

Improved Regret Bound for Safe Reinforcement Learning via Tighter Cost Pessimism and Reward Optimism

Anonymous authors

Paper under double-blind review

Keywords: Safe Reinforcement Learning, Constrained MDPs, Regret Analysis.

Summary

This paper studies the safe reinforcement learning problem formulated as an episodic finite-horizon tabular constrained Markov decision process with an unknown transition kernel and stochastic reward and cost functions. We propose a model-based algorithm based on novel cost and reward function estimators that provide tighter cost pessimism and reward optimism. While guaranteeing no constraint violation in every episode, our algorithm achieves a regret upper bound of $\tilde{O}((\bar{C} - \bar{C}_b)^{-1} H^{2.5} S \sqrt{AK})$ where \bar{C} is the cost budget for an episode, \bar{C}_b is the expected cost under a safe baseline policy over an episode, H is the horizon, and S , A and K are the number of states, actions, and episodes, respectively. This improves upon the best-known regret upper bound, and when $\bar{C} - \bar{C}_b = \Omega(H)$, the gap from the regret lower bound of $\Omega(H^{1.5} \sqrt{SAK})$ is $\tilde{O}(\sqrt{S})$. We deduce our cost and reward function estimators via a Bellman-type law of total variance to obtain tight bounds on the expected sum of the variances of value function estimates. This leads to a tighter dependence on the horizon in the function estimators. We also present numerical results to demonstrate the computational effectiveness of our proposed framework.

Contribution(s)

1. This paper presents an algorithm for episodic finite-horizon tabular constrained Markov decision processes with an improved regret upper bound of $\tilde{O}((\bar{C} - \bar{C}_b)^{-1} H^{2.5} S \sqrt{AK})$.
Context: The best-known regret upper bound is $\tilde{O}((\bar{C} - \bar{C}_b)^{-1} H^3 S \sqrt{AK})$ due to [Bura et al. \(2022\)](#), and our result improves it by a factor of $\tilde{O}(\sqrt{H})$.
2. Our algorithm ensures zero constraint violation for each episode, given the knowledge of a safe baseline policy.
Context: The guarantee is stronger than zero cumulative constraint violation, for which error cancellations are permitted across episodes. Hence, under our algorithm, there is no episode in which the constraint is violated. A safe baseline policy is necessary to enforce zero constraint violation, especially in the early stage ([Liu et al., 2021](#); [Bura et al., 2022](#)).
3. When $\bar{C} - \bar{C}_b = \Omega(H)$, the gap from the regret lower bound of $\Omega(H^{1.5} \sqrt{SAK})$ is $\tilde{O}(\sqrt{S})$, and our result nearly matches the lower bound in terms of H .
Context: The lower bound is originally derived for the unconstrained case ([Jin et al., 2020](#); [Domingues et al., 2021](#)), and it also works for the constrained case as we can take trivial cost functions.
4. The reduction in the regret upper bound is a consequence of our novel reward and cost function estimators. The key is to control the error of estimating the unknown transition kernel over each episode. In particular, we provide a tighter bound on the estimation error for each episode, based on a Bellman-type law of total variance. The bound is given by a function of the estimated transition kernel, whose choice can be optimized by the algorithm.
Context: Our Bellman-type law of total variance technique refines the analysis of [Bura et al. \(2022\)](#). The technique is inspired by [Chen & Luo \(2021\)](#), while they gave only a cumulative error bound across all episodes, and at the same time, the bound is expressed as a function of the true transition kernel which is unknown to the algorithm.

Improved Regret Bound for Safe Reinforcement Learning via Tighter Cost Pessimism and Reward Optimism

Anonymous authors

Paper under double-blind review

Abstract

1 This paper studies the safe reinforcement learning problem formulated as an episodic
 2 finite-horizon tabular constrained Markov decision process with an unknown transition
 3 kernel and stochastic reward and cost functions. We propose a model-based algorithm
 4 based on novel cost and reward function estimators that provide tighter cost pessimism
 5 and reward optimism. While guaranteeing no constraint violation in every episode, our
 6 algorithm achieves a regret upper bound of $\tilde{O}((\bar{C} - \bar{C}_b)^{-1} H^{2.5} S \sqrt{AK})$ where \bar{C} is
 7 the cost budget for an episode, \bar{C}_b is the expected cost under a safe baseline policy over
 8 an episode, H is the horizon, and S , A and K are the number of states, actions, and
 9 episodes, respectively. This improves upon the best-known regret upper bound, and
 10 when $\bar{C} - \bar{C}_b = \Omega(H)$, the gap from the regret lower bound of $\Omega(H^{1.5} \sqrt{SAK})$ is
 11 $\tilde{O}(\sqrt{S})$. The reduction in the regret upper bound is a consequence of our novel reward
 12 and cost function estimators. The key is to control the error of estimating the unknown
 13 transition kernel over each episode. In particular, we provide a tighter bound on the
 14 estimation error for each episode, based on a Bellman-type law of total variance to ana-
 15 lyze the expected sum of the variances of value function estimates. The bound is given
 16 by a function of the estimated transition kernel, whose choice can be optimized by the
 17 algorithm. This leads to a tighter dependence on the horizon in the function estimators.
 18 We also present numerical results to demonstrate the computational effectiveness of our
 19 proposed framework.

20 1 Introduction

21 Safe reinforcement learning (RL) aims to learn a policy that maximizes the cumulative reward and, at
 22 the same time, ensures that some safety requirements are satisfied during the learning process. Safe
 23 RL provides modeling frameworks for many practical scenarios where violating a safety constraint
 24 results in a critical situation. For example, it is crucial to enforce collision avoidance for autonomous
 25 driving (Isele et al., 2018; Krasowski et al., 2020) and robotics (Fisac et al., 2018; García & Shafie,
 26 2020). For financial planning, there exist legal and business regulations (Abe et al., 2010). For
 27 healthcare systems, service providers consider restrictions due to patients’ conditions (Coronato
 28 et al., 2020).

29 The standard approach is to formulate a safe RL problem as a constrained Markov decision process
 30 (CMDP), where the objective is to maximize the expected reward over a time horizon while there
 31 is a constraint that the expected cost should be under budget (Altman, 1999). The presence of con-
 32 straints, however, brings about challenges in developing solution methods for CMDPs. The Bellman
 33 optimality principle does not hold for CMDPs, and as a consequence, backward induction and the
 34 greedy operator cannot be directly applied to CMDPs (Altman, 1999). This makes online learning
 35 of CMDPs difficult, and we need significantly different frameworks and algorithms compared to the
 36 unconstrained setting (García et al., 2015; Efroni et al., 2020; Gu et al., 2024).

37 The first direction for online reinforcement learning of CMDPs is to consider *cumulative (or soft)*
 38 *constraint violation*, which sums up the constraint violations across episodes (Efroni et al., 2020).
 39 Here, the constraint violation in an episode is defined as the expected cost minus the budget. Then a
 40 policy can have a negative constraint violation, which means that a positive violation in one episode
 41 can be canceled out by a negative violation in another episode in the sum. This cancellation effect
 42 allows oscillating between such two cases, while still achieving zero cumulative constraint violation.
 43 This phenomenon can indeed be observed in practice (Stooke et al., 2020; Moskovitz et al., 2023).

44 The second direction attempts to remedy the issue of error cancellation with the notion of *hard*
 45 *constraint violation* (Efroni et al., 2020). It ignores episodes with a negative violation and takes
 46 the sum of only the positive constraint violations. Efroni et al. (2020) developed OptCMDP and its
 47 efficient variant, OptCMDP-bonus, that attain a regret upper bound and a hard constraint violation
 48 of $\tilde{O}(H^2\sqrt{S^2AK})$. Recently, Ghosh et al. (2024) proposed a model-free algorithm with the same
 49 asymptotic guarantees. However, as in the first setting, the algorithms cannot avoid episodes in
 50 which the constraint is violated. Thus, they are still not suitable for the aforementioned applications,
 51 where even a single incidence of violation can cause substantial problems.

52 The third approach seeks *zero (hard) constraint violation*, requiring that the constraint is satisfied
 53 in every episode (Simão et al., 2021). Satisfying constraints in the early stage is difficult when
 54 the model parameters, especially the transition kernel, are unknown. Simão et al. (2021) con-
 55 sidered some abstraction of the transition model under which they showed an algorithm with no
 56 constraint violation, but no regret upper bound was presented. Then Liu et al. (2021) came up
 57 with the first algorithm, OptPess-LP, that achieves a sublinear regret with no constraint violation,
 58 assuming the knowledge of a *safe baseline policy*. Here, a safe baseline policy is a policy under
 59 which the expected cost is lower than the budget. OptPess-LP guarantees a regret upper bound of
 60 $\tilde{O}((\bar{C} - \bar{C}_b)^{-1}H^3\sqrt{S^3AK})$ where \bar{C} is the budget, \bar{C}_b is the expected cost under the safe baseline
 61 policy, H is the length of the horizon, and S , A and K are the number of states, actions, and episodes,
 62 respectively. Bura et al. (2022) developed Doubly Optimistic Pessimistic Exploration (DOPE) with
 63 an improved regret upper bound of $\tilde{O}((\bar{C} - \bar{C}_b)^{-1}H^3\sqrt{S^2AK})$. DOPE is based on designing tight
 64 optimistic reward function estimators (reward optimism) and conservative cost function estimators
 65 (cost pessimism).

66 While DOPE establishes a tight regret upper bound with no constraint violation, there is still room
 67 for improvement. The regret lower bound of $\Omega(H^{1.5}\sqrt{SAK})$ for the unconstrained case (Jin et al.,
 68 2018; Domingues et al., 2021) also works as a lower bound for the constrained setting because we
 69 may take trivial cost functions. However, even when $\bar{C} - \bar{C}_b = \Omega(H)$, the regret upper bound
 70 of DOPE is as low as $\tilde{O}(H^2\sqrt{S^2AK})$ which has a gap of $\tilde{O}(\sqrt{HS})$ from the lower bound. This
 71 naturally motivates the following question.

72 *Is there an algorithm for learning CMDPs that guarantees no constraint violation during learning*
 73 *and achieves an improved regret upper bound?*

74 **Our Contributions** We answer this question affirmatively with an algorithm that improves upon
 75 DOPE via tighter reward optimism and cost pessimism. Our results are summarized in Table 1 and
 76 as follows.

- 77 • Our algorithm, DOPE+, achieves a regret upper bound of $\tilde{O}((\bar{C} - \bar{C}_b)^{-1}H^{2.5}\sqrt{S^2AK})$ and en-
 78 sures no constraint violation in every episode, with the knowledge of a safe baseline policy. This
 79 improves upon the best-known regret upper bound $\tilde{O}((\bar{C} - \bar{C}_b)^{-1}H^3\sqrt{S^2AK})$ attained by DOPE.
- 80 • When the gap $\bar{C} - \bar{C}_b$ between the budget and the expected cost under the safe baseline policy
 81 satisfies $\bar{C} - \bar{C}_b = \Omega(H)$, the regret upper bound becomes $\tilde{O}(H^{1.5}\sqrt{S^2AK})$. Then the gap from
 82 the regret lower bound of $\Omega(H^{1.5}\sqrt{SAK})$ is $\tilde{O}(\sqrt{S})$, which shows that the regret upper bound
 83 achieves the optimal dependence on the horizon H .
- 84 • The improvement comes from our novel reward and cost function estimators with tighter reward
 85 optimism and cost pessimism. We deduce the function estimators by providing a tighter upper

86 bound on the estimation error for each episode, based on a Bellman-type law of total variance to
 87 analyze the expected sum of the variances of value function estimates. The bound is given by a
 88 function of the estimated transition kernel, whose choice can be optimized by the algorithm. This
 89 leads to a tighter dependence on the horizon in the function estimators.

Table 1: Comparison of Safe RL algorithms for the Hard Constraint Violation Setting: OptCMDP, OptCMDP-bonus (Efroni et al., 2020), AlwaysSafe (Simão et al., 2021), OptPess-LP (Liu et al., 2021), DOPE (Bura et al., 2022), and DOPE+ (Algorithm 1).

Algorithms	Regret	Hard Constraint Violation
OptCMDP, OptCMDP-bonus	$\tilde{O}(H^2\sqrt{S^2AK})$	$\tilde{O}(H^2\sqrt{S^2AK})$
AlwaysSafe	Unknown	0
OptPess-LP	$\tilde{O}((\bar{C} - \bar{C}_b)^{-1}H^3\sqrt{S^3AK})$	0
DOPE	$\tilde{O}((\bar{C} - \bar{C}_b)^{-1}H^3\sqrt{S^2AK})$	0
DOPE+	$\tilde{O}((\bar{C} - \bar{C}_b)^{-1}H^{2.5}\sqrt{S^2AK})$	0

90 A more comprehensive literature review on online reinforcement learning of CMDPs is given in the
 91 supplementary material.

92 2 Preliminary

93 In this section, we introduce the problem setting and necessary definitions. In Section 2.1, we
 94 describe the episodic finite-horizon tabular CMDPs and its performance metrics. In Section 2.2, we
 95 define the confidence set for the transition kernel, and the confidence interval for the reward and cost
 96 functions, which are necessary for deriving our theoretical results.

97 2.1 Problem Setting

98 A finite-horizon tabular MDP is defined by a tuple $(\mathcal{S}, \mathcal{A}, H, \{P_h\}_{h=1}^{H-1}, p)$ where \mathcal{S} is the finite
 99 state space with $|\mathcal{S}| = S$, \mathcal{A} is the finite action space with $|\mathcal{A}| = A$, H is the finite-horizon,
 100 $P_h : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the transition kernel at step $h \in [H - 1]$, and p is the known initial
 101 distribution of the states. Here, $P_h(s' | s, a)$ is the probability of transitioning to state s' from state
 102 s when the chosen action is a at step $h \in [H - 1]$. Equivalently, we may define a single *non-*
 103 *stationary* transition kernel $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H] \rightarrow [0, 1]$ with $P(s' | s, a, h) = P_h(s' | s, a)$ and
 104 $P(s' | s, a, H) = p(s')$ for $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H - 1]$. We assume that $\{P_h\}_{h=1}^{H-1}$ and thus
 105 P are *unknown*.

106 Before an episode begins, the agent prepares a *stochastic policy* $\pi : \mathcal{S} \times [H] \times \mathcal{A} \rightarrow [0, 1]$ where
 107 $\pi(a | s, h)$ is the probability of taking action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$ at step h . Here, π can be viewed
 108 as a *non-stationary policy* as it may change over the horizon, and this is due to the non-stationarity
 109 of P over steps $h \in [H]$. Given a policy π_k for episode $k \in [K]$, the MDP proceeds with trajectory
 110 $\{s_h^{P, \pi_k}, a_h^{P, \pi_k}\}_{h \in [H]}$ generated by P .

111 The reward and cost functions are given by $f, g : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, 1]$, i.e., choosing action
 112 $a \in \mathcal{A}$ at state $s \in \mathcal{S}$ and step $h \in [H]$ generates a reward $f(s, a, h)$ and cost $g(s, a, h)$. Here,
 113 functions f and g are non-stationary over $h \in [H]$. However, the agent observes the noisy reward
 114 and cost. We denote the observed noisy reward and cost for episode $k \in [K]$ by $f_k(s, a, h)$ and
 115 $g_k(s, a, h)$, respectively. As in Liu et al. (2021), we assume that $f_k(s, a, h)$ and $g_k(s, a, h)$ are
 116 determined by independent¹ noisy random variables $\zeta_k^f(s, a, h)$ and $\zeta_k^g(s, a, h)$ following a zero-
 117 mean 1/2-sub-Gaussian distribution, i.e., $f_k(s, a, h) = f(s, a, h) + \zeta_k^f(s, a, h)$ and $g_k(s, a, h) =$

¹We may impose conditional independence.

118 $g(s, a, h) + \zeta_k^g(s, a, h)$. We note that $1/2$ -sub-Gaussian random variables ζ with zero mean satisfies
 119 $\mathbb{E}[\zeta] = 0$ and $\mathbb{E}[\exp(\lambda\zeta)] \leq \exp(\lambda^2/4)$. Then Hoeffding’s inequality implies the following.

Lemma 1. For any $\delta > 0$, with probability at least $1 - 4\delta$, it holds that for any $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$,

$$|f_k(s, a, h)|, |g_k(s, a, h)| \leq 1 + \sqrt{\ln(HSAK/\delta)}.$$

120 We define the value function $V_h^\pi(s; \ell, P)$ at state $s \in \mathcal{S}$ and step $h \in [H]$ for a given policy π ,
 121 function ℓ , and transition kernel P as

$$V_h^\pi(s; \ell, P) = \mathbb{E} \left[\sum_{j=h}^H \ell(s_j^{P,\pi}, a_j^{P,\pi}, j) \mid \ell, \pi, P, s_h^{P,\pi} = s \right].$$

122 Moreover, let $V_1^\pi(\ell, P) = \mathbb{E}_{s \sim p} [V_1^\pi(s; \ell, P) \mid \ell, \pi, P]$ where p is the known distribution of the
 123 initial state.

124 The goal of the constrained Markov decision process is to learn an optimal policy π^* defined as

$$\pi^* \in \operatorname{argmax}_{\pi \in \Pi} V_1^\pi(f, P) \quad \text{s.t.} \quad V_1^\pi(g, P) \leq \bar{C}$$

where \bar{C} is the budget on the expected cost over the horizon, and Π is the set of all policies. As the model parameters f, g, P are unknown, we develop a learning algorithm that computes policies over multiple episodes. For K episodes, we deduce policies π_1, \dots, π_K with the safety requirement that

$$V_1^{\pi_k}(g, P) \leq \bar{C} \quad \forall k \in [K]$$

holds with high probability. The safety requirement is equivalent to enforcing zero hard constraint violation where the hard constraint violation is defined as

$$\text{Violation}(\vec{\pi}) := \sum_{k=1}^K \max \{0, V_1^{\pi_k}(g, P) - \bar{C}\}$$

125 and $\vec{\pi} = (\pi_1, \dots, \pi_K)$ is a shorthand notation for the K policies. As a performance metric for a
 126 learning algorithm, we use the following notion of regret.

$$\text{Regret}(\vec{\pi}) := \sum_{k=1}^K \left(V_1^{\pi^*}(f, P) - V_1^{\pi_k}(f, P) \right).$$

127 To satisfy the safety constraint, we assume that a *strictly safe baseline policy* π_b is given to the agent.

128 **Assumption 1.** The agent knows a policy π_b and its expected cost $\bar{C}_b = V_1^{\pi_b}(g, P)$. We further
 129 assume that π_b is strictly feasible, i.e., $\bar{C}_b < \bar{C}$.

130 This assumption is necessary because the learning agent has no information about the underlying
 131 MDP at the beginning. Without a safe baseline policy, it is difficult to satisfy the constraint in the
 132 initial phase of learning. It is a commonly assumed condition for learning CMDPs (Simão et al.,
 133 2021; Liu et al., 2021; Bura et al., 2022). We also remark that strict feasibility of π_b is related to
 134 Slater’s condition in constrained optimization.

135 Lastly, we assume that the budget \bar{C} satisfies $\bar{C} \in (0, H)$. If $\bar{C} \geq H$, then as $V_1^\pi(g, P) \leq H$ for
 136 any policy π , the safety requirement is trivially satisfied. Moreover, we have \bar{C} is strictly positive
 137 because Assumption 1 imposes that $\bar{C} > \bar{C}_b$ and $\bar{C}_b = V_1^{\pi_b}(g, P) \geq 0$.

138 2.2 Confidence Sets and Intervals

139 We follow the standard Bernstein inequality-based confidence set construction for estimating the
 140 true transition kernel and use confidence intervals based on Hoeffding’s inequality for estimating
 141 reward and cost functions (Jin et al., 2020; Cohen et al., 2020).

142 As in [Efroni et al. \(2020\)](#); [Bura et al. \(2022\)](#), we maintain counters to keep track of the number of
 143 visits to each tuple (s, a, h) and tuple (s, a, s', h) . For each $k \in [K]$, we define $N_k(s, a, h)$ and
 144 $M_k(s, a, s', h)$ as the number of visits to tuple (s, a, h) and the number of visits to tuple (s, a, s', h)
 145 up to the first $k - 1$ episodes, respectively, for $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$. Given $N_k(s, a, h)$
 146 and $M_k(s, a, s', h)$, we define the empirical transition kernel \bar{P}_k for episode k as

$$\bar{P}_k(s' | s, a, h) = \frac{M_k(s, a, s', h)}{\max\{1, N_k(s, a, h)\}}.$$

147 Next, for some confidence parameter $\delta \in (0, 1)$, we define the confidence radius $\epsilon_k(s' | s, a, h)$ for
 148 $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$ and $k \in [K]$ as

$$\epsilon_k(s' | s, a, h) = 2\sqrt{\frac{\bar{P}_k(s' | s, a, h)(1 - \bar{P}_k(s' | s, a, h))L_\delta}{\max\{1, N_k(s, a, h) - 1\}}} + \frac{14 \ln(HSAK/\delta)}{3 \max\{1, N_k(s, a, h) - 1\}} \quad (1)$$

149 where $L_\delta = \ln(HSAK/\delta)$. Based on the empirical transition kernel and the radius, we define the
 150 confidence set \mathcal{P}_k for episode k as

$$\mathcal{P}_k = \left\{ \hat{P} : \left| \hat{P}(s' | s, a, h) - \bar{P}_k(s' | s, a, h) \right| \leq \epsilon_k(s' | s, a, h) \quad \forall (s, a, s', h) \right\}. \quad (2)$$

151 By the empirical Bernstein inequality due to [Maurer & Pontil \(2009\)](#), we can show the following.

152 **Lemma 2.** *For any $\delta > 0$, with probability at least $1 - 4\delta$, the true transition kernel P is contained*
 153 *in the confidence set \mathcal{P}_k for every episode $k \in [K]$.*

Next, for reward and cost functions, we define the confidence radius $R_k(s, a, h)$ for $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, $k \in [K]$ and $\delta \in (0, 1)$ as

$$R_k(s, a, h) = \sqrt{\frac{\ln(HSAK/\delta)}{\max\{1, N_k(s, a, h)\}}}.$$

154 We define empirical estimators \bar{f}_k and \bar{g}_k as

$$\bar{f}_k(s, a, h) = \frac{\sum_{j=1}^{k-1} f_j(s, a, h)n_j(s, a, h)}{\max\{1, N_k(s, a, h)\}}, \quad \bar{g}_k(s, a, h) = \frac{\sum_{j=1}^{k-1} g_j(s, a, h)n_j(s, a, h)}{\max\{1, N_k(s, a, h)\}}$$

155 where $f_j(s, a, h)$, $g_j(s, a, h)$ are the instantaneous reward and cost for episode $j \in [k - 1]$ and
 156 $n_j(s, a, h)$ is the indicator variable that returns 1 if the agent visited (s, a, h) in episode j and 0
 157 otherwise. Then we may deduce the following from Hoeffding's inequality.

158 **Lemma 3.** *For any $\delta > 0$, with probability at least $1 - 4\delta$, it holds that for any $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$*
 159 *and $k \in [K]$,*

$$\left| \bar{f}_k(s, a, h) - f(s, a, h) \right| \leq R_k(s, a, h), \quad \left| \bar{g}_k(s, a, h) - g(s, a, h) \right| \leq R_k(s, a, h).$$

160 3 Tighter Function Estimators

161 In this section, we introduce the tighter function estimators, which are crucial for achieving our
 162 theoretical results: (i) zero constraint violation and (ii) an improved regret upper bound. First, we
 163 show how to design the tighter pessimistic cost estimator \hat{g}_k , focusing on zero constraint violation.
 164 Accordingly, we present the reward estimator \hat{f}_k with an extra optimism to compensate for the
 165 pessimism of \hat{g}_k , which directly affects the regret upper bound.

166 **Remark 1.** The reason why we begin with designing \hat{g}_k is that a tighter \hat{g}_k can be translated to
 167 a tighter regret upper bound. To provide an intuition, let us consider the following optimization
 168 problem based on the estimated MDP: $\max_{\pi', P'} V_1^{\pi'}(\hat{f}_k, P')$ s.t. $V_1^{\pi'}(\hat{g}_k, P') \leq \bar{C}$. Once we
 169 take a tighter \hat{g}_k , the set of feasible solutions becomes larger. Then it leads to increase the optimal
 170 value $V_1^{\pi_k}(\hat{f}_k, P_k)$, where (π_k, P_k) is an optimal solution. Taking advantage of this, it allows us to
 171 have a tighter optimism for \hat{f}_k , which directly affects the regret upper bound. \square

172 Lemmas 2 and 3 motivate the following attempt to deduce feasible policies. For episode $k \in [K]$, we
 173 take a transition kernel P_k from the confidence set \mathcal{P}_k and $\bar{g}_k + R_k$ as a pessimistic (or conservative)
 174 estimator of the cost function g . Then we may compute a policy π_k that satisfies $V_1^{\pi_k}(\bar{g}_k + R_k, P_k) \leq$
 175 \bar{C} , which is an approximation of the constraint. However, even if $\bar{g}_k + R_k$ provides an upper bound
 176 on g , the issue is that $V_1^{\pi_k}(g, P) \not\leq V_1^{\pi_k}(\bar{g}_k + R_k, P_k)$. This is because the difference between the
 177 true transition kernel P and P_k can make $V_1^{\pi_k}(g, P)$ greater than $V_1^{\pi_k}(\bar{g}_k + R_k, P_k)$. That said, π_k
 178 does not necessarily satisfy the constraint, although it satisfies the approximate constraint.

179 Inspired by the challenge, the next question is as to whether we can design an approximate con-
 180 straint, satisfying which guarantees that the true constraint is also satisfied. Liu et al. (2021); Bura
 181 et al. (2022) considered this, and their idea was to add an extra pessimism to cost function estimators.
 182 Basically, we take functions of the form

$$\hat{g}_k(s, a, h) = \bar{g}_k(s, a, h) + R_k(s, a, h) + U_k(s, a, h) \quad (3)$$

183 for $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$ where U_k captures the error in estimating the true
 184 transition kernel P . In the above-discussed context, U_k considers the difference between P and
 185 P_k . Here, one needs to set U_k sufficiently large so that $V_1^{\pi_k}(g, P) \leq V_1^{\pi_k}(\hat{g}_k, P_k)$, in which case
 186 satisfying the corresponding approximate constraint $V_1^{\pi_k}(\hat{g}_k, P_k) \leq \bar{C}$ guarantees satisfaction of
 187 the true constraint.

188 On the other hand, choosing the right magnitude of U_k is important to control the regret function.
 189 When U_k is too large, \hat{g}_k is too conservative, and it prevents from getting a high reward. Indeed,
 190 Bura et al. (2022) improved upon Liu et al. (2021) by making U_k tighter. Our main contribution is
 191 to develop an even tighter U_k function than Bura et al. (2022).

Before we present our design of U_k , let us briefly discuss how to deduce the extra pessimism term U_k
 in general. As explained before, we want to guarantee $V_1^{\pi_k}(g, P) \leq V_1^{\pi_k}(\hat{g}_k, P_k)$ for any $P_k \in \mathcal{P}_k$.
 Then note that

$$V_1^{\pi_k}(g, P) \leq V_1^{\pi_k}(g, P_k) + |V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)|.$$

If the statement of Lemma 3 holds, then $V_1^{\pi_k}(g, P_k)$ is bounded above by $V_1^{\pi_k}(\bar{g}_k + R_k, P_k)$. There-
 fore, once we come up with some U_k such that $|V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)| \leq V_1^{\pi_k}(U_k, P_k)$, we get

$$V_1^{\pi_k}(g, P) \leq V_1^{\pi_k}(\bar{g}_k + R_k + U_k, P_k).$$

192 In this case, $\hat{g}_k = \bar{g}_k + R_k + U_k$ gives rise to a valid function estimator.

193 We devise our pessimism function U_k as follows.

194 **Theorem 1.** *Let π_k be any policy for episode k . Take*

$$U_k(s, a, h) = 8\sqrt{H}\varepsilon_k(s, a, h) + 4S\sqrt{HA/K} + \frac{2\ln(HSAK/\delta)\sqrt{HK/A} + \eta}{\max\{1, N_k(s, a, h) - 1\}} \quad (4)$$

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

$$\varepsilon_k(s, a, h) = 2\sqrt{\frac{S\ln(HSAK/\delta)}{\max\{1, N_k(s, a, h) - 1\}}} + \frac{14S\ln(HSAK/\delta)}{3\max\{1, N_k(s, a, h) - 1\}} \quad (5)$$

and $\eta = (19HS + 2H^{1.5}S + 10^4H^2S^2)\ln(HSAK/\delta)^2$. Then for any $\delta > 0$, it holds with proba-
 bility at least $1 - 14\delta$ that

$$|V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)| \leq V_1^{\pi_k}(U_k, P_k)$$

196 for any $P_k \in \mathcal{P}_k$ and $g : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, 1]$.

197 In the following remark, we demonstrate that our U_k indeed improves upon Bura et al. (2022).

198 **Remark 2.** [Bura et al. \(2022\)](#) set $U_k(s, a, h)$ as $2H\varepsilon_k(s, a, h)$, which has coefficient $2H$ in front of
 199 ε_k^2 . In contrast, our construction in [Theorem 1](#) has an improved coefficient of $8\sqrt{H}$. Although we
 200 have additional terms for U_k , the reduction of $\mathcal{O}(\sqrt{H})$ in the coefficient translates to the improve-
 201 ment of $\tilde{\mathcal{O}}(\sqrt{H})$ factor in the regret upper bound. \square

202 Next, we present our optimistic reward function estimator \hat{f}_k . We define the optimistic reward
 203 function estimator \hat{f}_k as

$$\hat{f}_k(s, a, h) = \min \left\{ B, \bar{f}_k(s, a, h) + \frac{3H}{\bar{C} - \bar{C}_b} R_k(s, a, h) + \frac{H}{\bar{C} - \bar{C}_b} U_k(s, a, h) \right\} \quad (6)$$

204 where $B = 1 + \sqrt{\ln(HSAK/\delta)}$. On top of $\bar{f}_k + R_k$, we take an additional optimistic term U_k
 205 for the reward function to compensate for U_k in \hat{g}_k , which reduces the search space of policies and
 206 hinders exploration. Furthermore, in \hat{f}_k , we multiply R_k and U_k by $\mathcal{O}(H/(\bar{C} - \bar{C}_b))$ to guarantee
 207 the extra optimism in \hat{f}_k truly promotes exploration. Nevertheless, taking the extra optimism can
 208 cause a substantial overestimation of the reward function. To avoid this, we take a truncation to B
 209 as in (6).

210 3.1 Proof Outline of [Theorem 1](#)

211 The value difference lemma ([Dann et al., 2017](#)) implies

$$V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k) = \mathbb{E} \left[\sum_{h=1}^H \ell(s_h^{P_k, \pi_k}, a_h^{P_k, \pi_k}, h) \mid \pi_k, P_k \right]$$

212 where $\ell(s, a, h)$ is given by

$$\sum_{s' \in \mathcal{S}} (P - P_k)(s' \mid s, a, h) V_{h+1}^{\pi_k}(s'; g, P) \quad (7)$$

213 with $V_{H+1}^{\pi_k} = 0$ and $(P - P_k)(s' \mid s, a, h) = P(s' \mid s, a, h) - P_k(s' \mid s, a, h)$. Here, [Bura et al.](#)
 214 (2022) used that $V_{h+1}^{\pi_k} \leq H$ and $|P - P_k| \leq |P - \bar{P}_k| + |\bar{P}_k - P_k| \leq 2\varepsilon_k$ by [Lemma 2](#). Then it
 215 follows that

$$|V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)| \leq \mathbb{E} \left[\sum_{h=1}^H 2H \sum_{s' \in \mathcal{S}} \varepsilon_k(s' \mid s_h^{P_k, \pi_k}, a_h^{P_k, \pi_k}, h) \mid \varepsilon_k, \pi_k, P_k \right]$$

216 whose right-hand side equals $V_1^{\pi_k}(U_k, P_k)$ where U_k is given by $2H\varepsilon_k$. This explains how [Bura](#)
 217 [et al. \(2022\)](#) deduced their pessimistic cost estimators.

218 To prove [Theorem 1](#) that establishes the validity of our choice of tighter U_k in (4), we need a more
 219 refined analysis of the difference term $|V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)|$. Note that $\ell(s, a, h)$ in (7) satisfies

$$|\ell(s, a, h)| \leq \underbrace{\left| \sum_{s' \in \mathcal{S}} (P - P_k)(s' \mid s, a, h) W_{h+1}^{\pi_k}(s'; g) \right|}_{I_1} + \underbrace{\left| \sum_{s' \in \mathcal{S}} (P - P_k)(s' \mid s, a, h) V_{h+1}^{\pi_k}(s'; g, P_k) \right|}_{I_2}$$

220 where $W_{h+1}^{\pi_k}(s'; g) = V_{h+1}^{\pi_k}(s'; g, P) - V_{h+1}^{\pi_k}(s'; g, P_k)$. We prove the following lemma to provide
 221 an upper bound on term I_1 .

²In fact, the original choice of [Bura et al. \(2022\)](#) was $U_k(s, a, h) = 2H \sum_{s' \in \mathcal{S}} \varepsilon_k(s' \mid s, a, h)$ where $\varepsilon_k(s' \mid s, a, h)$ is given in (1), but there is an issue with this choice. We need the property that U_k is nonincreasing in k to show [Lemma 6](#) and ([Proposition 4, Bura et al., 2022](#)), but their U_k can increase as $\bar{P}_k(s' \mid s, a, h)/N_k(s, a, h)$ can increase. As a fix, we may take $U_k(s, a, h) = 2H\varepsilon_k(s, a, h)$ where ε_k is given in (5). Note that ε_k is nonincreasing in k . At the same time, by the Cauchy-Schwarz inequality, $\varepsilon_k(s, a, h)$ is an upper bound on $\sum_{s' \in \mathcal{S}} \varepsilon_k(s' \mid s, a, h)$. As a result, our construction resolves the issue of [Bura et al. \(2022\)](#).

222 **Lemma 4.** Let π_k be any policy for episode $k \in [K]$, and let $g : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, 1]$ be an
 223 arbitrary cost function. Then for any $P, P_k \in \mathcal{P}_k$, we have

$$\mathbb{E} \left[\left| \sum_{s' \in \mathcal{S}} (P - P_k)(s' | s, a, h) (V_{h+1}^{\pi_k}(s'; g, P) - V_{h+1}^{\pi_k}(s'; g, P_k)) \right| \mid \pi_k, P_k \right] \leq V_1^{\pi_k}(U_{k,1}, P_k)$$

224 where

$$U_{k,1}(s, a, h) = \frac{10^4 H^2 S^2 \ln(HSAK/\delta)^2}{\max\{1, N_k(s, a, h)\}}.$$

225 The proof of this lemma is based on the value difference lemma to evaluate $V_{h+1}^{\pi_k}(s'; g, P) -$
 226 $V_{h+1}^{\pi_k}(s'; g, P_k)$. Here, the key part is to provide an upper bound that is represented as a value
 227 function of π_k and P_k . Hence, we have

$$\mathbb{E}[I_1 \mid \pi_k, P_k] \leq V_1^{\pi_k}(U_{k,1}, P_k).$$

228 Next, we consider term I_2 , which turns out to be the dominant one. Since P and P_k both define
 229 transition functions, I_2 equals

$$\left| \sum_{s' \in \mathcal{S}} (P - P_k)(s' | s, a, h) (V_{h+1}^{\pi_k}(s'; g, P_k) - \hat{\mu}_k(s, a, h)) \right|$$

where $\hat{\mu}_k(s, a, h) = \mathbb{E}_{s' \sim P_k(\cdot | s, a, h)} [V_{h+1}^{\pi_k}(s'; g, P_k)]$. Next, we observe that
 $|(P - P_k)(s' | s, a, h)| \leq 2\epsilon_k(s' | s, a, h)$ due to Lemma 2. Recall that $\epsilon_k(s' | s, a, h)$
 contains the term $\sqrt{P_k(s' | s, a, h)}$. As $P_k \in \mathcal{P}_k$ we deduce that $\sqrt{P_k(s' | s, a, h)} \leq$
 $\sqrt{P_k(s' | s, a, h)} + \epsilon_k(s' | s, a, h)$. As a result, by the Cauchy-Schwarz inequality, the anal-
 ysis boils down to providing an upper bound on the term

$$\sum_{s' \in \mathcal{S}} P_k(s' | s, a, h) (V_{h+1}^{\pi_k}(s'; g, P_k) - \hat{\mu}_k(s, a, h))^2,$$

which equals

$$\hat{\mathbb{V}}_k(s, a, h) := \text{Var}_{s' \sim P_k(\cdot | s, a, h)} [V_{h+1}^{\pi_k}(s'; g, P_k)].$$

230 Furthermore, our proof reveals that $V_1^{\pi_k}(\hat{\mathbb{V}}_k, P_k)$ is the important quantity to control. Applying a
 231 naïve upper bound on value functions gives $\hat{\mathbb{V}}_k \leq H^2$ and thus $V_1^{\pi_k}(\hat{\mathbb{V}}_k, P_k) \leq H^3$. However, this
 232 bound is not tight enough. Instead, we prove the following lemma based on a Bellman-type law of
 233 total variance (Azar et al., 2017; Chen & Luo, 2021).

234 **Lemma 5.** Let π_k be a policy for episode k . Then

$$V_1^{\pi_k}(\hat{\mathbb{V}}_k, P_k) \leq 2H^2$$

235 for any $P_k \in \mathcal{P}_k$ and $g : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, 1]$.

236 This improvement in the variance term leads to

$$\mathbb{E}[I_2 \mid \pi_k, P_k] \leq V_1^{\pi_k}(U_{k,2}, P_k)$$

237 where

$$U_{k,2}(s, a, h) = 8\sqrt{H}\epsilon_k(s, a, h) + 4S\sqrt{HA/K} + \frac{2L\sqrt{HK/A} + (19HS + 2H^{1.5}S)L_\delta^2}{\max\{1, N_k(s, a, h) - 1\}}$$

238 where $L_\delta = \ln(HSAK/\delta)$. Putting the pieces together, we complete the proof of Theorem 1, as we
 239 have $U_k(s, a, h) = U_{k,1}(s, a, h) + U_{k,2}(s, a, h)$. A complete proof is given in the supplementary
 240 material.

241 **3.2 Comparison with Previous Works**

242 Our main technical contribution is to provide a tighter upper bound on the term

$$|V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)| \quad (8)$$

243 over each episode $k \in [K]$. This improves upon the analysis of [Bura et al. \(2022\)](#), thereby providing
 244 tighter cost and reward function estimators. Recall that our upper bound given in Theorem 1 is in the
 245 form of $V_1^{\pi_k}(U_k, P_k)$ and the main technique is a Bellman-type law of total variance. While [Chen](#)
 246 [& Luo \(2021\)](#) applied a similar technique to control the error of estimating the unknown transition
 247 kernel, their result does not immediately translate to a proper function estimator for our setting. We
 248 elaborate on this below.

249 [Chen & Luo \(2021\)](#) gave an upper bound on the cumulative error given by

$$\sum_{k=1}^K |V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)| \leq C_1 \sum_{k=1}^K V_1^{\pi_k}(U_k, P) + C_2 \quad (9)$$

250 where $C_1 = 16\lambda S^2 A$, $C_2 = C_1 \tilde{\mathcal{O}}(H^3 \sqrt{K}) + 16 \ln^2(HSAK/\delta)/\lambda + \tilde{\mathcal{O}}(H^3 S^2 A)$ for any $\lambda > 0$,
 251 and $U_k = Hg$. However, the bound on the cumulative error does not lead to an upper bound on the
 252 error term (8) for each episode. Recall that to define \hat{f}_k, \hat{g}_k for each k , we need an upper bound on
 253 (8). Furthermore, the bound in (9) is written as a function of the true transition kernel P , which is
 254 not known to the agent. However, our algorithm as well as DOPE due to [Bura et al. \(2022\)](#) chooses
 255 an optimistic transition kernel, we require an upper bound on (8) that depends on the optimistic
 256 transition kernel to estimate the error caused by the choice.

257 Theorem 1 addresses these issues by providing an upper bound for (8) in the form of $V_1^{\pi_k}(U_k, P_k)$,
 258 thereby leading to our novel reward and cost function estimators \hat{f}_k, \hat{g}_k .

259 **4 Algorithm**

260 DOPE+, given by Algorithm 1, is a variant of DOPE by [Bura et al. \(2022\)](#) with our novel reward
 261 and cost function estimators from Section 3. Recall that our pessimistic cost estimator \hat{g}_k is given
 262 by (3) with the extra pessimism term U_k given in (4) and our optimistic reward estimator \hat{f}_k is given
 263 in (6).

264 As in [Efroni et al. \(2020\)](#); [Bura et al. \(2022\)](#), we compute our policy π_k for episode $k \in [K]$ by
 265 solving the following optimization problem.

$$(\pi_k, P_k) \in \operatorname{argmax}_{(\pi, Q) \in \Pi \times \mathcal{P}_k} \left\{ V_1^\pi(\hat{f}_k, Q) : V_1^\pi(\hat{g}_k, Q) \leq \bar{C} \right\} \quad (10)$$

266 where \mathcal{P}_k is the confidence set given by (2) and Π is the set of valid policies.

267 To solve (10) efficiently, we take the standard approach of using *occupancy measures* ([Altman,](#)
 268 [1999](#)). An occupancy measure is essentially a joint probability for the event that we observe the
 269 state-action pair (s, a) at step h and state s' at step $h + 1$. Introducing occupancy measure, we can
 270 reformulate (10) as an linear program in terms of an occupancy measure, which is referred to as the
 271 extended linear program ([Altman, 1999](#); [Efroni et al., 2020](#); [Bura et al., 2022](#)). By solving it, we
 272 obtain an optimal occupancy measure inducing an optimal solution to (10). We defer the formal
 273 description of the extended linear program to the supplementary material.

274 One issue, however, is that (10) can be infeasible at the beginning of the algorithm as \hat{g}_k can be
 275 too large to guarantee feasibility of (10). Hence, the algorithm executes the safe baseline policy π_b
 276 for the first few episodes until sufficient information is gathered so that (10) becomes feasible. The
 277 following lemma characterizes a sufficient number of episodes running the safe baseline policy to
 278 guarantee feasibility of (10).

Algorithm 1 Doubly Optimistic Pessimistic Exploration with Tighter Function Estimators (DOPE+)

Input: Safe baseline policy π_b and its expected cost for a single episode \bar{C}_b , and the number K_0 of episodes for the initial phase
Initialize: $N(s, a, h) = M(s, a, s', h) \leftarrow 0$ for $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$.
for $k = 1, \dots, K$ **do**
 Set counters $N_k \leftarrow N$ and $M_k \leftarrow M$.
 Compute \bar{P}_k, ϵ_k , and \mathcal{P}_k (Section 2.2).
 if $k \leq K_0$ **then**
 Set $\pi_k = \pi_b$.
 else
 Compute estimators \hat{f}_k and \hat{g}_k (Section 3).
 Deduce π_k, P_k from (10).
 end if
 Sample state s_1 from distribution p .
 for $h = 1, \dots, H$ **do**
 Sample a_h from $\pi_k(\cdot \mid s_h, h)$.
 Observe $f_k(s_h, a_h, h), g_k(s_h, a_h, h)$, and s_{h+1} determined by $P(\cdot \mid s_h, a_h, h)$.
 Update the counters N, M .
 end for
end for

279 **Lemma 6.** With probability at least $1 - 14\delta$, (π_b, P) is a feasible solution of (10) for any $k > K_0$
280 where

$$K_0 = \tilde{\mathcal{O}} \left(\frac{H^3 S^2 A}{(\bar{C} - \bar{C}_b)^2} \right) \quad (11)$$

281 where $\tilde{\mathcal{O}}(\cdot)$ hides factors polynomial in $\ln(HSAK/\delta)$.

282 5 Regret Analysis of DOPE+

283 Let us state our theoretical guarantees for DOPE+.

Theorem 2. Let $\vec{\pi} = (\pi_1, \dots, \pi_K)$ denote policies computed by DOPE+ with K_0 given in (11).
Then

$$\text{Violation}(\vec{\pi}) = 0$$

284 with probability at least $1 - 14\delta$.

285 Hence, DOPE+ achieves no constraint violation. The next theorem shows a regret upper bound for
286 DOPE+.

287 **Theorem 3.** Let $\vec{\pi} = (\pi_1, \dots, \pi_K)$ denote policies computed by DOPE+ with K_0 given in (11).
288 Then, with probability at least $1 - 16\delta$, we have

$$\text{Regret}(\vec{\pi}) = \tilde{\mathcal{O}} \left(\frac{H}{\bar{C} - \bar{C}_b} \left(H^{1.5} S \sqrt{AK} + \frac{H^4 S^3 A}{\bar{C} - \bar{C}_b} \right) \right)$$

289 where $\tilde{\mathcal{O}}(\cdot)$ hides factors polynomial in $\ln(HSAK/\delta)$.

290 **Remark 3.** Note that there is a gap of $\tilde{\mathcal{O}}((\bar{C} - \bar{C}_b)^{-1} H \sqrt{S})$ factor between our regret upper bound
291 and the lower bound $\Omega(H^{3/2} \sqrt{SAK})$ due to Jin et al. (2020); Domingues et al. (2021). In fact, the
292 instance from Domingues et al. (2021) is an unconstrained MDP. We observe that the $\mathcal{O}(H/(\bar{C} -$
293 $\bar{C}_b))$ factor in our regret upper bound is due to the constraint, which becomes a constant if $\bar{C} - \bar{C}_b =$
294 $\Omega(H)$. Hence, our regret upper bound nearly matches the regret lower bound in terms of H when
295 $\bar{C} - \bar{C}_b = \Omega(H)$. \square

296 **5.1 Constraint Violation Analysis**

297 We prove Theorem 2 as follows. For episode k with $k \leq K_0$, DOPE+ takes the safe baseline policy
 298 π_b , so no constraint violation is guaranteed. Then let us consider episode k with $k > K_0$. As
 299 explained in Section 3, we argue that

$$\begin{aligned} V_1^{\pi_k}(g, P) &\leq V_1^{\pi_k}(g, P_k) + |V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)| \\ &\leq V_1^{\pi_k}(\bar{g}_k + R_k, P_k) + V_1^{\pi_k}(U_k, P_k) \\ &= V_1^{\pi_k}(\hat{g}_k, P_k) \end{aligned}$$

300 where the second inequality is due to Lemma 3 and Theorem 1. Since (π_k, P_k) is a solution to (10),
 301 it holds that $V_1^{\pi_k}(\hat{g}_k, P_k) \leq \bar{C}$. Therefore, it follows that $V_1^{\pi_k}(g, P) \leq \bar{C}$ and thus the constraint is
 302 satisfied.

 303 **5.2 Regret Decomposition**

304 We provide an overview of the proof of Theorem 3. Since we execute the safe baseline policy π_b for
 305 the first K_0 episodes, we decompose the regret function as follows.

$$\begin{aligned} \text{Regret}(\bar{\pi}) &= \underbrace{\sum_{k=1}^{K_0} \left(V_1^{\pi^*}(f, P) - V_1^{\pi_b}(f, P) \right)}_{\text{(I)}} + \underbrace{\sum_{k=K_0+1}^K \left(V_1^{\pi^*}(f, P) - V_1^{\pi_k}(\hat{f}_k, P_k) \right)}_{\text{(II)}} \\ &\quad + \underbrace{\sum_{k=K_0+1}^K \left(V_1^{\pi_k}(\hat{f}_k, P_k) - V_1^{\pi_k}(\hat{f}_k, P) \right)}_{\text{(III)}} + \underbrace{\sum_{k=K_0+1}^K \left(V_1^{\pi_k}(\hat{f}_k, P) - V_1^{\pi_k}(f, P) \right)}_{\text{(IV)}}. \end{aligned}$$

306 Term (I) is due to executing π_b for K_0 episodes for feasibility. By Lemma 6, term (I) can be bounded
 307 by $\tilde{\mathcal{O}}((\bar{C} - \bar{C}_b)^{-2}(H^4 S^2 A))$ as $V_1^{\pi} \leq H$ for any policy π .

308 For term (II), we provide the following upper bound.

309 **Lemma 7.** *With probability at least $1 - 14\delta$,*

$$\sum_{k=K_0+1}^K \left(V_1^{\pi^*}(f, P) - V_1^{\pi_k}(\hat{f}_k, P_k) \right) = \tilde{\mathcal{O}} \left(\frac{H}{\bar{C} - \bar{C}_b} \left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) \right)$$

310 where $\tilde{\mathcal{O}}(\cdot)$ hides factor polynomial in $\ln(HSAK/\delta)$.

311 To prove the lemma, we define a new policy $\pi_k^{\alpha_k}$ for $k \in [K]$, which is induced by a con-
 312 vex combination of the occupancy measures associated with (π^*, P) and (π_b, P) with coeffi-
 313 cients $\alpha_k, 1 - \alpha_k \in (0, 1)$. We choose the value of α_k so that $(\pi_k^{\alpha_k}, P)$ is feasible to (10).

314 Then the optimality of (π_k, P_k) implies $V_1^{\pi_k^{\alpha_k}}(\hat{f}_k, P) \leq V_1^{\pi_k}(\hat{f}_k, P_k)$, which lets us to analyze
 315 $V_1^{\pi^*}(f, P) - V_1^{\pi_k^{\alpha_k}}(\hat{f}_k, P)$ with the same transition kernel P .

316 Term (III) comes from learning the unknown transition kernel. We apply a Bellman-type law of total
 317 variance to provide an upper bound on term (III).

318 **Lemma 8.** *With probability at least $1 - 16\delta$,*

$$\sum_{k=K_0+1}^K \left(V_1^{\pi_k}(\hat{f}_k, P_k) - V_1^{\pi_k}(\hat{f}_k, P) \right) = \tilde{\mathcal{O}} \left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right)$$

319 where $\tilde{\mathcal{O}}(\cdot)$ hides factor polynomial in $\ln(HSAK/\delta)$.

320 Term (IV) is due to the difference between f and our estimator \hat{f}_k .

321 **Lemma 9.** *With probability at least $1 - 14\delta$,*

$$\sum_{k=K_0+1}^K \left(V_1^{\pi_k}(\hat{f}_k, P) - V_1^{\pi_k}(f, P) \right) = \tilde{O} \left(\frac{H}{\bar{C} - \bar{C}_b} \left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) \right)$$

322 where $\tilde{O}(\cdot)$ hides factor polynomial in $\ln(HSAK/\delta)$.

323 6 Numerical Experiment

324 We evaluate DOPE+ on the three-state CMDP instance of [Zheng & Ratliff \(2020\)](#); [Simão et al. \(2021\)](#); [Bura et al. \(2022\)](#) with a few modifications. In Figure 1, we compare regret and constraint violation under DOPE+ and DOPE for 200,000 episodes when $H = 30$. We consider DOPE as a benchmark algorithm because it provides the best regret bound among the existing algorithms while ensuring zero constraint violation. Our results are averaged across 5 runs with different random seeds, and we display the 95% confidence interval with shaded regions. More details of the experiment setup can be found in the supplementary material including the MDP instance and algorithm parameters.

332 In Figure 1a, DOPE+ outperforms DOPE in terms of regret. This result demonstrates that DOPE+ improves upon DOPE computationally, in addition to our theoretical improvement. Figure 1b shows that both algorithms achieve zero constraint violation.

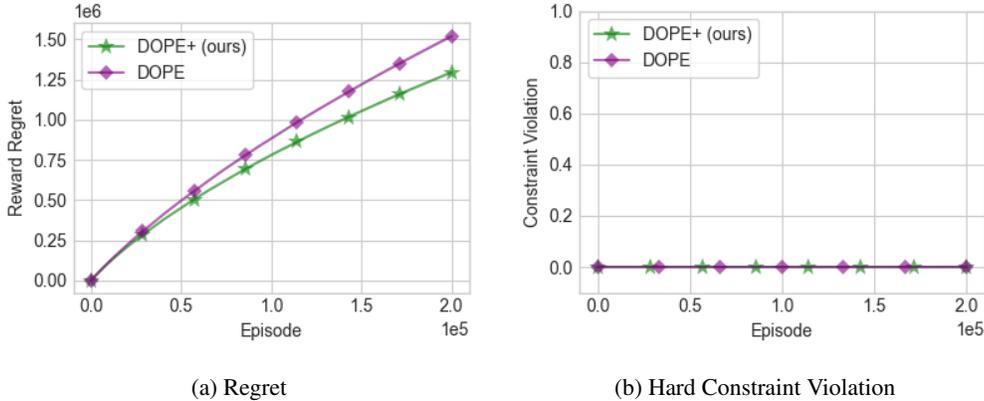


Figure 1: Comparison of DOPE+ and DOPE

335 7 Conclusion

336 In this paper, we investigate safe RL formulated as an episodic finite-horizon tabular CMDP. We
 337 propose novel reward and cost function estimators with tighter reward optimism and cost pessimism.
 338 Based on them, we develop DOPE+, which is a variant of DOPE due to [\(Bura et al., 2022\)](#). We prove
 339 that DOPE+ achieves regret upper bound $\tilde{O}((\bar{C} - \bar{C}_b)^{-1} H^{2.5} S \sqrt{AK})$ and zero hard constraint
 340 violation. The regret upper bound improves upon the best-known bound by a multiplicative factor of
 341 $\tilde{O}(\sqrt{H})$ factor. When $\bar{C} - \bar{C}_b = \Omega(H)$, the gap from the regret lower bound of $\Omega(H^{1.5} \sqrt{SAK})$ ([Jin](#)
 342 [et al., 2020](#); [Domingues et al., 2021](#)) is $\tilde{O}(\sqrt{S})$, and we would like to leave closing this gap as an
 343 open question in the zero hard constraint violation setting. We also present numerical results that
 344 demonstrate the computational effectiveness of DOPE+ compared to DOPE.

345 **References**

- 346 Naoki Abe, Prem Melville, Cezar Pendus, Chandan K. Reddy, David L. Jensen, Vince P. Thomas,
347 James J. Bennett, Gary F. Anderson, Brent R. Cooley, Melissa Kowalczyk, Mark Domick, and
348 Timothy Gardinier. Optimizing debt collections using constrained reinforcement learning. In
349 *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and*
350 *Data Mining*, KDD '10, pp. 75–84, New York, NY, USA, 2010. Association for Computing
351 Machinery. ISBN 9781450300551. DOI: 10.1145/1835804.1835817. URL <https://doi.org/10.1145/1835804.1835817>.
- 353 Eitan Altman. *Constrained Markov Decision Processes*, volume 7. CRC Press, 1999.
- 354 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforce-
355 ment learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International*
356 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*,
357 pp. 263–272. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/v70/](https://proceedings.mlr.press/v70/azar17a.html)
358 [azar17a.html](https://proceedings.mlr.press/v70/azar17a.html).
- 359 Alina Beygelzimer, John Langford, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual ban-
360 ddit algorithms with supervised learning guarantees. In Geoffrey Gordon, David Dunson, and
361 Miroslav Dudík (eds.), *Proceedings of the Fourteenth International Conference on Artificial In-*
362 *telligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pp. 19–26,
363 Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v15/beygelzimer11a.html)
364 [press/v15/beygelzimer11a.html](https://proceedings.mlr.press/v15/beygelzimer11a.html).
- 365 Kianté Brantley, Miroslav Dudík, Thodoris Lykouris, Sobhan Miryoosefi, Max Simchowitz, Alek-
366 sandrs Slivkins, and Wen Sun. Constrained episodic reinforcement learning in concave-convex
367 and knapsack settings. In *Proceedings of the 34th International Conference on Neural Informa-*
368 *tion Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN
369 9781713829546.
- 370 Archana Bura, Aria Hasanzadezonuzy, Dileep Kalathil, Srinivas Shakkottai, and Jean-Francois
371 Chamberland. DOPE: Doubly optimistic and pessimistic exploration for safe reinforcement learn-
372 ing. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in*
373 *Neural Information Processing Systems*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=U4BUMoVTrB2)
374 [id=U4BUMoVTrB2](https://openreview.net/forum?id=U4BUMoVTrB2).
- 375 Liyu Chen and Haipeng Luo. Finding the stochastic shortest path with low regret: the adversarial
376 cost and unknown transition case. In Marina Meila and Tong Zhang (eds.), *Proceedings of the*
377 *38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*
378 *Learning Research*, pp. 1651–1660. PMLR, 18–24 Jul 2021. URL [https://proceedings.](https://proceedings.mlr.press/v139/chen21l.html)
379 [mlr.press/v139/chen21l.html](https://proceedings.mlr.press/v139/chen21l.html).
- 380 Yi Chen, Jing Dong, and Zhaoran Wang. A primal-dual approach to constrained markov decision
381 processes, 2021. URL <https://arxiv.org/abs/2101.10895>.
- 382 Alon Cohen, Haim Kaplan, Yishay Mansour, and Aviv Rosenberg. Near-optimal regret bounds
383 for stochastic shortest path. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*
384 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*
385 *Research*, pp. 8210–8219. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v119/rosenberg20a.html)
386 [press/v119/rosenberg20a.html](https://proceedings.mlr.press/v119/rosenberg20a.html).
- 387 Antonio Coronato, Muddasar Naeem, Giuseppe De Pietro, and Giovanni Paragliola. Rein-
388 forcement learning for intelligent healthcare applications: A survey. *Artificial Intelligence*
389 *in Medicine*, 109:101964, 2020. ISSN 0933-3657. DOI: [https://doi.org/10.1016/j.artmed.](https://doi.org/10.1016/j.artmed.2020.101964)
390 [2020.101964](https://doi.org/10.1016/j.artmed.2020.101964). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S093336572031229X)
391 [S093336572031229X](https://www.sciencedirect.com/science/article/pii/S093336572031229X).

- 392 Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uni-
393 form pac bounds for episodic reinforcement learning. In I. Guyon, U. Von Luxburg,
394 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-
395 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
396 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
397 file/17d8da815fa21c57af9829fb0a869602-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/17d8da815fa21c57af9829fb0a869602-Paper.pdf).
- 398 Dongsheng Ding, Xiaohan Wei, Zhuoran Yang, Zhaoran Wang, and Mihailo Jovanovic. Provably
399 efficient safe exploration via primal-dual policy optimization. In *International conference on
400 artificial intelligence and statistics*, pp. 3304–3312. PMLR, 2021.
- 401 Omar Darwiche Domingues, Pierre Ménard, Emilie Kaufmann, and Michal Valko. Episodic re-
402 inforcement learning in finite mdps: Minimax lower bounds revisited. In Vitaly Feldman, Ka-
403 trina Ligett, and Sivan Sabato (eds.), *Proceedings of the 32nd International Conference on Al-
404 gorithmic Learning Theory*, volume 132 of *Proceedings of Machine Learning Research*, pp.
405 578–598. PMLR, 16–19 Mar 2021. URL [https://proceedings.mlr.press/v132/
406 domingues21a.html](https://proceedings.mlr.press/v132/domingues21a.html).
- 407 Yonathan Efroni, Shie Mannor, and Matteo Pirota. Exploration-exploitation in constrained mdps,
408 2020.
- 409 Jaime F Fisac, Anayo K Akametalu, Melanie N Zeilinger, Shahab Kaynama, Jeremy Gillula, and
410 Claire J Tomlin. A general safety framework for learning-based control in uncertain robotic
411 systems. *IEEE Transactions on Automatic Control*, 64(7):2737–2752, 2018.
- 412 Javier García and Diogo Shafie. Teaching a humanoid robot to walk faster through safe reinforce-
413 ment learning. *Engineering Applications of Artificial Intelligence*, 88:103360, 2020.
- 414 Javier García, Fern, and o Fernández. A comprehensive survey on safe reinforcement learning.
415 *Journal of Machine Learning Research*, 16(42):1437–1480, 2015. URL [http://jmlr.org/
416 papers/v16/garcial5a.html](http://jmlr.org/papers/v16/garcial5a.html).
- 417 Arnob Ghosh, Xingyu Zhou, and Ness Shroff. Provably efficient model-free con-
418 strained rl with linear function approximation. In S. Koyejo, S. Mohamed, A. Agar-
419 wal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Pro-
420 cessing Systems*, volume 35, pp. 13303–13315. Curran Associates, Inc., 2022. URL
421 [https://proceedings.neurips.cc/paper_files/paper/2022/file/
422 56b8f22d895c45f60eaac9580152afd9-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/56b8f22d895c45f60eaac9580152afd9-Paper-Conference.pdf).
- 423 Arnob Ghosh, Xingyu Zhou, and Ness Shroff. Towards achieving sub-linear regret and hard con-
424 straint violation in model-free rl. In *International Conference on Artificial Intelligence and Statis-
425 tics*, pp. 1054–1062. PMLR, 2024.
- 426 Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, and Alois Knoll. A
427 review of safe reinforcement learning: Methods, theory and applications, 2024. URL [https:
428 //arxiv.org/abs/2205.10330](https://arxiv.org/abs/2205.10330).
- 429 David Isele, Alireza Nakhaei, and Kikuo Fujimura. Safe reinforcement learning on autonomous
430 vehicles. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*,
431 pp. 1–6. IEEE, 2018.
- 432 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably
433 efficient? In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
434 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran
435 Associates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/
436 paper/2018/file/d3b1fb02964aa64e257f9f26a31f72cf-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/d3b1fb02964aa64e257f9f26a31f72cf-Paper.pdf).

- 437 Chi Jin, Tiancheng Jin, Haipeng Luo, Suvrit Sra, and Tiancheng Yu. Learning adversarial Markov
438 decision processes with bandit feedback and unknown transition. In Hal Daumé III and Aarti
439 Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume
440 119 of *Proceedings of Machine Learning Research*, pp. 4860–4869. PMLR, 13–18 Jul 2020.
441 URL <https://proceedings.mlr.press/v119/jin20c.html>.
- 442 Krishna C. Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic
443 finite-horizon mdp with constraints. *Proceedings of the AAAI Conference on Artificial Intelli-*
444 *gence*, 35(9):8030–8037, May 2021. DOI: 10.1609/aaai.v35i9.16979. URL <https://ojs.aaai.org/index.php/AAAI/article/view/16979>.
- 446 Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. Safe posterior sampling for constrained mdps
447 with bounded constraint violation, 2023. URL <https://arxiv.org/abs/2301.11547>.
- 448 Hanna Krasowski, Xiao Wang, and Matthias Althoff. Safe reinforcement learning for autonomous
449 lane changing using set-based prediction. In *2020 IEEE 23rd international conference on Intelli-*
450 *gent Transportation Systems (ITSC)*, pp. 1–7. IEEE, 2020.
- 451 Tao Liu, Ruida Zhou, Dileep Kalathil, Panganamala Kumar, and Chao Tian. Learning policies with
452 zero or bounded constraint violation for constrained MDPs. In A. Beygelzimer, Y. Dauphin,
453 P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*,
454 2021. URL https://openreview.net/forum?id=Nl7VO_Y7K4Q.
- 455 Andreas Maurer and Massimiliano Pontil. Empirical Bernstein bounds and sample variance penal-
456 ization. In *Proceedings of the 22nd Annual Conference on Learning Theory*, 2009.
- 457 Sobhan Miryoosefi and Chi Jin. A simple reward-free approach to constrained reinforcement learn-
458 ing. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan
459 Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume
460 162 of *Proceedings of Machine Learning Research*, pp. 15666–15698. PMLR, 17–23 Jul 2022.
461 URL <https://proceedings.mlr.press/v162/miryoosefi22a.html>.
- 462 Ted Moskowitz, Brendan O’Donoghue, Vivek Veeriah, Sebastian Flennerhag, Satinder Singh, and
463 Tom Zahavy. ReLOAD: Reinforcement learning with optimistic ascent-descent for last-iterate
464 convergence in constrained MDPs. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Bar-
465 bara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International*
466 *Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*,
467 pp. 25303–25336. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/moskovitz23a.html>.
- 469 Adrian Müller, Pragnya Alatur, Volkan Cevher, Giorgia Ramponi, and Niao He. Truly no-regret
470 learning in constrained MDPs. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian
471 Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st*
472 *International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning*
473 *Research*, pp. 36605–36653. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/muller24b.html>.
- 475 Shuang Qiu, Xiaohan Wei, Zhuoran Yang, Jieping Ye, and Zhaoran Wang. Upper confi-
476 dence primal-dual reinforcement learning for cmdp with adversarial loss. In H. Larochelle,
477 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural In-*
478 *formation Processing Systems*, volume 33, pp. 15277–15287. Curran Associates, Inc.,
479 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/ae95296e27d7f695f891cd26b4f37078-Paper.pdf)
480 [file/ae95296e27d7f695f891cd26b4f37078-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/ae95296e27d7f695f891cd26b4f37078-Paper.pdf).
- 481 Aviv Rosenberg and Yishay Mansour. Online convex optimization in adversarial Markov decision
482 processes. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th*
483 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*
484 *Research*, pp. 5478–5486. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/rosenberg19a.html>.
- 485

- 486 Thiago D. Simão, Nils Jansen, and Matthijs T. J. Spaan. Always safe: Reinforcement learning
487 without safety constraint violations during training. In *Proceedings of the 20th International*
488 *Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '21, pp. 1226–1235, Rich-
489 land, SC, 2021. International Foundation for Autonomous Agents and Multiagent Systems. ISBN
490 9781450383073.
- 491 Rahul Singh, Abhishek Gupta, and Ness B. Shroff. Learning in constrained markov decision
492 processes. *IEEE Transactions on Control of Network Systems*, 10(1):441–453, 2023. DOI:
493 10.1109/TCNS.2022.3203361.
- 494 Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning
495 by PID lagrangian methods. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*
496 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*
497 *Research*, pp. 9133–9143. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v119/stooke20a.html)
498 [press/v119/stooke20a.html](https://proceedings.mlr.press/v119/stooke20a.html).
- 499 Honghao Wei, Xin Liu, and Lei Ying. A provably-efficient model-free algorithm for infinite-horizon
500 average-reward constrained markov decision processes. *Proceedings of the AAAI Conference*
501 *on Artificial Intelligence*, 36(4):3868–3876, Jun. 2022a. DOI: 10.1609/aaai.v36i4.20302. URL
502 <https://ojs.aaai.org/index.php/AAAI/article/view/20302>.
- 503 Honghao Wei, Xin Liu, and Lei Ying. Triple-q: A model-free algorithm for constrained reinforce-
504 ment learning with sublinear regret and zero constraint violation. In *International Conference on*
505 *Artificial Intelligence and Statistics*, pp. 3274–3307. PMLR, 2022b.
- 506 Honghao Wei, Arnob Ghosh, Ness Shroff, Lei Ying, and Xingyu Zhou. Provably efficient model-
507 free algorithms for non-stationary cmdps. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de
508 Meent (eds.), *Proceedings of The 26th International Conference on Artificial Intelligence and*
509 *Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp. 6527–6570. PMLR,
510 25–27 Apr 2023. URL <https://proceedings.mlr.press/v206/wei23b.html>.
- 511 Tiancheng Yu, Yi Tian, Jingzhao Zhang, and Suvrit Sra. Provably efficient algorithms for multi-
512 objective competitive rl. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th In-*
513 *ternational Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning*
514 *Research*, pp. 12167–12176. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v139/yu21b.html)
515 [press/v139/yu21b.html](https://proceedings.mlr.press/v139/yu21b.html).
- 516 Liyuan Zheng and Lillian Ratliff. Constrained upper confidence reinforcement learning. In *Learning*
517 *for Dynamics and Control*, pp. 620–629. PMLR, 2020.

518
519
520

Supplementary Materials

The following content was not necessarily subject to peer review.

521 8 Related Work

522 In this section, we provide a more detailed discussion of related work to online learning of con-
523 strained Markov decision processes (CMDPs). As explained in the introduction, we review previous
524 works for the three frameworks, cumulative constraint violation, hard constraint violation, and zero
525 constraint violation.

526 **Cumulative Constraint Violation** Starting with the work of [Efroni et al. \(2020\)](#), online learn-
527 ing of CMDPs has been an active area of research in reinforcement learning, especially with the
528 framework of cumulative (or soft) constraint violation ([Brantley et al., 2020](#); [Qiu et al., 2020](#); [Zheng
529 & Ratliff, 2020](#); [Kalagarla et al., 2021](#); [Ding et al., 2021](#); [Chen et al., 2021](#); [Yu et al., 2021](#); [Liu
530 et al., 2021](#); [Wei et al., 2022a;b](#); [Singh et al., 2023](#); [Miryoosefi & Jin, 2022](#); [Ghosh et al., 2022](#); [Wei
531 et al., 2023](#); [Kalagarla et al., 2023](#)). Among these works, [Brantley et al. \(2020\)](#) studied a knapsack
532 constrained formulation, and [Qiu et al. \(2020\)](#) studied the setting where the reward functions are
533 adversarially given and the cost functions are sampled from a fixed but unknown distribution. More-
534 over, [Zheng & Ratliff \(2020\)](#) considered the case where the transition kernel is known to the agent,
535 and [Kalagarla et al. \(2021\)](#) studied a PAC bound for learning CMDPs. [Ding et al. \(2021\)](#); [Chen et al.
536 \(2021\)](#) developed model-free algorithms for CMDPs, although these approaches require access to
537 simulators, while [Yu et al. \(2021\)](#) studied vector-valued Markov games for a variant of constrained
538 MDPs. [Liu et al. \(2021\)](#) introduced the first algorithm that achieves zero cumulative constraint vio-
539 lation. [Wei et al. \(2022a\)](#) and [Singh et al. \(2023\)](#) considered the infinite-horizon average-reward set-
540 ting. Moreover, [Wei et al. \(2022b\)](#) came up with a model-free algorithm for finite-horizon episodic
541 tabular CMDPs. [Miryoosefi & Jin \(2022\)](#) studied the reward-free setting, and [Ghosh et al. \(2022\)](#)
542 proposed an algorithm for the linear MDP setting, which leads to a model-free algorithm for tabular
543 CMDPs. Lastly, [Wei et al. \(2023\)](#) considered non-stationary CMDPs, while [Kalagarla et al. \(2023\)](#)
544 developed a posterior sampling-based algorithm that guarantees a Bayesian regret upper bound.

545 [Wei et al. \(2022b\)](#) introduced model-free and simulator-free algorithms to solve tabular CMDPs.
546 These algorithms were analyzed under soft constraint violations, thus they do not guarantee safety
547 in all episodes. In contrast, [Müller et al. \(2024\)](#); [Ghosh et al. \(2024\)](#) presented PD-based algorithms
548 with hard constraint violations, though these suffer from high regret and constraint violations. On
549 the other hand, [Liu et al. \(2021\)](#) proposed the LP-based algorithm OptPess-LP, which achieves
550 zero hard constraint violations with sublinear regret by employing *optimistic pessimism in the face
551 of uncertainty (OPFU)*. The pessimism in the cost function estimator ensures safety but hampers
552 exploration. To address this, [Bura et al. \(2022\)](#) recently proposed DOPE, incorporating optimism
553 for the transition kernel to improve the regret bound.

554 **Hard Constraint Violation** The notion of hard constraint violation was introduced by [Efroni et al.
555 \(2020\)](#). [Efroni et al. \(2020\)](#) developed an LP-based algorithm for controlling hard constraint vio-
556 lation and raised an open question of whether there exists a primal-dual algorithm for the setting.
557 Recently, [Ghosh et al. \(2024\)](#) established an algorithm that guarantees a sublinear regret upper bound
558 and a sublinear upper bound on hard constraint violation. Their algorithm is for the linear MDP set-
559 ting, and it provides a model-free algorithm for the tabular setting. In fact, their analysis shows that
560 for the tabular case, one may get a tighter performance guarantees. [Müller et al. \(2024\)](#) developed a
561 simpler primal-dual algorithm that guarantees a sublinear regret upper bound and a sublinear upper
562 bound on hard constraint violation, answering the question of [Efroni et al. \(2020\)](#).

563 **Zero Constraint Violation** [Simão et al. \(2021\)](#) considered the importance of achieving no con-
564 straint violation, which is equivalent to zero hard constraint violation. They showed an algorithm

565 that guarantees no constraint violation, but their result relies on the assumption of some abstrac-
 566 tion of the transition model, and moreover, there is no regret upper bound given for the algorithm.
 567 Liu et al. (2021) established the first algorithm that achieves a sublinear regret while guaranteeing
 568 zero hard constraint violation. After Liu et al. (2021), (Bura et al., 2022) proposed their algorithm,
 569 DOPE, which improves upon Liu et al. (2021) to show a smaller regret upper bound.

570 9 Auxiliary Measures and Notations

571 In this section, we first summarize notations in Table 2. Next, we define some auxiliary measures
 572 and notations that are useful for the analysis of DOPE+.

Table 2: Summary of Notations

Notation	Definition
K	The number of episodes
H	The finite horizon
$[H]$	The set $\{1, 2, \dots, H\}$
\mathcal{S}, S	The finite state space \mathcal{S} and the number of states $S = \mathcal{S} $
\mathcal{A}, A	The finite action space \mathcal{A} and the number of actions $A = \mathcal{A} $
P	The true transition kernel $P(s, a, s', h) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H] \rightarrow [0, 1]$
p	The initial distribution of the states
\mathcal{P}_k	The confidence set of the transition kernel for episode $k \in [K]$
P_k	The transition kernel obtained from DOPE+ for episode $k \in [K]$, $P_k \in \mathcal{P}_k$
f, g	The reward and cost function
f_k, g_k	The instantaneous reward and cost for episode $k \in [K]$
\hat{f}_k, \bar{g}_k	The empirical estimators of f, g for episode $k \in [K]$
\hat{f}_k, \hat{g}_k	The optimistic/pessimistic estimators of f, g for episode $k \in [K]$
L_δ	$\ln(HSAK/\delta)$ for some confidence parameter $\delta \in (0, 1)$
$V_h^\pi(s; f, P)$	The value function at state s and step h under f and P
$Q_h^\pi(s, a; f, P)$	The action-value function at state s and step h for action a under f and P
$N_k(s, a, h)$	The number of visits (s, a, h) up to the first $k - 1$ episodes
$M_k(s, a, s', h)$	The number of visits (s, a, s', h) up to the first $k - 1$ episodes
$n_k(s, a, h)$	The indicator variable for visits (s, a, h) for episode $k \in [K]$
π^*	The benchmark policy
π_k	The policy obtained from DOPE+ for episode $k \in [K]$
π_b	The safe baseline policy
\bar{C}_b	The expected cost of π_b for a single episode
\bar{C}	The budget on the expected cost
$q^{P, \pi}$	The occupancy measure with respect to policy π and transition kernel P
q^*	The occupancy measure q^{P, π^*}
q_b	The occupancy measure q^{P, π_b}
q_k	The occupancy measure q^{P, π_k}
\hat{q}_k	The occupancy measure q^{P_k, π_k}
$\Delta(P)$	The set of occupancy measures inducing P
$\Delta(P, k)$	The set of occupancy measures inducing $P_k \in \mathcal{P}_k$

573 We define the *state-action value function* for $(s, a) \in \mathcal{S} \times \mathcal{A}$ at step h with a function $\ell : \mathcal{S} \times \mathcal{A} \times$
 574 $[H] \rightarrow [0, 1]$ and transition kernel P as follows.

$$Q_h^\pi(s, a; \ell, P) = \mathbb{E} \left[\sum_{j=h}^H \ell \left(s_j^{P, \pi}, a_j^{P, \pi}, j \right) \mid \ell, \pi, P, s_h^{P, \pi} = s, a_h^{P, \pi} = a \right].$$

Let $\mathbf{Q}^{P,\pi,\ell}$ denote the $(S \times A \times H)$ -dimensional vector whose coordinates are for $(s, a, h) \in S \times \mathcal{A} \times [H]$,

$$(\mathbf{Q}^{P,\pi,\ell})_{(s,a,h)} = Q_h^\pi(s, a; \ell, P).$$

575 Given a policy π and transition kernel P , we define $q^{P,\pi}(s, a, h | s', m)$ as for $(s, a, s') \in S \times \mathcal{A} \times S$
 576 and $1 \leq m \leq h \leq H$,

$$q^{P,\pi}(s, a, h | s', m) = \mathbb{P}\left[s_h^{P,\pi} = s, a_h^{P,\pi} = a \mid \pi, P, s_m^{P,\pi} = s'\right].$$

Given two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{S \times A \times H}$, let $\mathbf{u} \odot \mathbf{v}$, $\mathbf{u} \wedge \mathbf{v}$ be defined as the vector obtained from coordinate-wise products and coordinate-wise minimization of \mathbf{u} and \mathbf{v} , respectively, i.e., for $(s, a, h) \in S \times \mathcal{A} \times [H]$,

$$(\mathbf{u} \odot \mathbf{v})_{(s,a,h)} = \mathbf{u}_{(s,a,h)} \times \mathbf{v}_{(s,a,h)}, \quad (\mathbf{u} \wedge \mathbf{v})_{(s,a,h)} = \min\{\mathbf{u}_{(s,a,h)}, \mathbf{v}_{(s,a,h)}\}.$$

Let $\vec{\mathbf{h}}$ and $\vec{\mathbf{B}}$ be $(S \times A \times H)$ -dimensional vectors all of whose coordinates are h and $1 + \sqrt{L\delta}$, respectively, i.e., for $(s, a, j) \in S \times \mathcal{A} \times [H]$,

$$\vec{\mathbf{h}}_{(s,a,j)} = j, \quad \vec{\mathbf{B}}_{(s,a,j)} = 1 + \sqrt{L\delta}.$$

577 10 Extended Linear Program

578 In this section, we provide a formal definition of occupancy measures for a finite-horizon MDP.
 579 Then we provide a reformulation of (10) using occupancy measures, which is called the extended
 580 linear program (Efroni et al., 2020; Bura et al., 2022).

581 Given a policy π and a transition kernel P , let $\bar{q}^{P,\pi} : S \times \mathcal{A} \times S \times [H] \rightarrow [0, 1]$ be defined
 582 as $\bar{q}^{P,\pi}(s, a, s', h) = \mathbb{P}[(s_h^{P,\pi}, a_h^{P,\pi}, s_{h+1}^{P,\pi}) = (s, a, s') \mid \pi, P]$ for $(s, a, s', h) \in S \times \mathcal{A} \times$
 583 $S \times [H]$. Note that any \bar{q} defined as the above equation has the following properties. (C1)
 584 $\sum_{(s,a,s') \in S \times \mathcal{A} \times S} \bar{q}(s, a, s', h) = 1$, (C2) $\sum_{(s',a) \in S \times \mathcal{A}} \bar{q}(s, a, s', h) = \sum_{(s',a) \in S \times \mathcal{A}} \bar{q}(s', a, s, h -$
 585 $1)$, $s \in S$, $h = 2, \dots, H$. The *occupancy measure* $q^{P,\pi} : S \times \mathcal{A} \times [H] \rightarrow [0, 1]$ associated with
 586 policy π and transition kernel P is defined as (C3) $q^{P,\pi}(s, a, h) = \sum_{s' \in S} \bar{q}^{P,\pi}(s, a, s', h)$. Then it
 587 follows that $q^{P,\pi}(s, a, h) = \mathbb{P}[(s_h^{P,\pi}, a_h^{P,\pi}) = (s, a) \mid \pi, P]$. Hence, if a policy π is chosen, then the
 588 occupancy measure for a finite-horizon MDP with transition kernel P is determined. Conversely,
 589 any $q \in S \times \mathcal{A} \times [H] \rightarrow [0, 1]$ with $\bar{q} : S \times \mathcal{A} \times S \times [H] \rightarrow [0, 1]$ satisfying (C1), (C2), and (C3)
 590 induces a transition kernel P^q and a policy π^q given as follows.

$$P^q(s' | s, a, h) = \frac{\bar{q}(s, a, s', h)}{\sum_{s'' \in S} \bar{q}(s, a, s'', h)}, \quad (12)$$

$$\pi^q(a | s, h) = \frac{q(s, a, h)}{\sum_{b \in \mathcal{A}} q(s, b, h)}.$$

591 Next, we provide a lemma that characterizes valid occupancy measures for a finite-horizon MDP.

592 **Lemma 10.** *Let $q : S \times \mathcal{A} \times [H] \rightarrow [0, 1]$. Then q is a valid occupancy measure that induces*
 593 *transition kernel P if and only if there exists $\bar{q} : S \times \mathcal{A} \times S \times [H] \rightarrow [0, 1]$ that satisfies (C1), (C2),*
 594 *(C3), and $P^q = P$.*

Proof. Given the finite-horizon MDP associated with transition kernel P , we may define a loop-free MDP as follows. We define its state space as $S' := S \times [H + 1]$, which can be viewed as $H + 1$ layers $S \times \{h\}$ for $h \in [H + 1]$. Its transition kernel P' is given by $P'((s', h + 1) | (s, h), a) = P(s' | s, a, h)$ for $(s, a, s', h) \in S \times \mathcal{A} \times S \times [H]$. Next, given \bar{q} , we may define an occupancy measure q' for the loop-free MDP as $q'((s, h), a, (s', h + 1)) = \bar{q}(s, a, s', h)$ for $(s, a, s', h) \in S \times \mathcal{A} \times S \times [H]$. Then

it follows from (Rosenberg & Mansour, 2019, Lemma 3.1) that q' is a valid occupancy measure for the loop-free MDP with transition kernel P' if and only if q' satisfies

$$\sum_{(s,a,s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}} q'((s,h), a, (s', h+1)) = 1 \quad \text{for } h = 1, \dots, H, \quad (\text{C1}')$$

$$\sum_{(s',a) \in \mathcal{S} \times \mathcal{A}} q'((s,h), a, (s', h+1)) = \sum_{(s',a) \in \mathcal{S} \times \mathcal{A}} q'((s', h-1), a, (s,h)) \quad \begin{array}{l} \forall s \in \mathcal{S}, \\ h \in \{2, \dots, H\} \end{array} \quad (\text{C2}')$$

and $P^{q'} = P'$ where $P^{q'}$ is given by

$$P^{q'}((s', h+1) | (s, h), a) = \frac{q'((s, h), a, (s', h+1))}{\sum_{s'' \in \mathcal{S}} q'((s, h), a, (s'', h+1))} = \frac{\bar{q}(s, a, s', h)}{\sum_{s'' \in \mathcal{S}} \bar{q}(s, a, s'', h)}.$$

595 Here, the conditions are equivalent to (C1), (C2), and $P^{\bar{q}} = P$. Moreover, q' is a valid occupancy
596 measure with P' if and only if q is a valid occupancy measure with P , as required. \square

Therefore, there is a one-to-one correspondence between the set of policies and the set of occupancy measures that give rise to transition kernel P . We define $\Delta(P)$ as the set of occupancy measures inducing the true transition kernel P .

$$\Delta(P) = \{q : \exists \bar{q} \text{ satisfying (C1),(C2),(C3), } P^q = P\}.$$

597 Moreover, the value function for reward function f , policy π_k , and transition kernel P can be writ-
598 ten in terms of occupancy measure q^{P, π_k} as $V_1^{\pi_k}(f, P) = \sum_{(s,a,h)} q^{P, \pi_k}(s, a, h) f(s, a, h)$. Let
599 $q^{P, \pi}, \mathbf{f}$ denote $(\mathcal{S} \times \mathcal{A} \times H)$ -dimensional vector representations for $q^{P, \pi}, f$, respectively. Then
600 it follows that $V_1^{\pi_k}(f, P) = \langle \mathbf{f}, \mathbf{q}^{P, \pi_k} \rangle$ where $\langle \cdot, \cdot \rangle$ is the inner product. Consequently, (10) is
601 equivalent to

$$\max_{q \in \Delta(P, k)} \left\{ \langle \hat{\mathbf{f}}_k, \mathbf{q} \rangle : \langle \hat{\mathbf{g}}_k, \mathbf{q} \rangle \leq \bar{C} \right\} \quad (13)$$

where $\hat{\mathbf{f}}_k, \hat{\mathbf{g}}_k$ are the vector representations of \hat{f}_k, \hat{g}_k , respectively, and

$$\Delta(P, k) = \{q : \exists \bar{q} \text{ satisfying (C1),(C2),(C3), } P^q \in \mathcal{P}_k\}.$$

602 Next, we reformulate (10) as an extended linear program. Due to the definition of $\Delta(P, k)$, (13) is
603 equivalent to the following linear program. Given $\hat{f}_k(s, a, h), \hat{g}_k(s, a, h), \bar{P}_k(s' | s, a, h), \epsilon_k(s' |$
604 $s, a, h), p(s)$ for $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$,

$$\begin{aligned} \max \quad & \sum_{(s,a,s',h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \hat{f}_k(s, a, h) \bar{q}(s, a, s', h) \\ \text{s.t.} \quad & \sum_{(s,a,s',h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \hat{g}_k(s, a, h) \bar{q}(s, a, s', h) \leq \bar{C}, \\ & \sum_{(a,s') \in \mathcal{A} \times \mathcal{S}} \bar{q}(s, a, s', h) = \sum_{(a,s') \in \mathcal{A} \times \mathcal{S}} \bar{q}(s', a, s, h-1) \quad \forall s \in \mathcal{S}, h = 2, \dots, H, \\ & \sum_{(a,s') \in \mathcal{A} \times \mathcal{S}} \bar{q}(s, a, s', 1) = p(s) \quad \forall s \in \mathcal{S}, \\ & \bar{q}(s, a, s', h) \leq (\bar{P}_k(s' | s, a, h) + \epsilon_k(s' | s, a, h)) \sum_{s'' \in \mathcal{S}} \bar{q}(s, a, s'', h) \quad \forall (s, a, s', h), \\ & \bar{q}(s, a, s', h) \geq (\bar{P}_k(s' | s, a, h) - \epsilon_k(s' | s, a, h)) \sum_{s'' \in \mathcal{S}} \bar{q}(s, a, s'', h) \quad \forall (s, a, s', h), \\ & 0 \leq \bar{q}(s, a, s', h) \quad \forall (s, a, s', h). \end{aligned} \quad (14)$$

605 In fact, the constraint $\sum_{(s,a,s')} \bar{q}(s, a, s', h) = 1$ for $h \in [H]$ corresponding to (C1) is not necessary,
606 because we can derive it from other constraints. To be more specific, the third constraint implies

607 that $\sum_{(s,a,s')} \bar{q}(s,a,s',1) = 1$ as $\sum_s p(s) = 1$. Then we can deduce from the second constraint
 608 that $\sum_{(s,a,s')} \bar{q}(s,a,s',h) = 1$ for $h \in [H]$. Additionally, we call the above linear program as an
 609 extended linear program due to the fifth and sixth constraints.

610 One natural question to the extended LP defined in (14) is how hard it is to solve. Indeed, we can
 611 easily observe that the dimension of the decision variable \bar{q} is S^2AH , and the number of constraints
 612 is $\mathcal{O}(S^2AH)$. Hence, the computational complexity for solving (14) is equivalent to solving an LP
 613 with a S^2AH -dimensional decision variable and $\mathcal{O}(S^2AH)$ constraints.

614 11 Good Event

615 In this section, we first prove Lemma 1 which ensures that all instantaneous reward and cost values
 616 are bounded. Then we prove Lemma 2 that describes important properties of the confidence sets
 617 estimating the true transition kernel. Next, we show Lemma 3 which delineates the accuracy of our
 618 estimators of the reward function f and the cost function g .

619 Furthermore, we prove Lemma 11 that is useful to bound value functions with respect to estimated
 620 reward and cost functions. Then we define the notion of the *good event* \mathcal{E} that the statements of
 621 Lemmas 1 to 3 and 11 hold. Taking the union bound, we deduce that the good event \mathcal{E} holds with
 622 probability at least $1 - 14\delta$ (Lemma 12).

623 Lastly, we prove Lemma 13 which considers the difference between the true transition kernel and
 624 any \bar{P} contained in the confidence set \mathcal{P}_k .

Proof of Lemma 1. It follows from Hoeffding's inequality (Lemma 21) and the union bound that for any $(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$,

$$\mathbb{P}\left(|f_k(s,a,h) - f(s,a,h)| \geq \sqrt{L_\delta}\right) \leq 2 \cdot \exp(-L_\delta) = \frac{2\delta}{HSAK}.$$

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

$$\mathbb{P}\left(|g_k(s,a,h) - g(s,a,h)| \geq \sqrt{L_\delta}\right) \leq 2 \cdot \exp(-L_\delta) = \frac{2\delta}{HSAK}.$$

Taking the union bound, it follows that with probability at least $1 - 4\delta$,

$$|f_k(s,a,h) - f(s,a,h)|, |g_k(s,a,h) - g(s,a,h)| \leq \sqrt{L_\delta}$$

holds for all $(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$. Since $f(s,a,h), g(s,a,h) \in [0,1]$ for any $(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]$, it holds with probability at least $1 - 4\delta$ that

$$|f_k(s,a,h)|, |g_k(s,a,h)| \leq 1 + \sqrt{L_\delta},$$

625 as required. □

626 The following lemma is a modification of (Jin et al., 2020, Lemma 8) to our finite-horizon MDP
 627 setting.

628 **Proof of Lemma 2.** We will show that with probability at least $1 - 4\delta$,

$$|P(s' | s, a, h) - \bar{P}_k(s' | s, a, h)| \leq \epsilon_k(s' | s, a, h) \tag{15}$$

629 where

$$\epsilon_k(s' | s, a, h) = 2\sqrt{\frac{\bar{P}_k(s' | s, a, h)(1 - \bar{P}_k(s' | s, a, h))L_\delta}{\max\{1, N_k(s, a, h) - 1\}}} + \frac{14L_\delta}{3 \max\{1, N_k(s, a, h) - 1\}}$$

630 holds for every $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$ and every episode $k \in [K]$.

Let us first consider the case $N_k(s, a, h) \leq 1$. As we may assume that $HS^2AK \geq 2$, it follows that

$$\epsilon_k(s' | s, a, h) = \frac{14L\delta}{3 \max\{1, N_k(s, a, h) - 1\}} \geq \frac{14}{3} \ln 2 > 1.$$

631 Then (15) holds because $0 \leq P(s' | s, a, h), \bar{P}_k(s' | s, a, h) \leq 1$.

Assume that $n = N_k(s, a, h) \geq 2$. Then we define Z_1, \dots, Z_n as follows.

$$Z_j = \begin{cases} 1, & \text{if the transition after the } j\text{th visit to } (s, a, h) \text{ is } s', \\ 0, & \text{otherwise.} \end{cases}$$

Then Z_1, \dots, Z_n are i.i.d. with mean $P(s' | s, a, h)$, and we have

$$\sum_{j=1}^n Z_j = M_k(s, a, s', h).$$

632 Moreover, the sample variance V_n of Z_1, \dots, Z_n is given by

$$\begin{aligned} V_n &= \frac{1}{N_k(s, a, h)(N_k(s, a, h) - 1)} M_k(s, a, s', h) (N_k(s, a, h) - M_k(s, a, s', h)) \\ &= \frac{N_k(s, a, h)}{(N_k(s, a, h) - 1)} \bar{P}_k(s' | s, a, h) (1 - \bar{P}_k(s' | s, a, h)). \end{aligned} \quad (16)$$

633 Then it follows from Lemma 22 that with probability at least $1 - 2\delta/(HS^2AK)$,

$$\begin{aligned} &P(s' | s, a, h) - \bar{P}_k(s' | s, a, h) \\ &\leq \sqrt{\frac{2\bar{P}_k(s' | s, a, h) (1 - \bar{P}_k(s' | s, a, h)) \ln(HS^2AK/\delta)}{N_k(s, a, h) - 1}} + \frac{7 \ln(HS^2AK/\delta)}{3(N_k(s, a, h) - 1)}. \end{aligned} \quad (17)$$

634 Here, as we assumed that $N_k(s, a, h) \geq 2$, we have $N_k(s, a, h) - 1 = \max\{1, N_k(s, a, h) - 1\}$.

635 In addition, we know that $\ln(HS^2AK/\delta) \leq 2L\delta$. Then (17) implies that with probability at least

636 $1 - 2\delta/(HS^2AK)$,

$$P(s' | s, a, h) - \bar{P}_k(s' | s, a, h) \leq \epsilon_k(s' | s, a, h). \quad (18)$$

637 Next, we apply Lemma 22 to variables $1 - Z_1, \dots, 1 - Z_n$ that are i.i.d. and have mean $1 - \bar{P}_k(s' | s, a, h)$.

638 Moreover, the sample variance of $1 - Z_1, \dots, 1 - Z_n$ is also equal to V_n defined as in (16).

639 Therefore, based on the same argument, we deduce that with probability at least $1 - 2\delta/(HS^2AK)$,

$$-P(s' | s, a, h) + \bar{P}_k(s' | s, a, h) \leq \epsilon_k(s' | s, a, h). \quad (19)$$

640 By applying union bound to (18) and (19), with probability at least $1 - 4\delta/(HS^2AK)$, (15) holds

641 for (s, a, s', h) . Furthermore, by applying union bound over all $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$, it

642 follows that with probability at least $1 - 4\delta$, (15) holds for every $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$,

643 as required. \square

644 Next, we state the proof of Lemma 3 based on Hoeffding's inequality.

645 **Proof of Lemma 3.** If $N_k(s, a, h) = \sum_{j=1}^{k-1} n_j(s, a, h) = 0$, then $\bar{f}_k(s, a, h) = \bar{g}_k(s, a, h) =$

646 0 while $R_k(s, a, h) \geq 1$ when we may assume that $HS^2AK \geq 4$. In this case, the statements

647 trivially hold. Now we consider when $\sum_{j=1}^{k-1} n_j(s, a, h) \geq 1$. Note that $f_k(s, a, h) = f(s, a, h) +$

648 $\zeta_k^f(s, a, h)$ and $g_k(s, a, h) = g(s, a, h) + \zeta_k^g(s, a, h)$ where $\zeta_k^f(s, a, h)$ and $\zeta_k^g(s, a, h)$ are i.i.d. $1/2$ -

649 sub-Gaussian random variables with zero mean for each $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$.

650 Then it follows from the Hoeffding's inequality provided in Lemma 21 that for a given $(s, a, h) \in$
 651 $\mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$,

$$|\bar{f}_k(s, a, h) - f(s, a, h)| \leq R_k(s, a, h) \quad (20)$$

with probability at least $1 - 2\delta/(HSAK)$. By applying union bound, (20) holds with probability at least $1 - 2\delta$ for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$. Likewise, we deduce for g that with probability at least $1 - 2\delta$,

$$|\bar{g}_k(s, a, h) - g(s, a, h)| \leq R_k(s, a, h)$$

652 for $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$ as desired. \square

653 Next, using Lemma 23 that states the Bernstein-type concentration inequality for a martingale dif-
 654 ference sequence, we prove the following lemma that is useful for our analysis. Lemma 11 is a mod-
 655 ification of (Jin et al., 2020, Lemma 10) and (Chen & Luo, 2021, Lemma 8) to our finite-horizon
 656 MDP setting.

657 **Lemma 11.** *With probability at least $1 - 2\delta$, we have*

$$\sum_{k=1}^K \sum_{(s,a,h)} \frac{q_k(s, a, h)}{\max\{1, N_k(s, a, h)\}} \leq 2HSA \ln K + 2HSA + 4H \ln(H/\delta) \quad (21)$$

$$\sum_{k=1}^K \sum_{(s,a,h)} \frac{q_k(s, a, h)}{\sqrt{\max\{1, N_k(s, a, h)\}}} \leq 2H\sqrt{SAK} + 2HSA \ln K + 3HSA + 5H \ln(H/\delta) \quad (22)$$

658 *Proof.* We define ξ_1 as $\xi_1 = \emptyset$ and for $k \geq 2$, we define ξ_k as

$$\left\{ s_h^{P, \pi_{k-1}}, a_h^{P, \pi_{k-1}}, f_{k-1}(s_h^{P, \pi_{k-1}}, a_h^{P, \pi_{k-1}}, h), g_{k-1}(s_h^{P, \pi_{k-1}}, a_h^{P, \pi_{k-1}}, h) \right\}_{h=1}^H$$

where π_{k-1} denotes the policy for episode $k - 1$ and

$$\left(s_1^{P, \pi_{k-1}}, a_1^{P, \pi_{k-1}}, \dots, s_h^{P, \pi_{k-1}}, a_h^{P, \pi_{k-1}} \right)$$

659 is the trajectory generated under policy π_{k-1} and transition kernel P . Then for $k \in [K]$, let \mathcal{F}_k be
 660 defined as the σ -algebra generated by the random variables in $\xi_1 \cup \dots \cup \xi_k$. Then it follows that
 661 $\mathcal{F}_1, \dots, \mathcal{F}_k$ give rise to a filtration.

662 Note that

$$\sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s, a, h)}{\max\{1, N_k(s, a, h)\}} = \sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{n_k(s, a, h)}{\max\{1, N_k(s, a, h)\}} + \sum_{k=1}^K Y_k \quad (23)$$

where

$$Y_k = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{-n_k(s, a, h) + q_k(s, a, h)}{\max\{1, N_k(s, a, h)\}}.$$

663 As $\mathbb{E}[n_k(s, a, h) | \pi_k, P] = q_k(s, a, h)$ holds for every $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, we know that
 664 Y_1, \dots, Y_K is a martingale difference sequence. We know that $Y_k \leq 1$ for each $k \in [K]$. Let $\mathbb{E}_k[\cdot]$
 665 denote $\mathbb{E}[\cdot | \mathcal{F}_k, P]$. Since π_k is \mathcal{F}_k -measurable, we have $\mathbb{E}_k[n_k(s, a, h)] = q_k(s, a, h)$. Then we

666 deduce

$$\begin{aligned}
\mathbb{E}_k [Y_k^2] &= \sum_{(s,a),(s',a') \in \mathcal{S} \times \mathcal{A}} \frac{\mathbb{E}_k [(n_k(s,a,h) - q_k(s,a,h))(n_k(s',a',h) - q_k(s',a',h))]}{\max\{1, N_k(s,a,h)\} \cdot \max\{1, N_k(s',a',h)\}} \\
&= \sum_{(s,a),(s',a') \in \mathcal{S} \times \mathcal{A}} \frac{\mathbb{E}_k [n_k(s,a,h)n_k(s',a',h) - q_k(s,a,h)q_k(s',a',h)]}{\max\{1, N_k(s,a,h)\} \cdot \max\{1, N_k(s',a',h)\}} \\
&\leq \sum_{(s,a),(s',a') \in \mathcal{S} \times \mathcal{A}} \frac{\mathbb{E}_k [n_k(s,a,h)n_k(s',a',h)]}{\max\{1, N_k(s,a,h)\} \cdot \max\{1, N_k(s',a',h)\}} \\
&\leq \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{\mathbb{E}_k [n_k(s,a,h)]}{\max\{1, N_k(s,a,h)\}} \\
&= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}
\end{aligned}$$

where the second equality holds because it follows from $\mathbb{E}_k [n_k(s,a,h)] = q_k(s,a,h)$ for $(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]$ that

$$\mathbb{E}_k [q_k(s,a,h)n_k(s',a',h)] = \mathbb{E}_k [q_k(s',a',h)n_k(s,a,h)] = q_k(s,a,h)q_k(s',a',h),$$

the second inequality holds because $n_k(s,a,h)n_k(s',a',h) = 0$ if $(s,a) \neq (s',a')$, and the last equality holds true because $\mathbb{E}_k [n_k(s,a,h)] = q_k(s,a,h)$ for any $(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]$. Then we may apply Lemma 23 with $\lambda = 1/2$, and we deduce that with probability at least $1 - \delta/H$,

$$\sum_{k=1}^K Y_k \leq \frac{1}{2} \sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} + 2 \ln(H/\delta).$$

Plugging this inequality to (23), it follows that

$$\sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} = 2 \sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{n_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} + 4 \ln(H/\delta).$$

667 Here, the first term on the right-hand side can be bounded as follows. We have

$$\begin{aligned}
&\sum_{k=1}^K \frac{n_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} \\
&= \sum_{k=1}^K \frac{n_k(s,a,h)}{\max\{1, N_{k+1}(s,a,h)\}} + \sum_{k=1}^K \left(\frac{n_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} - \frac{n_k(s,a,h)}{\max\{1, N_{k+1}(s,a,h)\}} \right) \\
&\leq \sum_{k=1}^K \frac{n_k(s,a,h)}{\max\{1, N_{k+1}(s,a,h)\}} + \sum_{k=1}^K \left(\frac{1}{\max\{1, N_k(s,a,h)\}} - \frac{1}{\max\{1, N_{k+1}(s,a,h)\}} \right) \\
&\leq \sum_{k=1}^K \frac{n_k(s,a,h)}{\max\{1, N_{k+1}(s,a,h)\}} + 1 \\
&\leq \ln K + 1.
\end{aligned}$$

where the first inequality is due to $n_k(s,a,h) \leq 1$ and the last inequality holds because

$$n_k(s,a,h) = N_{k+1}(s,a,h) - N_k(s,a,h) \quad \text{and} \quad N_K(s,a,h) + n_K(s,a,h) \leq K.$$

668 Therefore, it follows that

$$\sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{n_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{k=1}^K \frac{n_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} = SA \ln K + SA.$$

As a result, for any fixed $h \in [H]$,

$$\sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s, a, h)}{\max\{1, N_k(s, a, h)\}} \leq 2SA \ln K + 2SA + 4 \ln(H/\delta)$$

669 holds with probability at least $1 - \delta/H$. By union bound, (21) holds with probability at least $1 - \delta$.

670 Next, we will show that (22) holds.

$$\sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s, a, h)}{\sqrt{\max\{1, N_k(s, a, h)\}}} = \sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{n_k(s, a, h)}{\sqrt{\max\{1, N_k(s, a, h)\}}} + \sum_{k=1}^K Z_k \quad (24)$$

where

$$Z_k = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{-n_k(s, a, h) + q_k(s, a, h)}{\sqrt{\max\{1, N_k(s, a, h)\}}}.$$

671 As $\mathbb{E}_k[n_k(s, a, h)] = q_k(s, a, h)$ holds for every $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, we know that Z_1, \dots, Z_K

672 is a martingale difference sequence. We know that $Z_k \leq 1$ for each $k \in [K]$. Then we deduce

$$\begin{aligned} \mathbb{E}_k[Z_k^2] &\leq \sum_{(s,a),(s',a') \in \mathcal{S} \times \mathcal{A}} \frac{\mathbb{E}_k[n_k(s, a, h)n_k(s', a', h)]}{\sqrt{\max\{1, N_k(s, a, h)\}} \cdot \sqrt{\max\{1, N_k(s', a', h)\}}} \\ &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{\mathbb{E}_k[n_k(s, a, h)]}{\max\{1, N_k(s, a, h)\}} \\ &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s, a, h)}{\max\{1, N_k(s, a, h)\}} \end{aligned}$$

where the first inequality is derived by the same argument when bounding $\mathbb{E}_k[Y_k^2]$, the first equality holds because $n_k(s, a, h)n_k(s', a', h) = 0$ if $(s, a) \neq (s', a')$, and the last equality holds true because $\mathbb{E}_k[n_k(s, a, h)] = q_k(s, a, h)$ for any $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$. Then we may apply Lemma 23 with $\lambda = 1$, and we deduce that with probability at least $1 - \delta/H$,

$$\sum_{k=1}^K Z_k \leq \sum_{k=1}^K \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{q_k(s, a, h)}{\max\{1, N_k(s, a, h)\}} + \ln(H/\delta).$$

673 Then with probability at least $1 - \delta$, (21) holds and

$$\begin{aligned} \sum_{h \in [H]} \sum_{k=1}^K Z_k &\leq \sum_{k=1}^K \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} \frac{q_k(s, a, h)}{\max\{1, N_k(s, a, h)\}} + H \ln(H/\delta) \\ &= 2HSA \ln K + 2HSA + 5H \ln(H/\delta). \end{aligned} \quad (25)$$

674 holds. Moreover, we have

$$\begin{aligned} &\sum_{k=1}^K \frac{n_k(s, a, h)}{\sqrt{\max\{1, N_k(s, a, h)\}}} \\ &= \sum_{k=1}^K \frac{n_k(s, a, h)}{\sqrt{\max\{1, N_{k+1}(s, a, h)\}}} + \sum_{k=1}^K \left(\frac{n_k(s, a, h)}{\sqrt{\max\{1, N_k(s, a, h)\}}} - \frac{n_k(s, a, h)}{\sqrt{\max\{1, N_{k+1}(s, a, h)\}}} \right) \\ &\leq \sum_{k=1}^K \frac{n_k(s, a, h)}{\sqrt{\max\{1, N_{k+1}(s, a, h)\}}} + \sum_{k=1}^K \left(\frac{1}{\sqrt{\max\{1, N_k(s, a, h)\}}} - \frac{1}{\sqrt{\max\{1, N_{k+1}(s, a, h)\}}} \right) \\ &\leq \sum_{k=1}^K \frac{n_k(s, a, h)}{\sqrt{\max\{1, N_{k+1}(s, a, h)\}}} + 1 \\ &\leq 2\sqrt{N_{K+1}(s, a, h)} + 1. \end{aligned}$$

675 where the last equality holds because $n_k(s, a, h) = N_{k+1}(s, a, h) - N_k(s, a, h)$. Then

$$\begin{aligned} \sum_{k=1}^K \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} \frac{n_k(s, a, h)}{\sqrt{\max\{1, N_k(s, a, h)\}}} &\leq \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} 2\sqrt{N_{K+1}(s, a, h)} + HSA \\ &\leq 2\sqrt{HSA \sum_{(s,a,h)} N_{K+1}(s, a, h)} + HSA \\ &\leq 2H\sqrt{SA\bar{K}} + HSA \end{aligned}$$

676 where the second equality is due to the Cauchy-Schwarz inequality. Then it follows from (24)
677 and (25) that (22) holds. \square

678 Recall that the good event \mathcal{E} is the event that the statements of Lemmas 1 to 3 and 11 hold.

679 **Lemma 12.** *The good event \mathcal{E} holds with probability at least $1 - 14\delta$, i.e., $\mathbb{P}[\mathcal{E}] \geq 1 - 14\delta$.*

680 *Proof.* The proof follows from the union bound. \square

681 Lemma 2 bounds the difference between the true transition kernel P and the empirical transition
682 kernel \bar{P}_k . Based on Lemma 2, the next lemma bounds the difference between the true transition
683 kernel and any \hat{P} contained in the confidence set \mathcal{P}_k . Lemma 13 is a modification of (Jin et al.,
684 2020, Lemma 8) to our finite-horizon MDP setting.

685 **Lemma 13.** *Under the good event \mathcal{E} , we have*

$$\left| \hat{P}(s' | s, a, h) - P(s' | s, a, h) \right| \leq \epsilon_k^*(s' | s, a, h) \quad (26)$$

686 where

$$\epsilon_k^*(s' | s, a, h) = 6\sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + 94\frac{L_\delta}{\max\{1, N_k(s, a, h)\}}$$

687 for every $\hat{P} \in \mathcal{P}_k$ and every $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$.

Proof. We follow the proof of (Cohen et al., 2020, Lemma B.13). Note that

$$\max\{1, N_k(s, a, h) - 1\} \geq \frac{1}{2} \cdot \max\{1, N_k(s, a, h)\}$$

holds for any value of $N_k(s, a, h)$. We know that $1 - \bar{P}_k(s' | s, a) \leq 1$. Furthermore, as we assumed
that $P \in \mathcal{P}_k$, we have that

$$\bar{P}_k(s' | s, a, h) \leq P(s' | s, a, h) + \sqrt{\frac{8\bar{P}_k(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + \frac{28L_\delta}{3\max\{1, N_k(s, a, h)\}}.$$

688 We may view this as a quadratic inequality in terms of $x = \sqrt{\bar{P}_k(s' | s, a, h)}$. Note that $x^2 \leq$
689 $ax + b + c$ for any $a, b, c \geq 0$ implies that $x \leq a + \sqrt{b} + \sqrt{c}$. Therefore, we deduce that

$$\begin{aligned} \sqrt{\bar{P}_k(s' | s, a, h)} &\leq \sqrt{P(s' | s, a, h)} + \left(2\sqrt{2} + \sqrt{\frac{28}{3}}\right) \sqrt{\frac{L_\delta}{\max\{1, N_k(s, a, h)\}}} \\ &\leq \sqrt{P(s' | s, a, h)} + 13\sqrt{\frac{L_\delta}{\max\{1, N_k(s, a, h)\}}}. \end{aligned}$$

690 Using this bound on $\sqrt{\bar{P}_k(s' | s, a, h)}$, we obtain the following.

$$\begin{aligned}
 \epsilon_k(s' | s, a, h) &\leq \sqrt{\frac{8\bar{P}_k(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + \frac{28L_\delta}{3\max\{1, N_k(s, a, h)\}} \\
 &\leq \sqrt{\frac{8P(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + \left(13\sqrt{8} + \frac{28}{3}\right) \frac{L_\delta}{\max\{1, N_k(s, a, h)\}} \\
 &\leq 3\sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + 47\frac{L_\delta}{\max\{1, N_k(s, a, h)\}} \\
 &= \frac{1}{2} \cdot \epsilon_k^*(s' | s, a, h)
 \end{aligned} \tag{27}$$

Since we assumed that $P \in \mathcal{P}_k$,

$$|P(s' | s, a, h) - \bar{P}_k(s' | s, a, h)| \leq \frac{1}{2} \cdot \epsilon_k^*(s' | s, a, h).$$

Moreover, for any $\hat{P} \in \mathcal{P}_k$, we have

$$|\hat{P}(s' | s, a, h) - \bar{P}_k(s' | s, a, h)| \leq \epsilon_k(s' | s, a, h) \leq \frac{1}{2} \cdot \epsilon_k^*(s' | s, a, h).$$

By the triangle inequality, it follows that

$$|\hat{P}(s' | s, a, h) - P(s' | s, a, h)| \leq \epsilon_k^*(s' | s, a, h),$$

691 as required. □

692 We note that the above lemma holds when we replace $P(s' | s, a, h)$ of $\epsilon_k^*(s' | s, a, h)$ into $\hat{P}(s' |$
 693 $s, a, h)$ for any $\hat{P} \in \mathcal{P}_k$. Specifically, under the good event \mathcal{E} , we have for $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times$
 694 $\mathcal{S} \times [H]$,

$$|\hat{P}(s' | s, a, h) - P(s' | s, a, h)| \leq 6\sqrt{\frac{\hat{P}(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + 94\frac{L_\delta}{\max\{1, N_k(s, a, h)\}}. \tag{28}$$

It can be obtained by applying

$$\bar{P}_k(s' | s, a, h) \leq \hat{P}(s' | s, a, h) + \sqrt{\frac{8\bar{P}_k(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + \frac{28L_\delta}{3\max\{1, N_k(s, a, h)\}}$$

695 with the same argument for the remaining part of the proof.

696 12 Missing Proofs for Section 3: Tighter Function Estimators

697 **Proof of Lemma 4.** The proof is based on Lemma 10 of [Chen & Luo \(2021\)](#) with further so-
 698 phisticated evaluations. We consider an arbitrary cost function $g : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [-B, B]$
 699 for some boundedness constant $B > 0$. Let $\mathbf{q}_{(s', h+1)}^{P_k, \pi_k}, \mathbf{q}_{(s', h+1)}^{P, \pi_k}, \mathbf{g}$ be the vector representa-
 700 tions of $q^{P_k, \pi_k}(\cdot | s', h+1), q^{P, \pi_k}(\cdot | s', h+1) : \mathcal{S} \times \mathcal{A} \times \{h+1, \dots, H\} \rightarrow [0, 1]$, and

701 $g_{(h+1)} : \mathcal{S} \times \mathcal{A} \times \{h+1, \dots, H\} \rightarrow [-B, B]$ respectively. Note that

$$\begin{aligned}
& \left| \sum_{(s,a,s',h)} q_k(s,a,h) ((P - P_k)(s' | s,a,h)) (V_{h+1}^{\pi_k}(s';g,P_k) - V_{h+1}^{\pi_k}(s';g,P)) \right| \\
& \leq \sum_{(s,a,s',h)} q_k(s,a,h) \epsilon_k^*(s' | s,a,h) |(V_{h+1}^{\pi_k}(s';g,P_k) - V_{h+1}^{\pi_k}(s';g,P))| \\
& = \sum_{(s,a,s',h)} q_k(s,a,h) \epsilon_k^*(s' | s,a,h) \left| \langle \mathbf{q}_{(s',h+1)}^{P_k, \pi_k} - \mathbf{q}_{(s',h+1)}^{P, \pi_k}, \mathbf{g}_{(h+1)} \rangle \right| \\
& \leq BH \sum_{(s,a,s',h)} q_k(s,a,h) \epsilon_k^*(s' | s,a,h) \sum_{\substack{(s'',a'',s'''), \\ m \geq h+1}} q_k(s'',a'',m | s',h+1) \epsilon_k^*(s''' | s'',a'',m)
\end{aligned}$$

where the first inequality is from Lemma 13, the first equality holds because $V_{h+1}^{\pi_k}(s';g,P_k) = \langle \mathbf{q}_{(s',h+1)}^{P_k, \pi_k}, \mathbf{g}_{(h+1)} \rangle$ and $V_{h+1}^{\pi_k}(s';g,P) = \langle \mathbf{q}_{(s',h+1)}^{P, \pi_k}, \mathbf{g}_{(h+1)} \rangle$, the second inequality is due to Lemma 18. Remember that the definition of ϵ_k^* is given by

$$\epsilon_k^*(s' | s,a,h) = 6 \sqrt{\frac{P(s' | s,a,h)L_\delta}{\max\{1, N_k(s,a,h)\}}} + 94 \frac{L_\delta}{\max\{1, N_k(s,a,h)\}}.$$

702 Then it follows that

$$\begin{aligned}
& L_\delta^{-2} \sum_{(s,a,s',h)} q_k(s,a,h) \epsilon_k^*(s' | s,a,h) \sum_{(s'',a'',s'''), m \geq h+1} q_k(s'',a'',m | s',h+1) \epsilon_k^*(s''' | s'',a'',m) \\
& \leq 36 \underbrace{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \sqrt{\frac{q_k(s,a,h)^2 P(s' | s,a,h)}{\max\{1, N_k(s,a,h)\}}} \sqrt{\frac{q_k(s'',a'',m | s',h+1)^2 P(s''' | s'',a'',m)}{\max\{1, N_k(s'',a'',m)\}}}}_{\text{Term 1}} \\
& \quad + 564 \underbrace{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \sqrt{\frac{q_k(s,a,h)^2 P(s' | s,a,h)}{\max\{1, N_k(s,a,h)\}}} \frac{q_k(s'',a'',m | s',h+1)}{\max\{1, N_k(s'',a'',m)\}}}_{\text{Term 2}} \\
& \quad + 564 \underbrace{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} \sqrt{\frac{q_k(s'',a'',m | s',h+1)^2 P(s''' | s'',a'',m)}{\max\{1, N_k(s'',a'',m)\}}}}_{\text{Term 3}} \\
& \quad + 8836 \underbrace{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} \frac{q_k(s'',a'',m | s',h+1)}{\max\{1, N_k(s'',a'',m)\}}}_{\text{Term 4}}.
\end{aligned}$$

703 Term 1 can be bounded as follows.

$$\begin{aligned}
 \text{Term 1} &\leq \sqrt{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \frac{q_k(s,a,h)P(s''' | s'',a'',m)q_k(s'',a'',m | s',h+1)}{\max\{1, N_k(s,a,h)\}}} \\
 &\quad \times \sqrt{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \frac{q_k(s'',a'',m | s',h+1)P(s' | s,a,h)q_k(s,a,h)}{\max\{1, N_k(s'',a'',m)\}}} \\
 &\leq \sqrt{HS \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}} \sqrt{HS \sum_{(s'',a'',m)} \frac{q_k(s'',a'',m)}{\max\{1, N_k(s'',a'',m)\}}} \\
 &= HS \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}
 \end{aligned}$$

704 where the first inequality is from the Cauchy-Schwarz inequality. We can bound Term 2 as the
 705 following argument.

$$\begin{aligned}
 \text{Term 2} &\leq \sqrt{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \frac{q_k(s,a,h)q_k(s'',a'',m | s',h+1)}{\max\{1, N_k(s,a,h)\} \max\{1, N_k(s'',a'',m)\}}} \\
 &\quad \times \sqrt{\sum_{\substack{(s,a,s',h), \\ (s'',a'',s'''), \\ m \geq h+1}} \frac{q_k(s'',a'',m | s',h+1)P(s' | s,a,h)q_k(s,a,h)}{\max\{1, N_k(s'',a'',m)\}}} \\
 &\leq \sqrt{HS^2 \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}} \sqrt{HS \sum_{(s'',a'',m)} \frac{q_k(s'',a'',m)}{\max\{1, N_k(s'',a'',m)\}}} \\
 &= HS^{1.5} \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}.
 \end{aligned}$$

Similar to Term 2, we have an upper bound on Term 3 as follows.

$$\text{Term 3} = HS^{1.5} \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}.$$

706 Since $1/\max\{1, N_k(s,a,h)\} \leq 1$, we bound Term 4 in the following way.

$$\text{Term 4} \leq HS^2 \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}.$$

707 Finally, we deduce that

$$\begin{aligned}
 &\left| \sum_{(s,a,s',h)} q_k(s,a,h) (P - P_k)(s' | s,a,h) (V_{h+1}^{\pi_k}(s'; g, P_k) - V_{h+1}^{\pi_k}(s'; g, P)) \right| \\
 &\leq 10^4 BH^2 S^2 L_\delta^2 \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}}
 \end{aligned}$$

708 as desired.

709

□

Proof of Lemma 5. Let π_k be a policy for episode k . Moreover, let $P_k \in \mathcal{P}_k$, and let $g : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, 1]$ be an arbitrary cost function. Then we may define the occupancy measure $\hat{q}_k = q^{P_k, \pi_k}$ associated with policy π_k and transitional kernel P_k . Then we know that $V_1^{\pi_k}(\hat{V}_k, P_k) = \langle \hat{q}_k, \hat{V}_k \rangle$. Moreover, it follows from Lemma 19 that

$$\langle \hat{q}_k, \hat{V}_k \rangle \leq \text{Var} [\langle \hat{n}_k, \mathbf{g} \rangle \mid g, \pi_k, P_k]$$

710 where \hat{n}_k is a vector representation of $\hat{n}_k = n^{P_k, \pi_k}$. Furthermore, by Lemma 15 with $B = 1$, we
711 have

$$\begin{aligned} \text{Var} [\langle \hat{n}_k, \mathbf{g} \rangle \mid g, \pi_k, P_k] &\leq \mathbb{E}[\langle \hat{n}_k, \mathbf{g} \rangle^2 \mid g, \pi_k, P_k] \\ &\leq 2\langle \hat{q}_k, \vec{h} \odot \mathbf{g} \rangle \\ &\leq 2H^2 \end{aligned}$$

712 as desired. □

713 Having proved Lemmas lemma 4 and 5, we are ready to prove Theorem 1 which is the crucial part
714 of deducing our tighter function estimators.

715 **Proof of Theorem 1.** We assume that the good event \mathcal{E} holds, which holds with probability at least
716 $1 - 14\delta$ according to Lemma 12. We observe that $|V_1^{\pi_k}(g, P) - V_1^{\pi_k}(g, P_k)|$ can be rewritten by
717 $|\langle \mathbf{g}, \mathbf{q}_k - \hat{q}_k \rangle|$ using occupancy measures. By Lemma 17, it follows that

$$\begin{aligned} &|\langle \mathbf{g}, \mathbf{q}_k - \hat{q}_k \rangle| \\ &= \left| \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \hat{q}_k(s, a, h) (P - P_k)(s' \mid s, a, h) V_{h+1}^{\pi_k}(s'; g, P) \right| \\ &\leq \underbrace{\left| \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \hat{q}_k(s, a, h) (P - P_k)(s' \mid s, a, h) V_{h+1}^{\pi_k}(s'; g, P_k) \right|}_{\text{Term 1}} \\ &\quad + \underbrace{\left| \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \hat{q}_k(s, a, h) (P - P_k)(s' \mid s, a, h) (V_{h+1}^{\pi_k}(s'; g, P) - V_{h+1}^{\pi_k}(s'; g, P_k)) \right|}_{\text{Term 2}} \end{aligned}$$

718 where $(P - P_k)(s' \mid s, a, h) = P(s' \mid s, a, h) - P_k(s' \mid s, a, h)$.

719 To bound Term 2, we use bound

$$P(s' \mid s, a, h) - P_k(s' \mid s, a, h) \leq 6\sqrt{\frac{P_k(s' \mid s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} + 94\frac{L_\delta}{\max\{1, N_k(s, a, h)\}}$$

as explained in (28). This is because $\hat{q}_k = q^{P_k, \pi_k}$ is an occupancy measure with respect to $P_k \in \mathcal{P}_k$, not P . Then we can apply Lemma 4 and obtain

$$\text{Term 2} \leq 10^4 H^2 S^2 L_\delta^2 \sum_{(s, a, h)} \frac{\hat{q}_k(s, a, h)}{\max\{1, N_k(s, a, h)\}}.$$

720 Next, we bound Term 1. Note that $\sum_{s'} (P(s' | s, a, h) - P_k(s' | s, a, h)) = 0$. Then it follows that

$$\begin{aligned}
 \text{Term 1} &= \left| \sum_{(s,a,s',h)} \widehat{q}_k(s, a, h) (P - P_k)(s' | s, a, h) (V_{h+1}^{\pi_k}(g, P_k) - \widehat{\mu}_k(s, a, h)) \right| \\
 &\leq 2 \sum_{(s,a,s',h)} \widehat{q}_k(s, a, h) \epsilon_k(s' | s, a, h) |V_{h+1}^{\pi_k}(g, P_k) - \widehat{\mu}_k(s, a, h)| \\
 &= 4 \underbrace{\sum_{(s,a,s',h)} \widehat{q}_k(s, a, h) \sqrt{\frac{\bar{P}_k(s' | s, a, h) L_\delta}{\max\{1, N_k(s, a, h) - 1\}}} |V_{h+1}^{\pi_k}(s'; g, P_k) - \widehat{\mu}_k(s, a, h)|}_{\text{Term 3}} \\
 &\quad + \underbrace{\frac{28}{3} \sum_{(s,a,s',h)} \widehat{q}_k(s, a, h) \frac{L_\delta}{\max\{1, N_k(s, a, h) - 1\}} |V_{h+1}^{\pi_k}(s'; g, P) - \widehat{\mu}_k(s, a, h)|}_{\text{Term 4}}
 \end{aligned}$$

where $\widehat{\mu}_k(s, a, h) = \mathbb{E}_{s' \sim P_k(\cdot | s, a, h)} [V_{h+1}^{\pi_k}(s'; g, P_k)]$. The first inequality is from $|(P - P_k)(s' | s, a, h)| \leq |(P - \bar{P}_k)(s' | s, a, h)| + |(\bar{P}_k - P_k)(s' | s, a, h)| \leq 2\epsilon_k(s' | s, a, h)$ for any $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$ under the good event \mathcal{E} . We note that $\bar{P}_k(s' | s, a, h) \leq P_k(s' | s, a, h) + \epsilon_k(s' | s, a, h)$ and define

$$\widehat{V}_k(s, a, h) = \sum_{s'} P_k(s' | s, a, h) |V_{h+1}^{\pi_k}(s'; g, P_k) - \widehat{\mu}_k(s, a, h)|^2.$$

721 Then we can bound Term 3 as the following.

$$\begin{aligned}
 &\text{Term 3} \\
 &\leq \sqrt{L_\delta} \sum_{(s,a,s',h)} \widehat{q}_k(s, a, h) \sqrt{\frac{(P_k + \epsilon_k)(s' | s, a, h)}{\max\{1, N_k(s, a, h) - 1\}}} |V_{h+1}^{\pi_k}(s'; g, P_k) - \widehat{\mu}_k(s, a, h)| \\
 &\leq \sqrt{L_\delta} \sqrt{\sum_{(s,a,s',h)} \widehat{q}_k(s, a, h) (P_k + \epsilon_k)(s' | s, a, h) |V_{h+1}^{\pi_k}(s'; g, P_k) - \widehat{\mu}_k(s, a, h)|^2} \\
 &\quad \times \sqrt{\sum_{(s,a,s',h)} \frac{\widehat{q}_k(s, a, h)}{\max\{1, N_k(s, a, h) - 1\}}} \\
 &\leq \sqrt{L_\delta} \sqrt{\sum_{(s,a,h)} \widehat{q}_k(s, a, h) \widehat{V}_k(s, a, h) + 4H^2 \sum_{(s,a,s',h)} \widehat{q}_k(s, a, h) \epsilon_k(s' | s, a, h)} \\
 &\quad \times \sqrt{\sum_{(s,a,s',h)} \frac{\widehat{q}_k(s, a, h)}{\max\{1, N_k(s, a, h) - 1\}}}
 \end{aligned}$$

722 where the second inequality follows from the Cauchy-Schwarz inequality and the last inequality is
 723 due to $|V_{h+1}^{\pi_k}(s'; g, P_k) - \widehat{\mu}_k(s, a, h)| \leq 2H$.

By Lemma 5, we deduce that

$$\sum_{(s,a,h)} \widehat{q}_k(s, a, h) \widehat{V}_k(s, a, h) \leq 2H^2.$$

724 Due to the AM-GM inequality, we have

$$\begin{aligned}
& \sqrt{2H^2 + 4H^2 \sum_{(s,a,s',h)} \widehat{q}_k(s,a,h) \epsilon_k(s' | s,a,h)} \sqrt{\sum_{(s,a,s',h)} \frac{\widehat{q}_k(s,a,h)}{\max\{1, N_k(s,a,h) - 1\}}} \\
& \leq \left(\sqrt{2H^2} + \sqrt{4H^2 \sum_{(s,a,s',h)} \widehat{q}_k(s,a,h) \epsilon_k(s' | s,a,h)} \right) \sqrt{\sum_{(s,a,s',h)} \frac{\widehat{q}_k(s,a,h)}{\max\{1, N_k(s,a,h) - 1\}}} \\
& \leq \frac{H^2}{\alpha_1} + \frac{2H^2}{\alpha_2} \sum_{(s,a,s',h)} \widehat{q}_k(s,a,h) \epsilon_k(s' | s,a,h) + \frac{\alpha_1 + \alpha_2}{2} \sum_{(s,a,h)} \frac{S \cdot \widehat{q}_k(s,a,h)}{\max\{1, N_k(s,a,h) - 1\}}
\end{aligned}$$

725 for any $\alpha_1, \alpha_2 > 0$. By taking $\alpha_1 = \sqrt{HKL_\delta}/(S\sqrt{A})$, $\alpha_2 = \sqrt{H^3L_\delta}$, we obtain

Term 3

$$\begin{aligned}
& \leq \sum_{(s,a,h)} \widehat{q}_k(s,a,h) \left(\frac{S\sqrt{HA}}{\sqrt{K}} + 2\sqrt{H} \sum_{s'} \epsilon_k(s' | s,a,h) + \frac{\sqrt{HK} + \sqrt{H^3S^2A}}{2\sqrt{A}} \frac{L_\delta}{\max\{1, N_k(s,a,h) - 1\}} \right) \\
& \leq \sum_{(s,a,h)} \widehat{q}_k(s,a,h) \left(\frac{S\sqrt{HA}}{\sqrt{K}} + 2\sqrt{H} \epsilon_k(s,a,h) + \frac{\sqrt{HK} + \sqrt{H^3S^2A}}{2\sqrt{A}} \frac{L_\delta}{\max\{1, N_k(s,a,h) - 1\}} \right).
\end{aligned}$$

726 Note that the last inequality follows from

$$\begin{aligned}
\sum_{s'} \epsilon_k(s' | s,a,h) &= \sum_{s'} \left(\sqrt{\frac{4\bar{P}_k(s' | s,a,h)L_\delta}{\max\{1, N_k(s,a,h) - 1\}}} + \frac{14L_\delta}{3\max\{1, N_k(s,a,h) - 1\}} \right) \\
&\leq \sqrt{S} \sqrt{\frac{4\sum_{s'} \bar{P}_k(s' | s,a,h)L_\delta}{\max\{1, N_k(s,a,h) - 1\}}} + \frac{14SL_\delta}{3\max\{1, N_k(s,a,h) - 1\}} \\
&= \sqrt{\frac{4SL_\delta}{\max\{1, N_k(s,a,h) - 1\}}} + \frac{14SL_\delta}{3\max\{1, N_k(s,a,h) - 1\}} \\
&= \epsilon_k(s,a,h)
\end{aligned}$$

727 where the inequality is due to the Cauchy-Schwarz inequality and the second equality is due to

728 $\sum_{s'} \bar{P}_k(s' | s,a,h) \leq 1$.

Since $|V_{h+1}^{\pi_k}(s'; g, P) - \widehat{\mu}_k(s,a,h)| \leq 2H$, Term 4 can be bounded as follows.

$$\text{Term 4} \leq 2HSL_\delta \sum_{(s,a,h)} \frac{\widehat{q}_k(s,a,h)}{\max\{1, N_k(s,a,h) - 1\}}.$$

729 Finally, we proved that

$$\begin{aligned}
& |\langle \mathbf{g}, \mathbf{q}_k - \widehat{\mathbf{q}}_k \rangle| \\
& \leq 4 \cdot (\text{Term 3}) + \frac{28}{3} \cdot (\text{Term 4}) + (\text{Term 2}) \\
& \leq \sum_{(s,a,h)} \widehat{q}_k(s,a,h) \left(\frac{4S\sqrt{HA}}{\sqrt{K}} + 8\sqrt{H} \epsilon_k(s,a,h) + \frac{2\sqrt{HK}L_\delta}{\sqrt{A}\max\{1, N_k(s,a,h) - 1\}} \right) \\
& \quad + \left(\left(\frac{56}{3}HS + 2H^{1.5}S \right) L_\delta + 10^4H^2S^2L_\delta^2 \right) \sum_{(s,a,h)} \frac{\widehat{q}_k(s,a,h)}{\max\{1, N_k(s,a,h) - 1\}}
\end{aligned}$$

730 as required. \square

731 **13 Missing Proofs for Section 4: Safe Exploration**

 732 In this section, we prove Lemma 6 that provides an asymptotic upper bound on a sufficient number
 733 of episodes executing π_b , which is denoted by K_0 , for feasibility of (10).

Lemma 14. *Assume that the good event \mathcal{E} holds. Let π_k be any policy for episode k , and let P be the true transition kernel. Let q_k denote the occupancy measure q^{P, π_k} associated with π_k and P . For R_k, U_k , we have*

$$\sum_{k=1}^K \langle R_k + U_k, q_k \rangle = \mathcal{O} \left(\left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) L_\delta^3 \right).$$

 734 *Proof.* Note that $\sum_{k=1}^K \langle R_k + U_k, q_k \rangle$ can be rewritten as

$$\begin{aligned} & \sum_{k=1}^K \langle R_k + U_k, q_k \rangle \\ &= \sum_{k=1}^K \sum_{(s,a,h)} q_k(s,a,h) \sqrt{\frac{L_\delta}{\max\{1, N_k(s,a,h)\}}} \\ & \quad + \sum_{k=1}^K \sum_{(s,a,h)} q_k(s,a,h) \left(\frac{4S\sqrt{HA}}{\sqrt{K}} + 8\sqrt{H}\varepsilon_k(s,a,h) + \frac{2(\sqrt{HK} + \sqrt{H^3 S^2 A})L_\delta}{\sqrt{A} \max\{1, N_k(s,a,h) - 1\}} \right) \\ & \quad + \left(\frac{56}{3} HSL_\delta + 10^4 H^2 S^2 L_\delta^2 \right) \sum_{k=1}^K \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h) - 1\}}. \end{aligned}$$

 Since $\sum_{(s,a,h)} \hat{q}_k(s,a,h) = H$, we have

$$\sum_{k=1}^K \sum_{(s,a,h)} q_k(s,a,h) \cdot \frac{4S\sqrt{HA}}{\sqrt{K}} = \mathcal{O}(H^{1.5} S \sqrt{AK}).$$

735 Furthermore, Lemma 11 implies that

$$\begin{aligned} & \sum_{k=1}^K \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h)\}} = \mathcal{O}(HSA \ln K + H \ln(H/\delta)), \\ & \sum_{k=1}^K \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\sqrt{\max\{1, N_k(s,a,h)\}}} = \mathcal{O}(H\sqrt{SAK} + HSA \ln K + H \ln(H/\delta)). \end{aligned}$$

736 Then it follows that

$$\sum_{k=1}^K \sum_{(s,a,h)} q_k(s,a,h) \sqrt{\frac{L_\delta}{\max\{1, N_k(s,a,h)\}}} = \mathcal{O} \left((H\sqrt{SAK} + HSA) L_\delta^2 \right).$$

 737 Since $\max\{1, N_k(s,a,h) - 1\} \geq \frac{1}{2} \max\{1, N_k(s,a,h)\}$, we have

$$\sum_{k=1}^K \sum_{(s,a,h)} q_k(s,a,h) \frac{(\sqrt{HK} + \sqrt{H^3 S^2 A})L_\delta}{\sqrt{A} \max\{1, N_k(s,a,h) - 1\}} = \mathcal{O} \left((H^{1.5} S \sqrt{AK} + H^{2.5} S^2 A) L_\delta^2 \right),$$

738 and moreover,

$$(HSL_\delta + H^2 S^2 L_\delta^2) \sum_{k=1}^K \sum_{(s,a,h)} \frac{q_k(s,a,h)}{\max\{1, N_k(s,a,h) - 1\}} = \mathcal{O} \left(H^3 S^3 A L_\delta^3 \right).$$

739 Next, by Lemma 11, $\sum_{k=1}^K \sum_{(s,a,h)} q_k(s, a, h) \left(\sqrt{H} \varepsilon_k(s, a, h) \right)$ can be bounded as follows.

$$\begin{aligned} & \sum_{k=1}^K \sum_{(s,a,h)} q_k(s, a, h) \left(\sqrt{H} \varepsilon_k(s, a, h) \right) \\ &= \sqrt{H} \sum_{k=1}^K \sum_{(s,a,h)} q_k(s, a, h) \left(\sqrt{\frac{4SL_\delta}{\max\{1, N_k(s, a, h) - 1\}}} + \frac{14SL_\delta}{3 \max\{1, N_k(s, a, h) - 1\}} \right) \\ &= \mathcal{O} \left(\left(H^{1.5} S \sqrt{AK} + H^{1.5} S^2 A \right) L_\delta^2 \right). \end{aligned}$$

As a result, we have proved that

$$\sum_{k=1}^K \langle \mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_k \rangle = \mathcal{O} \left(\left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) L_\delta^3 \right),$$

740 as required. □

741 We are ready to prove Lemma 6 based on Lemma 14.

742 **Proof of Lemma 6.** We closely follow the proof of (Bura et al., 2022, Proposition 4). We assume
743 that the good event \mathcal{E} holds, which holds with probability at least $1 - 14\delta$. Let $q_b = q^{P, \pi_b}$ be the
744 occupancy measure associated with the safe baseline policy π_b and the true transition kernel P . Then
745 q_b is a feasible solution of (13) if $\langle \hat{\mathbf{g}}_k, \mathbf{q}_b \rangle \leq \bar{C}$ holds. To find a sufficient condition, we deduce that

$$\begin{aligned} \langle \hat{\mathbf{g}}_k, \mathbf{q}_b \rangle &= \langle \bar{\mathbf{g}}_k + \mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle \\ &\leq \langle \mathbf{g} + 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle \\ &= \bar{C}_b + \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle \end{aligned}$$

746 where the first equality is from the definition of $\hat{\mathbf{g}}_k$, the inequality is from Lemma 3, and the last
747 equality follows from $\langle \mathbf{g}, \mathbf{q}_b \rangle = \bar{C}_b$. This implies that a sufficient condition for $\langle \hat{\mathbf{g}}_k, \mathbf{q}_b \rangle \leq \bar{C}$ is
748 given by

$$\langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle < \bar{C} - \bar{C}_b. \quad (29)$$

Note that $\langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle$ decreases as k increases because

$$\frac{1}{\max\{1, N_k(s, a, h)\}}, \quad \frac{1}{\sqrt{\max\{1, N_k(s, a, h)\}}}$$

can only decrease as k increases. Then suppose that K_0 is the last episode where (29) does not hold. By definition, $K_0 + 1$ is the first episode satisfying $\langle \hat{\mathbf{g}}_k, \mathbf{q}_b \rangle < \bar{C}$. Due to the strict inequality, occupancy measures other than q_b can be potentially feasible to (13). This implies that DOPE+ can sufficiently explore policies other than π_b after K_0 episodes. Then we have

$$K_0(\bar{C} - \bar{C}_b) < \sum_{k=1}^{K_0} \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle.$$

Since q_b induces the true transition kernel, we can apply Lemma 14. Then the right-hand side is bounded as follows.

$$\sum_{k=1}^{K_0} \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle = \tilde{\mathcal{O}} \left(H^{1.5} S \sqrt{AK_0} \right).$$

Hence, K_0 satisfies

$$K_0 = \tilde{\mathcal{O}} \left(\frac{H^3 S^2 A}{(\bar{C} - \bar{C}_b)^2} \right).$$

Then we have

$$\langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle \leq \langle 2\mathbf{R}_{K_0+1} + \mathbf{U}_{K_0+1}, \mathbf{q}_b \rangle \leq \bar{C} - \bar{C}_b \quad \forall k = K_0 + 1, \dots, K.$$

749 This implies that (10) is feasible after episode K_0 when (π_b, P) becomes a feasible solution in
 750 episode K_0 . \square

751 14 Detailed Proofs for the Regret Analysis

752 In this section, we prove Theorem 2 that guarantees zero constraint violation for DOPE+. Next, we
 753 provide the proofs of Lemmas 7, 8 and 9. Lastly, we show Theorem 3 that gives us the regret upper
 754 bound.

755 14.1 Details of Constraint Violation Analysis

756 **Proof of Theorem 2.** We assume that the good event \mathcal{E} holds, which is the case with probability
 757 at least $1 - 14\delta$. Let π_k, P_k denote the policy and the transition kernel obtained from DOPE+ for
 758 episode k , respectively. Let $q_k = q^{P, \pi_k}, \hat{q}_k = q^{P_k, \pi_k}$. We know that the constraint is satisfied if
 759 $V_1^{\pi_k}(g, P) = \langle \mathbf{g}, \mathbf{q}_k \rangle \leq \bar{C}$ for each $k \in [K]$. For $k \leq K_0$, there is no constraint violation because
 760 we take $\pi_k = \pi_b$. Now we consider the case when $k > K_0$. We have

$$\begin{aligned} \langle \mathbf{g}, \mathbf{q}_k \rangle &= \langle \mathbf{g}, \hat{\mathbf{q}}_k \rangle + \langle \mathbf{g}, \mathbf{q}_k - \hat{\mathbf{q}}_k \rangle \\ &\leq \langle \bar{\mathbf{g}}_k + \mathbf{R}_k, \hat{\mathbf{q}}_k \rangle + \langle \mathbf{g}, \mathbf{q}_k - \hat{\mathbf{q}}_k \rangle \\ &\leq \langle \bar{\mathbf{g}}_k + \mathbf{R}_k, \hat{\mathbf{q}}_k \rangle + \langle \mathbf{U}_k, \hat{\mathbf{q}}_k \rangle \\ &= \langle \hat{\mathbf{g}}_k, \hat{\mathbf{q}}_k \rangle \\ &\leq \bar{C} \end{aligned}$$

761 where the first inequality follows from Lemma 3, the second inequality is from Theorem 1, and the
 762 last inequality is due to the update rule of DOPE+. This implies that π_k holds $\langle \mathbf{g}, \mathbf{q}_k \rangle \leq \bar{C}$ for
 763 $k > K_0$. Thus, we showed that $\text{Violation}(\bar{\pi}) = 0$ with probability at least $1 - 14\delta$. \square

764 14.2 Details of Regret Analysis

Proof of Lemma 7. We closely follow the proof of (Bura et al., 2022, Lemma 18). We assume that
 the good event \mathcal{E} holds, which is the case with probability at least $1 - 14\delta$. We observe that

$$\sum_{k=K_0+1}^K \left(V_1^{\pi^*}(f, P) - V_1^{\pi_k}(f_k, P_k) \right) = \sum_{k=K_0+1}^K \langle \mathbf{f}, \mathbf{q}^* \rangle - \sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k \rangle.$$

765 By Lemma 10, there exist $\bar{q}_b(s, a, s', h)$ and $\bar{q}^*(s, a, s', h)$ such that $q_b(s, a, h) =$
 766 $\sum_{s' \in \mathcal{S}} \bar{q}_b(s, a, s', h)$ and $q^*(s, a, h) = \sum_{s' \in \mathcal{S}} \bar{q}^*(s, a, s', h)$, respectively. Then we define the
 767 new occupancy measure $q_{\alpha_k}(s, a, h)$ satisfying $q_{\alpha_k}(s, a, h) = \sum_{s' \in \mathcal{S}} \bar{q}_{\alpha_k}(s, a, s', h)$ where

$$\bar{q}_{\alpha_k}(s, a, s', h) = (1 - \alpha_k) \bar{q}_b(s, a, s', h) + \alpha_k \bar{q}^*(s, a, s', h) \quad (30)$$

768 for $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$ and $\alpha_k \in [0, 1]$. Now we verify (C1),(C2) and (C3) in Lemma 10
 769 to say q_{α_k} is a valid occupancy measure. Since \bar{q}_{α_k} is a convex combination of \bar{q}_b and \bar{q}^* , (C1),(C2)
 770 hold. For (C3), we can show that q_{α_k} induces the true transition kernel P as follows. Since we
 771 know q_b and q^* induce P , it follows that $\bar{q}_b(s, a, s', h) = P(s' | s, a, h) \sum_{s'' \in \mathcal{S}} \bar{q}_b(s, a, s'', h)$ and
 772 $\bar{q}^*(s, a, s', h) = P(s' | s, a, h) \sum_{s'' \in \mathcal{S}} \bar{q}^*(s, a, s'', h)$ for $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$. Then
 773 $\bar{q}_{\alpha_k}(s, a, s', h) = P(s' | s, a, h) \sum_{s'' \in \mathcal{S}} \bar{q}_{\alpha_k}(s, a, s'', h)$ can be derived from (30), which implies
 774 that q_{α_k} induces the true transition kernel P . Hence, q_{α_k} is a valid occupancy measure inducing the
 775 true transition kernel P .

776 To use the optimality of \hat{q}_k in our analysis, we expect that q_{α_k} is a feasible solution for (13). Under
 777 the good event \mathcal{E} , we know that $q_{\alpha_k} \in \Delta(P, k)$ due to $P \in \mathcal{P}_k$. Then it is sufficient to find a
 778 condition for α_k satisfying $\langle \hat{g}_k, \mathbf{q}_{\alpha_k} \rangle \leq \bar{C}$. We deduce that

$$\begin{aligned} \langle \hat{g}_k, \mathbf{q}_{\alpha_k} \rangle &= \langle \bar{g}_k + \mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle \\ &\leq \langle \mathbf{g} + 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle \\ &= (1 - \alpha_k) \langle \mathbf{g} + 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle + \alpha_k \langle \mathbf{g} + 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}^* \rangle \\ &\leq (1 - \alpha_k) (\bar{C}_b + \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle) + \alpha_k (\bar{C} + \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}^* \rangle) \end{aligned}$$

779 where the first inequality is from Lemma 3 and the last inequality is from $\langle \mathbf{g}, \mathbf{q}_b \rangle = \bar{C}_b$ and
 780 $\langle \mathbf{g}, \mathbf{q}^* \rangle \leq \bar{C}$. Furthermore, the second equality is true because (30) implies that $q_{\alpha_k}(s, a, h) =$
 781 $(1 - \alpha_k)q_b(s, a, h) + \alpha_k q^*(s, a, h)$. Hence, a sufficient condition of α_k for $\langle \hat{g}_k, \mathbf{q}_{\alpha_k} \rangle \leq \bar{C}$ is given
 782 by

$$\alpha_k \leq \frac{\bar{C} - \bar{C}_b - \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle}{\bar{C} - \bar{C}_b + \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}^* \rangle - \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle}.$$

783 Remember that, in the proof of Lemma 6, we defined K_0 so that $K_0 + 1$ is the first episode satisfying
 784 $\langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle \leq \bar{C} - \bar{C}_b$. This guarantees that there exists some $\alpha_k \in [0, 1]$ satisfying the above
 785 inequality for $k > K_0$.

786 Now, for some α_k , we claim that

$$\langle \mathbf{f}, \mathbf{q}^* \rangle \leq \langle \bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle. \quad (31)$$

To show (31), we first take for $\beta \geq 1$,

$$\mathbf{f}_\beta = \bar{\mathbf{f}}_k + 3\beta \mathbf{R}_k + \beta \mathbf{U}_k.$$

787 Then we find α_k, β satisfying $\langle \mathbf{f}, \mathbf{q}^* \rangle \leq \langle \mathbf{f}_\beta, \mathbf{q}_{\alpha_k} \rangle$. By Lemma 3, we have

$$\begin{aligned} \langle \mathbf{f}_\beta, \mathbf{q}_{\alpha_k} \rangle &= \langle \bar{\mathbf{f}}_k + 3\beta \mathbf{R}_k + \beta \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle \\ &\geq \langle \mathbf{f} + 2\beta \mathbf{R}_k + \beta \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle \\ &= (1 - \alpha_k) \langle \mathbf{f} + 2\beta \mathbf{R}_k + \beta \mathbf{U}_k, \mathbf{q}_b \rangle + \alpha_k \langle \mathbf{f} + 2\beta \mathbf{R}_k + \beta \mathbf{U}_k, \mathbf{q}^* \rangle. \end{aligned}$$

788 We have $\langle \mathbf{f}, \mathbf{q}^* \rangle \leq \langle \mathbf{f}_\beta, \mathbf{q}_{\alpha_k} \rangle$ if β satisfies

$$\beta \geq \frac{(1 - \alpha_k) (\langle \mathbf{f}, \mathbf{q}^* \rangle - \langle \mathbf{f}, \mathbf{q}_b \rangle)}{(1 - \alpha_k) \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle + \alpha_k \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}^* \rangle}.$$

789 By taking

$$\alpha_k = \frac{\bar{C} - \bar{C}_b - \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle}{\bar{C} - \bar{C}_b + \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}^* \rangle - \langle 2\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_b \rangle}, \quad (32)$$

790 it follows that

$$\beta \geq \frac{\langle \mathbf{f}, \mathbf{q}^* \rangle - \langle \mathbf{f}, \mathbf{q}_b \rangle}{\bar{C} - \bar{C}_b}.$$

791 Since $\langle \mathbf{f}, \mathbf{q}^* \rangle - \langle \mathbf{f}, \mathbf{q}_b \rangle \leq H$, it is sufficient to take

$$\beta = \frac{H}{\bar{C} - \bar{C}_b}. \quad (33)$$

792 For α_k satisfying (32), we showed that q_{α_k} is a feasible solution for (13). Then it follows
 793 $\langle \hat{\mathbf{f}}_k, \mathbf{q}_{\alpha_k} \rangle \leq \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k \rangle$ due to optimality of \hat{q}_k . Furthermore, for β satisfying (33), we have (31).

794 Hence, we deduce that

$$\begin{aligned}
 & \langle \mathbf{f}, \mathbf{q}^* \rangle - \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k \rangle \\
 & \leq \langle \mathbf{f}, \mathbf{q}^* \rangle - \langle \hat{\mathbf{f}}_k, \mathbf{q}_{\alpha_k} \rangle \\
 & = \langle \mathbf{f}, \mathbf{q}^* \rangle - \langle \bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle \\
 & \quad + \langle \bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle - \langle \vec{\mathbf{B}} \wedge (\bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k), \mathbf{q}_{\alpha_k} \rangle \\
 & \leq \langle \bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle - \langle \vec{\mathbf{B}} \wedge (\bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k), \mathbf{q}_{\alpha_k} \rangle
 \end{aligned}$$

where the last inequality is from (31). Furthermore, under the good event \mathcal{E} , we know that $f_k(s, a, h) \leq B$ for $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \in [K]$, where $B = 1 + \sqrt{L_\delta}$. This implies that $\bar{f}_k(s, a, h) \leq B$. Thus, we have

$$\langle \bar{\mathbf{f}}_k, \mathbf{q}_{\alpha_k} \rangle \leq \langle \vec{\mathbf{B}} \wedge (\bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k), \mathbf{q}_{\alpha_k} \rangle.$$

795 Then it follows that

$$\begin{aligned}
 & \langle \bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle - \langle \vec{\mathbf{B}} \wedge (\bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k), \mathbf{q}_{\alpha_k} \rangle \\
 & \leq \langle \bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle - \langle \bar{\mathbf{f}}_k, \mathbf{q}_{\alpha_k} \rangle \\
 & = \langle \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle.
 \end{aligned}$$

796 Finally, we proved that

$$\langle \mathbf{f}, \mathbf{q}^* \rangle - \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k \rangle \leq \langle \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle.$$

797 By Lemma 14, we have

$$\begin{aligned}
 \sum_{k=K_0+1}^K \langle \mathbf{f}, \mathbf{q}^* \rangle - \sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k \rangle & \leq \sum_{k=K_0+1}^K \langle \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k, \mathbf{q}_{\alpha_k} \rangle \\
 & = \mathcal{O} \left(\left(\frac{H^{2.5}}{\bar{C} - \bar{C}_b} S \sqrt{AK} + \frac{H^4}{\bar{C} - \bar{C}_b} S^3 A \right) L_\delta^3 \right)
 \end{aligned}$$

798 as desired. \square

Proof of Lemma 8. The lemma is a direct consequence of Lemma 20 with $B = \mathcal{O}(L_\delta)$. Hence, we have

$$\sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k - \mathbf{q}_k \rangle = \mathcal{O} \left(\left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) L_\delta^4 \right)$$

799 with probability at least $1 - 2\delta$ under the good event \mathcal{E} . By taking the union bound, the statement
800 holds with probability at least $1 - 16\delta$. \square

Proof of Lemma 9. We assume that the good event \mathcal{E} holds, which is the case with probability at least $1 - 14\delta$. The left-hand side of Lemma 9 can be rewritten as

$$\sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k - \mathbf{f}, \mathbf{q}_k \rangle.$$

801 Under the good event \mathcal{E} , we have $\bar{f}_k(s, a, h) \leq f(s, a, h) + R_k(s, a, h)$ for $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$
 802 and $k \in [K]$. Furthermore, $H/(\bar{C} - \bar{C}_b) \geq 1$ due to $\bar{C} - \bar{C}_b \leq H$. Then it follows that

$$\begin{aligned} \sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k - \mathbf{f}, \mathbf{q}_k \rangle &= \sum_{k=K_0+1}^K \langle \bar{\mathbf{B}} \wedge (\bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k) - \mathbf{f}, \mathbf{q}_k \rangle \\ &\leq \sum_{k=K_0+1}^K \langle \bar{\mathbf{f}}_k + \frac{3H}{\bar{C} - \bar{C}_b} \mathbf{R}_k + \frac{H}{\bar{C} - \bar{C}_b} \mathbf{U}_k - \mathbf{f}, \mathbf{q}_k \rangle \\ &\leq \frac{H}{\bar{C} - \bar{C}_b} \sum_{k=K_0+1}^K \langle 4\mathbf{R}_k + \mathbf{U}_k, \mathbf{q}_k \rangle \\ &= \mathcal{O} \left(\left(\frac{H^{2.5}}{\bar{C} - \bar{C}_b} S \sqrt{AK} + \frac{H^4}{\bar{C} - \bar{C}_b} S^3 A \right) L_\delta^3 \right) \end{aligned}$$

803 where the last equality is due to Lemma 14. \square

804 **Proof of Theorem 3.** We assume that the good event \mathcal{E} holds, which is the case with probability at
 805 least $1 - 14\delta$. We decompose the regret as follows using occupancy measures.

$$\begin{aligned} \text{Regret}(\bar{\pi}) &= \underbrace{\sum_{k=1}^{K_0} \langle \mathbf{f}, \mathbf{q}^* \rangle - \sum_{k=1}^{K_0} \langle \mathbf{f}, \mathbf{q}_k \rangle}_{\text{(I)}} + \underbrace{\sum_{k=K_0+1}^K \langle \mathbf{f}, \mathbf{q}^* \rangle - \sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k \rangle}_{\text{(II)}} \\ &\quad + \underbrace{\sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k, \hat{\mathbf{q}}_k - \mathbf{q}_k \rangle}_{\text{(III)}} + \underbrace{\sum_{k=K_0+1}^K \langle \hat{\mathbf{f}}_k - \mathbf{f}, \mathbf{q}_k \rangle}_{\text{(IV)}}. \end{aligned}$$

806 As explained in Section 5.2, we can upper bound term (I) as

$$\tilde{\mathcal{O}} \left(\frac{H^4 S^2 A}{(\bar{C} - \bar{C}_b)^2} \right).$$

807 because $K_0 = \tilde{\mathcal{O}} \left(\frac{H^3 S^2 A}{(\bar{C} - \bar{C}_b)^2} \right)$ due to Lemma 6 and $\langle \mathbf{f}, \mathbf{q}^* \rangle \leq H$.

808 By Lemma 7, we have

$$\text{Term (II)} = \mathcal{O} \left(\left(\frac{H^{2.5}}{\bar{C} - \bar{C}_b} S \sqrt{AK} + \frac{H^4}{\bar{C} - \bar{C}_b} S^3 A \right) L_\delta^3 \right).$$

809 By Lemma 8, with probability at least $1 - 2\delta$, it follows that

$$\text{Term (III)} = \mathcal{O} \left(\left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) L_\delta^4 \right).$$

810 Moreover, it follows from Lemma 9 that

$$\text{Term (IV)} = \mathcal{O} \left(\left(\frac{H^{2.5}}{\bar{C} - \bar{C}_b} S \sqrt{AK} + \frac{H^4}{\bar{C} - \bar{C}_b} S^3 A \right) L_\delta^3 \right).$$

Hence, by taking the union bound,

$$\text{Regret}(\bar{\pi}) = \tilde{\mathcal{O}} \left(\frac{H}{\bar{C} - \bar{C}_b} \left(H^{1.5} S \sqrt{AK} + \frac{H^4 S^3 A}{\bar{C} - \bar{C}_b} \right) \right)$$

811 with probability at least $1 - 16\delta$. \square

812 **15 Technical Lemmas**

813 In this section, we provide technical lemmas that are crucial for our regret and constraint violation
 814 analysis. The following lemma is from (Chen & Luo, 2021) with a few modifications, and it is useful
 815 to bound the variance of $\langle \mathbf{n}_k, \mathbf{f}_k \rangle$.

Lemma 15. (Chen & Luo, 2021, Lemma 2) *Let π_k be any policy for episode k , and let q_k denote the occupancy measure q^{P, π_k} . Let $\ell : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [-B, B]$ be an arbitrary function, and let P be an arbitrary transition kernel. Then*

$$\mathbb{E} [\langle \mathbf{n}_k, \ell \rangle^2 \mid \ell, \pi_k, P] \leq 2B \langle \mathbf{q}_k, \vec{\mathbf{h}} \odot \ell \rangle$$

816 where $\mathbf{q}_k, \mathbf{n}_k, \ell$ are the vector representations of q_k, n_k, ℓ .

817 *Proof.* For ease of notation, let $\mathbb{E}_k [\cdot]$ denotes $\mathbb{E} [\cdot \mid \ell, \pi_k, P]$, and let s_h and a_h denote s_h^{P, π_k} and
 818 a_h^{P, π_k} , respectively for $h \in [H]$. Note that

$$\begin{aligned} \mathbb{E}_k [\langle \mathbf{n}_k, \ell \rangle^2] &= \mathbb{E}_k \left[\left(\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} n_k(s, a, h) \ell(s, a, h) \right)^2 \right] \\ &= \mathbb{E}_k \left[\left(\sum_{h=1}^H \ell(s_h, a_h, h) \right)^2 \right] \\ &\leq 2 \mathbb{E}_k \left[\sum_{h=1}^H \ell(s_h, a_h, h) \left(\sum_{m=h}^H \ell(s_m, a_m, m) \right) \right] \\ &= 2 \mathbb{E}_k \left[\sum_{h=1}^H \mathbb{E}_k \left[\ell(s_h, a_h, h) \left(\sum_{m=h}^H \ell(s_m, a_m, m) \right) \mid s_h, a_h \right] \right] \\ &= 2 \mathbb{E}_k \left[\sum_{h=1}^H \ell(s_h, a_h, h) \mathbb{E}_k \left[\sum_{m=h}^H \ell(s_m, a_m, m) \mid s_h, a_h \right] \right] \\ &= 2 \mathbb{E}_k \left[\sum_{h=1}^H \ell(s_h, a_h, h) Q_h^{\pi_k}(s_h, a_h; \ell, P) \right] \\ &= 2 \mathbb{E}_k \left[\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} n_k(s, a, h) \ell(s, a, h) Q_h^{\pi_k}(s, a; \ell, P) \right] \end{aligned}$$

819 where the first inequality holds because $(\sum_{h=1}^H x_h)^2 \leq 2 \sum_{h=1}^H x_h (\sum_{m=h}^H x_m)$. Moreover,

$$\begin{aligned} &\mathbb{E}_k \left[\sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} n_k(s, a, h) \ell(s, a, h) Q_h^{\pi_k}(s, a; \ell, P) \right] \\ &= \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \ell(s, a, h) Q_h^{\pi_k}(s, a; \ell, P) \mathbb{E}_k [n_k(s, a, h)] \\ &= \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \ell(s, a, h) Q_h^{\pi_k}(s, a; \ell, P) q_k(s, a, h) \\ &= \langle \mathbf{q}_k, \ell \odot \mathbf{Q}^{P, \pi_k, \ell} \rangle. \end{aligned}$$

Therefore, it follows that

$$\mathbb{E}_k [\langle \mathbf{n}_k, \ell \rangle^2] \leq 2 \langle \mathbf{q}_k, \ell \odot \mathbf{Q}^{P, \pi_k, \ell} \rangle.$$

820 Next, observe that

$$\begin{aligned}
& \langle \mathbf{q}_k, \ell \odot \mathbf{Q}^{P, \pi_k, \ell} \rangle \\
& \leq B \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} Q_h^{\pi_k}(s, a; \ell, P) q_k(s, a, h) \\
& = B \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \pi_k(a | s, h) Q_h^{\pi_k}(s, a; \ell, P) \left(\sum_{a' \in \mathcal{A}} q_k(s, a', h) \right) \\
& = B \sum_{h=1}^H \sum_{s \in \mathcal{S}} V_h^{\pi_k}(s; \ell, P) \left(\sum_{a' \in \mathcal{A}} q_k(s, a', h) \right) \\
& = B \sum_{h=1}^H \sum_{s \in \mathcal{S}} \left(\sum_{m=h}^H \sum_{(s'', a'') \in \mathcal{S} \times \mathcal{A}} q_k(s'', a'', m | s, h) \ell(s'', a'', m) \right) \left(\sum_{a' \in \mathcal{A}} q_k(s, a', h) \right) \\
& = B \sum_{h=1}^H \sum_{m=h}^H \sum_{(s'', a'') \in \mathcal{S} \times \mathcal{A}} \sum_{s \in \mathcal{S}} q_k(s'', a'', m | s, h) \left(\sum_{a' \in \mathcal{A}} q_k(s, a', h) \right) \ell(s'', a'', m) \\
& = B \sum_{h=1}^H \sum_{m=h}^H \sum_{(s'', a'') \in \mathcal{S} \times \mathcal{A}} q_k(s'', a'', m) \ell(s'', a'', m) \\
& = B \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} h \cdot q_k(s, a, h) \ell(s, a, h) \\
& = B \langle \mathbf{q}_k, \vec{h} \odot \ell \rangle
\end{aligned}$$

where the first inequality holds because $\ell(s, a, h) \leq B$ for any (s, a, h) , the first equality holds because

$$q_k(s, a, h) = \pi_k(a | s, h) \sum_{a' \in \mathcal{A}} q_k(s, a', h),$$

the fifth equality follows from

$$\sum_{s \in \mathcal{S}} q_k(s'', a'', m | s, h) \left(\sum_{a' \in \mathcal{A}} q_k(s, a', h) \right) = q_k(s'', a'', m).$$

821 Therefore, we get that $\langle \mathbf{q}_k, \ell \odot \mathbf{Q}^{P, \pi_k, \ell} \rangle \leq B \langle \mathbf{q}_k, \vec{h} \odot \ell \rangle$ as required. \square

822 The following lemma is from the first statement of (Chen & Luo, 2021, Lemma 7) with a few
823 modifications to adapt the proof to our setting.

824 **Lemma 16.** (Chen & Luo, 2021, Lemma 7) *Let π be a policy, and let \tilde{P}, \hat{P} be two different transition
825 kernels. We denote by \tilde{q} the occupancy measure $q^{\tilde{P}, \pi}$ associated with \tilde{P} and π , and we denote by \hat{q}
826 the occupancy measure $q^{\hat{P}, \pi}$ associated with \hat{P} and π . Then*

$$\begin{aligned}
& \hat{q}(s, a, h) - \tilde{q}(s, a, h) \\
& = \sum_{(s', a', s'')} \sum_{m=1}^{h-1} \tilde{q}(s', a', m) \left(\hat{P}(s'' | s', a', m) - \tilde{P}(s'' | s', a', m) \right) \hat{q}(s, a, h | s'', m+1).
\end{aligned}$$

Proof. We prove the first statement by induction on h . When $h = 1$, note that

$$\hat{q}(s, a, h) = \tilde{q}(s, a, h) = \pi(a | s, 1) \cdot p(s).$$

827 Hence, both the left-hand side and right-hand side are equal to 0. Next, assume that the equality
 828 holds with $h - 1 \geq 1$. Then we consider h . By the definition of occupancy measure,

$$\begin{aligned}
 & \widehat{q}(s, a, h) - \widetilde{q}(s, a, h) \\
 &= \pi(a \mid s, h) \sum_{(s', a')} (\widehat{P}(s \mid s', a', h - 1) \widehat{q}(s', a', h - 1) - \widetilde{P}(s \mid s', a', h - 1) \widetilde{q}(s', a', h - 1)) \\
 &= \pi(a \mid s, h) \underbrace{\sum_{(s', a')} \widehat{P}(s \mid s', a', h - 1) (\widehat{q}(s', a', h - 1) - \widetilde{q}(s', a', h - 1))}_{\text{Term 1}} \\
 &\quad + \underbrace{\pi(a \mid s, h) \sum_{(s', a')} \widetilde{q}(s', a', h - 1) (\widehat{P}(s \mid s', a', h - 1) - \widetilde{P}(s \mid s', a', h - 1))}_{\text{Term 2}}.
 \end{aligned}$$

829 To provide an upper bound on Term 1, we use the induction hypothesis for $h - 1$:

$$\begin{aligned}
 & \widehat{q}(s', a', h - 1) - \widetilde{q}(s', a', h - 1) \\
 &= \sum_{(s'', a'', s''')} \sum_{m=1}^{h-2} \widetilde{q}(s'', a'', m) \left((\widehat{P} - \widetilde{P})(s''' \mid s'', a'', m) \right) \widehat{q}(s', a', h - 1 \mid s''', m + 1)
 \end{aligned}$$

where

$$(\widehat{P} - \widetilde{P})(s''' \mid s'', a'', m) = \widehat{P}(s''' \mid s'', a'', m) - \widetilde{P}(s''' \mid s'', a'', m).$$

In addition, observe that

$$\pi(a \mid s, h) \sum_{(s', a')} \widehat{P}(s \mid s', a', h - 1) \widehat{q}(s', a', h - 1 \mid s''', m + 1) = \widehat{q}(s, a, h \mid s''', m + 1).$$

830 Therefore, it follows that Term 1 is equal to

$$\begin{aligned}
 & \sum_{(s'', a'', s''')} \sum_{m=1}^{h-2} \widetilde{q}(s'', a'', m) \left((\widehat{P} - \widetilde{P})(s''' \mid s'', a'', m) \right) \widehat{q}(s, a, h \mid s''', m + 1) \\
 &= \sum_{(s', a', s'')} \sum_{m=1}^{h-2} \widetilde{q}(s', a', m) \left(\widehat{P}(s'' \mid s', a', m) - \widetilde{P}(s'' \mid s', a', m) \right) \widehat{q}(s, a, h \mid s'', m + 1).
 \end{aligned}$$

Next, we upper bound Term 2. Note that

$$\widehat{q}(s, a, h \mid s'', h) = \pi(a \mid s'', h) \cdot \mathbf{1}[s'' = s].$$

831 Then it follows that

$$\begin{aligned}
 & \pi(a \mid s, h) (\widehat{P}(s \mid s', a', h - 1) - \widetilde{P}(s \mid s', a', h - 1)) \\
 &= \sum_{s'' \in \mathcal{S}} \mathbf{1}[s'' = s] \cdot \pi(a \mid s'', h) (\widehat{P}(s'' \mid s', a', h - 1) - \widetilde{P}(s'' \mid s', a', h - 1)) \\
 &= \sum_{s'' \in \mathcal{S}} \widehat{q}(s, a, h \mid s'', h) (\widehat{P}(s'' \mid s', a', h - 1) - \widetilde{P}(s'' \mid s', a', h - 1)),
 \end{aligned}$$

implying in turn that Term 2 equals

$$\sum_{(s', a', s'') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}} \widetilde{q}(s', a', h - 1) (\widehat{P}(s'' \mid s', a', h - 1) - \widetilde{P}(s'' \mid s', a', h - 1)) \widehat{q}(s, a, h \mid s'', h).$$

832 Adding the equivalent expression of Term 1 and that of Term 2 that we have obtained, we get the
 833 right-hand side of the statement. \square

834 The following lemma is called value difference lemma (Dann et al., 2017). Based on Lemma 13
 835 and Lemma 16, we show the following lemma, which is a modification of (Chen & Luo, 2021,
 836 Lemma 7, the second statement).

837 **Lemma 17.** *Let π be a policy, and let \tilde{P}, \hat{P} be two different transition kernels. We denote by \tilde{q} the*
 838 *occupancy measure $q^{\tilde{P}, \pi}$ associated with \tilde{P} and π , and we denote by \hat{q} the occupancy measure $q^{\hat{P}, \pi}$*
 839 *associated with \hat{P} and π . Let $\ell : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [-B, B]$ be an arbitrary function. If $\tilde{P}, \hat{P} \in \mathcal{P}_k$,*
 840 *then we have*

$$\begin{aligned} |\langle \ell, \hat{q} - \tilde{q} \rangle| &= \left| \sum_{(s,a,s',h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \tilde{q}(s, a, h) \left(\hat{P}(s' | s, a, h) - \tilde{P}(s' | s, a, h) \right) V_{h+1}^\pi(s'; \ell, \hat{P}) \right| \\ &\leq BH \sum_{(s,a,s',h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \tilde{q}(s, a, h) \epsilon_k^*(s' | s, a, h) \end{aligned}$$

841 where \hat{q}, \tilde{q}, ℓ are the vector representations of \hat{q}, \tilde{q}, ℓ .

842 *Proof.* First, observe that

$$\langle \ell, \hat{q} - \tilde{q} \rangle = \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} (\hat{q}(s, a, h) - \tilde{q}(s, a, h)) \ell(s, a, h).$$

843 By Lemma 16, the right-hand side can be rewritten so that we obtain the following.

$$\begin{aligned} &\langle \ell, \hat{q} - \tilde{q} \rangle \\ &= \sum_{(s,a,h)} \sum_{(s',a',s'')} \sum_{m=1}^{h-1} \tilde{q}(s', a', m) \left((\hat{P} - \tilde{P})(s'' | s', a', m) \right) \hat{q}(s, a, h | s'', m+1) \ell(s, a, h) \\ &= \sum_{m=1}^H \sum_{(s',a',s'')} \tilde{q}(s', a', m) \left((\hat{P} - \tilde{P})(s'' | s', a', m) \right) \sum_{\substack{(s,a,h), \\ h>m}} \hat{q}(s, a, h | s'', m+1) \ell(s, a, h) \\ &= \sum_{m=1}^H \sum_{(s',a',s'')} \tilde{q}(s', a', m) \left((\hat{P} - \tilde{P})(s'' | s', a', m) \right) V_{m+1}^\pi(s''; \ell, \hat{P}) \\ &= \sum_{h=1}^H \sum_{(s',a',s'')} \tilde{q}(s', a', h) \left(\hat{P}(s'' | s', a', h) - \tilde{P}(s'' | s', a', h) \right) V_{h+1}^\pi(s''; \ell, \hat{P}). \end{aligned}$$

844 Since $\tilde{P}, \hat{P} \in \mathcal{P}_k$, Lemma 13 implies that

$$\begin{aligned} |\langle \ell, \hat{q} - \tilde{q} \rangle| &\leq \sum_{h=1}^H \sum_{(s',a',s'')} \tilde{q}(s', a', h) \left| \hat{P}(s'' | s', a', h) - \tilde{P}(s'' | s', a', h) \right| V_{h+1}^\pi(s''; \ell, \hat{P}) \\ &\leq \sum_{h=1}^H \sum_{(s',a',s'')} \tilde{q}(s', a', h) (2\epsilon_k(s'' | s', a', h)) V_{h+1}^\pi(s''; \ell, \hat{P}) \\ &\leq \sum_{h=1}^H \sum_{(s',a',s'')} \tilde{q}(s', a', h) \epsilon_k^*(s'' | s', a', h) V_{h+1}^\pi(s''; \ell, \hat{P}) \\ &\leq BH \sum_{h=1}^H \sum_{(s',a',s'')} \tilde{q}(s', a', h) \epsilon_k^*(s'' | s', a', h) \\ &= BH \sum_{(s,a,s',h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} \tilde{q}(s, a, h) \epsilon_k^*(s' | s, a, h) \end{aligned}$$

845 where the third inequality holds because $V_{h+1}^\pi(s''; \ell, \hat{P}) \leq BH$, as required. \square

Lemma 18. Let π be a policy, and let \tilde{P}, \hat{P} be two different transition kernels. We denote by \tilde{q} the occupancy measure $q^{\tilde{P}, \pi}$ associated with \tilde{P} and π , and we denote by \hat{q} the occupancy measure $q^{\hat{P}, \pi}$ associated with \hat{P} and π . Let $(s, h) \in \mathcal{S} \times [H]$, and consider $\tilde{q}(\cdot | s, h), \hat{q}(\cdot | s, h) : \mathcal{S} \times \mathcal{A} \times \{h, \dots, H\}$. If $\tilde{P}, \hat{P} \in \mathcal{P}_k$, then we have

$$|\langle \ell_{(h)}, \hat{q}_{(s,h)} - \tilde{q}_{(s,h)} \rangle| \leq BH \sum_{(s', a', s'', m) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \{h, \dots, H\}} \tilde{q}(s', a', m | s, h) \epsilon_k^*(s'' | s', a', m)$$

846 where $\tilde{q}_{(s,h)}, \hat{q}_{(s,h)}, \ell_{(h)}$ are the vector representations of $\tilde{q}(\cdot | s, h), \hat{q}(\cdot | s, h) : \mathcal{S} \times \mathcal{A} \times$
 847 $\{h, \dots, H\} \rightarrow [0, 1]$ and $\ell_{(h)} : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [-B, B]$.

848 *Proof.* The proof follows the same argument used to prove Lemmas 16 and 17. \square

849 The following lemma is called a Bellman-type law of total variance lemma (Azar et al., 2017; Chen
 850 & Luo, 2021). We follow the proof of (Chen & Luo, 2021, Lemma 4) after some changes to adapt
 851 to our setting.

Lemma 19. (Chen & Luo, 2021, Lemma 4) Let π_k be the policy for episode k , P be an arbitrary transition kernel, and let q_k denote the occupancy measure q^{P, π_k} . Let $\ell : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [-B, B]$ be an arbitrary reward function, and define $\mathbb{V}_k(s, a, h) = \text{Var}_{s' \sim P(\cdot | s, a, h)} [V_{h+1}^{\pi_k}(s'; \ell, P)]$. Then

$$\langle \mathbf{q}_k, \mathbb{V}_k \rangle \leq \text{Var}[\langle \mathbf{n}_k, \ell \rangle | \ell, \pi_k, P]$$

852 where $\mathbf{q}_k, \mathbb{V}_k, \mathbf{n}_k, \ell$ are the vector representations of $q_k, \mathbb{V}_k, n_k, \ell$.

853 *Proof.* For ease of notation, let s_h and a_h denote s_h^{P, π_k} and a_h^{P, π_k} , respectively for $h \in [H]$. More-
 854 over, let $V(s, h)$ denote $V_h^\pi(s; \ell, P)$ for $(s, h) \in \mathcal{S} \times [H]$. Note that

$$\langle \mathbf{n}_k, \ell \rangle = \sum_{(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]} \ell(s, a, h) n_k(s, a, h) = \sum_{h=1}^H \ell(s_h, a_h, h).$$

855 For ease of notation, let $\mathbb{E}_k[\cdot]$ and $\text{Var}_k[\cdot]$ denote $\mathbb{E}[\cdot | \ell, \pi_k, P]$ and $\text{Var}[\cdot | \ell, \pi_k, P]$, respectively.
 856 Then

$$\mathbb{E}_k[\langle \mathbf{n}_k, \ell \rangle] = \mathbb{E}_k \left[\sum_{h=1}^H \ell(s_h, a_h, h) \right] = \mathbb{E}_k \left[\mathbb{E} \left[\sum_{h=1}^H \ell(s_h, a_h, h) | \ell, \pi_k, P, s_1 \right] \right] = \mathbb{E}_k[V(s_1, 1)].$$

857 Moreover,

$$\begin{aligned} \text{Var}_k[\langle \mathbf{n}_k, \ell \rangle] &= \mathbb{E}_k \left[\left(\sum_{h=1}^H \ell(s_h, a_h, h) - \mathbb{E}_k[V(s_1, 1)] \right)^2 \right] \\ &= \mathbb{E}_k \left[\left(\sum_{h=1}^H \ell(s_h, a_h, h) - V(s_1, 1) + V(s_1, 1) - \mathbb{E}_k[V(s_1, 1)] \right)^2 \right] \\ &= \mathbb{E}_k \left[\left(\sum_{h=1}^H \ell(s_h, a_h, h) - V(s_1, 1) \right)^2 \right] + \mathbb{E}_k \left[(V(s_1, 1) - \mathbb{E}_k[V(s_1, 1)])^2 \right] \\ &\quad + 2\mathbb{E}_k \left[\left(\sum_{h=1}^H \ell(s_h, a_h, h) - V(s_1, 1) \right) (V(s_1, 1) - \mathbb{E}_k[V(s_1, 1)]) \right] \\ &\geq \mathbb{E}_k \left[\left(\sum_{h=1}^H \ell(s_h, a_h, h) - V(s_1, 1) \right)^2 \right] \end{aligned}$$

858 where the inequality is by $\mathbb{E}_k [V(s_1, 1) - \mathbb{E}_k [V(s_1, 1) | s_1]] = 0$ and
 859 $(V(s_1, 1) - \mathbb{E}_k [V(s_1, 1)])^2 \geq 0$. Therefore,

$$\text{Var}_k [\langle \mathbf{n}_k, \ell \rangle] \geq \mathbb{E}_k \left[\left(\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) + \ell(s_1, a_1, 1) + V(s_2, 2) - V(s_1, 1) \right)^2 \right].$$

860 Note that

$$861 \quad \mathbb{E}_k \left[\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) \mid s_1, a_1, s_2 \right] = \mathbb{E}_k \left[\sum_{h=2}^H \ell(s_h, a_h, h) \mid s_2 \right] - V(s_2, 2) = 0. \quad (34)$$

862 Then

$$\begin{aligned} & \text{Var}_k [\langle \mathbf{n}_k, \ell \rangle] \\ & \geq \mathbb{E}_k \left[\left(\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) \right)^2 \right] + \mathbb{E}_k \left[(\ell(s_1, a_1, 1) + V(s_2, 2) - V(s_1, 1))^2 \right] \\ & \quad + 2\mathbb{E}_k \left[\mathbb{E}_k \left[\left(\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) \right) (\ell(s_1, a_1, 1) + V(s_2, 2) - V(s_1, 1)) \mid s_1, a_1, s_2 \right] \right] \\ & = \mathbb{E}_k \left[\left(\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) \right)^2 \right] + \mathbb{E}_k \left[(\ell(s_1, a_1, 1) + V(s_2, 2) - V(s_1, 1))^2 \right] \\ & \quad + 2\mathbb{E}_k \left[(\ell(s_1, a_1, 1) + V(s_2, 2) - V(s_1, 1)) \mathbb{E}_k \left[\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) \mid s_1, a_1, s_2 \right] \right] \\ & = \mathbb{E}_k \left[\left(\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) \right)^2 \right] + \mathbb{E}_k \left[(\ell(s_1, a_1, 1) + V(s_2, 2) - V(s_1, 1))^2 \right] \end{aligned}$$

863 where the last equality follows from (34). Here, the second term from the right-most side can be
 864 bounded from below as follows.

$$\begin{aligned} & \mathbb{E}_k \left[(\ell(s_1, a_1, 1) + V(s_2, 2) - V(s_1, 1))^2 \right] \\ & = \mathbb{E}_k \left[\left(\ell(s_1, a_1, 1) + \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) - V(s_1, 1) + V(s_2, 2) - \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) \right)^2 \right] \\ & = \mathbb{E}_k \left[\left(\ell(s_1, a_1, 1) + \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) - V(s_1, 1) \right)^2 \right] \\ & \quad + \mathbb{E}_k \left[\left(V(s_2, 2) - \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) \right)^2 \right] \\ & \quad + 2\mathbb{E}_k \left[\left(\ell(s_1, a_1, 1) + \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) - V(s_1, 1) \right) \left(V(s_2, 2) - \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) \right) \right] \\ & = \mathbb{E}_k \left[\left(\ell(s_1, a_1, 1) + \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) - V(s_1, 1) \right)^2 \right] \\ & \quad + \mathbb{E}_k \left[\left(V(s_2, 2) - \sum_{s' \in \mathcal{S}} P(s' \mid s_1, a_1, 1) V(s', 2) \right)^2 \right] \\ & \geq \mathbb{E}_k [V_k(s_1, a_1, 1)] \end{aligned}$$

866 where third equality holds because

$$\begin{aligned}
 & \mathbb{E}_k \left[\left(\ell(s_1, a_1, 1) + \sum_{s' \in \mathcal{S}} P(s' | s_1, a_1, 1) V(s', 2) - V(s_1, 1) \right) \left(V(s_2, 2) - \sum_{s' \in \mathcal{S}} P(s' | s_1, a_1, 1) V(s', 2) \right) \middle| s_1, a_1 \right] \\
 867 \quad &= \left(\ell(s_1, a_1, 1) + \sum_{s' \in \mathcal{S}} P(s' | s_1, a_1, 1) V(s', 2) - V(s_1, 1) \right) \mathbb{E}_k \left[V(s_2, 2) - \sum_{s' \in \mathcal{S}} P(s' | s_1, a_1, 1) V(s', 2) \middle| s_1, a_1 \right] \\
 &= \left(\ell(s_1, a_1, 1) + \sum_{s' \in \mathcal{S}} P(s' | s_1, a_1, 1) V(s', 2) - V(s_1, 1) \right) \times 0
 \end{aligned}$$

868 and the last inequality holds because

$$\mathbb{E}_k \left[\left(V(s_2, 2) - \sum_{s' \in \mathcal{S}} P(s' | s_1, a_1, 1) V(s', 2) \right)^2 \right] = \mathbb{E}_k [\mathbb{V}_k(s_1, a_1, 1)].$$

869 Then it follows that

$$\begin{aligned}
 \text{Var}_k [\langle \mathbf{n}_k, \ell \rangle] &\geq \mathbb{E}_k \left[\left(\sum_{h=1}^H \ell(s_h, a_h, h) - V(s_1, 1) \right)^2 \right] \\
 &\geq \mathbb{E}_k \left[\left(\sum_{h=2}^H \ell(s_h, a_h, h) - V(s_2, 2) \right)^2 \right] + \mathbb{E}_k [\mathbb{V}_k(s_1, a_1, 1)].
 \end{aligned}$$

870 Repeating the same argument, we deduce that

$$\text{Var}_k [\langle \mathbf{n}_k, \ell \rangle] \geq \sum_{h=1}^H \mathbb{E}_k [\mathbb{V}_k(s_h, a_h, h)] = \sum_{(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]} q_k(s, a, h) \mathbb{V}_k(s, a, h) = \langle \mathbf{q}_k, \mathbb{V}_k \rangle,$$

871 as required. \square

872 Next, we provide Lemma 20, which is a modification of (Chen & Luo, 2021, Lemma 9) to our
873 finite-horizon MDP setting.

874 **Lemma 20.** *Assume that the good event \mathcal{E} holds. Let π_k be any policy for episode k , let P_k be any
875 transition kernel from \mathcal{P}_k for episode k , and let P be the true transition kernel. Let q_k, \hat{q}_k denote the
876 occupancy measures $q^{P, \pi_k}, \hat{q}^{P_k, \pi_k}$, respectively. Let $\ell_k : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [-B, B]$ be an arbitrary
877 reward function for episode k . With probability at least $1 - 2\delta$,*

$$\sum_{k=1}^K |\langle \ell_k, \mathbf{q}_k - \hat{\mathbf{q}}_k \rangle| = \mathcal{O} \left(B \left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) L_\delta^3 \right).$$

878 where $\mathbf{q}_k, \hat{\mathbf{q}}_k, \ell_k$ are the vector representations of q_k, \hat{q}_k, ℓ_k .

Proof. We define ξ_1 as $\xi_1 = \{\ell_1, \pi_1\}$ and for $k \geq 2$, we define ξ_k as

$$\left\{ s_1^{P, \pi_{k-1}}, a_1^{P, \pi_{k-1}}, \dots, s_h^{P, \pi_{k-1}}, a_h^{P, \pi_{k-1}}, \ell_k, \pi_k \right\}$$

where π_{k-1} and π_k denote the policies for episode $k-1$ and episode k , respectively, and

$$\left(s_1^{P, \pi_{k-1}}, a_1^{P, \pi_{k-1}}, \dots, s_h^{P, \pi_{k-1}}, a_h^{P, \pi_{k-1}} \right)$$

879 is the trajectory generated under policy π_{k-1} and transition kernel P . Then for $k \in [K]$, let \mathcal{H}_k be
880 defined as the σ -algebra generated by the random variables in $\xi_1 \cup \dots \cup \xi_k$. Then it follows that
881 $\mathcal{H}_1, \dots, \mathcal{H}_k$ give rise to a filtration.

Let us define

$$\mu_k(s, a, h) = \mathbb{E}_{s' \sim P(\cdot | s, a, h)} [V_{h+1}^{\pi_k}(s'; \ell_k, P)].$$

882 Note that

$$\begin{aligned} & \sum_{k=1}^K |(\ell_k, \mathbf{q}_k - \widehat{\mathbf{q}}_k)| \\ &= \sum_{k=1}^K \left| \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} q_k(s, a, h) (P(s' | s, a, h) - P_k(s' | s, a, h)) V_{h+1}^{\pi_k}(s'; \ell_k, P_k) \right| \\ &\leq \sum_{k=1}^K \left| \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} q_k(s, a, h) (P(s' | s, a, h) - P_k(s' | s, a, h)) V_{h+1}^{\pi_k}(s'; \ell_k, P) \right| \\ &\quad + \mathcal{O}(BH^3 S^3 AL_\delta^3) \end{aligned}$$

883 where the equality is due to Lemma 17 and the inequality is due to Lemmas 4 and 11.

884 Moreover,

$$\begin{aligned} & \sum_{k=1}^K \left| \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} q_k(s, a, h) (P(s' | s, a, h) - P_k(s' | s, a, h)) V_{h+1}^{\pi_k}(s'; \ell_k, P) \right| \\ &= \sum_{k=1}^K \left| \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} q_k(s, a, h) ((P - P_k)(s' | s, a, h)) (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h)) \right| \\ &\leq \sum_{k=1}^K \sum_{(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]} q_k(s, a, h) \epsilon_k^*(s' | s, a, h) |V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h)| \\ &\leq \mathcal{O} \left(\sum_{k=1}^K \sum_{\substack{(s, a, s', h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} q_k(s, a, h) \sqrt{\frac{P(s' | s, a, h) L_\delta}{\max\{1, N_k(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \right) \\ &\quad + \mathcal{O} \left(BHS \sum_{k=1}^K \sum_{(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]} \frac{q_k(s, a, h) L_\delta}{\max\{1, N_k(s, a, h)\}} \right) \\ &\leq \mathcal{O} \left(\sum_{k=1}^K \sum_{\substack{(s, a, s', h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} q_k(s, a, h) \sqrt{\frac{P(s' | s, a, h) L_\delta}{\max\{1, N_k(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \right) \\ &\quad + \mathcal{O}(BH^2 S^2 AL_\delta^2) \end{aligned}$$

885 where the first equality holds because $\sum_{s' \in \mathcal{S}} (P - P_k)(s' | s, a, h) = 0$ and $\mu_k(s, a, h)$
 886 is independent of s' , the first inequality is due to Lemma 13, the second inequality is from
 887 $|V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h)| \leq 2BH$, and the last inequality is from Lemma 11. Recall that
 888 $q_k(s, a, h) = \mathbb{E}[n_k(s, a, h) | \pi_k, P]$, which implies that

$$\begin{aligned} & \sum_{k=1}^K \mathbb{E}[X_k | \mathcal{H}_k, P] \\ &= \sum_{k=1}^K \sum_{\substack{(s, a, s', h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} q_k(s, a, h) \sqrt{\frac{P(s' | s, a, h) L_\delta}{\max\{1, N_k(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \end{aligned}$$

where

$$X_k = \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2.$$

Here, we have

$$0 \leq X_k \leq \mathcal{O} \left(BH S \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} n_k(s, a, h) \sqrt{L_\delta} \right) = \mathcal{O}(BH^2 S \sqrt{L_\delta}).$$

889 Then it follows from Lemma 26 that with probability at least $1 - \delta$,

$$\begin{aligned} & \sum_{k=1}^K \mathbb{E}[X_k | \mathcal{H}_k, P] \\ & \leq 2 \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \\ & \quad + \mathcal{O}(BH^2 S L_\delta^{1.5}). \end{aligned}$$

890 Note that

$$\begin{aligned} & \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \\ & \leq \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_{k+1}(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \\ & \quad + BH \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \left(\sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_k(s, a, h)\}}} - \sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_{k+1}(s, a, h)\}}} \right) \\ & \leq \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_{k+1}(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \\ & \quad + BH \sqrt{S} \sum_{k=1}^K \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} \left(\sqrt{\frac{L_\delta}{\max\{1, N_k(s, a, h)\}}} - \sqrt{\frac{L_\delta}{\max\{1, N_{k+1}(s, a, h)\}}} \right) \\ & \leq \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \sqrt{\frac{P(s' | s, a, h)L_\delta}{\max\{1, N_{k+1}(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \\ & \quad + \mathcal{O}(BH^2 S^{1.5} A \sqrt{L_\delta}). \end{aligned}$$

where the first inequality holds because $|V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h)| \leq BH$, the second inequality holds because $n_k(s, a, h) \leq 1$ and the Cauchy-Schwarz inequality implies that

$$\sum_{s' \in \mathcal{S}} \sqrt{P(s' | s, a, h)} \leq \sqrt{S \sum_{s' \in \mathcal{S}} P(s' | s, a, h)} = \sqrt{S},$$

and the third inequality follows from

$$\sum_{k=1}^K \left(\sqrt{\frac{1}{\max\{1, N_k(s, a, h)\}}} - \sqrt{\frac{1}{\max\{1, N_{k+1}(s, a, h)\}}} \right) \leq \sqrt{\frac{1}{\max\{1, N_1(s, a, h)\}}} = 1.$$

891 Next, the Cauchy-Schwarz inequality implies the following.

$$\begin{aligned} & \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \sqrt{\frac{P(s' | s, a, h) L_\delta}{\max\{1, N_{k+1}(s, a, h)\}}} (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \\ & \leq \sqrt{\sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) P(s' | s, a, h) (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2} \\ & \quad \times \sqrt{\sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) \frac{L_\delta}{\max\{1, N_{k+1}(s, a, h)\}}} \end{aligned}$$

892 Here, the second term can be bounded as follows.

$$\begin{aligned} \sum_{k=1}^K \sum_{(s,a,s',h)} n_k(s, a, h) \frac{L_\delta}{\max\{1, N_{k+1}(s, a, h)\}} &= SL_\delta \sum_{k=1}^K \sum_{(s,a,h)} \frac{n_k(s, a, h)}{\max\{1, N_{k+1}(s, a, h)\}} \\ &= SL_\delta \sum_{(s,a,h)} \sum_{k=1}^K \frac{n_k(s, a, h)}{\max\{1, N_{k+1}(s, a, h)\}} \\ &= \mathcal{O}(HS^2 AL_\delta^2). \end{aligned}$$

For $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$, we define

$$\mathbb{V}_k(s, a, h) = \text{Var}_{s' \sim P(\cdot | s, a, h)} [V_{h+1}^{\pi_k}(s'; \ell_k, P)].$$

893 Then

$$\begin{aligned} \mathbb{V}_k(s, a, h) &= \mathbb{E}_{s' \sim P(\cdot | s, a, h)} \left[(V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \right] \\ &= \sum_{s' \in \mathcal{S}} P(s' | s, a, h) (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \end{aligned}$$

894 Furthermore, with probability at least $1 - \delta$,

$$\begin{aligned} & \sum_{k=1}^K \sum_{\substack{(s,a,s',h) \in \\ \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]}} n_k(s, a, h) P(s' | s, a, h) (V_{h+1}^{\pi_k}(s'; \ell_k, P) - \mu_k(s, a, h))^2 \\ &= \sum_{k=1}^K \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} n_k(s, a, h) \mathbb{V}_k(s, a, h) \\ &= \sum_{k=1}^K \langle \mathbf{q}_k, \mathbb{V}_k \rangle + \sum_{k=1}^K \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} (n_k(s, a, h) - q_k(s, a, h)) \mathbb{V}_k(s, a, h) \\ &\leq \sum_{k=1}^K \text{Var}[\langle n_k, \ell_k \rangle | \ell_k, \pi_k, P] + \mathcal{O}\left(B^2 H^3 \sqrt{K \ln(1/\delta)}\right) \end{aligned}$$

895 where $\mathbb{V}_k \in \mathbb{R}^{SAH}$ is the vector representation of \mathbb{V}_k and the inequality follows from Lemma 19,
 896 $\mathbb{V}_k(s, a, h) \leq B^2 H^2$,

$$\begin{aligned} & \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} (n_k(s, a, h) - q_k(s, a, h)) \mathbb{V}_k(s, a, h) \\ & \leq \sum_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} (n_k(s, a, h) + q_k(s, a, h)) B^2 H^2 \\ & \leq 2B^2 H^3, \end{aligned}$$

897 and Lemma 24. Therefore, we finally have proved that

$$\begin{aligned} \sum_{k=1}^K |\langle \ell_k, \mathbf{q}_k - \hat{\mathbf{q}}_k \rangle| &= \mathcal{O} \left(\sqrt{HS^2 AL_\delta^2 \left(\sum_{k=1}^K \text{Var} [\langle n_k, \ell_k \rangle \mid \ell_k, \pi_k, P] + B^2 H^3 \sqrt{K \ln \frac{1}{\delta}} \right)} \right) \\ &+ \mathcal{O}(BH^3 S^3 AL_\delta^3). \end{aligned}$$

Moreover, we know from Lemma 15 that

$$\text{Var} [\langle n_k, \ell_k \rangle \mid \ell_k, \pi_k, P] \leq \mathbb{E} [\langle n_k, \ell_k \rangle^2 \mid \ell_k, \pi_k, P] \leq 2B \langle \mathbf{q}_k, \vec{\mathbf{h}} \odot \ell_k \rangle,$$

898 and therefore, it follows that

$$\begin{aligned} \sum_{k=1}^K |\langle \ell_k, \mathbf{q}_k - \hat{\mathbf{q}}_k \rangle| &= \mathcal{O} \left(\left(\sqrt{HS^2 A \left(B \sum_{k=1}^K \langle \mathbf{q}_k, \vec{\mathbf{h}} \odot \ell_k \rangle + B^2 H^3 \sqrt{K} \right)} + BH^3 S^3 A \right) L_\delta^3 \right) \\ &= \mathcal{O} \left(\left(\sqrt{B^2 H^3 S^2 AK + B^2 H^4 S^2 A \sqrt{K}} + BH^3 S^3 A \right) L_\delta^3 \right) \\ &= \mathcal{O} \left(\left(\sqrt{B^2 H^3 S^2 AK + B^2 H^3 S^2 AK + B^2 H^5 S^2 A} + BH^3 S^3 A \right) L_\delta^3 \right) \\ &= \mathcal{O} \left(B \left(H^{1.5} S \sqrt{AK} + H^3 S^3 A \right) L_\delta^3 \right) \end{aligned}$$

899 where the second equality holds because $\langle \mathbf{q}_k, \vec{\mathbf{h}} \odot \ell_k \rangle = \mathcal{O}(BH^2)$ and the third equality holds
 900 because $B^2 H^4 S^2 A \sqrt{K} = \mathcal{O}(B^2 (H^3 S^2 AK + H^5 S^2 A))$. \square

901 16 Concentration Inequalities

902 **Lemma 21.** (Hoeffding's inequality) For i.i.d. random variables Z_1, \dots, Z_n following 1/2-sub-
 903 Gaussian with zero mean,

$$\begin{aligned} \mathbb{P} \left(\frac{1}{n} \sum_{j=1}^n Z_j \geq \epsilon \right) &\leq \exp(-n\epsilon^2), \\ \mathbb{P} \left(\frac{1}{n} \sum_{j=1}^n Z_j \leq -\epsilon \right) &\leq \exp(-n\epsilon^2). \end{aligned}$$

Lemma 22. (Maurer & Pontil, 2009, Theorem 4) Let $Z_1, \dots, Z_n \in [0, 1]$ be i.i.d. random variables with mean z , and let $\delta > 0$. Then with probability at least $1 - \delta$,

$$z - \frac{1}{n} \sum_{j=1}^n Z_j \leq \sqrt{\frac{2V_n \ln(2/\delta)}{n}} + \frac{7 \ln(2/\delta)}{3(n-1)}$$

where V_n is the sample variance given by

$$V_n = \frac{1}{n(n-1)} \sum_{1 \leq j < k \leq n} (Z_j - Z_k)^2.$$

904 Next, we need the following Bernstein-type concentration inequality for martingales due to [Beygelzimer et al. \(2011\)](#). We take the version used in ([Jin et al., 2020](#), Lemma 9).

Lemma 23. ([Beygelzimer et al., 2011](#), Theorem 1) *Let Y_1, \dots, Y_n be a martingale difference sequence with respect to a filtration $\mathcal{F}_1, \dots, \mathcal{F}_n$. Assume that $Y_j \leq R$ almost surely for all $j \in [n]$. Then for any $\delta \in (0, 1)$ and $\lambda \in (0, 1/R]$, with probability at least $1 - \delta$, we have*

$$\sum_{j=1}^n Y_j \leq \lambda \sum_{j=1}^n \mathbb{E}[Y_j^2 | \mathcal{F}_j] + \frac{\ln(1/\delta)}{\lambda}.$$

Lemma 24 (Azuma’s inequality). *Let Y_1, \dots, Y_n be a martingale difference sequence with respect to a filtration $\mathcal{F}_1, \dots, \mathcal{F}_n$. Assume that $|Y_j| \leq B$ for $j \in [n]$. Then with probability at least $1 - \delta$, we have*

$$\left| \sum_{j=1}^n Y_j \right| \leq B \sqrt{2n \ln(2/\delta)}.$$

906 Next, we need the following concentration inequalities due to [Cohen et al. \(2020\)](#).

907 **Lemma 25.** ([Cohen et al., 2020](#), Theorem D.3) *Let $\{X_n\}_{n=1}^\infty$ be a sequence of i.i.d. random variables with expectation μ . Suppose that $0 \leq X_n \leq B$ holds almost surely for all n . Then with*
908 *probability at least $1 - \delta$, the following holds for all $n \geq 1$ simultaneously:*

$$\left| \sum_{i=1}^n (X_i - \mu) \right| \leq 2\sqrt{B\mu n \ln \frac{2n}{\delta}} + B \ln \frac{2n}{\delta},$$

$$\left| \sum_{i=1}^n (X_i - \mu) \right| \leq 2\sqrt{B \sum_{i=1}^n X_i \ln \frac{2n}{\delta}} + 7B \ln \frac{2n}{\delta}.$$

Lemma 26. ([Cohen et al., 2020](#), Lemma D.4) *Let $\{X_n\}_{n=1}^\infty$ be a sequence of random variables adapted to the filtration $\{\mathcal{F}_n\}_{n=1}^\infty$. Suppose that $0 \leq X_n \leq B$ holds almost surely for all n . Then with probability at least $1 - \delta$, the following holds for all $n \geq 1$ simultaneously:*

$$\sum_{i=1}^n \mathbb{E}[X_i | \mathcal{F}_i] \leq 2 \sum_{i=1}^n X_i + 4B \ln(2n/\delta).$$

910 17 Experimental Setup Details

911 We evaluate DOPE+ via the following numerical experiment. We first explain the details of our
912 CMDP setting, which is a modification of the three-state CMDP instances of [Zheng & Ratliff \(2020\)](#);
913 [Simão et al. \(2021\)](#); [Bura et al. \(2022\)](#). We define the state space $\{s_1, s_2, s_3\}$ and the action space
914 $\{a_1, a_2\}$. In Figure 2, we illustrate the transition probability. For taking a_1 at s_1 , the agent remains
915 in s_1 with probability 0.8, and moves to s_2 with probability 0.2. For taking a_2 at s_1 , the agent moves
916 to s_2 with probability 0.8, and remains in s_2 with probability 0.2. Furthermore, the same transition
917 rule is applied to s_2 and s_3 .

918 Next, we present the reward function f and the cost function g . When the agent takes a_1 , no reward
919 or cost occurs. Then it can be written as $f(s, a_1) = g(s, a_1) = 0$ for $s = s_1, s_2, s_3$. When
920 a_2 is taken, the reward occurs depending on the current state. Specifically, we set $f(s_1, a_2) =$
921 $1/3$, $f(s_2, a_2) = 2/3$, and $f(s_3, a_2) = 1$. On the other hand, for any state, the same amount of
922 cost is incurred for a_2 , i.e., $g(s_1, a_2) = g(s_2, a_2) = g(s_3, a_2) = 1$. Hence, a_2 is an action with a
923 high reward and a high cost while a_1 is an action with zero reward and zero cost. Furthermore, for
924 taking action a at state s , the agent can observe the noisy reward $f(s, a) + \zeta_1$ and the noisy cost
925 $g(s, a) + \zeta_2$, where ζ_1, ζ_2 are independently drawn from a zero-mean $1/2$ -sub-Gaussian distribution.

