} (1 - \eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty},$$

where we choose $\alpha = t^{-3/8}$. Firstly, when t obeys $\lfloor (1-\alpha)t \rfloor = \lfloor t-t^{5/8} \rfloor \geq 50$ (satisfied by $t \geq 100$), for any i such that $1 \leq i \leq \lfloor (1-\alpha)t \rfloor$, we have

$$\prod_{j=i}^{t} (1 - \eta_j) \leq \prod_{j=\lfloor (1-\alpha)t \rfloor}^{t} (1 - \eta_j) \leq \exp\left(-\sum_{j=\lfloor (1-\alpha)t \rfloor}^{t} \eta_j\right)
\stackrel{\text{(i)}}{\leq} \exp\left(-\sum_{j=\lfloor (1-\alpha)t \rfloor}^{t} \frac{(\log j)^2}{(2 + c_1)j^{5/8}}\right) \stackrel{\text{(ii)}}{\leq} \exp\left(-\frac{\alpha t (\log t)^2}{(2 + c_1)t^{5/8}}\right),$$

where (i) follows from (105), and (ii) follows from (106) for $j \ge \lfloor (1 - \alpha)t \rfloor \ge 50$. Together with the fact that $\|\boldsymbol{d}_i\|_{\infty} \le 1$, the above bound implies

$$\sum_{i=1}^{\lfloor (1-\alpha)t\rfloor} \prod_{j=i}^{t} (1-\eta_j) \|\boldsymbol{d}_i\|_{\infty} \le t \cdot \exp\left(-\frac{\alpha t^{3/8} (\log t)^2}{2+c_1}\right) = \exp\left(\log t - (\log t)^2/(2+c_1)\right). \tag{40}$$

On the other hand, (30) implies

$$\|\boldsymbol{d}_i\|_{\infty} \le \frac{(i-1)^{-1/8} - i^{-1/5}}{(i-1)^{-1/8}} \le 1 - (1-1/i)^{1/8} \le \frac{1}{4i} \le \frac{1}{4(1-\alpha)t},$$

for all $i \ge 1 + \lfloor (1 - \alpha)t \rfloor$. The second inequality above is implied by

$$1 - (1 - 1/i)^{1/8} = 1 - \exp\left(\frac{\log(1 - 1/i)}{8}\right) \le 1 - \exp\left(-2/(8i)\right) \le 2/(8i),$$

where for $i \ge 2$, we use the fact that $\log(1-x) \ge -2x$ for $x \le \log 2/2$ and $e^x \ge 1+x$ for $x \in \mathbb{R}$. Therefore, noting that $\alpha \le 0.25$ for $t \ge 150$, we have

$$\sum_{1+\lfloor (1-\alpha)t\rfloor}^{t} \prod_{j=i}^{t} (1-\eta_{j}) \|\boldsymbol{d}_{i}\|_{\infty} \leq \sum_{1+\lfloor (1-\alpha)t\rfloor}^{t} \|\boldsymbol{d}_{i}\|_{\infty} \leq \frac{\alpha t}{4(1-\alpha)t} \leq \alpha/3 = \frac{1}{3t^{3/8}}.$$
 (41)

Combining (40) and (41), once $T/\log T \ge 100$, we have

$$\begin{split} \left\| \sum_{i=1}^{t} \eta_{i}^{(t)} \boldsymbol{d}_{i} / \eta_{i} \right\|_{\infty} &\leq \exp\left(\log t - (\log t)^{2} / (2 + c_{1})\right) + \frac{1}{3t^{3/8}} \\ &\leq \exp\left(\log T - \frac{(\log T - \log \log T)^{2}}{2 + c_{1}}\right) + \frac{(\log T)^{3/8}}{3T^{3/8}} \\ &\leq \exp\left(\log T - \frac{(\log T)^{2}}{4(2 + c_{1})}\right) + \frac{(\log T)^{3/8}}{3T^{3/8}} \leq \frac{2(\log T)^{3/8}}{T^{3/8}} \end{split}$$

for all $T/\log T \le t \le T$, where the third inequality follows from $\log T \ge 2\log\log T$, and the last inequality holds as long as T obeys $(\log T)^3/(8+4c_1) \ge \frac{11}{8}\log T$.

B.2 Bounds on ζ_t

Similar to the proof of Lemma 4.2, the ℓ_{∞} -norm of ζ_t can be bounded as

$$\|\boldsymbol{\zeta}_t\|_{\infty} \leq \eta_0^{(t)} + 3 \sum_{i=1}^{\lfloor (1-\beta)t\rfloor} \eta_i^{(t)} \gamma_i.$$

By definition of $\{\eta_t\}$, as long as $\beta < 1/2$ and $t \ge 100$,

$$\eta_0^{(t)} = \prod_{j=1}^t (1 - \eta_j) \le \exp\left(-\sum_{j=1}^t \eta_j\right)
\le \exp\left(-\sum_{j=50}^t \frac{(\log j)^2}{(2 + c_1)j^{5/8}}\right) \le \exp\left(-\frac{(t - 50)(\log t)^2}{(2 + c_1)t^{5/8}}\right), \tag{42}$$

where the second inequality follows from (105), and the last inequality follows from the monotonicity (106).

Similarly, when $\lfloor (1-\beta)t \rfloor \geq 50$ (which holds as long as $\beta < 1/2$ and $t \geq 100$), for any $1 \leq i \leq \lfloor (1-\beta)t \rfloor$, it holds that

$$\eta_i^{(t)} \le \prod_{j=1+\lfloor (1-\beta)t\rfloor}^t (1-\eta_j) \le \exp\left(-\sum_{j=1+\lfloor (1-\beta)t\rfloor}^t \eta_j\right) \le \exp\left(-\frac{\beta t (\log t)^2}{(2+c_1)t^{5/8}}\right). \tag{43}$$

Combining (42) and (43), once $\beta < 1/2$, we have

$$\|\zeta_t\|_{\infty} \le \exp\left(-\frac{(t-50)(\log t)^2}{(2+c_1)t^{5/8}}\right) + 3\exp\left(\log t - \frac{\beta t(\log t)^2}{(2+c_1)t^{5/8}}\right)$$

$$\le \exp\left(-\frac{t(\log t)^2}{2(2+c_1)t^{5/8}}\right) + 3\exp\left(\log t - \frac{\beta t(\log t)^3}{(2+c_1)t^{5/8}}\right)$$

for all $t \ge 100$. Note that when $T/\log T \ge 100$, for all t such that $T/\log T \le t \le T$,

$$t^{3/8}(\log t)^2 \geq (T/\log T)^{3/8}(\log T - \log\log T)^2 \geq (T/\log T)^{3/8} \cdot (\log T)^2/4 \geq T^{3/8}(\log T)^2/4,$$

where the second inequality follows from $\log T \geq 2 \log \log T$. Meanwhile, for $T/\log T \leq t \leq T$,

$$\log t - \frac{\beta t (\log t)^2}{(2+c_1)t^{5/8}} \le \log T - \frac{c_2}{T^{1/8}(\log T)^2} \frac{T^{3/8}}{(\log T)^{3/8}} \frac{(\log T)^2}{4(2+c_1)} \le \log T - \frac{c_2 T^{1/4}}{4(2+c_1)(\log T)^{3/8}}$$

Therefore, for any $T/\log T \le t \le \log T$,

$$\|\zeta_t\|_{\infty} \le \exp\left(-\frac{T^{3/8}(\log T)^2}{2(2+c_1)}\right) + \exp\left(\log(3T) - \frac{c_2T^{1/4}}{4(2+c_1)(\log T)^{3/8}}\right) \le \frac{4}{T}$$

if T obeys

$$T^{3/8} \log T \ge 2(2+c_1)$$
 and $c_2 T^{1/4} \ge 4(2+c_1)(\log T)^{11/8}$.

B.3 Bound on ξ_t

Similar to the proof of Lemma 4.3, we write

$$\boldsymbol{\xi}_{t} = \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \boldsymbol{z}_{i}, \quad \text{where} \quad \boldsymbol{z}_{i} = \eta_{i}^{(t)} \gamma_{i} (\boldsymbol{P}_{i} - \boldsymbol{P}) \boldsymbol{V}_{i-1}. \tag{44}$$

Then (entries of) $\{z_i\}$ are martingale differences, and we shall obtain high-probability bound on $\|\xi_t\|_{\infty}$ via Freedman's inequality as usual.

Firstly, parallel to (36), when $\lfloor (1-\beta)t \rfloor \geq 50$ and $1-\beta \geq 1/2$,

$$\max_{\lfloor (1-\beta)t \rfloor < i \le t} \| \boldsymbol{z}_i \|_{\infty} \le \max_{\lfloor (1-\beta)t \rfloor < i \le t} 2\eta_i \le \frac{4(\log t)^2}{c_1 t^{5/8}} =: R, \tag{45}$$

where the third inequality follows from (107).

We denote the entrywise conditional variance of z_i as $Var(z_i | V_{i-1}, ..., V_0) \in \mathbb{R}^{|S||A|}$, whose j-th element is the variance of $[z_i]_j$ conditional on $V_{i-1}, ..., V_0$. By the definition of z_i in (44),

$$\boldsymbol{W}_{t} = \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \operatorname{Var}(\boldsymbol{z}_{i} \mid \boldsymbol{V}_{i-1}, \dots, \boldsymbol{V}_{0}) = \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \gamma_{i}^{2} (\eta_{i}^{(t)})^{2} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}).$$

Similar to the situation in Lemma 4.3 (c.f. Appendix A.6), we have the crude entrywise upper bounded

$$W_{t} \leq \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \eta_{i}^{(t)}\right) \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \gamma_{i}^{2} \eta_{i}^{(t)} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i-1})$$

$$\leq \frac{2(\log t)^{2}}{c_{1}t^{5/8}} \max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i-1}),$$

where the last inequality follows from (107) and the fact that $\sum_{i=1+\lfloor (1-\beta)t\rfloor}^t \eta_i^{(t)} \leq 1$ due to (101). In addition, by the boundedness of $\|V_{i-1}\|_{\infty} \leq 1$, we also have a deterministic upper bound of W_t that

$$\|\boldsymbol{W}_t\|_{\infty} \le \frac{2(\log t)^2}{c_1 t^{5/8}} =: \sigma^2.$$
 (46)

To apply Freedman's inequality, we further choose the (smallest) positive integer K such that for some constant $c_2 \ge 1$,

$$\frac{\sigma^2}{2^K} \le \frac{2(\log t)^2}{c_1 t},$$

which leads to

$$\frac{3\log t}{8\log 2} \le K \le \frac{3\log t}{8\log 2} + 1.$$

In view of (45) and (46) together with a union bound over all $|\mathcal{S}||\mathcal{A}|$ entries, Freedman's inequality (Lemma D.6) implies that for any fixed t obeying $\lfloor (1-\beta)t \rfloor \geq 50$ and $\beta < 1/2$,

$$|\xi_{t}| \leq \frac{8(\log t)^{2}}{3c_{1}t^{5/8}} \log \frac{2|\mathcal{S}||\mathcal{A}| \log t}{\delta \log 2} \mathbf{1} + \sqrt{8\left(\mathbf{W}_{t} + \frac{2(\log t)^{2}}{c_{1}t}\mathbf{1}\right) \log \frac{2|\mathcal{S}||\mathcal{A}| \log t}{\delta \log 2}}$$

$$\leq \frac{8(\log t)^{2} \log \frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{3c_{1}t^{5/8}} \mathbf{1} + \sqrt{\frac{16(\log t)^{2} \log \frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{c_{1}t^{5/8}}} \left(\max_{|(1-\beta)t| < i \leq t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i-1}) + t^{-3/8}\mathbf{1}\right)}$$
(47)

holds with probability at least $1 - \delta$, where the second line follows from the fact that $\log t/t \le \log 2/2$ for $t \ge 4$. Finally, for any fixed t obeying $T/\log T \le t \le \log T$, we have

$$\frac{(\log t)^2 \log \frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{t^{5/8}} \leq \frac{(\log T)^{\frac{21}{8}} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{5/8}}, \quad \frac{(\log t)^2 \log \frac{|\mathcal{S}||\mathcal{A}|t}{\delta}}{t} \leq \frac{(\log T)^3 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T}.$$

The bound in (47) thus translates to

$$|\boldsymbol{\xi}_t| \leq 5\sqrt{\frac{(\log T)^3 \log \frac{|\mathcal{A}||\mathcal{S}|T}{\delta}}{c_1 T^{5/8}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-3/8}\mathbf{1}\right)},$$

as long as t additionally satisfies

$$T^{5/8} \ge \frac{64(\log T)^3}{9c_1} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}.$$

B.4 Recursive bound

Combining the three bounds in preceding subsections, for sufficiently large T satisfying all the mentioned conditions and any fixed t obeying $T/\log T \le t \le T$, it holds with probability at least $1-\delta$ that

$$\Delta_{t} \leq \frac{2(\log T)^{3/8}}{T^{3/8}} + \frac{4}{T} + 5\sqrt{\frac{(\log T)^{3} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{c_{1}T^{5/8}}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-3/8} \mathbf{1} \right)
+ \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \eta_{i}^{(t)} \gamma_{i} \boldsymbol{P}^{\pi_{t-1}} \Delta_{i-1}
\leq \left(1 + \frac{5}{\sqrt{c_{1}}}\right) \sqrt{\frac{(\log T)^{3} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{5/8}}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i-1}) + T^{-3/8} \mathbf{1} \right)
+ \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \eta_{i}^{(t)} \gamma_{i} \boldsymbol{P}^{\pi_{i-1}} \Delta_{i-1}, \tag{48}$$

where the second inequality follows from the fact that

$$\frac{2(\log T)^{3/8}}{T^{3/8}} + \frac{4}{T} \le \sqrt{\frac{(\log T)^3 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{5/8}}}$$

for T so large that $T/\log T \ge 100$. Once $(1-\beta) \ge 3/4$, for any fixed t obeying $3T/(2\log T) \le t \le T$, we apply (48) with a union bound over $\{k : \frac{2t}{3} \le k \le t\}$ — as a result, with probability at least $1-\delta$, one has

$$\Delta_k \le \sqrt{\varphi_t} + \sum_{i=1+\lfloor (1-\beta)k\rfloor}^k \eta_i^{(k)} \gamma_i \mathbf{P}^{\pi_{i-1}} \Delta_{i-1}, \quad \text{for all } \frac{2t}{3} \le k \le t,$$
(49)

where we define

$$\varphi_{t} = 2\left(1 + \frac{5}{\sqrt{c_{1}}}\right)^{2} \cdot \frac{(\log T)^{3} \log \frac{|S||\mathcal{A}|T}{\delta}}{T^{5/8}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i-1}) + T^{-3/8}\mathbf{1}\right) \\
\geq \left(1 + \frac{5}{\sqrt{c_{1}}}\right)^{2} \cdot \frac{(\log T)^{3} \log \frac{|S||\mathcal{A}|T}{\delta/T}}{T^{5/8}} \left(\max_{\lfloor (1-\beta)t\rfloor < i \leq t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i-1}) + T^{-3/8}\mathbf{1}\right).$$
(50)

In the following, we are to follow exactly the same recipe as in the proof of Proposition 4.4 (c.f. Appendix A.3), where we correspondingly define the $\{\alpha_i^{(t)}\}$ according this set of learning rates $\{\eta_j\}$ here. We also set

$$H := T^{1/8} (\log T)^2$$

so that $(1-\beta)^H \ge \exp(-2\beta H) \ge 2/3$ as long as $c_2 \le \log(3/2)/2$. Thus, when T is sufficiently large so that $1-\beta \ge 2/3$, in parallel with (24), we have

$$\boldsymbol{\Delta}_{t} \leq \max_{(i_{1},\dots,i_{H})\in\mathcal{I}_{t}} \left\{ \underbrace{\left(\boldsymbol{I} + \sum_{h=1}^{H-1} \prod_{k=1}^{h} (\gamma_{i_{k}+1} \boldsymbol{P}^{\pi_{i_{k}}})\right) \sqrt{\varphi_{t}}}_{=:\boldsymbol{\beta}_{1}} + \underbrace{\prod_{h=1}^{H} (\gamma_{i_{h}+1} \boldsymbol{P}^{\pi_{i_{h}}}) |\boldsymbol{\Delta}_{i_{H}}|}_{=:\boldsymbol{\beta}_{2}} \right\}.$$
(51)

In the following, we are to bound β_1 and β_2 in (51) separately. By definition of the discount factor $\{\gamma_t\}$, we have $\gamma_j \leq \gamma_T$ for all $2t/3 \leq j \leq T$, indicating

$$|\beta_2| \le \gamma_t^H \prod_{h=1}^H \mathbf{P}^{\pi_{i_h}} |\mathbf{\Delta}_H| \le \gamma_T^H \left\| \prod_{h=1}^H \mathbf{P}^{\pi_{i_h}} \right\|_1 \|\mathbf{\Delta}_H\|_{\infty} \le \gamma_T^H \le \frac{1}{T}, \tag{52}$$

where (i) follows from the bounded magnitude $\|\Delta_H\|_{\infty} \leq 1$ and the fact that $\prod_{h=1}^H P^{\pi_{i_h}}$ is a probability transition matrix; (ii) follows from

$$\gamma_T^H = \left(1 - T^{-1/8}\right)^{T^{1/8}(\log T)^2} \le \exp\left(-(\log T)^2\right) \le \frac{1}{T}.$$

Again, for β_1 , parallel to (26), we have the entrywise upper bound

$$|\beta_{1}|^{2} \leq \frac{1}{1 - \gamma_{T}} \sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \frac{2(1 + \frac{5}{\sqrt{c_{1}}})^{2} (\log T)^{3} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta/T}}{T^{5/8}} \left(\max_{\lfloor (1-\beta)k \rfloor < i \leq k} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i-1}) + T^{-3/8} \mathbf{1} \right)$$

$$\leq \frac{4(1 + \frac{5}{\sqrt{c_{1}}})^{2} (\log T)^{3} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/2}} \left(\sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i}) + \mathbf{1} \right),$$
(53)

We employ the following lemma to bound the first term of (53); it is parallel to Lemma A.1, and the proof follows exactly the same arguments hence we omit here.

Lemma B.1. Suppose $T \ge 100$, then it holds for all $T/\log T \le t \le T$ that

$$\sum_{h=0}^{H-1} \gamma_T^h \prod_{k=1}^h \mathbf{P}^{\pi_{i_k}} \max_{\lfloor t/2 \rfloor \le i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_i) \le \left(7 + 72t_{mix} (\log T)^{1/8} + 8T^{1/8} \max_{\lfloor t/2 \rfloor \le i < t} \|\mathbf{\Delta}_i\|_{\infty}\right) \mathbf{1}.$$
 (54)

Invoking Lemma B.1, we obtain an upper bound

$$|\beta_1|^2 \le \frac{4(1+\frac{5}{\sqrt{c_1}})^2(\log T)^3 \log \frac{|S||A|T}{\delta}}{T^{1/2}} \left(8+72t_{\text{mix}}(\log T)^{1/8}+8T^{1/8} \max_{\lfloor t/2\rfloor \le i < t} \|\boldsymbol{\Delta}_i\|_{\infty}\right) \mathbf{1}. \tag{55}$$

Recalling (51), by the upper bound (52) on $|\beta_2|$ and (55) on $|\beta_1|$, for any fixed t such that $3T/(2\log T) \le t \le T$,

$$\Delta_{t} \leq \frac{1}{T} + \sqrt{\frac{4(1 + \frac{5}{\sqrt{c_{1}}})^{2}(\log T)^{3} \log \frac{|S||A|T}{\delta}}{T^{1/2}} \left(8 + 72t_{\text{mix}}(\log T)^{1/8} + 8T^{1/8} \max_{\lfloor t/2 \rfloor \leq i < t} \|\Delta_{i}\|_{\infty}\right) \mathbf{1}}$$

$$\leq 4\left(1 + \frac{5}{\sqrt{c_{1}}}\right) \sqrt{\frac{(\log T)^{4} \log \frac{|S||A|T}{\delta}}{T^{1/2}}} \left(8(\log T)^{-1/8} + 72t_{\text{mix}} + 8(T/\log T)^{1/8} \max_{\lfloor t/2 \rfloor \leq i < t} \|\Delta_{i}\|_{\infty}\right) \mathbf{1}} \tag{56}$$

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

We now further take a union bound for (56) over $\{t: 3T/(2\log T) \le t \le T\}$, which leads to the simultaneous $(1-\delta)$ high-probability bound

$$\Delta_t \leq 16 \left(1 + \frac{5}{\sqrt{c_1}}\right) \sqrt{\frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/2}} \left((\log T)^{-1/8} + 9t_{\text{mix}} + (T/\log T)^{1/8} \max_{\lfloor t/2 \rfloor \leq i < t} \|\Delta_i\|_{\infty}\right)} \mathbf{1}$$

for all t such that $3T/(2\log T) \le t \le T$, since $\log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta/T} \le 2\log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}$

B.5 Solving the recursive bound

We solve the recursive bounds to obtain the final high-probability bound. For any positive integer k, we define

$$u_k = \max \left\{ \|\boldsymbol{\Delta}_i\|_{\infty} \colon 2^k \frac{3T}{2\log T} \le i \le T \right\}.$$

A naive upper bound is $u_k \leq 1$ for all k. Furthermore, by (15) and the definition of u_k , with probability at least $1 - \delta$, it holds simultaneously for all $k \geq 1$ that

$$u_{k+1} \le \sqrt{\frac{c_3(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/2}}} \Big((\log T)^{-1/8} + 9t_{\text{mix}} + (T/\log T)^{1/8} u_k \Big)$$

for some absolute constant $c_3 > 0$. If there exists some k such that $2^{k+1} \le 2 \log T/3$ and

$$u_k \le ((\log T)^{-1/8} + 9t_{\text{mix}})(T/\log T)^{-1/8},$$
(57)

then $2^{k+1} \frac{3T}{2 \log T} \le T$, and

$$\|\boldsymbol{\Delta}_T\|_{\infty} \le u_{k+1} \le \sqrt{\frac{c_3(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/2}}} (2(\log T)^{-1/8} + 18t_{\text{mix}}).$$

On the other hand, suppose $u_j > ((\log T)^{-1/8} + 9t_{\text{mix}})(T/\log T)^{-1/8}$ for any $1 \le j \le k$ for which $2^{k+1} \le 2\log T/3$. Then

$$u_{j+1} \le \sqrt{\frac{2c_3(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{3/8}} u_j}, \quad \text{for all } 1 \le j \le k,$$

which implies

$$\log u_{j+1} \le \frac{1}{2} \log u_j + \frac{1}{2} \log \theta, \quad \text{where} \quad \theta = 2c_3 \cdot T^{-3/8} (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}.$$

Recursively applying this relation for $1 \le j < k$ yields

$$2^k \log u_k \le \log u_0 + \sum_{j=0}^{k-1} 2^j \log \theta \le (2^k - 1) \log \theta,$$

where we used the fact that $u_0 \leq 1$. Hence

$$u_k \le \theta^{1-1/2^k} = \left(2c_3 \cdot T^{-3/8} (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}\right)^{1-1/2^k}$$

$$= \left(2c_3 \cdot (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}\right)^{1-1/2^k} \cdot T^{-3/8} \cdot T^{\frac{3}{8\cdot 2^k}}$$

$$\le 2c_3 \cdot \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{3/8}} \cdot \exp\left(3 \log T \cdot 2^{-k}/8\right).$$

Now we can set $2^k \ge 3 \log T/16$, so that $2^k 3T/(2 \log T) \le T$, and the above upper bound translates to

$$\|\Delta_T\|_{\infty} \le u_k \le 16c_3 \cdot \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{3/8}}.$$

Combining the two cases above, we arrive at

$$\|\mathbf{\Delta}_T\|_{\infty} \le \left(c_4\sqrt{(\log T)^{-1/8} + 9t_{\text{mix}}} + 16c_3T^{-1/8}\right) \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/4}}$$

with probability at least $1 - \delta$ for T obeying the conditions of Proposition 4.4.

We now take a moment to collect all the conditions we impose on T, which are

$$(\log T)^3/(8+4c_1) \ge \frac{11}{8} \log T, \quad T/\log T \ge 100, \quad T^{3/8} \log T \ge 2(2+c_1),$$

$$c_2 T^{1/4} \ge 4(2+c_1)(\log T)^{11/8}, \quad T^{5/8} \ge \frac{64(\log T)^3}{9c_1} \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}.$$

They can be further simplified to

$$(\log T)^2 \ge 11(2+c_1)/2$$
, $T/\log T \ge 100$,
 $c_2 T^{1/4} \ge 4(2+c_1)(\log T)^{11/8}$, $T^{5/8} \ge \frac{64(\log T)^3}{9c_1}\log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}$.

B.6 From estimation to policy learning

Finally, we are to leverage the performance difference lemma for discounted reward MDP to obtain the bounds on average reward performance.

Invoking Lemma D.9 with $\gamma = \gamma_T$, $\pi' = \pi_{\gamma_T}^*$ and $\pi = \pi_T$, the greedy policy obtained in the T-th iteration, we know that

$$V_{\gamma_T}^{\pi_T}(s) - V_{\gamma_T}^*(s) \ge -T^{1/8} \cdot \sup_{s' \in \mathcal{S}} |Q_{\gamma_T}^*(s', \pi_T(s')) - V_{\gamma_T}^*(s')|.$$

Since π_T is the greedy policy with respect to Q_T , for any $s' \in \mathcal{S}$, we abve

$$\begin{aligned} Q_{\gamma_T}^*(s', \pi_T(s')) &\geq Q_T(s', \pi_T(s')) - \|\mathbf{\Delta}_T\|_{\infty} \\ &\geq Q_T(s', \pi_{\gamma_T}^*(s')) - \|\mathbf{\Delta}_T\|_{\infty} \\ &\geq Q_{\gamma_T}^*(s', \pi_{\gamma_T}^*(s')) - 2\|\mathbf{\Delta}_T\|_{\infty} \geq V_{\gamma_T}^*(s') - 2\|\mathbf{\Delta}_T\|_{\infty}, \end{aligned}$$

and also $Q_{\gamma_T}^*(s', \pi_T(s')) \leq \max_{a \in \mathcal{A}} Q_{\gamma_T}^*(s', a) = V_{\gamma_T}^*(s')$. Therefore, we have $0 \geq V_{\gamma_T}^{\pi_T}(s) - V_{\gamma_T}^*(s) \geq -2\|\mathbf{\Delta}_T\|_{\infty}$. Further invoking Lemma D.1, and recalling that π^* is the optimal policy for average reward, we have

$$V^{\pi_{T}}(s) \geq V_{\gamma_{T}}^{\pi_{T}}(s) - (1 - \gamma_{T})t_{\text{mix}}$$

$$\geq V_{\gamma_{T}}^{*}(s) - 2\|\mathbf{\Delta}_{T}\|_{\infty} - 3(1 - \gamma_{T})t_{\text{mix}}$$

$$\geq V_{\gamma_{T}}^{\pi^{*}}(s) - 2\|\mathbf{\Delta}_{T}\|_{\infty} - 3(1 - \gamma_{T})t_{\text{mix}}$$

$$\geq V^{\pi^{*}}(s) - 2\|\mathbf{\Delta}_{T}\|_{\infty} - 6(1 - \gamma_{T})t_{\text{mix}},$$

where the first line follows from Lemma D.1 for π_T , the second line uses the previous result, the third line uses the optimality of $\pi_{\gamma T}^*$ for γ_T -discounted reward, and the last line uses Lemma D.1 for π^* . Therefore, on the event that $\|\mathbf{\Delta}_T\|_{\infty} \leq \left(c_4\sqrt{(\log T)^{-1/8} + 3t_{\text{mix}}} + 16c_3T^{-1/8}\right) \frac{(\log T)^4 \log \frac{|S||A|T}{\delta}}{T^{1/4}}$, we have

$$V^*(s) - V^{\pi_T}(s) \le 2\left(c_4\sqrt{(\log T)^{-1/8} + 3t_{\text{mix}}} + 16c_3T^{-1/8}\right) \frac{(\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}}{T^{1/8}} + 6T^{-1/8}t_{\text{mix}}.$$

Since $t_{\text{mix}} \geq 1$, the above bound translates to

$$V^*(s) - V^{\pi_T}(s) \le \frac{ct_{\text{mix}}}{T^{1/8}} (\log T)^4 \log \frac{|\mathcal{S}||\mathcal{A}|T}{\delta}$$

for some absolute constant c > 0 that only depends on c_1, c_2 . This completes the proof of Theorem 3.3.

C Technical proofs regarding Theorem 5.2

In this section, we provide detailed proofs for results regarding Algorithm 2. Appendix C.1 provides a proof sketch for Theorem 5.2, while the remaining of this section provides detailed proofs for supportive results.

C.1 Sketch of analysis for Algorithm 2

In this section, we provide a sketch of analysis for the estimation error of Algorithm 2, which forms the basis for proving Theorem 5.2. The general approach is similar to Section 4.1 with slightly different bounds. Let

$$\boldsymbol{\delta}_t = \mathbf{q}_t - \mathbf{q}_{\alpha_t}^*,$$

for $t \geq 0$, which is the estimation error of Algorithm 2 in the t-th iteration.

Our updating rule (17) in the t-th iteration satisfies

$$\mathbf{q}_t = (1 - \theta_t)\mathbf{q}_{t-1} + \theta_t [\mathbf{r} + \alpha_t \mathbf{P}_t \mathbf{v}_{t-1}].$$

Employing the Bellman equation $\mathbf{q}_{\alpha_t}^* = \mathbf{r} + \alpha_t \mathbf{P} \mathbf{v}_{\alpha_t}^*$, we obtain the decomposition

$$\delta_{t} = \mathbf{q}_{t} - \mathbf{q}_{\alpha_{t}}^{*} = (1 - \theta_{t})\mathbf{q}_{t-1} + \theta_{t} \left[\mathbf{r} + \alpha_{t} \mathbf{P}_{t} \mathbf{v}_{t-1} \right] - \mathbf{q}_{\alpha_{t}}^{*}
= (1 - \theta_{t})\delta_{t-1} + (1 - \theta_{t})(\mathbf{q}_{\alpha_{t-1}}^{*} - \mathbf{q}_{\alpha_{t}}^{*}) + \theta_{t} \left[\mathbf{r} + \alpha_{t} \mathbf{P}_{t} \mathbf{v}_{t-1} - \mathbf{q}_{\alpha_{t}}^{*} \right]
= (1 - \theta_{t})\delta_{t-1} + (1 - \theta_{t})(\mathbf{q}_{\alpha_{t-1}}^{*} - \mathbf{q}_{\alpha_{t}}^{*}) + \theta_{t} \left[\alpha_{t} \mathbf{P}_{t} \mathbf{v}_{t-1} - \alpha_{t} \mathbf{P} \mathbf{v}_{\alpha_{t}}^{*} \right]
= (1 - \theta_{t})\delta_{t-1} + \theta_{t}\alpha_{t} \left[\mathbf{P}_{t} \mathbf{v}_{t-1} - \mathbf{P} \mathbf{v}_{\alpha_{t-1}}^{*} \right] + (1 - \theta_{t})(\mathbf{q}_{\alpha_{t-1}}^{*} - \mathbf{q}_{\alpha_{t}}^{*}) + \theta_{t} \mathbf{P}(\mathbf{q}_{\alpha_{t-1}}^{*} - \mathbf{q}_{\alpha_{t}}^{*})$$
(58)

Similar as the arguments in the analysis of Algorithm 1 (see Section 4), we have

$$\boldsymbol{P}_t \mathbf{v}_{t-1} - \boldsymbol{P} \mathbf{v}^*_{\alpha_{t-1}} = (\boldsymbol{P}_t - \boldsymbol{P}) \mathbf{v}_{t-1} + \boldsymbol{P} (\mathbf{v}_{t-1} - \mathbf{v}^*_{\alpha_{t-1}}),$$

where the term $P(\mathbf{v}_{t-1} - \mathbf{v}_{\alpha_{t-1}}^*)$ is linked to $\boldsymbol{\delta}_{t-1}$ via

$$P(\mathbf{v}_{t-1} - \mathbf{v}_{\alpha_{t-1}}^*) \le P^{\pi_{t-1}} \delta_{t-1},\tag{59a}$$

$$P(\mathbf{v}_{t-1} - \mathbf{v}_{\alpha_{t-1}}^*) \ge P^{\pi_{\alpha_{t-1}}^*} \delta_{t-1}. \tag{59b}$$

Here with some abuse of notations, we define π_t as the greedy policy with respect to q_t . Plugging (59a) and (59b) back into (58) leads to the recursive relation

$$\delta_t \le (1 - \theta_t)\delta_t + \theta_t \alpha_t \left[\mathbf{P}^{\pi_{t-1}} \delta_{t-1} + (\mathbf{P}_t - \mathbf{P}) \mathbf{v}_{t-1} \right] + \bar{\mathbf{d}}_t; \tag{60a}$$

$$\delta_t \ge (1 - \theta_t)\delta_t + \theta_t \alpha_t \left[\boldsymbol{P}^{\pi_{\alpha_{t-1}}^*} \delta_{t-1} + (\boldsymbol{P}_t - \boldsymbol{P}) \mathbf{v}_{t-1} \right] + \bar{\boldsymbol{d}}_t, \tag{60b}$$

where we define the switching error in the t-th iteration as

$$\bar{\boldsymbol{d}}_t = (1 - \theta_t)(\mathbf{q}_{\alpha_{t-1}}^* - \mathbf{q}_{\alpha_t}^*) + \theta_t \boldsymbol{P}(\mathbf{q}_{\alpha_{t-1}}^* - \mathbf{q}_{\alpha_t}^*).$$

Applying (60a) and (60b) recursively, we arrive at

$$\delta_{t} \leq \theta_{0}^{(t)} \delta_{0} + \sum_{i=1}^{t} \theta_{i}^{(t)} \mathbf{d}_{i} / \theta_{i} + \sum_{i=1}^{t} \theta_{i}^{(t)} \alpha_{i} \left[(\mathbf{P}_{i} - \mathbf{P}) \mathbf{v}_{i-1} + \mathbf{P}^{\pi_{t-1}} \delta_{i-1} \right],$$

$$\delta_{t} \geq \theta_{0}^{(t)} \delta_{0} + \sum_{i=1}^{t} \theta_{i}^{(t)} \mathbf{d}_{i} / \theta_{i} + \sum_{i=1}^{t} \theta_{i}^{(t)} \alpha_{i} \left[(\mathbf{P}_{i} - \mathbf{P}) \mathbf{v}_{i-1} + \mathbf{P}^{\pi_{\alpha_{i-1}}^{*}} \delta_{i-1} \right].$$

$$(61)$$

where we define $\theta_t^{(t)} = 1$,

$$\theta_0^{(t)} = \prod_{j=1}^t (1 - \theta_j), \text{ and } \theta_i^{(t)} = \theta_i \cdot \prod_{j=i+1}^t (1 - \theta_j), \forall i \ge 1.$$

Now we set (with some abuse of notation)

$$\beta = \frac{c_2}{T^{1/3}(\log T)^2}$$

for some constant $c_2 > 0$. The upper bound of (61) can be decomposed as

$$\delta_{t} \leq \underbrace{\theta_{0}^{(t)} \delta_{0} + \sum_{i=1}^{\lfloor (1-\beta)t \rfloor} \theta_{i}^{(t)} \alpha_{i} [(\boldsymbol{P}_{i} - \boldsymbol{P}) \mathbf{v}_{i-1} + \boldsymbol{P}^{\pi_{t-1}} \delta_{i-1}]}_{=:\boldsymbol{\zeta}_{t}} + \underbrace{\sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \theta_{i}^{(t)} \alpha_{i} (\boldsymbol{P}_{i} - \boldsymbol{P}) \mathbf{v}_{i-1}}_{=:\boldsymbol{\xi}_{t}} + \underbrace{\sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \theta_{i}^{(t)} \alpha_{i} \boldsymbol{P}^{\pi_{i-1}} \delta_{i-1} + \underbrace{\sum_{i=1}^{t} \theta_{i}^{(t)} \bar{\boldsymbol{d}}_{i} / \theta_{i}}_{=:\boldsymbol{\omega}_{t}}$$

The three terms, similar to before, are bounded as follows.

Lemma C.1. For all sufficiently large t such that $t \ge \max\{e^{2+c_1}, 100\}$, it holds that

$$\left\| \sum_{i=1}^t \theta_i^{(t)} \boldsymbol{d}_i / \theta_i \right\|_{\infty} \le 3t^{-2/9}.$$

Proof of Lemma C.1. See Appendix C.5 for a detailed proof.

Lemma C.2. Suppose T is sufficiently large such that $c_2 \log T \ge 5(2+c_1)(\log T)^{1/3} + \log \log T$, $c_1T^3 \ge 3(\log T)^5$ and $(T/\log T)^{1/3} \ge 4(2+c_1)$. Then

$$\left\| \theta_0^{(t)} \boldsymbol{\delta}_0 + \sum_{i=1}^{\lfloor (1-\beta)t \rfloor} \theta_i^{(t)} \theta_i^{(t)} \alpha_i (\boldsymbol{P}_i - \boldsymbol{P}) \mathbf{v}_{i-1} \right\|_{\infty} \le 2/t$$

for any t such that $T/\log T \le t \le T$.

Proof of Lemma C.2. See Appendix C.6 for a detailed proof.

We utilize the Freedman's inequality to control the term ξ_t .

Lemma C.3. When $\beta < 1/2$ and $T/\log T \ge \max\{e^{c_1}, 100\}$, for any fixed t such that $T/\log T \le t \le T$,

$$|\boldsymbol{\xi}_t| \leq \sqrt{\frac{c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^2}{T^{2/3}} \left(\max_{1+\lfloor (1-\beta)t\rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) + 2(\log T)^6\right)}$$

holds with probability at least $1 - \delta$, where $c_3 = 16/c_1 + 144/c_1^2$ is an absolute constant.

Proof of Lemma C.3. See Appendix C.7 for a detailed proof.

Based on the above three lemmas, we have the following recursive relation, which form the base of proving Theorem 5.2.

Proposition C.4. Suppose T is sufficiently large such that

$$T/\log T \ge \max\{e^{2+c_1}, 100, 64(2+c_1)^3\},\$$
 $c_3(\log T)^5 \ge 5, \quad c_1 T^3 \ge 3(\log T)^5,\$
 $2c_3 T \ge (1+3/c_1)^2(\log T)^2$
 $c_2 \log T \ge 5(2+c_1)(\log T)^{1/3} + \log\log T$

for some constant $c_2 > 0$ and $c_3 = 16/c_1 + 144/c_1^2$. Then with probability at least $1 - \delta$,

$$\|\boldsymbol{\delta}_t\|_{\infty} \leq \sqrt{\frac{4c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}} \left(1 + c_4(\log T)^2 \max_{|t/2| \leq i < t} \|\boldsymbol{\delta}_i\|_{\infty}\right)}$$

holds simultaneously for all $3T/(2\log T) \le t \le T$, where we define the constant $c_4 = 2(1+3/c_1)$.

Proof of Proposition C.4. See Appendix C.4 for a detailed proof.

C.2 Proof of Theorem 5.2

Proof of Theorem 5.2. We solve the recursive bound of Propositions C.4 that

$$\|\boldsymbol{\delta}_{t}\|_{\infty} \leq \sqrt{\frac{4c_{3}(\log T)^{2}(\log\frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^{2}}{T^{2/9}}\left(1 + c_{4}(\log T)^{2} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\delta}_{i}\|_{\infty}\right)}$$
(62)

with probability at least $1 - \delta$ for all $3T/(2 \log T) \le t \le T$. For any positive integer k, we define

$$w_k = \max \left\{ \|\boldsymbol{\delta}_i\|_{\infty} \colon 2^k \frac{3T}{2\log T} \le i \le T \right\}.$$

By (62) and the definition of w_k , with probability at least $1 - \delta$, it holds simultaneously for all $k \ge 1$ that

$$u_{k+1} \le \sqrt{\frac{4c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}} \left(1 + c_4(\log T)^2 u_k\right)}.$$

If there exists some k such that $u_k \leq 1$ and $2^{k+1} \leq 2 \log T/3$, then

$$\|\boldsymbol{\delta}_T\|_{\infty} \le u_{k+1} \le \sqrt{\frac{4c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}} \left(1 + c_4(\log T)^2\right)} \le \frac{c_5(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{1/9}}$$

for some sufficiently large constant $c_5 > 0$. Otherwise, suppose $u_j > 1$ for any $1 \le j \le k$ where $k \le |\log_2(2\log T/3)| - 1$, then

$$u_{j+1} \le \sqrt{\frac{4c_3(1+c_4)(\log T)^4(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}}u_j}, \text{ for all } 1 \le j \le k,$$

which implies

$$\log u_{j+1} \le \frac{1}{2} \log u_j + \frac{1}{2} \log \theta$$
, where $\theta = \frac{4c_3(1+c_4)(\log T)^4(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}}$.

Recursively applying this relation for $1 \le j < k$ yields

$$2^k \log u_k \le \log u_0 + \sum_{j=0}^{k-1} 2^j \log \theta \le \log T + (2^k - 1) \log \theta,$$

where we used the fact that $u_0 \leq \max_{1 \leq i \leq T} \{ \|\mathbf{v}_{\alpha_i}^*\|_{\infty} + \|\mathbf{v}_i\|_{\infty} \} \leq 3(\log T)^2 T^{1/3}/c_1 + T^{1/9}$, as long as T satisfies $3(\log T)^2 T^{1/3}/c_1 + T^{1/9} \leq T$. Hence

$$u_k \le \theta^{1-1/2^k} (\log T)^{1/2^k} \le \frac{4c_3(1+c_4)(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}} \cdot T^{2^{-k+1}}$$

as long as $T^{7/9} \ge 4c_3(1+c_4)(\log T)^4(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2$. Now we can set $k = \log_2(T/6)$, so that $2^k \cdot 3T/\log T < T$ and $T^{2^{-k+1}} < e^{12}$, hence

$$\|\boldsymbol{\delta}\|_T \le \frac{c_6(\log T)^4(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}}$$

for the constant $c_6 = 4c_3(1 + c_4)$. Combining the two cases above, for any fixed T obeying the given conditions,

$$\|\boldsymbol{\delta}\|_T \le \frac{c_5(\log T)^4(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{1/9}}$$

holds with probability at least $1 - \delta$ for some constant $c_5 > 0$. Finally, by Lemma 5.1, we have

$$\|\mathbf{q}_{T} - \mathbf{q}^{*} - J^{*}/(1 - \alpha_{T})\mathbf{1}\|_{\infty} \leq \|\boldsymbol{\delta}_{T}\|_{\infty} + \|\mathbf{q}_{\alpha_{T}}^{*} - \mathbf{q}^{*} - J^{*}/(1 - \alpha_{T})\mathbf{1}\|_{\infty}$$

$$\leq \frac{B + c_{5}(\log T)^{4}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^{2}}{T^{1/9}}.$$

Put it another way, given any $\varepsilon > 0$, as long as T satisfies

$$T^{1/9} \ge \frac{B + c_5(\log T)^4(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{\varepsilon},$$

we have $\|\mathbf{q}_T - \mathbf{q}^* - J^*/(1 - \alpha_T)\mathbf{1}\|_{\infty} \le \varepsilon$. On the same event of probability at least $1 - \delta$, we also have

$$\|\mathbf{v}_{T} - \mathbf{v}^{*} - J^{*}/(1 - \alpha_{T})\mathbf{1}\|_{\infty} \leq \|\mathbf{v}_{T} - \mathbf{v}_{\alpha_{T}}^{*}\|_{\infty} + \|\mathbf{v}_{\alpha_{T}}^{*} - \mathbf{v}^{*} - J^{*}/(1 - \alpha_{T})\mathbf{1}\|_{\infty}$$

$$\leq \frac{B + c_{5}(\log T)^{4}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^{2}}{T^{1/9}}.$$

We thus have $\|\mathbf{v}_T - \mathbf{v}^* - J^*/(1 - \alpha_T)\mathbf{1}\|_{\infty} \le \varepsilon$ as well. Still by Lemma 5.1, we also obtain the form of \mathbf{v}^* and \mathbf{q}^* , which by standard MDP theory (Puterman, 2014) satisfies the Bellman equation (1). Therefore, we complete the proof of Theorem 5.2.

C.3 Proof of Corollary 5.3

Proof of Corollary 5.3. Recall that $\pi_T(s) = \operatorname{argmax}_{a \in \mathcal{A}} q_T(s, a)$ for all $s \in \mathcal{S}$ is the greedy policy, and π^* is the optimal policy for average reward. Then under the conditions of Theorem 5.2, we have

$$\|q_T(s,a) - q^{\pi^*}(s,a) - J^*/(1-\alpha_T)\|_{\infty} \le \varepsilon,$$
 (63)

where $\varepsilon = \frac{B + c(\log T)^4(\log \frac{T|S||A|}{\delta})^2}{T^{1/9}}$, since the function q^* defined in Theorem 5.2 equals q^{π^*} . Assume without loss of generality that π^* and π_T are both deterministic policies. Invoking Lemma D.8 with $\pi' = \pi_T$ and $\pi = \pi^*$, we know that

$$J^{\pi_T} - J^{\pi^*} = \mathbb{E}_{s \sim d_{\pi'}} \left[\sum_{a \in \mathcal{A}} (\pi_T(a \mid s) - \pi^*(a \mid s)) q^*(s, a) \right]$$

$$\geq -\max_{s \in \mathcal{S}} \left| q^*(s, \pi_T(s)) - q^*(s, \pi^*(s)) \right|.$$

Here by the optimality of π^* and (63), we have

$$q^*(s, \pi^*(s)) \ge q^*(s, \pi_T(s))$$

$$\ge q_T(s, \pi_T(s)) - J^*/(1 - \alpha_T) - \varepsilon$$

$$\ge q_T(s, \pi^*(s)) - J^*/(1 - \alpha_T) - \varepsilon$$

$$\ge q^*(s, \pi^*(s)) - 2\varepsilon,$$

where the second and fourth lines uses the uniform bound (63), and the third line uses the property of greedy policy π_T . This leads to

$$J^{\pi_T} - J^{\pi^*} > -2\varepsilon$$

as desired. \Box

C.4 Proof of Proposition C.4

Proof of Proposition C.4. Putting Lemmas C.1, C.2 and C.3 together, we arrive at

$$\begin{aligned} \boldsymbol{\delta}_t &\leq \left(3t^{-2/9} + 2/t\right)\mathbf{1} + \sum_{i=1+\lfloor (1-\beta)t\rfloor}^t \theta_i^{(t)} \alpha_i \boldsymbol{P}^{\pi_{i-1}} \boldsymbol{\delta}_{i-1} \\ &+ \sqrt{\frac{c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^2}{T^{2/3}}} \left(\max_{1+\lfloor (1-\beta)t\rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) + 2(\log T)^6 \mathbf{1}\right) \end{aligned}$$

with probability at least $1 - \delta$ for any fixed t within $T/\log T \le t \le T$, as long as T satisfies

$$T/\log T \ge \max\{e^{2+c_1}, 100\}, \quad (T/\log T)^{1/3} \ge 4(2+c_1),$$

 $c_2 \log T \ge 5(2+c_1)(\log T)^{1/3} + \log\log T, \quad c_1 T^3 \ge 3(\log T)^5.$

Further simplifying this inequality leads to

$$\boldsymbol{\delta}_{t} \leq \sqrt{\frac{2c_{3}(\log T)^{2}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^{2}}{T^{2/3}} \left(\max_{1+\lfloor (1-\beta)t\rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) + 2T^{2/9}\mathbf{1}\right)} + \sum_{i=1+\lfloor (1-\beta)t\rfloor}^{t} \theta_{i}^{(t)} \alpha_{i} \boldsymbol{P}^{\pi_{i-1}} \boldsymbol{\delta}_{i-1}, \tag{64}$$

as long as T additionally satisfy $(\log T)^{50/9} \ge \frac{9}{2c_3}$. Once $(1-\beta) \ge 3/4$, for any fixed t obeying $3T/(2\log T) \le t \le T$, we apply (64) with a union bound over $\{k \colon \frac{2t}{3} \le k \le t\}$ — as a result, with probability at least $1-\delta$,

$$\boldsymbol{\delta}_{k} \leq \sqrt{\boldsymbol{\phi}_{t}} + \sum_{i=1+\lfloor (1-\beta)k \rfloor}^{k} \theta_{i}^{(k)} \alpha_{i} \boldsymbol{P}^{\pi_{i-1}} \boldsymbol{\delta}_{i-1}, \quad \text{for all } \frac{2t}{3} \leq k \leq t,$$
 (65)

where we define

$$\phi_t = \frac{2c_3(\log T)^2(\log \frac{T|S||A|}{\delta})^2}{T^{2/3}} \left(\max_{1+\lfloor (1-\beta)k \rfloor \le i \le k} \text{Var}_{\mathbf{P}}(\mathbf{v}_{i-1}) + 2T^{2/9} \mathbf{1} \right).$$
(66)

Following exactly the same arguments as in the proof of Proposition 4.4, we define

$$\lambda_i^{(t)} := \frac{\theta_{i+1}^{(t)}}{\sum_{j=\lfloor (1-\beta)t\rfloor}^{t-1} \theta_{j+1}^{(t)}}, \quad \lfloor (1-\beta)t \rfloor \le i \le t-1, \tag{67}$$

and arrive at the decomposition

$$\delta_{t} \leq \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \left(\lambda_{i_{1}}^{(t)}\sqrt{\phi_{t}} + \theta_{i_{1}+1}^{(t)}\alpha_{i_{1}+1}P^{\pi_{i_{1}}}\delta_{i_{1}}\right) \\
\leq \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \left[\lambda_{i_{1}}^{(t)}\sqrt{\phi_{t}} + \theta_{i_{1}+1}^{(t)}\alpha_{i_{1}+1}P^{\pi_{i_{1}}}\sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \left(\lambda_{i_{2}}^{(i_{1})}\sqrt{\phi_{t}} + \theta_{i_{2}+1}^{(i_{1})}\alpha_{i_{2}+1}P^{\pi_{i_{2}}}\delta_{i_{2}}\right)\right] \\
\leq \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \lambda_{i_{1}}^{(t)}\sqrt{\phi_{t}} + \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \lambda_{i_{1}}^{(t)}\lambda_{i_{2}}^{(i_{1})}\left(\alpha_{i_{1}+1}P^{\pi_{i_{1}}}\right)\sqrt{\phi_{t}} \\
+ \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \theta_{i_{1}+1}^{(t)}\theta_{i_{2}+1}^{(i_{1})} \prod_{k=1}^{2} (\alpha_{i_{k}+1}P^{\pi_{i_{k}}})\delta_{i_{2}}\right) \\
\leq \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \lambda_{i_{1}}^{(t)}\lambda_{i_{2}}^{(i_{1})}\left\{I + \alpha_{i_{1}+1}P^{\pi_{i_{1}}}\right\}\sqrt{\phi_{t}} \\
+ \sum_{i_{1}=\lfloor(1-\beta)t\rfloor}^{t-1} \sum_{i_{2}=\lfloor(1-\beta)i_{1}\rfloor}^{i_{1}-1} \theta_{i_{1}+1}^{(t)}\theta_{i_{2}+1}^{(i_{1})} \prod_{k=1}^{2} (\alpha_{i_{k}+1}P^{\pi_{i_{k}}})\delta_{i_{2}}. \tag{68}$$

To further apply such recursion, we define

$$G = T^{1/9} \log T$$
, and $\lambda_{\{i_k\}_{k=1}^G} := \lambda_{i_1}^{(t)} \lambda_{i_2}^{(i_1)} \dots \lambda_{i_G}^{(i_{G-1})} \ge 0$

for any $t > i_1 > i_2 > \cdots > i_H$, which satisfies

$$\alpha_{\{i_k\}_{k=1}^G} \ge \theta_{i_1+1}^{(t)} \theta_{i_2+1}^{(i_1)} \cdots \theta_{i_G+1}^{(i_{G-1})}.$$

We define the index set

$$\mathcal{I}_t = \{(i_1, \dots, i_G) : \lfloor (1 - \beta)t \rfloor \le i_1 \le t - 1, \ \lfloor (1 - \beta)i_{j-1} \rfloor \le i_j \le i_{j-1} - 1, \ \forall 1 \le j < G \},\$$

which satisfies

$$\sum_{(i_1,\dots,i_G)\in\mathcal{I}_t} \lambda_{\{i_k\}_{k=1}^G} = 1.$$

Note that once $T^{2/9} \log T \ge 6$ (implied by a stronger condition $T^{2/9} \ge (\log T)^5$ imposed before), we have

$$(1-\beta)^G = (1-c_2T^{-1/3}(\log T)^{-2})^{T^{1/9}\log T} \ge \exp\left(-2c_2T^{-1/3}(\log T)^{-2} \cdot T^{1/9}\log T\right) \ge \frac{2}{3},$$

which implies

$$i_1 > i_2 > \cdots > i_G \ge 2t/3$$
, for all $(i_1, \ldots, i_G) \in \mathcal{I}_t$.

Recursively invoking (68), we obtain

$$\delta_{t} \leq \sum_{(i_{1},\dots,i_{G})\in\mathcal{I}_{t}} \lambda_{\{i_{k}\}_{k=1}^{H}} \left\{ \left(\boldsymbol{I} + \sum_{h=1}^{G-1} \prod_{k=1}^{h} (\alpha_{i_{k}+1} \boldsymbol{P}^{\pi_{i_{k}}}) \right) \sqrt{\phi_{t}} + \prod_{h=1}^{G} (\alpha_{i_{h}+1} \boldsymbol{P}^{\pi_{i_{h}}}) |\delta_{i_{G}}| \right\}$$

$$\leq \max_{(i_{1},\dots,i_{G})\in\mathcal{I}_{t}} \left\{ \underbrace{\left(\boldsymbol{I} + \sum_{h=1}^{G-1} \prod_{k=1}^{G} (\alpha_{i_{k}+1} \boldsymbol{P}^{\pi_{i_{k}}}) \right) \sqrt{\phi_{t}}}_{=:\beta_{1}} + \underbrace{\prod_{h=1}^{G} (\alpha_{i_{h}+1} \boldsymbol{P}^{\pi_{i_{h}}}) |\delta_{i_{G}}|}_{=:\beta_{2}} \right\}. \tag{69}$$

We now treat the two terms β_1 and β_2 separately. The easier part is

$$\begin{split} \prod_{h=1}^{G} (\alpha_{i_h+1} \boldsymbol{P}^{\pi_{i_h}}) |\boldsymbol{\delta}_{i_G}| &\leq \prod_{h=1}^{G} \alpha_{i_h+1} \cdot \left\| \prod_{h=1}^{G} \boldsymbol{P}^{\pi_{i_h}} \right\|_{1} \cdot \|\boldsymbol{\delta}_{i_G}\|_{\infty} \mathbf{1} \\ &\leq \alpha_T^{G} \cdot \max_{2t/3 \leq j \leq t-1} \|\boldsymbol{\delta}_{j}\|_{\infty} \mathbf{1} \\ &\leq (1 - T^{-1/9})^{T^{1/9} \log T} \cdot (T^{1/9} + 3T^{1/3} (\log T)^3 / c_1) \leq (1 + 3/c_1) T^{-2/3} (\log T)^3, \end{split}$$

where the second inequality uses the fact that $\prod_{h=1}^{G} \boldsymbol{P}^{\pi_{i_h}}$ is a probability matrix, and the third inequality uses $\|\boldsymbol{\delta}_j\|_{\infty} \leq \|\mathbf{q}_j\|_{\infty} + \|\mathbf{q}_{\alpha_j}^*\|_{\infty} \leq 3j^{1/3}(\log j)^2/c_1 + j^{1/9}$ by Lemma E.2. On the other hand, following exactly the same argument as 26, we have

$$|\beta_{1}|^{2} \leq \sum_{h=0}^{G-1} \alpha_{T}^{h} \cdot \sum_{h=0}^{G-1} \alpha_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \phi_{t}$$

$$\leq \frac{1}{1 - \alpha_{T}} \sum_{h=0}^{G-1} \alpha_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \frac{2c_{3}(\log T)^{2}(\log \frac{T|S||A|}{\delta})^{2}}{T^{2/3}} \Big(\max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{v}_{i-1}) + 2T^{2/9} \mathbf{1} \Big)$$

$$\leq \frac{2c_{3}(\log T)^{2}(\log \frac{T|S||A|}{\delta})^{2}}{T^{2/3}} \Big(2T^{4/9} \mathbf{1} + \sum_{h=0}^{G-1} \alpha_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{v}_{i}) \Big). \tag{70}$$

Applying Lemma D.7, we know that

$$\max_{\lfloor t/2 \rfloor \le i < t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_i) \le \max_{\lfloor t/2 \rfloor \le i < t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{\alpha_i}^*) + 2(1 + 3/c_1)(\log T)^2 T^{1/3} \max_{\lfloor t/2 \rfloor \le i < t} \|\boldsymbol{\delta}_i\|_{\infty},$$

which, by the boundedness that $\|\operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{\alpha_i}^*)\|_{\infty} \leq \|\mathbf{v}_{\alpha_i}^*\|_{\infty}^2 \leq \frac{1}{(1-\alpha_i)^2} \leq T^{2/9}$, is further bounded as

$$\max_{\lfloor t/2 \rfloor \le i < t} \text{Var}_{\mathbf{P}}(\mathbf{v}_i) \le T^{2/9} + 2(1 + 3/c_1)(\log T)^2 T^{1/3} \max_{\lfloor t/2 \rfloor \le i < t} \|\boldsymbol{\delta}_i\|_{\infty}.$$

Plug back into (70), we have

$$|\beta_{1}|^{2} \leq \frac{2c_{3}(\log T)^{2}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^{2}}{T^{2/3}} \left(2T^{4/9}\mathbf{1} + \frac{T^{2/9}}{1 - \alpha_{T}}\mathbf{1} + \frac{2(1 + 3/c_{1})(\log T)^{2}T^{1/3}}{1 - \alpha_{T}} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\delta}_{i}\|_{\infty}\right) \\ \leq \frac{2c_{3}(\log T)^{2}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^{2}}{T^{2/3}} \left(3T^{4/9}\mathbf{1} + c_{4}(\log T)^{2}T^{4/9} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\delta}_{i}\|_{\infty}\right), \tag{71}$$

where we define the constant $c_4 = 2(1 + 3/c_1) > 0$. Therefore, combining with (69), we have

$$\delta_{t} \leq (1+3/c_{1})T^{-2/3}(\log T)^{3} + \sqrt{\frac{2c_{3}(\log T)^{2}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^{2}}{T^{2/3}}} \left(3T^{4/9}\mathbf{1} + c_{4}(\log T)^{2}T^{4/9} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\delta}_{i}\|_{\infty}\right) \\
\leq \sqrt{\frac{4c_{3}(\log T)^{2}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^{2}}{T^{2/3}}} \left(4T^{4/9}\mathbf{1} + c_{4}(\log T)^{2}T^{4/9} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\delta}_{i}\|_{\infty}\right)} \tag{72}$$

$$= \sqrt{\frac{4c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^2}{T^{2/9}}} \left(\mathbf{1} + c_4(\log T)^2 \max_{\lfloor t/2 \rfloor \le i < t} \|\boldsymbol{\delta}_i\|_{\infty} \right)^2$$

that holds with probability at least $1-\delta$ for any fixed t obeying $3T/(2\log T) \le t \le T$, as long as T satisfies $T^{10/9}(\log \frac{T|S||A|}{\delta})^2 \ge \frac{(1+3/c_1)^2}{2c_3}(\log T)^4$ to ensure (72). Taking a union bound over such t and noting that $\log \frac{T|S||A|}{\delta/T} \le 2\log \frac{T|S||A|}{\delta}$, we know that

$$\delta_t \le \sqrt{\frac{4c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}} \left(1 + c_4(\log T)^2 \max_{|t/2| \le i \le t} \|\delta_i\|_{\infty}\right)}$$

holds simultaneously for all $3T/(2 \log T) \le t \le T$ with probability at least $1 - \delta/2$. On the other hand, following exactly the same arguments, starting from the lower bound in (61) leads to

$$\delta_t \ge -\sqrt{\frac{4c_3(\log T)^2(\log \frac{T|\mathcal{S}||\mathcal{A}|}{2\delta})^2}{T^{2/9}} \left(1 + c_4(\log T)^2 \max_{\lfloor t/2 \rfloor \le i < t} \|\delta_i\|_{\infty}\right)}$$

simultaneously for all $3T/(2\log T) \le t \le T$ with probability at least $1 - \delta/2$. Finally, we summarize all the conditions as

$$T/\log T \ge \max\{e^{2+c_1}, 100\}, \quad (T/\log T)^{1/3} \ge 4(2+c_1),$$

$$c_2 \log T \ge 5(2+c_1)(\log T)^{1/3} + \log\log T, \quad c_1 T^3 \ge 3(\log T)^5,$$

$$(\log T)^{50/9} \ge \frac{9}{2c_3}, \quad T^{10/9}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^2 \ge \frac{(1+3/c_1)^2}{2c_3}(\log T)^4,$$

which can be simplified to

$$T/\log T \ge \max\{e^{2+c_1}, 100, 64(2+c_1)^3\}, \quad c_3(\log T)^5 \ge 5, \quad c_1T^3 \ge 3(\log T)^5, c_2\log T \ge 5(2+c_1)(\log T)^{1/3} + \log\log T, \quad 2c_3T \ge (1+3/c_1)^2(\log T)^2.$$

We thus complete the proof of Proposition C.4 by taking a union bound.

C.5 Proof of Lemma C.1

Proof of Lemma C.1. Firstly, Lemma D.3 implies that each term d_i can be bounded as

$$\|\boldsymbol{d}_{i}\|_{\infty} = \|(1-\theta_{i})(\mathbf{q}_{\alpha_{i-1}}^{*} - \mathbf{q}_{\alpha_{i}}^{*}) + \theta_{i}\boldsymbol{P}(\mathbf{q}_{\alpha_{i-1}}^{*} - \mathbf{q}_{\alpha_{i}}^{*})\|_{\infty}$$

$$\leq \|\mathbf{q}_{\alpha_{i-1}}^{*} - \mathbf{q}_{\alpha_{i}}^{*}\|_{\infty} \leq \frac{1}{1-\alpha_{i}} - \frac{1}{1-\alpha_{i-1}} = i^{1/9} - (i-1)^{1/9} \leq (i-1)^{-8/9}/9 \leq i^{-8/9},$$

where the second last inequality follows from the fact that $x^{1/9} - y^{1/9} = (x - y)z^{-8/9}/9$ for some z lying between x and y.

For any fixed $t \ge 1$, let $\lambda = t^{-1/3}$, then the switching error can be decomposed as

$$\sum_{i=1}^{t} \theta_i^{(t)} \boldsymbol{d}_i / \theta_i = \sum_{i=1}^{\lfloor (1-\lambda)t \rfloor} \theta_i^{(t)} \boldsymbol{d}_i / \theta_i + \sum_{\lfloor (1-\lambda)t \rfloor+1}^{t} \theta_i^{(t)} \boldsymbol{d}_i / \theta_i, \tag{73}$$

where the first summation can be bounded as

$$\sum_{i=1}^{\lfloor (1-\lambda)t\rfloor} \theta_i^{(t)} \boldsymbol{d}_i / \theta_i \le \sum_{i=1}^{\lfloor (1-\lambda)t\rfloor} \theta_i^{(t)} \|\boldsymbol{d}_i\|_{\infty} / \theta_i \le \sum_{i=1}^{\lfloor (1-\lambda)t\rfloor} \prod_{j=i+1}^{t} (1-\theta_j) \cdot i^{-8/9}, \tag{74}$$

where invoking (109) we know that once $(1 - \lambda)t = t - t^{2/3} \ge 50$,

$$\prod_{j=i+1}^{t} (1 - \theta_j) \le \exp\left(-\sum_{j=\lfloor (1-\lambda)t \rfloor + 1}^{t} \theta_j\right) \le \exp\left(-\lambda t \cdot c'(\log t)^2 t^{-2/3}\right) \le t^{-1}$$

as long as $\log t \ge 1/c' = 2 + c_1$. When $\lambda < 1/2$, the second summation is bounded as

$$\sum_{\lfloor (1-\lambda)t\rfloor+1}^t \theta_i^{(t)} \boldsymbol{d}_i/\theta_i \leq \sum_{\lfloor (1-\lambda)t\rfloor+1}^t \|\boldsymbol{d}_i\|_{\infty} \leq \sum_{\lfloor (1-\lambda)t\rfloor+1}^t i^{-8/9} \leq \lambda t \cdot (t/2)^{-8/9} \leq 2t^{-2/9}.$$

Putting them together, we arrive at

$$\left\| \sum_{i=1}^{t} \theta_i^{(t)} \boldsymbol{d}_i / \theta_i \right\|_{\infty} \le 3t^{-2/9}$$

as long as $\log t \ge 2 + c_1$ and $t - t^{2/3} > 50$ (which is satisfied when $t \ge 100$).

C.6 Proof of Lemma C.2

Proof of Lemma C.2. By the initialization, $\|\boldsymbol{\delta}_0\| \leq 20$. Once $t \geq 50$, by the monotonicity in (109),

$$\begin{aligned} \|\theta_0^{(t)} \boldsymbol{\delta}_0\|_{\infty} &\leq \prod_{j=1}^t (1 - \theta_j) \|\boldsymbol{\delta}_0\|_{\infty} \leq 20 \exp\left(-\sum_{j=1}^t \theta_j\right) \leq 20 \exp\left(-\sum_{j=50}^t \theta_j\right) \\ &\leq 20 \exp\left(-\sum_{j=50}^t (\log t)^2 t^{-2/3} / (2 + c_1)\right) \leq 20 \exp\left(-(\log t)^2 (t - 50) t^{-2/3} / (2 + c_1)\right) \leq t^{-1} \end{aligned}$$

as long as $t^{1/3} \log t \ge 4(2+c_1)$. By the boundedness of $\|\mathbf{v}_i\|_{\infty} \le 3(\log i)^2 i^{1/3}/c_1$, invoking (110) we know that when $|(1-\beta)t| \ge 50$,

$$\left\| \sum_{i=1}^{\lfloor (1-\beta)t \rfloor} \theta_i^{(t)} \theta_i^{(t)} \alpha_i (\mathbf{P}_i - \mathbf{P}) \mathbf{v}_{i-1} \right\|_{\infty} \leq \sum_{i=1}^{\lfloor (1-\beta)t \rfloor} \theta_i^{(t)} \| \mathbf{v}_{i-1} \|_{\infty}$$

$$\leq t \cdot \max_{0 \leq i \leq t} \| \mathbf{v}_i \|_{\infty} \cdot \exp\left(-\beta t \cdot \frac{(\log t)^2}{(2 + c_2)t^{2/3}} \right)$$

$$\leq 3(\log t)^2 t^{4/3} / c_1 \cdot \exp\left(-\beta t \cdot \frac{(\log t)^2}{(2 + c_2)t^{2/3}} \right)$$

Here for any t such that $T/\log T \le t \le T$, we have

$$\beta t \cdot \frac{(\log t)^2}{(2+c_1)t^{2/3}} = \frac{c_2 t^{1/3} (\log t)^2}{(2+c_1)T^{1/3}} \ge \frac{c_2 (\log t)^2}{(2+c_1)(\log T)^{1/3}} \ge 5 \log t,$$

as long as $c_2 \log T \geq 5(2+c_1)(\log T)^{1/3} + \log \log T$, which leads to

$$\left\| \sum_{i=1}^{\lfloor (1-\beta)t\rfloor} \theta_i^{(t)} \theta_i^{(t)} \alpha_i (\boldsymbol{P}_i - \boldsymbol{P}) \mathbf{v}_{i-1} \right\|_{\infty} \le t^{-1}$$

for any t such that $T/\log T \le t \le T$, if T satisfies $c_1T^3 \ge 3(\log T)^5$. We thus complete the proof of Lemma C.2.

C.7 Proof of Lemma C.3

Proof of Lemma C.3. We write

$$\boldsymbol{x}_i := \theta_i^{(t)} \alpha_i (\boldsymbol{P}_i - \boldsymbol{P}) \mathbf{v}_{i-1},$$

so that $\boldsymbol{\xi}_t = \sum_{i=1+\lfloor (1-\beta)t\rfloor}^t \boldsymbol{x}_i$, where $\{\boldsymbol{x}_i\}_{i=1+\lfloor (1-\beta)t\rfloor}^t$ is a martingale difference squence with the coarse deterministic bound

$$\|\boldsymbol{x}_i\|_{\infty} \le \theta_i^{(t)} \alpha_i \|\mathbf{v}_{i-1}\|_{\infty} \le \theta_i \|\mathbf{v}_{i-1}\|_{\infty} \le \frac{(\log i)^2}{c_1 i^{2/3}} \cdot 3(\log i)^2 i^{1/3} / c_1 \le 3c_1^{-2} (\log t)^4 t^{-1/3} =: R.$$

as long as $\beta < 1/2$. Furthermore, we define the sum of conditional variances as

$$\begin{aligned} \boldsymbol{W}_{t} &= \sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \operatorname{Var}(\boldsymbol{x}_{i} \mid \mathbf{v}_{1}, \dots, \mathbf{v}_{i-1}) \\ &= \sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \left(\theta_{i}^{(t)}\right)^{2} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) \\ &\leq \sum_{i=1+\lfloor (1-\beta)t \rfloor}^{t} \theta_{i}^{(t)} \cdot \left(\max_{1+\lfloor (1-\beta)t \rfloor \leq i \leq t} \theta_{i}^{(t)}\right) \cdot \left(\max_{1+\lfloor (1-\beta)t \rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1})\right) \\ &\leq \frac{(\log t)^{2}}{c_{1}(1-\beta)^{2/3}t^{2/3}} \cdot \max_{1+\lfloor (1-\beta)t \rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) \leq \frac{2(\log t)^{2}}{c_{1}t^{2/3}} \cdot \max_{1+\lfloor (1-\beta)t \rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) \end{aligned}$$

as long as $\beta < 1/2$, where the last line follows from the fact that $\sum_{i=1+\lfloor (1-\beta)t\rfloor}^t \theta_i^{(t)} \leq \sum_{i=0}^t \theta_i^{(t)} = 1$. Since $\|\mathbf{v}_i\|_{\infty} \leq 3(\log i)^2 i^{1/3}/c_1$, we have a coarse deterministic upper bound that

$$\|\boldsymbol{W}_t\|_{\infty} \le \frac{18(\log t)^2}{c_1^3 t^{2/3}} \cdot (\log t)^4 t^{2/3} = 18(\log t)^6 / c_1^3 =: \sigma^2.$$

We now choose the positive integer K such that

$$\frac{(\log t)^2}{c_1 t} \le \frac{\sigma^2}{2^K} \le \frac{2(\log t)^2}{c_1 t}, \quad K \le \frac{\log \frac{2c_1^2 t}{(\log t)^4}}{\log 2} \le 2\log \frac{2t}{(\log t)^2} \le t/2,$$

as long as $t \ge \max\{e^{c_1}, 100\}$. Applying the Freedman's inequality (c.f. Lemma D.6) with a union bound over all $|\mathcal{S}||\mathcal{A}|$ entries, we know that with probability at least $1 - \delta$, for any fixed t such that $T/\log T \le t \le T$,

$$\begin{split} |\boldsymbol{\xi}_{t}| &\leq \sqrt{8 \max\left\{\boldsymbol{W}_{t}, \frac{\sigma^{2}}{2^{K}}\right\} \log \frac{2K|\mathcal{S}||\mathcal{A}|}{\delta}} + \frac{4R}{3} \log \frac{2K|\mathcal{S}||\mathcal{A}|}{\delta} \\ &\leq \sqrt{8 \left(\frac{2(\log t)^{2}}{c_{1}t^{2/3}} \cdot \max_{1+\lfloor (1-\beta)t\rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) + \frac{\sigma^{2}}{2^{K}}\right) \log \frac{2K|\mathcal{S}||\mathcal{A}|}{\delta}} + \frac{4R}{3} \log \frac{2K|\mathcal{S}||\mathcal{A}|}{\delta} \\ &\leq \sqrt{\frac{16(\log t)^{2} \log(t|\mathcal{S}||\mathcal{A}|/\delta)}{c_{1}t^{2/3}} \left(\max_{1+\lfloor (1-\beta)t\rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) + \frac{1}{t^{1/3}}\right)} + \frac{12(\log t)^{4}}{c_{1}^{2}t^{1/3}} \log \frac{t|\mathcal{S}||\mathcal{A}|}{\delta}}{\delta} \\ &\leq \sqrt{\frac{c_{3}(\log T)^{2}(\log \frac{T|\mathcal{S}||\mathcal{A}|}{\delta})^{2}}{T^{2/3}} \left(\max_{1+\lfloor (1-\beta)t\rfloor \leq i \leq t} \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i-1}) + 2(\log T)^{6}}\right)} \end{split}$$

where $c_3 = 32/c_1 + 144/c_1^4$. This completes the proof of Lemma C.3.

D Supporting lemmas

D.1 Relations of value functions

We quote the following lemma adapted from Jin and Sidford (2021, Lemma 2). It relates the average reward J^{π} and the rescaled discounted ones Q^{π}_{γ} , V^{π}_{γ} of any policy π . The proof the lemma is similar to that of Jin and Sidford (2021) and is omitted here.

Lemma D.1 (Lemma 2 of Jin and Sidford (2021)). For any policy $\pi: \mathcal{S} \to \Delta(\mathcal{A})$ and any discount factor $\gamma \in (0,1]$, it holds that $\|J^{\pi}\mathbf{1} - \mathbf{Q}_{\gamma}^{\pi}\|_{\infty} \leq 3(1-\gamma)t_{mix}$ and $\|J^{\pi}\mathbf{1} - \mathbf{V}_{\gamma}^{\pi}\|_{\infty} \leq (1-\gamma)t_{mix}$.

The following lemma is adapted from Dong et al. (2021) and De Farias and Van Roy (2006), showing that our framework is applicable when the finite mixing time condition is replaced by the reward averaging time τ (see Remark 2.2). The proof the lemma is similar to that of De Farias and Van Roy (2006) and is omitted here.

Lemma D.2 (Lemma 2 of Dong et al. (2021), Theorem 4.1 of De Farias and Van Roy (2006)). For any policy $\pi: \mathcal{S} \to \Delta(\mathcal{A})$ and any discount factor $\gamma \in (0,1]$, it holds that $\|J^{\pi}\mathbf{1} - \mathbf{Q}_{\gamma}^{\pi}\|_{\infty} \leq (1-\gamma)\tau_{\pi}$ and $\|J^{\pi}\mathbf{1} - \mathbf{V}_{\gamma}^{\pi}\|_{\infty} \leq (1-\gamma)\tau_{\pi}$

The following lemma bounds the difference between rescaled optimal value and Q-functions with two discount factors.

Lemma D.3. For any γ_1, γ_2 such that $0 < \gamma_1 < \gamma_2 < 1$, it holds that

$$\mathbf{0} \le \frac{\mathbf{Q}_{\gamma_2}^*}{1 - \gamma_2} - \frac{\mathbf{Q}_{\gamma_1}^*}{1 - \gamma_1} \le \left(\frac{1}{1 - \gamma_2} - \frac{1}{1 - \gamma_1}\right) \mathbf{1};\tag{75a}$$

$$\mathbf{0} \le \frac{V_{\gamma_2}^*}{1 - \gamma_2} - \frac{V_{\gamma_1}^*}{1 - \gamma_1} \le \left(\frac{1}{1 - \gamma_2} - \frac{1}{1 - \gamma_1}\right) \mathbf{1},\tag{75b}$$

where 0 (resp. 1) is a vector or matrix with all entries equal to 0 (resp. 1).

Proof of Lemma D.3. We first show (75a). The upper bound follows directly from Lemma 17 of Dong et al. (2021) with a rescaling. For any policy π , we note that

$$\frac{\boldsymbol{Q}_{\gamma_2}^{\pi}}{1 - \gamma_2} = \mathbb{E}_{\pi} \left[\sum_{k=1}^{\infty} \gamma_2^{t-1} r(s_k, a_k) \, \middle| \, s_1 = s \right] \ge \mathbb{E}_{\pi} \left[\sum_{k=1}^{\infty} \gamma_1^{t-1} r(s_k, a_k) \, \middle| \, s_1 = s \right] = \frac{\boldsymbol{Q}_{\gamma_1}^{\pi}}{1 - \gamma_1},$$

which follows from $\gamma_2 > \gamma_1$ and the non-negativeness of the reward function r. Applying the above relation to $\pi_{\gamma_1}^*$, we obtain

$$\frac{\bm{Q}_{\gamma_2}^*}{1-\gamma_2} \geq \frac{\bm{Q}_{\gamma_2}^{\pi_{\gamma_1}^*}}{1-\gamma_2} \geq \frac{\bm{Q}_{\gamma_1}^{\pi_{\gamma_1}^*}}{1-\gamma_1} = \frac{\bm{Q}_{\gamma_1}^*}{1-\gamma_1},$$

which completes the proof of lower bound in (75a). Further noting $V_{\gamma}^{*}(s) = \max_{a' \in \mathcal{A}} Q_{\gamma}^{*}(s, a')$ for any $s \in \mathcal{S}$, the fact that maximum is a contraction map leads to (75b) and completes the proof of Lemma D.3.

Lemma D.4. For any discount factors γ_1, γ_2 such that $0 < \gamma_1 < \gamma_2 < 1$, it holds that

$$\|\boldsymbol{Q}_{\gamma_1}^* - \boldsymbol{Q}_{\gamma_2}^*\|_{\infty} \le \frac{\gamma_2 - \gamma_1}{1 - \gamma_2} \quad and \quad \|\boldsymbol{V}_{\gamma_1}^* - \boldsymbol{V}_{\gamma_2}^*\|_{\infty} \le \frac{\gamma_2 - \gamma_1}{1 - \gamma_2}.$$
 (76)

Proof of Lemma D.4. The boundedness of reward function $\mathbf{0} \leq \mathbf{r} \leq \mathbf{1}$ implies $\mathbf{0} \leq \mathbf{Q}_{\gamma_i}^* \leq \mathbf{1}$ for i = 1, 2. For any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, applying the lower bound in Lemma D.3 yields

$$Q_{\gamma_2}^*(s,a) \ge \frac{1-\gamma_2}{1-\gamma_1} Q_{\gamma_1}^*(s,a) = Q_{\gamma_1}^*(s,a) + \frac{\gamma_1-\gamma_2}{1-\gamma_1} Q_{\gamma_1}^*(s,a) \ge Q_{\gamma_1}^*(s,a) - \frac{\gamma_2-\gamma_1}{1-\gamma_1} Q_{\gamma_2}^*(s,a) = Q_{\gamma_1}^*(s,a) - Q_{\gamma_2}^*(s,a) = Q_{\gamma_1}^*(s,a) - Q_{\gamma_2}^*(s,a) = Q_{\gamma_2}^*(s,a) - Q_{\gamma_2}^*(s,a) = Q_{\gamma_2}^*(s,a) - Q_{$$

where the second inequality follows from $Q_{\gamma_1}^*(s,a) \leq 1$. On the other hand, the upper bound in Lemma D.3 together with $Q_{\gamma_1}^*(s,a) \geq 0$ implies

$$Q_{\gamma_2}^*(s,a) \le \frac{1-\gamma_2}{1-\gamma_1} Q_{\gamma_1}^*(s,a) + \frac{\gamma_2-\gamma_1}{1-\gamma_1} \le Q_{\gamma_1}^*(s,a) + \frac{\gamma_2-\gamma_1}{1-\gamma_1},$$

which leads to $\|Q_{\gamma_1}^* - Q_{\gamma_2}^*\|_{\infty} \le \frac{\gamma_2 - \gamma_1}{1 - \gamma_2}$. The second inequality in (76) follows from the same arguments.

Lemma D.5. For any $i \geq 1$, it holds that

$$\operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_i) - \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_i}^*) \le 4\|\boldsymbol{\Delta}_i\|_{\infty}. \tag{77}$$

Proof of Lemma D.5. We start by bounding $V_i - V_{\gamma_i}^*$ for all $i \ge 1$. Recalling that $\Delta_i = Q_i - Q_{\gamma_i}^*$, we have

$$\boldsymbol{V}_{i} - \boldsymbol{V}_{\gamma_{i}}^{*} = \boldsymbol{P}^{\pi_{i}} \boldsymbol{Q}_{i} - \boldsymbol{P}^{\pi_{\gamma_{i}}^{*}} \boldsymbol{Q}_{\gamma_{i}}^{*} \overset{\text{(i)}}{\leq} \boldsymbol{P}^{\pi_{i}} \boldsymbol{Q}_{i} - \boldsymbol{P}^{\pi_{i}} \boldsymbol{Q}_{\gamma_{i}}^{*} \leq \|\boldsymbol{P}^{\pi_{i}}\|_{1} \|\boldsymbol{Q}_{i} - \boldsymbol{Q}_{\gamma_{i}}^{*}\|_{\infty} \overset{\text{(ii)}}{\leq} \|\boldsymbol{\Delta}_{i}\|_{\infty},$$

where (i) follows from the optimality of $\pi_{\gamma_i}^*$ with respect to $Q_{\gamma_i}^*$, and (ii) holds since P^{π_i} is a probability transition matrix. Similarly,

$$V_i - V_{\gamma_i}^* \geq P^{\pi_{\gamma_i}^*} Q_i - P^{\pi_{\gamma_i}^*} Q_{\gamma_i}^* \geq - \|P^{\pi_{\gamma_i}^*}\|_1 \|Q_i - Q_{\gamma_i}^*\|_\infty \geq - \|\Delta_i\|_\infty,$$

which leads to

$$\|V_i - V_{\gamma_i}^*\|_{\infty} \le \|\Delta_i\|_{\infty} \tag{78}$$

for all $i \geq 1$. By the definition of entrywise variance, we have

$$\operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{i}) - \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_{i}}^{*}) = \left(\boldsymbol{P}(\boldsymbol{V}_{i} \circ \boldsymbol{V}_{i}) - (\boldsymbol{P}\boldsymbol{V}_{i}) \circ (\boldsymbol{P}\boldsymbol{V}_{i})\right) - \left(\boldsymbol{P}(\boldsymbol{V}_{\gamma_{i}}^{*} \circ \boldsymbol{V}_{\gamma_{i}}^{*}) - (\boldsymbol{P}\boldsymbol{V}_{\gamma_{i}}^{*}) \circ (\boldsymbol{P}\boldsymbol{V}_{\gamma_{i}}^{*})\right)$$

$$= \boldsymbol{P}(\boldsymbol{V}_{i} \circ \boldsymbol{V}_{i} - \boldsymbol{V}_{\gamma_{i}}^{*} \circ \boldsymbol{V}_{\gamma_{i}}^{*}) + (\boldsymbol{P}\boldsymbol{V}_{\gamma_{i}}^{*}) \circ (\boldsymbol{P}\boldsymbol{V}_{\gamma_{i}}^{*}) - (\boldsymbol{P}\boldsymbol{V}_{i}) \circ (\boldsymbol{P}\boldsymbol{V}_{i})\right)$$

$$= \boldsymbol{P}((\boldsymbol{V}_{i} - \boldsymbol{V}_{\gamma_{i}}^{*}) \circ (\boldsymbol{V}_{i} + \boldsymbol{V}_{\gamma_{i}}^{*})) - (\boldsymbol{P}\boldsymbol{V}_{i} - \boldsymbol{P}\boldsymbol{V}_{\gamma_{i}}^{*}) \circ (\boldsymbol{P}\boldsymbol{V}_{i} + \boldsymbol{P}\boldsymbol{V}_{\gamma_{i}}^{*}). \tag{79}$$

Furthermore, we have

$$\begin{aligned} \left\| \boldsymbol{P} \big((\boldsymbol{V}_i - \boldsymbol{V}_{\gamma_i}^*) \circ (\boldsymbol{V}_i + \boldsymbol{V}_{\gamma_i}^*) \big) \right\|_{\infty} &\leq \| \boldsymbol{P} \|_1 \left\| (\boldsymbol{V}_i - \boldsymbol{V}_{\gamma_i}^*) \circ (\boldsymbol{V}_i + \boldsymbol{V}_{\gamma_i}^*) \right\|_{\infty} \\ &\leq \| \boldsymbol{P} \|_1 \left\| \boldsymbol{V}_i - \boldsymbol{V}_{\gamma_i}^* \right\|_{\infty} \| \boldsymbol{V}_i + \boldsymbol{V}_{\gamma_i}^* \right\|_{\infty} \leq 2 \| \boldsymbol{\Delta}_i \|_{\infty}, \end{aligned}$$

where the last inequality follows from $\|P\|_1 = 1$ and $\|V_i + V_{\gamma_i}^*\|_{\infty} \le \|V_i\|_{\infty} + \|V_{\gamma_i}^*\|_{\infty} \le 2$ by the bounded magnitude, and (78). Meanwhile, similar arguments yield

$$\begin{aligned} \left\| (PV_i - PV_{\gamma_i}^*) \circ (PV_i + PV_{\gamma_i}^*) \right\|_{\infty} &\leq \|PV_i - PV_{\gamma_i}^*\|_{\infty} \|PV_i + PV_{\gamma_i}^*\|_{\infty} \\ &\leq \|P\|_1 \|V_i - V_{\gamma_i}^*\|_{\infty} \|P\|_1 \|V_i + V_{\gamma_i}^*\|_{\infty} \leq 2 \|\Delta_i\|_{\infty}. \end{aligned}$$

Plugging the two bounds back in (79) completes the proof of (77).

D.2 Proof of Lemma 5.1

Proof of Lemma 5.1. By Feinberg and Shwartz (2012, Theorem 8.1), there exists an optimal policy π^* (namely, a Blackwell optimal policy) for the average reward and some $\alpha^* \in (0,1)$, such that for all $\alpha \in (0,\alpha^*)$, π^* is also the optimal policy for α -discounted reward. Since π^* is the optimal policy for average reward, its long-term average reward equals J^* . We define $\mathbf{v}^* \in \mathbb{R}^{|\mathcal{S}|}$ by

$$v^*(s) = \mathbb{E}_{\pi^*} \left[\sum_{k=0}^{\infty} (r_k - J^*) \mid s_0 = s \right], \quad \forall s \in \mathcal{S},$$

namely, the bias function of π^* . We also define $\mathbf{P}^* \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ as the transition matrix induced by π^* , and $\mathbf{P}^{**} = \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} (\mathbf{P}^*)^{t-1}$ as the long-term average transition matrix. Furthermore, let

$$H_P = (\boldsymbol{I} - \boldsymbol{P}^* + \boldsymbol{P}^{**})^{-1}(\boldsymbol{I} - \boldsymbol{P}^*) \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$$

By Puterman (2014, Theorem 8.2.3), letting $\rho = (1 - \alpha)/\alpha$, for α that is sufficiently close to 1, or equivalently for ρ that is sufficiently close to 0, we have the Laurent series expansion

$$\mathbf{v}_{\alpha}^{\pi^*} = (1+\rho) \left[\rho^{-1} J^* \mathbf{1} + \mathbf{v}^* + \sum_{n=1}^{\infty} \rho^n \mathbf{y}_n \right],$$

where J^* is the long-term average reward of π^* , $\mathbf{y}_n = (-1)^n H_P^{n+1} \mathbf{f}$, where $\mathbf{f} \in \mathbb{R}^{|\mathcal{S}|}$ is the vectorization of $f(s) = r(s, \pi^*(s))$ that also satisfies $\mathbf{v}^* = H_P \mathbf{f}$. Thus, recalling the Blackwell optimality of π^* and taking α sufficiently close to 1, we have

$$\mathbf{v}_{\alpha}^* = \mathbf{v}_{\alpha}^{\pi^*} = (1+\rho) \left[\rho^{-1} J^* \mathbf{1} + \mathbf{v}^* + \sum_{n=1}^{\infty} \rho^n \mathbf{y}_n \right].$$

It is clear to see $\lim_{\rho\to 0}\sum_{n=0}^{\infty}\rho^n\|\mathbf{y}_{n+1}\|_{\infty}<\infty$ from the definition of \mathbf{y}_n . Thus, since $\rho=(1-\alpha)/\alpha$,

$$\sum_{n=1}^{\infty} \rho^n \mathbf{y}_n = \alpha \frac{\mathbf{v}_{\alpha}^* - \frac{J^*}{1-\alpha} \mathbf{1} - \frac{\mathbf{v}^*}{\alpha}}{1-\alpha}$$
$$= \alpha \frac{\mathbf{v}_{\alpha}^* - \frac{J^*}{1-\alpha} \mathbf{1} - \mathbf{v}^*}{1-\alpha} - \mathbf{v}^*,$$

which means

$$\lim_{\alpha \to 1} \frac{\left\| \mathbf{v}_{\alpha}^* - \frac{J^*}{1 - \alpha} \mathbf{1} - \mathbf{v}^* \right\|_{\infty}}{1 - \alpha} < \infty.$$

Therefore, there exists a constant $B_1 > 0$ which only depends on the underlying MDP that

$$\sup_{\alpha \in (0,1)} \frac{\left\| \mathbf{v}_{\alpha}^* - \frac{J^*}{1-\alpha} \mathbf{1} - \mathbf{v}^* \right\|_{\infty}}{1-\alpha} \le B_1.$$
(80)

We now define the function

$$q^*(s, a) = r(s, a) - J^* + \mathbb{E}_{s' \sim P(\cdot \mid s, a)} [v^*(s')]$$

for all state-action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$, so that (v^*, q^*) satisfies the Bellman equation (1) by standard MDP theory (Puterman, 2014). The equivalent vectorized representation is

$$\mathbf{q}^* = \mathbf{r} - J^* \mathbf{1} + \mathbf{P}^{\pi^*} \mathbf{v}^*.$$

Then the Bellman equation for optimal discounted value functions implies

$$\mathbf{q}_{\alpha}^{*} = \mathbf{r} + \alpha \mathbf{P}^{\pi_{\alpha}^{*}} \mathbf{v}_{\alpha}^{*}$$

$$\geq \mathbf{r} + \alpha \mathbf{P}^{\pi^{*}} \mathbf{v}_{\alpha}^{*}$$

$$\geq \mathbf{r} + \alpha \mathbf{P}^{\pi^{*}} \left(\frac{J^{*}}{1-\alpha} \mathbf{1} + \mathbf{v}^{*} - (1-\alpha)B_{1} \mathbf{1} \right)$$

$$\geq \mathbf{r} + \mathbf{P}^{\pi^{*}} \mathbf{v}^{*} + \left(\frac{\alpha}{1-\alpha} J^{*} - (1-\alpha) \|\mathbf{v}^{*}\|_{\infty} - \alpha(1-\alpha)B_{1} \right) \mathbf{1}$$

$$= \mathbf{q}^{*} + J^{*} \mathbf{1} + \left(\frac{\alpha}{1-\alpha} J^{*} - (1-\alpha) \|\mathbf{v}^{*}\|_{\infty} - \alpha(1-\alpha)B_{1} \right) \mathbf{1}$$

$$= \mathbf{q}^{*} + \left(\frac{J^{*}}{1-\alpha} - (1-\alpha) \|\mathbf{v}^{*}\|_{\infty} - \alpha(1-\alpha)B_{1} \right) \mathbf{1},$$

where the second line follows from the optimality of π_{α}^* with respect to \mathbf{v}_{α}^* , the third line follows from (80), and the fifth line follows from the Bellman equation (1) for (q^*, v^*) . Similarly, we have

$$\mathbf{q}_{\alpha}^{*} = \mathbf{r} + \alpha \mathbf{P}^{\pi_{\alpha}^{*}} \mathbf{v}_{\alpha}^{*}$$

$$\leq \mathbf{r} + \alpha \mathbf{P}^{\pi_{\alpha}^{*}} \left(\frac{J^{*}}{1-\alpha} \mathbf{1} + \mathbf{v}^{*} + (1-\alpha)B_{1} \mathbf{1} \right)$$

$$\leq \mathbf{r} + \alpha \mathbf{P}^{\pi^{*}} \mathbf{v}^{*} + \left(\frac{\alpha}{1-\alpha} J^{*} + \alpha(1-\alpha)B_{1} \right) \mathbf{1}$$

$$\leq \mathbf{q}^{*} + J^{*} \mathbf{1} + \left(\frac{\alpha}{1-\alpha} J^{*} + (1-\alpha) \|\mathbf{v}^{*}\|_{\infty} + \alpha(1-\alpha)B_{1} \right) \mathbf{1}$$

$$= \mathbf{q}^{*} + \left(\frac{J^{*}}{1-\alpha} - (1-\alpha) \|\mathbf{v}^{*}\|_{\infty} - \alpha(1-\alpha)B_{1} \right) \mathbf{1}.$$

Combining the above two inequalities, we know that there exists some constant $B_2 > 0$ such that

$$\sup_{\alpha \in (0,1)} \frac{\left\| \mathbf{q}_{\alpha}^* - \frac{J^*}{1-\alpha} \mathbf{1} - \mathbf{q}^* \right\|_{\infty}}{1-\alpha} \le B_2.$$

Taking $B = \max\{B_1, B_2\}$ completes the proof of Lemma 5.1.

D.3 Freedman's inequality

The following lemma is adapted by Li et al. (2021) from Freedman's inequality Freedman (1975).

Lemma D.6 (Theorem 5 of Li et al. (2021)). Suppose $Y_n = \sum_{k=1}^n X_k \in \mathbb{R}$, where $\{X_k\}$ is a real-valued scalar sequence obeying

$$|X_k| \le R$$
 and $\mathbb{E}[X_k \mid \{X_j\}_{j < k}] = 0, \quad \forall \ k \ge 1.$

Define

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

where we write \mathbb{E}_{k-1} for the expectation conditional on $\{X_j\}_{j < k}$. Then for any given $\sigma^2 \geq 0$, one has

$$\mathbb{P}(|Y_n| \ge y \text{ and } W_n \le \sigma^2) \le 2 \exp\left(-\frac{y^2/2}{\sigma^2 + Ry/3}\right).$$

In addition, suppose $W_n \leq \sigma^2$ holds deterministically. For any positive integer K, with probability at least $1 - \delta$ one has

$$|Y_n| \le \sqrt{8 \max\left\{W_n, \frac{\sigma^2}{2^K}\right\} \log \frac{2K}{\delta}} + \frac{4R}{3} \log \frac{2K}{\delta}. \tag{81}$$

D.4 Proof of Lemma A.1

Proof of Lemma A.1. Employing (77) of Lemma D.5, we have the entrywise bound

$$\sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{i})$$

$$\leq \sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \left\{ \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{\gamma_{i}}^{*}) + \|\mathbf{\Delta}_{i}\|_{\infty} \mathbf{1} \right\}$$

$$\leq \sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{\gamma_{i}}^{*}) + 4 \sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \|\mathbf{\Delta}_{i}\|_{\infty} \mathbf{1}$$

$$= \sum_{h=0}^{H-1} \gamma_{T}^{h} \prod_{k=1}^{h} \mathbf{P}^{\pi_{i_{k}}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\mathbf{P}}(\mathbf{V}_{\gamma_{i}}^{*}) + \frac{4}{1 - \gamma_{T}} \max_{\lfloor t/2 \rfloor \leq i < t} \|\mathbf{\Delta}_{i}\|_{\infty} \mathbf{1}, \tag{82}$$

where the third line follows from $\max_i(\boldsymbol{a}_i + \boldsymbol{b}_i) \leq \max_i \boldsymbol{a}_i + \max_i \boldsymbol{b}_i$, and the last line follows from the fact that $\prod_{k=1}^h \boldsymbol{P}^{\pi_{i_k}} \mathbf{1} = \mathbf{1}$ since the product is still a probability transition matrix. Following the same arguments as the way we bound (79), it holds for all $i \geq 1$ that

$$\|\operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_i}^*) - \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_t}^*)\|_{\infty} \le 4\|\boldsymbol{Q}_{\gamma_i}^* - \boldsymbol{Q}_{\gamma_t}^*\|_{\infty}.$$
 (83)

On the other hand,

$$\operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_{t}}^{*}) = \boldsymbol{P}(\boldsymbol{V}_{\gamma_{t}}^{*} \circ \boldsymbol{V}_{\gamma_{t}}^{*}) - (\boldsymbol{P}\boldsymbol{V}_{\gamma_{t}}^{*}) \circ (\boldsymbol{P}\boldsymbol{V}_{\gamma_{t}}^{*})$$

$$= \boldsymbol{P}^{\pi_{i_{h+1}}}(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}) + \boldsymbol{P}(\boldsymbol{V}_{\gamma_{t}}^{*} \circ \boldsymbol{V}_{\gamma_{t}}^{*}) - \boldsymbol{P}^{\pi_{i_{h+1}}}(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}) - (\boldsymbol{P}\boldsymbol{V}_{\gamma_{t}}^{*}) \circ (\boldsymbol{P}\boldsymbol{V}_{\gamma_{t}}^{*})$$

$$= \boldsymbol{P}^{\pi_{i_{h+1}}}(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}) + (\boldsymbol{P}^{\pi_{\gamma_{t}}^{*}}(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}) - \boldsymbol{P}^{\pi_{i_{h+1}}}(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}))$$

$$-\frac{1}{\gamma_{t}^{2}}(\boldsymbol{Q}_{\gamma_{t}}^{*} - (1 - \gamma_{t})\boldsymbol{r}) \circ (\boldsymbol{Q}_{\gamma_{t}}^{*} - (1 - \gamma_{t})\boldsymbol{r}), \tag{84}$$

where the last equality follows from the Bellman equation $\mathbf{Q}_{\gamma}^* = (1-\gamma)\mathbf{r} + \gamma \mathbf{P}\mathbf{V}_{\gamma}^*$ for all $\gamma \in (0,1)$. Using similar arguments as Li et al. (2021), the second term in (84) can be bounded as

$$\left\| \boldsymbol{P}^{\pi_{\gamma_t}^*} (\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{Q}_{\gamma_t}^*) - \boldsymbol{P}^{\pi_{i_{h+1}}} (\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{Q}_{\gamma_t}^*) \right\|_{\infty}$$
(85)

$$= \left\| P \Pi^{\pi_{i_{h+1}}} (Q_{\gamma_t}^* \circ Q_{\gamma_t}^*) - P \Pi^{\pi_{\gamma_t}^*} (Q_{\gamma_t}^* \circ Q_{\gamma_t}^*) \right\|_{\infty}$$

$$(86)$$

$$\leq \|P\|_{1} \|\Pi^{\pi_{i_{h+1}}}(Q_{\gamma_{t}}^{*} \circ Q_{\gamma_{t}}^{*}) - P\Pi^{\pi_{\gamma_{t}}^{*}}(Q_{\gamma_{t}}^{*} \circ Q_{\gamma_{t}}^{*})\|_{\infty}$$
(87)

$$\stackrel{\text{(i)}}{\leq} \left\| (\boldsymbol{\Pi}^{\pi_{i_{h+1}}} \boldsymbol{Q}_{\gamma_{t}}^{*} - \boldsymbol{\Pi}^{\pi_{\gamma_{t}}^{*}} \boldsymbol{Q}_{\gamma_{t}}^{*}) \circ (\boldsymbol{\Pi}^{\pi_{i_{h+1}}} \boldsymbol{Q}_{\gamma_{t}}^{*} + \boldsymbol{\Pi}^{\pi_{\gamma_{t}}^{*}} \boldsymbol{Q}_{\gamma_{t}}^{*}) \right\|_{\infty}$$

$$(88)$$

$$\leq \|\boldsymbol{\Pi}^{\pi_{i_{h+1}}} \boldsymbol{Q}_{\gamma_{t}}^{*} - \boldsymbol{V}_{\gamma_{t}}^{*}\|_{\infty} \|\boldsymbol{\Pi}^{\pi_{i_{h+1}}} \boldsymbol{Q}_{\gamma_{t}}^{*} + \boldsymbol{\Pi}^{\pi_{\gamma_{t}}^{*}} \boldsymbol{Q}_{\gamma_{t}}^{*}\|_{\infty}$$

$$(89)$$

$$\stackrel{\text{(ii)}}{\leq} 2 \| \mathbf{\Pi}^{\pi_{i_{h+1}}} \mathbf{Q}_{\gamma_{t}}^* - \mathbf{\Pi}^{\pi_{i_{h+1}}} \mathbf{Q}_{i_{h+1}} \|_{\infty} + 2 \| \mathbf{\Pi}^{\pi_{i_{h+1}}} \mathbf{Q}_{i_{h+1}} - \mathbf{V}_{\gamma_{t}}^* \|_{\infty}$$

$$(90)$$

$$\stackrel{\text{(iii)}}{\leq} 2 \| \boldsymbol{Q}_{\gamma_{t}}^{*} - \boldsymbol{Q}_{i_{h+1}}^{*} \|_{\infty} + 2 \| \boldsymbol{V}_{i_{h+1}} - \boldsymbol{V}_{\gamma_{t}}^{*} \|_{\infty}$$
(91)

$$\stackrel{\text{(iv)}}{\leq} 4 \max_{\lfloor t/2 \rfloor \leq i < t} \left\| \boldsymbol{Q}_{\gamma_t}^* - \boldsymbol{Q}_i \right\|_{\infty} \tag{92}$$

$$\stackrel{\text{(v)}}{\leq} 4 \max_{\lfloor t/2 \rfloor \leq i < t} \left\| \boldsymbol{Q}_{\gamma_t}^* - \boldsymbol{Q}_{\gamma_i}^* \right\|_{\infty} + 4 \max_{\lfloor t/2 \rfloor \leq i < t} \| \boldsymbol{\Delta}_i \|_{\infty}, \tag{93}$$

where (i) follows from the fact that $\|P\|_1 = 1$, (ii) is due to the boundedness of $\|Q_{\gamma_t}^*\| \le 1$, (iii) follows from $\|\Pi^{\pi_{i_{h+1}}}\|_1 = 1$ and $V_{i_{h+1}} = \Pi^{\pi_{i_{h+1}}}Q_{i_{h+1}}$, (iv) follows from $\|V_{i_{h+1}} - V_{\gamma_t}^*\|_{\infty} \le \|Q_{\gamma_t}^* - Q_{i_{h+1}}^*\|_{\infty}$, and (v) follows from the property of entrywise maximum. We thus obtain the entrywise upper bound (uniform over all $\lfloor t/2 \rfloor \le i < t$)

$$\max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_{i}}^{*}) \leq \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_{t}}^{*}) + 4 \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{Q}_{\gamma_{i}}^{*} - \boldsymbol{Q}_{\gamma_{t}}^{*}\|_{\infty} \mathbf{1}$$

$$\leq \boldsymbol{P}^{\pi_{i_{h+1}}}(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}) + 8 \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{Q}_{\gamma_{t}}^{*} - \boldsymbol{Q}_{\gamma_{i}}^{*}\|_{\infty} \mathbf{1} + 4 \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\Delta}_{i}\|_{\infty} \mathbf{1}$$

$$- \frac{1}{\gamma_{t}^{2}}(\boldsymbol{Q}_{\gamma_{t}}^{*} - (1 - \gamma_{t})\boldsymbol{r}) \circ (\boldsymbol{Q}_{\gamma_{t}}^{*} - (1 - \gamma_{t})\boldsymbol{r})$$

$$\leq \frac{1}{\gamma_{t}}(\gamma_{t}\boldsymbol{P}^{\pi_{i_{h+1}}} - \boldsymbol{I})(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}) + 8 \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{Q}_{\gamma_{t}}^{*} - \boldsymbol{Q}_{\gamma_{i}}^{*}\|_{\infty} \mathbf{1} + 4 \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\Delta}_{i}\|_{\infty} \mathbf{1}$$

$$+ \frac{2(1 - \gamma_{t})}{\gamma_{t}^{2}} \boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{r} - \frac{(1 - \gamma_{t})^{2}}{\gamma_{t}^{2}} \boldsymbol{r} \circ \boldsymbol{r}$$

$$\leq \frac{1}{\gamma_{T}} (\gamma_{T} \boldsymbol{P}^{\pi_{i_{h+1}}} - \boldsymbol{I})(\boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{Q}_{\gamma_{t}}^{*}) + 8 \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{Q}_{\gamma_{t}}^{*} - \boldsymbol{Q}_{\gamma_{i}}^{*}\|_{\infty} \mathbf{1} + 4 \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\Delta}_{i}\|_{\infty} \mathbf{1}$$

$$+ \frac{2(1 - \gamma_{t})}{\gamma_{t}^{2}} \boldsymbol{Q}_{\gamma_{t}}^{*} \circ \boldsymbol{r} - \frac{(1 - \gamma_{t})^{2}}{\gamma_{t}^{2}} \boldsymbol{r} \circ \boldsymbol{r}, \tag{94}$$

where the first inequality follows from (83) and the second inequality follows from (84) and (93), and the last inequality follows from the monotonicity of γ_t . Applying (94) to (82) yields the telescoping sum

$$\begin{split} &\sum_{h=0}^{H-1} \gamma_T^h \prod_{k=1}^h \boldsymbol{P}^{\pi_{i_k}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_i) \\ &\leq \sum_{h=0}^{H-1} \gamma_T^h \prod_{k=1}^h \boldsymbol{P}^{\pi_{i_k}} \max_{\lfloor t/2 \rfloor \leq i < t} \operatorname{Var}_{\boldsymbol{P}}(\boldsymbol{V}_{\gamma_i}^*) + \frac{4}{1 - \gamma_T} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\Delta}_i\|_{\infty} \mathbf{1} \\ &\leq \frac{1}{\gamma_T} \sum_{h=0}^{H-1} \gamma_T^h \prod_{k=1}^h \boldsymbol{P}^{\pi_{i_k}} \left(\gamma_T \boldsymbol{P}^{\pi_{i_{h+1}}} - \boldsymbol{I} \right) (\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{Q}_{\gamma_t}^*) \\ &\quad + \left(\frac{8}{1 - \gamma_t} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{Q}_{\gamma_t}^* - \boldsymbol{Q}_{\gamma_i}^*\|_{\infty} + \frac{8}{1 - \gamma_t} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\Delta}_i\|_{\infty} + \frac{2(1 - \gamma_t)}{\gamma_t^2 (1 - \gamma_T)} \|\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{r}\|_{\infty} \right) \mathbf{1} \\ &= \frac{1}{\gamma_T} \left(\gamma_T^H \prod_{k=1}^H \boldsymbol{P}^{\pi_{i_{h+1}}} - \boldsymbol{I} \right) (\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{Q}_{\gamma_t}^*) \\ &\quad + \left(\frac{8}{1 - \gamma_T} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{Q}_{\gamma_t}^* - \boldsymbol{Q}_{\gamma_t}^*\|_{\infty} + \frac{8}{1 - \gamma_T} \max_{\lfloor t/2 \rfloor \leq i < t} \|\boldsymbol{\Delta}_i\|_{\infty} + \frac{2(1 - \gamma_t)}{\gamma_t^2 (1 - \gamma_T)} \|\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{r}\|_{\infty} \right) \mathbf{1}, \end{split}$$

where the first inequality follows from (82), and the second inequality uses $\prod_{k=1}^{h} P^{\pi_{i_k}} \mathbf{1} = \mathbf{1}$ as the product is a probability transition matrix. As the rows of the probability transition matrix $\prod_{k=1}^{H} P^{\pi_{i_{k+1}}}$ all sum up to one, we have

$$\left\|\frac{1}{\gamma_T}\left(\gamma_T^H\prod_{k=1}^H \boldsymbol{P}^{\pi_{i_{h+1}}} - \boldsymbol{I}\right)(\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{Q}_{\gamma_t}^*)\right\|_{\infty} \leq \frac{2}{\gamma_T}\|\boldsymbol{Q}_{\gamma_t}^*\|_{\infty}^2 \leq \frac{2}{\gamma_T} \leq 4$$

by the boundedness of $\|Q_{\gamma}^*\|_{\infty} \le 1$ for all $\gamma \in (0,1)$, as well as the fact that $\gamma_t \ge 1/2$ for $t \ge 160$. On the other hand, for any $|t/2| \le i < t$, invoking Lemma D.1, we obtain

$$\begin{aligned} \|\boldsymbol{Q}_{\gamma_{t}}^{*} - \boldsymbol{Q}_{\gamma_{i}}^{*}\|_{\infty} &\leq \|\boldsymbol{Q}_{\gamma_{t}}^{*} - J^{*}\boldsymbol{1}\|_{\infty} + \|J^{*}\boldsymbol{1} - \boldsymbol{Q}_{\gamma_{i}}^{*}\|_{\infty} \\ &\leq 3(1 - \gamma_{t})t_{\text{mix}} + 3(1 - \gamma_{i})t_{\text{mix}} \leq 3(1 + 2^{1/5})t^{-1/5}t_{\text{mix}} \leq 9T^{-1/5}(\log T)^{1/5}t_{\text{mix}} \end{aligned}$$

for $T/\log T \le t \le \log T$. In addition, for $T \ge 160$, we also have

$$\frac{2(1-\gamma_t)}{\gamma_t^2(1-\gamma_T)} \|\boldsymbol{Q}_{\gamma_t}^* \circ \boldsymbol{r}\|_{\infty} \leq 3 \|\boldsymbol{Q}_{\gamma_t}^*\|_{\infty} \|\boldsymbol{r}\|_{\infty} \leq 3,$$

which, together with the above pieces, leads to (27) and completes the proof of Lemma A.1.

Lemma D.7. For any $i \geq 1$, it holds that

$$\|\operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_i) - \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{\alpha_i}^*)\|_{\infty} \le 3(\log i)^2 i^{1/3} \|\boldsymbol{\delta}_i\|_{\infty}.$$
(95)

Proof of Lemma D.7. We first bound $\mathbf{v}_i - \mathbf{v}_{\alpha_i}^*$ for all $i \geq 1$. Recalling that $\delta_i = \mathbf{q}_i - \mathbf{q}_{\alpha_i}^*$, we have

$$\mathbf{v}_i - \mathbf{v}_{lpha_i}^* = oldsymbol{P}^{\pi_i} \mathbf{q}_i - oldsymbol{P}^{\pi_{lpha_i}^*} \mathbf{q}_{lpha_i}^* \overset{ ext{(i)}}{\leq} oldsymbol{P}^{\pi_i} \mathbf{q}_i - oldsymbol{P}^{\pi_i} \mathbf{q}_{lpha_i}^* \leq \|oldsymbol{P}^{\pi_i}\|_1 \|\mathbf{q}_i - \mathbf{q}_{lpha_i}^*\|_{\infty} \overset{ ext{(ii)}}{\leq} \|oldsymbol{\delta}_i\|_{\infty},$$

where (i) follows from the optimality of $\pi_{\alpha_i}^*$ with respect to $\mathbf{q}_{\alpha_i}^*$, and (ii) holds since \mathbf{P}^{π_i} is a probability transition matrix. Similarly,

$$\mathbf{v}_i - \mathbf{v}_{\alpha_i}^* \geq \boldsymbol{P}^{\pi_{\alpha_i}^*} \mathbf{q}_i - \boldsymbol{P}^{\pi_{\alpha_i}^*} \mathbf{q}_{\alpha_i}^* \geq - \|\boldsymbol{P}^{\pi_{\alpha_i}^*}\|_1 \|\mathbf{q}_i - \mathbf{q}_{\alpha_i}^*\|_{\infty} \geq - \|\boldsymbol{\delta}_i\|_{\infty},$$

which leads to

$$\|\mathbf{v}_i - \mathbf{v}_{\alpha_i}^*\|_{\infty} \le \|\boldsymbol{\delta}_i\|_{\infty} \tag{96}$$

for all $i \geq 1$. By the definition of entrywise variance, we have

$$\operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i}) - \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{\alpha_{i}}^{*}) = (\boldsymbol{P}(\mathbf{v}_{i} \circ \mathbf{v}_{i}) - (\boldsymbol{P}\mathbf{v}_{i}) \circ (\boldsymbol{P}\mathbf{v}_{i})) - (\boldsymbol{P}(\mathbf{v}_{\alpha_{i}}^{*} \circ \mathbf{v}_{\alpha_{i}}^{*}) - (\boldsymbol{P}\mathbf{v}_{\alpha_{i}}^{*}) \circ (\boldsymbol{P}\mathbf{v}_{\alpha_{i}}^{*}))$$

$$= \boldsymbol{P}(\mathbf{v}_{i} \circ \mathbf{v}_{i} - \mathbf{v}_{\alpha_{i}}^{*} \circ \mathbf{v}_{\alpha_{i}}^{*}) + (\boldsymbol{P}\mathbf{v}_{\alpha_{i}}^{*}) \circ (\boldsymbol{P}\mathbf{v}_{\alpha_{i}}^{*}) - (\boldsymbol{P}\mathbf{v}_{i}) \circ (\boldsymbol{P}\mathbf{v}_{i})$$

$$= \boldsymbol{P}((\mathbf{v}_{i} - \mathbf{v}_{\alpha_{i}}^{*}) \circ (\mathbf{v}_{i} + \mathbf{v}_{\alpha_{i}}^{*})) - (\boldsymbol{P}\mathbf{v}_{i} - \boldsymbol{P}\mathbf{v}_{\alpha_{i}}^{*}) \circ (\boldsymbol{P}\mathbf{v}_{i} + \boldsymbol{P}\mathbf{v}_{\alpha_{i}}^{*}). \tag{97}$$

Furthermore, we have

$$\begin{aligned} \left\| \boldsymbol{P} \big((\mathbf{v}_i - \mathbf{v}_{\alpha_i}^*) \circ (\mathbf{v}_i + \mathbf{v}_{\alpha_i}^*) \big) \right\|_{\infty} &\leq \| \boldsymbol{P} \|_1 \left\| (\mathbf{v}_i - \mathbf{v}_{\alpha_i}^*) \circ (\mathbf{v}_i + \mathbf{v}_{\alpha_i}^*) \right\|_{\infty} \\ &\leq \| \boldsymbol{P} \|_1 \left\| \mathbf{v}_i - \mathbf{v}_{\alpha_i}^* \right\|_{\infty} \| \mathbf{v}_i + \mathbf{v}_{\alpha_i}^* \right\|_{\infty} \\ &\leq \left(\| \mathbf{v}_i \|_{\infty} + \| \mathbf{v}_{\alpha_i}^* \|_{\infty} \right) \| \boldsymbol{\delta}_i \|_{\infty}, \end{aligned}$$

where the last inequality follows from $\|P\|_1 = 1$, the triangular inequality, and (96). Meanwhile, similar arguments yield

$$\begin{aligned} \left\| (\boldsymbol{P} \mathbf{v}_i - \boldsymbol{P} \mathbf{v}_{\alpha_i}^*) \circ (\boldsymbol{P} \mathbf{v}_i + \boldsymbol{P} \mathbf{v}_{\alpha_i}^*) \right\|_{\infty} &\leq \| \boldsymbol{P} \mathbf{v}_i - \boldsymbol{P} \mathbf{v}_{\alpha_i}^* \|_{\infty} \| \boldsymbol{P} \mathbf{v}_i + \boldsymbol{P} \mathbf{v}_{\alpha_i}^* \|_{\infty} \\ &\leq \| \boldsymbol{P} \|_1 \| \mathbf{v}_i - \mathbf{v}_{\alpha_i}^* \|_{\infty} \| \boldsymbol{P} \|_1 \| \mathbf{v}_i + \mathbf{v}_{\alpha_i}^* \|_{\infty} \leq \left(\| \mathbf{v}_i \|_{\infty} + \| \mathbf{v}_{\alpha_i}^* \|_{\infty} \right) \| \boldsymbol{\delta}_i \|_{\infty}. \end{aligned}$$

Plugging the two bounds back in (97) yields

$$\begin{aligned} \| \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{i}) - \operatorname{Var}_{\boldsymbol{P}}(\mathbf{v}_{\alpha_{i}}^{*}) \|_{\infty} &\leq 2 (\|\mathbf{v}_{i}\|_{\infty} + \|\mathbf{v}_{\alpha_{i}}^{*}\|_{\infty}) \|\boldsymbol{\delta}_{i}\|_{\infty} \\ &\leq 2 (i^{1/9} + 3(\log i)^{2} i^{1/3} / c_{1}) \|\boldsymbol{\delta}_{i}\|_{\infty} \leq 2 (1 + 3 / c_{1}) (\log i)^{2} i^{1/3} \|\boldsymbol{\delta}_{i}\|_{\infty}. \end{aligned}$$

D.5 Performance difference lemmas

The following performance difference lemma for average reward is adapted from Even-Dar et al. (2009, Lemma 4.1).

Lemma D.8 (Performance difference lemma). For any policy π whose long-term average reward is J^{π} , we define its bias function as

$$q^{\pi}(s,a) = \mathbb{E}_{\pi} \Big[\sum_{k=1}^{\infty} (r(s_k, a_k) - J^{\pi}) \mid s_1 = s, a_1 = a \Big],$$

Then for any two policies π and π' , the difference between their long-term average reward is

$$J^{\pi'} - J^{\pi} = \mathbb{E}_{s \sim d_{\pi'}} \Big[\sum_{a \in \mathcal{A}} \big(\pi'(a \mid s) - \pi(a \mid s) \big) q^{\pi}(s, a) \Big],$$

where d_{π} is the long-term average visit probability of s under policy π .

We quote without proof the performance difference lemma for discounted reward (Kakade and Langford, 2002). **Lemma D.9** (Performance difference lemma). Recall that Q_{γ}^{π} and V_{γ}^{π} is the (scaled) Q- and value functions for any policy π and discount factor γ . For any policies π and π' , it holds that

$$V_{\gamma}^{\pi}(s) - V_{\gamma}^{\pi'}(s) = \frac{1}{1 - \gamma} \mathbb{E}_{s' \sim d^{\pi}} \mathbb{E}_{a' \sim \pi(\cdot \mid s')} \left[Q_{\gamma}^{\pi'}(s', a') - V_{\gamma}^{\pi'}(s') \right], \quad \text{for all } s \in \mathcal{S},$$

where $d^{\pi}(s') = (1 - \gamma) \sum_{t=1}^{\infty} \mathbb{P}_{\pi}(S_t = s' \mid s_0 = s)$.

E Useful facts

In this section, we collect some useful facts that are useful for the technical proofs.

E.1 Magnitude of estimates

Lemma E.1. Suppose the initialization of Algorithm 1 satisfies $0 \le Q_0 \le 1$ (and thus $0 \le V_0 \le 1$), then $0 \le Q_t \le 1$ and $\|\Delta_t\|_{\infty} \le 1$ for all $t \ge 1$.

Proof of Lemma E.1. Note that if $0 \le Q_{t-1} \le 1$ and $0 \le V_{t-1} \le 1$, then by the updating rule (6), since $0 \le r \le 1$ and each row of P_t adds up to 1, we have $Q_t \ge 0$ and

$$Q_t \leq (1 - \eta_t)\mathbf{1} + \eta_t [(1 - \gamma_t)\mathbf{1} + \gamma_t \mathbf{1}] = \mathbf{1}.$$

The desired result then follows from an induction argument. As a consequence, we have $\|\Delta_t\|_{\infty} \leq 1$ by the boundedness of optimal value functions that $\mathbf{0} \leq Q_{\gamma_t}^* \leq \mathbf{1}$ for all $t \geq 1$.

Lemma E.2. With $\mathbf{q}_0 = \mathbf{0}$ and $\theta_t = (1 + \frac{c_1 t^{2/3}}{(\log t)^2})^{-1}$ for $t \geq 2$ and $\theta_1 = 0$, Algorithm 2 obeys $\mathbf{0} \leq \mathbf{q}_t \leq 3(\log t)^2 t^{1/3}/c_1 \mathbf{1}$ and $\mathbf{0} \leq \mathbf{v}_t \leq 3(\log t)^2 t^{1/3}/c_1 \mathbf{1}$ for all $t \geq 2$.

Proof of Lemma E.2. Noting the induction that

$$\|\mathbf{q}_{t}\|_{\infty} \leq (1 - \theta_{t}) \|\mathbf{q}_{t-1}\|_{\infty} + \theta_{t} (\|\mathbf{r}\|_{\infty} + \alpha_{t} \|\mathbf{q}_{t-1}\|_{\infty})$$

$$\leq (1 - \theta_{t} + \theta_{t}\alpha_{t}) \|\mathbf{q}_{t-1}\|_{\infty} + \theta_{t} \leq \|\mathbf{q}_{t-1}\|_{\infty} + \theta_{t}.$$

we have

$$\|\mathbf{q}_t\|_{\infty} \le \sum_{i=1}^t \theta_i \le (\log t)^2 \sum_{i=1}^t i^{-2/3}/c_1 \le 3(\log t)^2 t^{1/3}/c_1,$$

where the second inequality follows from (108). On the other hand, as $0 \le r \le 1$, we have

$$\mathbf{q}_t \ge (1 - \theta)\mathbf{q}_{t-1} + \theta_t \alpha_t \mathbf{q}_{t-1},$$

which, by induction, leads to $\mathbf{q}_t \geq \mathbf{0}$ for all $t \geq 1$. The definition of \mathbf{v}_t implies $\mathbf{0} \leq \mathbf{v}_t \leq 3(\log t)^2 t^{1/3}/c_1 \mathbf{1}$ for all $t \geq 2$, which complestes the proof of Lemma E.2.

E.2 Learning rates for estimating optimal value in Theorem 3.1

In this part, we collect some useful facts about the learning rates speficied for Theorem 3.1 that would be used repeatedly. Recall that we define the learning rates as $\eta_t = (1 + \frac{c_1 t^{3/5}}{(\log t)^3})^{-1}$, $\forall t \geq 2$ for some constant $c_1 > 0$.

Bound on η_i . Firstly, we have the naive upper bound

$$\eta_j = \frac{1}{1 + \frac{c_1 j^{3/5}}{(\log j)^3}} \le \frac{(\log j)^3}{c_1 j^{3/5}}, \quad \text{for all } j \ge 2.$$
(98)

One can also check that $j^{3/5}/(\log j)^3 \ge 0.1$ for all $j \ge 1$ with the convention that $1/0 = \infty$, which leads to the lower bounds

$$\eta_j = \frac{1}{1 + \frac{c_1 j^{3/5}}{(\log j)^3}} \ge \frac{1}{\frac{10 j^{3/5}}{(\log j)^3} + \frac{c_1 j^{3/5}}{(\log j)^3}} \ge \frac{(\log j)^3}{(10 + c_1)j^{3/5}}, \quad \text{for all } j \ge 2.$$
(99)

Furthermore, it can be easily checked that $(\log j)^3 j^{-3/5}$ is decreasing in j for $j \ge 150$. Hence

$$\eta_j \ge \eta_t \text{ and } \frac{(\log j)^3}{j^{3/5}} \ge \frac{(\log t)^3}{t^{3/5}}, \quad \text{for all } t > 150 \text{ and } 150 \le j \le t.$$

Compound learning rates. Some associated quantities in the analysis include $\eta_t^{(t)} = \eta_t$,

$$\eta_0^{(t)} = \prod_{j=1}^t (1 - \eta_j), \text{ and } \eta_i^{(t)} = \eta_i \cdot \prod_{j=i+1}^t (1 - \eta_j), \quad \forall 1 \le i < t.$$

The sum of $\eta_i^{(t)}$ over *i* obeys

$$\sum_{i=0}^{t} \eta_i^{(t)} = \prod_{j=1}^{t} (1 - \eta_j) + \eta_1 \prod_{j=2}^{t} (1 - \eta_j) + \dots + \eta_{t-1} (1 - \eta_t) + \eta_t = 1.$$
 (101)

Moreover, we consider a generic $\beta \in (0,1)$ and bound the compound learning rates for $i \leq \lfloor (1-\beta)t \rfloor$ and $i > \lfloor (1-\beta)t \rfloor$ separately.

• When $\beta < 1/2$, for any $t \ge 300$ (so that $\lfloor (1-\beta)t \rfloor \ge 150$) and any $i \le \lfloor (1-\beta)t \rfloor$,

$$\eta_i^{(t)} \le \prod_{j=i+1}^t (1 - \eta_j) \le \exp\left(-\sum_{j=\lfloor (1-\beta)t\rfloor}^t \eta_j\right)
\le \exp\left(-\sum_{j=\lfloor (1-\beta)t\rfloor}^t \eta_t\right) \le \exp\left(-\beta t \cdot \frac{(\log t)^3}{c_1 t^{3/5}}\right),$$
(102)

where the third inequality follows from (100), and the fourth inequality follows from (98).

• When $\beta < 1/2$, for any $t \ge 300$ and any $i > |(1 - \beta)t| \ge t/2 \ge 150$,

$$\eta_i^{(t)} \le \eta_i \le \frac{(\log i)^3}{c_1 i^{3/5}} \le \frac{2(\log t)^3}{c_1 t^{3/5}}$$
(103)

following (98) and (100) and the fact that $\log i \leq \log t$ and $i^{-3/5} \leq 2t^{-3/5}$ for any $i \geq t/2$.

E.3 Learning rates for optimal policy in Theorem 3.3

We now present some useful facts about the learning rates that are particular to Theorem 3.3. Recall that in Theorem 3.3, we set $\gamma_t = 1 - t^{-1/8}$ and

$$\eta_t = \frac{1}{1 + \frac{c_1 t^{5/8}}{(\log t)^2}}, \quad t \ge 2$$

for some constant $c_1 > 0$ in Algorithm 1.

Bound on η_{j} . Firstly, we have the naive upper bound

$$\eta_j = \frac{1}{1 + \frac{c_1 j^{5/8}}{(\log j)^2}} \le \frac{(\log j)^2}{c_1 j^{5/8}}, \quad \text{for all } j \ge 2.$$
(104)

One can check that $j^{5/8}/(\log j)^2 \ge 0.5$ for all $j \ge 1$ with the convention that $1/0 = \infty$, which leads to the lower bounds

$$\eta_j = \frac{1}{1 + \frac{c_1 j^{5/8}}{(\log j)^2}} \ge \frac{1}{\frac{2j^{5/8}}{(\log j)^2} + \frac{c_1 j^{5/8}}{(\log j)^2}} \ge \frac{(\log j)^2}{(2 + c_1)j^{5/8}}, \quad \text{for all } j \ge 2.$$
(105)

Furthermore, it can be easily checked that $(\log j)^2 j^{-5/8}$ is decreasing in j for $j \geq 50$. Hence

$$\eta_j \ge \eta_t \text{ and } \frac{(\log j)^2}{j^{5/8}} \ge \frac{(\log t)^2}{t^{5/8}}, \quad \text{for all } t > 50 \text{ and } 50 \le j \le t.$$
(106)

Compound learning rates. Some associated quantities in the analysis include $\eta_t^{(t)} = \eta_t$,

$$\eta_0^{(t)} = \prod_{j=1}^t (1 - \eta_j), \text{ and } \eta_i^{(t)} = \eta_i \cdot \prod_{j=i+1}^t (1 - \eta_j), \quad \forall 1 \le i < t.$$

The sum of $\eta_i^{(t)}$ over *i* obeys

$$\sum_{i=0}^{t} \eta_i^{(t)} = \prod_{j=1}^{t} (1 - \eta_j) + \eta_1 \prod_{j=2}^{t} (1 - \eta_j) + \dots + \eta_{t-1} (1 - \eta_t) + \eta_t = 1.$$

Moreover, we consider a generic $\beta \in (0,1)$ and bound the compound learning rates for $i \leq \lfloor (1-\beta)t \rfloor$ and $i > \lfloor (1-\beta)t \rfloor$ separately.

• When $\beta < 1/2$, for any $t \ge 100$ (so that $\lfloor (1-\beta)t \rfloor \ge 50$) and any $i \le \lfloor (1-\beta)t \rfloor$,

$$\eta_i^{(t)} \le \prod_{j=i+1}^t (1 - \eta_j) \le \exp\left(-\sum_{j=\lfloor (1-\beta)t \rfloor}^t \eta_j\right)
\le \exp\left(-\sum_{j=\lfloor (1-\beta)t \rfloor}^t \eta_t\right) \le \exp\left(-\beta t \cdot \frac{(\log t)^2}{c_1 t^{5/8}}\right),$$

where the third inequality follows from (100), and the fourth inequality follows from (98).

• When $\beta < 1/2$, for any $t \ge 100$ and any $i > \lfloor (1-\beta)t \rfloor \ge t/2 \ge 50$,

$$\eta_i^{(t)} \le \eta_i \le \frac{(\log i)^3}{c_1 i^{5/8}} \le \frac{2(\log t)^2}{c_1 t^{5/8}}$$
(107)

following (104) and (106) and the fact that $\log i \le \log t$ and $i^{-5/8} \le 2t^{-5/8}$ for any $i \ge t/2$.

E.4 Learning rates for bias function in Theorem 5.2

We now proceed to some basic facts on the learning rates of Algorithm 2. Recall that we define the learning rates as

$$\theta_t = \frac{1}{1 + \frac{c_1 t^{2/3}}{(\log t)^2}}, \quad \forall \ t \ge 2.$$

Bound on θ_j . Firstly, we have the naive upper bound

$$\theta_j = \frac{1}{1 + \frac{c_1 j^{2/3}}{(\log j)^2}} \le \frac{(\log j)^2}{c_1 j^{2/3}}, \quad \text{for all } j \ge 1.$$
 (108)

Also, one can check that $j^{2/3}/(\log j)^2 \ge 1/2$, which leads to

$$\theta_j = \frac{1}{1 + \frac{c_1 j^{2/3}}{(\log j)^2}} \ge \frac{(\log j)^2}{(2 + c_1)j^{2/3}}, \quad \text{for all } j \ge 1.$$

Furthermore, it can be easily checked that $(\log j)^2 j^{-2/3}$ is decreasing in j for $j \geq 30$. Hence

$$\theta_j \ge \theta_t \text{ and } \frac{(\log j)^2}{j^{2/3}} \ge \frac{(\log t)^2}{t^{2/3}}, \quad \text{for all } t > 50 \text{ and } 50 \le j \le t.$$
 (109)

Compound learning rates. Some associated quantities in the analysis include $\theta_t^{(t)} = 1$,

$$\theta_0^{(t)} = \prod_{j=1}^t (1 - \theta_j), \text{ and } \theta_i^{(t)} = \theta_i \cdot \prod_{j=i+1}^t (1 - \theta_j), \forall 1 \le i < t.$$

The sum of $\theta_i^{(t)}$ over *i* obeys

$$\sum_{i=0}^{t} \theta_i^{(t)} = \prod_{j=1}^{t} (1 - \theta_j) + \theta_1 \prod_{j=2}^{t} (1 - \theta_j) + \dots + \theta_{t-1} (1 - \theta_t) + \theta_t = 1.$$

Moreover, we also define

$$\beta = \frac{c_2}{T^{1/3} \log T},$$

and bound the compound learning rates for $i \leq \lfloor (1-\beta)t \rfloor$ and $i > \lfloor (1-\beta)t \rfloor$ separately.

• For any t such that $\lfloor (1-\beta)t \rfloor \geq 50$ (so that $t \geq 50$) and any $i \leq \lfloor (1-\beta)t \rfloor$,

$$\theta_i^{(t)} \le \prod_{j=i+1}^t (1 - \theta_j) \le \exp\left(-\sum_{j=\lfloor (1-\beta)t\rfloor}^t \theta_j\right)$$

$$\le \exp\left(-\sum_{j=\lfloor (1-\beta)t\rfloor}^t \theta_t\right) \le \exp\left(-\beta t \cdot \frac{(\log t)^2}{(2+c_1)t^{2/3}}\right),\tag{110}$$

where the third inequality follows from (109), and the fourth inequality follows from (108).

• For any t such that $\lfloor (1-\beta)t \rfloor \geq 50$ (so that $t \geq 50$) and any $i > \lfloor (1-\beta)t \rfloor$,

$$\theta_i^{(t)} \le \theta_i \le \frac{c'(\log i)^3}{i^{3/5}} \le \frac{c'(\log t)^3}{t^{3/5}}$$

following (108) and (109).