

Policy Optimization for Stochastic Shortest Path

Liyu Chen

University of Southern California

LIYUC@USC.EDU

Haipeng Luo

University of Southern California

HAIPENGL@USC.EDU

Aviv Rosenberg

Tel-Aviv University

AVIVROS007@GMAIL.COM

Editors: Po-Ling Loh and Maxim Raginsky

Abstract

Policy optimization is among the most popular and successful reinforcement learning algorithms, and there is increasing interest in understanding its theoretical guarantees. In this work, we initiate the study of policy optimization for the stochastic shortest path (SSP) problem, a goal-oriented reinforcement learning model that strictly generalizes the finite-horizon model and better captures many applications. We consider a wide range of settings, including stochastic and adversarial environments under full information or bandit feedback, and propose a policy optimization algorithm for each setting that makes use of novel correction terms and/or variants of dilated bonuses (Luo et al., 2021). For most settings, our algorithm is shown to achieve a near-optimal regret bound.

One key technical contribution of this work is a new approximation scheme to tackle SSP problems that we call *stacked discounted approximation* and use in all our proposed algorithms. Unlike the finite-horizon approximation that is heavily used in recent SSP algorithms, our new approximation enables us to learn a near-stationary policy with only logarithmic changes during an episode and could lead to an exponential improvement in space complexity.

1. Introduction

Stochastic Shortest Path (SSP) is a goal-oriented reinforcement learning setting, where a learner tries to reach a goal state with minimum total cost. Compared to the heavily studied finite-horizon setting, SSP is often a better model for capturing many real-world applications such as games, car navigation, robotic manipulations, and others. We study the online learning problem in SSP, where the learner interacts with an environment with unknown cost and transition function for multiple episodes. In each episode, the learner starts from an initial state, sequentially takes an action, incurs a cost, and transits to the next state until the goal state is reached. The goal of the learner is to achieve low regret, defined as the difference between her total cost and the expected cost of the optimal policy. A unique challenge of learning SSP is to trade off between two objectives: reaching the goal state and minimizing the cost. Indeed, neither reaching the goal as fast as possible nor minimizing the cost alone solves the problem.

Policy Optimization (PO) is among the most popular methods in reinforcement learning due to its strong empirical performance and favorable theoretical properties. Unlike value-based approaches such as Q learning, PO-type methods directly optimize the policy in an incremental manner. Many widely used practical algorithms fall into this category, such as REINFORCE (Williams,

1992), NPG (Kakade, 2001), and TRPO (Schulman et al., 2015). They are also easy to implement and computationally efficient compared to other methods such as those operating over the occupancy measure space (e.g., (Zimin and Neu, 2013)). From a theoretical perspective, PO is a general framework that works for different types of environments, including stochastic costs or even adversarial costs (Shani et al., 2020), function approximation (Cai et al., 2020), and non-stationary environments (Fei et al., 2020). Despite its popularity in applications, most theoretical works on PO focus on simple models such as finite-horizon models (Cai et al., 2020; Shani et al., 2020; Luo et al., 2021) and discounted models (Liu et al., 2019; Wang et al., 2020; Agarwal et al., 2021), which are often oversimplifications of real-life applications. In particular, PO methods have not been applied to regret minimization in SSP as far as we know.

Motivated by this gap, in this work, we systematically study policy optimization in SSP. We consider a wide range of different settings and for each of them discuss how to design a policy optimization algorithm with a strong regret bound. Specifically, our main results are as follows:

- In Section 3, we first propose an important technique used in all our algorithms: *stacked discounted approximation*. It reduces any SSP instance to a special Markov Decision Process (MDP) with a stack of $\mathcal{O}(\ln K)$ layers (K is the total number of episodes), each of which contains a discounted MDP (hence the name) such that the learner stays in the same layer with a certain probability γ and proceeds to the next layer with probability $1 - \gamma$. This approximation not only resolves the difficulty of having dynamic and potentially unbounded episode lengths in the PO analysis, but more importantly leads to a near-stationary policies with only $\mathcal{O}(\ln K)$ changes within an episode. Compared to the commonly used finite-horizon approximation (Chen et al., 2021d; Chen and Luo, 2021; Cohen et al., 2021) which changes the policy at every step of an episode, our approach could lead to an exponential improvement in space complexity and is also more natural since the optimal policy for SSP is indeed stationary.
- Building on the stacked discounted approximation, in Section 4, we design PO algorithms for two types of stochastic environments considered in the literature. In the first type (called stochastic costs), the cost for each visit of a state-action pair is an i.i.d. sample of an unknown distribution and is revealed to the learner immediately after the visit. Our algorithm achieves $\tilde{\mathcal{O}}(B_* S \sqrt{AK})$ regret in this case, close to the minimax bound $\tilde{\mathcal{O}}(B_* \sqrt{SAK})$ (Cohen et al., 2021), where S is the number of states, A is the number of actions, and B_* is the maximum expected cost of the optimal policy starting from any states. In the second type (called stochastic adversary following (Chen and Luo, 2021)), the cost function for each episode is fixed and an i.i.d. sample of an unknown distribution, and only at the end of the episode, the learner observes the entire cost function (full-information feedback) or the costs for all visited state-action pairs (bandit feedback). Our algorithm achieves $\tilde{\mathcal{O}}(\sqrt{DT_* K} + DS\sqrt{AK})$ regret with full information and $\tilde{\mathcal{O}}(\sqrt{DT_* SAK} + DS\sqrt{AK})$ regret with bandit feedback, where D is the diameter of the MDP and T_* is the expected hitting time of the optimal policy starting from the initial state. These bounds match the best existing results from (Chen and Luo, 2021) (and exhibit a \sqrt{S} gap in the second term $DS\sqrt{AK}$ compared to their lower bounds).
- Finally, in Section 5, we further study SSP with adversarial costs and design PO algorithms that achieve $\tilde{\mathcal{O}}(T_* \sqrt{DK} + \sqrt{DT_* S^2 AK})$ regret with full information and $\tilde{\mathcal{O}}(\sqrt{T_{\max}^5 S^2 AK})$ regret with bandit feedback, where T_{\max} is the maximum expected hitting time of the optimal policy over all states. The best existing bounds for these settings are $\tilde{\mathcal{O}}(\sqrt{DT_* S^2 AK})$ and $\tilde{\mathcal{O}}(\sqrt{DT_* S^3 A^2 K})$ respectively (Chen and Luo, 2021).

Table 1: Comparison of regret bound, time complexity, and space complexity of different SSP algorithms. We consider five feedback types: SC (stochastic costs), SAF (Stochastic Adversary, Full information), SAB (Stochastic Adversary, Bandit feedback), AF (Adversarial, Full information), and AB (Adversarial, Bandit feedback). Operator $\tilde{O}(\cdot)$ is hidden for simplicity. Time complexity of $\text{poly}(S, A, T_{\max})$ is due to optimization in the occupancy measure space.

	Regret	Time	Space	Feedback
Cohen et al. (2021)	$B_\star \sqrt{SAK}$	$S^3 A^2 T_{\max}$	$S^2 AT_{\max}$	SC
Our work	$B_\star S \sqrt{AK}$	$S^2 AT_{\max} K$	$S^2 A$	
Chen and Luo (2021)	$\sqrt{DT_\star K} + DS\sqrt{AK}$	$\text{poly}(S, A, T_{\max}) \cdot K$	$S^2 AT_{\max}$	SAF
Our work	$\sqrt{DT_\star K} + DS\sqrt{AK}$	$S^2 AT_{\max} K$	$S^2 A$	
Chen and Luo (2021)	$\sqrt{SADT_\star K} + DS\sqrt{AK}$	$\text{poly}(S, A, T_{\max}) \cdot K$	$S^2 AT_{\max}$	SAB
Our work	$\sqrt{SADT_\star K} + DS\sqrt{AK}$	$S^2 AT_{\max} K$	$S^2 A$	
Chen and Luo (2021)	$\sqrt{S^2 ADT_\star K}$	$\text{poly}(S, A, T_{\max}) \cdot K$	$S^2 AT_{\max}$	AF
Our work	$\sqrt{(S^2 A + T_\star) DT_\star K}$	$S^2 AT_{\max} K$	$S^2 A$	
Chen and Luo (2021)	$\sqrt{S^3 A^2 DT_\star K}$	$\text{poly}(S, A, T_{\max}) \cdot K$	$S^2 AT_{\max}$	AB
Our work	$\sqrt{S^2 AT_{\max}^5 K}$	$\text{poly}(S, A, T_{\max}) \cdot K$	$S^2 A$	

We also include Table 1 with a comprehensive comparison between SSP algorithms for a better understanding of our contributions. While our regret bounds do not always match the state-of-the-art, we emphasize again that our algorithms are more space-efficient due to the stacked discounted approximation ($S^2 A$ versus $S^2 AT_{\max}$ in Table 1). It is also more time-efficient in some cases (for feedback types SAF, SAB, and AF in Table 1). We also note that in the analysis of stacked discounted approximation, a regret bound starting from any state (not just the initial state) is important, and PO indeed provides such a guarantee while other methods based on occupancy measure do not. In other words, PO is especially compatible with our stacked discounted approximation. Moreover, our results also significantly improve our theoretical understanding on PO, and pave the way for future study on more challenging problems such as SSP with function approximation, where in some cases PO is the only method known to be computationally and statistically efficient (Luo et al., 2021).

Other Techniques To achieve our results for stochastic environments, we make two other technical contributions. First, in order to control the cost estimation error optimally, we derive a set of novel correction terms fed to the PO algorithm, which resolves some technical difficulties brought by PO due to its lack of optimism and also greatly simplifies the analysis. Second, due to the soft policy updates, the standard PO analysis leads to an undesirable dominating term related to T_\star or even T_{\max} in the regret, and we develop a refined analysis on the value difference between learner’s policies and the optimal policy to reduce this to a lower order term.

To achieve our results for adversarial environments, we develop a tighter variance-aware bound for the stability term in the PO analysis, which plays a key role in removing the T_{\max} dependency in the dominating term of the regret bound in the full information setting. We further extend the dilated bonuses of (Luo et al., 2021) (for the finite-horizon setting) to the stacked discounted MDPs, which is essential for both the full information setting and the bandit feedback setting.

Related Work Regret minimization in SSP has received much attention recently for both stochastic environment (Tarbouriech et al., 2020; Cohen et al., 2020, 2021; Tarbouriech et al., 2021; Chen et al., 2021a,b; Jafarnia-Jahromi et al., 2021) and adversarial environment (Rosenberg and Mansour, 2021; Chen et al., 2021d; Chen and Luo, 2021). All previous approaches are either value-based (e.g. Q learning) or occupancy-measure-based, while we take the first step in studying the more practical and versatile PO methods. Among numerous studies on PO, the closest to our work are the recent ones by Shani et al. (2020) and Luo et al. (2021) for the special case of finite-horizon MDPs.

The use of variance information (Lattimore and Hutter, 2012; Azar et al., 2017; Zhou et al., 2021; Zhang et al., 2021; Kim et al., 2021) and correction terms (Steinhardt and Liang, 2014; Wei and Luo, 2018; Chen et al., 2021c) is crucial for achieving optimal and adaptive regret bound in online learning. In this work we heavily make use of these ideas as mentioned.

2. Preliminaries

An SSP instance is defined by a Markov Decision Process (MDP) $\mathcal{M} = (\mathcal{S}, s_{\text{init}}, g, \mathcal{A}, P)$. Here, \mathcal{S} is the state space, $s_{\text{init}} \in \mathcal{S}$ is the initial state, $g \notin \mathcal{S}$ is the goal state, \mathcal{A} is the action space, and $P = \{P_{s,a}\}_{(s,a) \in \mathcal{S} \times \mathcal{A}}$ with $P_{s,a} \in \Delta_{\mathcal{S}_+}$ is the transition function, where $\mathcal{S}_+ = \mathcal{S} \cup \{g\}$ and $\Delta_{\mathcal{S}_+}$ is the simplex over \mathcal{S}_+ .

The learning protocol is as follows: the learner interacts with the environment for K episodes. In episode k , the learner starts in initial state s_{init} , sequentially takes an action, incurs a cost (which might not be observed immediately), and transits to the next state until the goal state g is reached. Formally, at the i -th step of episode k , the learner observes state s_i^k (with $s_1^k = s_{\text{init}}$), takes action a_i^k , suffers cost c_i^k , and transits to the next state $s_{i+1}^k \sim P_{s_i^k, a_i^k}$. Denote by I_k the length of episode k , such that $s_{I_k+1}^k = g$ when I_k is finite. Note that the heavily studied finite-horizon setting is a special case of SSP where I_k is always guaranteed to be some fixed number.

Proper Policies and Related Concepts At a high level, the learner’s goal is to reach the goal state with minimum cost. Thus, we focus on *proper policies*: a stationary policy $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$ is a mapping that assigns to each state a distribution over actions, and it is proper if following π from any initial state reaches the goal state with probability 1. Denote by Π the set of proper policies (assumed to be non-empty). Given a proper policy π , a transition function P , and a cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, we define its value function and action-value function as follows: $V^{\pi, P, c}(s) = \mathbb{E} \left[\sum_{i=1}^I c(s_i, a_i) \mid \pi, P, s_1 = s \right]$ and $Q^{\pi, P, c}(s, a) = c(s, a) + \mathbb{E}_{s' \sim P_{s,a}} [V^{\pi, P, c}(s')]$, where the expectation in $V^{\pi, P, c}$ is over the randomness of action $a_i \sim \pi(\cdot | s_i)$, next state $s_{i+1} \sim P_{s_i, a_i}$, and the number of steps I before reaching g . Also define the *advantage function* $A^{\pi, P, c}(s, a) = Q^{\pi, P, c}(s, a) - V^{\pi, P, c}(s)$.

We consider two types of environments: stochastic environments and adversarial environments, which differ in the way costs are generated (and revealed), discussed in detail below.

Stochastic Environments We start with the simpler environment with a fixed “ground truth” cost: there exists an unknown mean cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow [c_{\min}, 1]$, and the costs incurred by the learner are i.i.d samples from some distribution with support $[c_{\min}, 1]$ and mean c . Here, $c_{\min} \in [0, 1]$ is a global lower bound.¹ We consider the following three types of cost feedback.

1. Unlike many previous works for stochastic costs that require $c_{\min} > 0$ in their analysis, our methods allow $c_{\min} = 0$.

1. **Stochastic costs:** whenever the learner visits state-action pair (s, a) , she immediately observes (and incurs) an i.i.d cost sampled from some unknown distribution with mean $c(s, a)$.
2. **Stochastic adversary, full information:** before learning starts, an adversary samples K i.i.d. cost functions $\{c_k\}_{k=1}^K$ from some unknown distribution with mean c . At the i -th step of episode k , the learner incurs cost $c_i^k = c_k(s_i^k, a_i^k)$. Only at the end of this episode (after the goal state is reached), the learner observes the entire cost function c_k .
3. **Stochastic adversary, bandit feedback:** this is the same as above, except that at the end of episode k , the learner only observes the costs of all visited state-action pairs: $\{c_k(s_i^k, a_i^k)\}_{i=1}^{I_k}$.

The learner’s objective is to minimize her regret, defined as the difference between her total incurred cost and the total expected cost of the best proper policy: $R_K = \sum_{k=1}^K \sum_{i=1}^{I_k} c_i^k - K \cdot V^{\pi^*, P, c}(s_{\text{init}})$, where π^* is the optimal proper policy satisfying $\pi^* \in \arg\min_{\pi \in \Pi} V^{\pi, P, c}(s)$ for all $s \in \mathcal{S}$.

Adversarial Environments We also consider the more challenging environment that adapts to learner’s behavior in a possibly malicious manner. Specifically, in episode k , the environment decides an arbitrary cost function $c_k : \mathcal{S} \times \mathcal{A} \rightarrow [c_{\min}, 1]$ which could depend on the learner’s algorithm as well as her randomness before episode k . The learner then suffers cost $c_i^k = c_k(s_i^k, a_i^k)$ at the i -th step of episode k . Similarly to the stochastic adversary case, the learner observes information on c_k only after she reaches the goal state in episode k , and she observes the entire c_k in the full-information setting or just the cost of visited state-action pairs $\{c_k(s_i^k, a_i^k)\}_{i=1}^{I_k}$ in the bandit setting. The objective is again to minimize her regret against the optimal proper policy in hindsight: $R_K = \sum_{k=1}^K \left(\sum_{i=1}^{I_k} c_i^k - V^{\pi^*, P, c_k}(s_{\text{init}}) \right)$, where we overload the notation π^* to denote the overall optimal proper policy such that $\pi^* \in \arg\min_{\pi \in \Pi} \sum_{k=1}^K V^{\pi, P, c_k}(s)$ for all $s \in \mathcal{S}$.

Key Parameters and Notations Let $T^\pi(s)$ be one plus the expected number of steps to reach the goal if one follows policy π starting from state s . Four parameters play a key role in our analysis and regret bounds: $B_\star = \max_s V^{\pi^*, P, c}(s)$, the maximum expected cost of the optimal policy starting from any state; $T_\star = T^{\pi^*}(s_{\text{init}})$, the hitting time of the optimal policy starting from the initial state; $T_{\max} = \max_s T^{\pi^*}(s)$, the maximum hitting time of the optimal policy starting from any state; and $D = \max_s \min_\pi T^\pi(s)$, the SSP-diameter. We also define the *fast policy* π_f such that $\pi_f \in \arg\min_\pi T^\pi(s)$ for all state s . Similarly to previous works, in most discussions we assume the knowledge of all four parameters and the fast policy, and defer to [Appendix E](#) what we can achieve when some of these are unknown. We also assume $B_\star \geq 1$ for simplicity.

For $n \in \mathbb{N}_+$, we define $[n] = \{1, \dots, n\}$. $\mathbb{E}_k[\cdot]$ denotes the conditional expectation given everything before episode k . The notation $\tilde{\mathcal{O}}(\cdot)$ hides all logarithmic terms including $\ln K$ and $\ln \frac{1}{\delta}$ for some confidence level $\delta \in (0, 1)$. For a distribution $\tilde{P} \in \Delta_{\mathcal{S}_+}$ and a function $V : \mathcal{S}_+ \rightarrow \mathbb{R}$, define $\tilde{P}V = \mathbb{E}_{s \sim \tilde{P}}[V(s)]$.

3. Stacked Discounted Approximation and Algorithm Template

Policy optimization algorithm have been naturally derived in many MDP models. In the finite-horizon setting, one can update the policy at the end of each episode using the cost for this episode that is always bounded. In the discounted setting or average reward setting with some ergodic assumption, one can also update the policy after a certain fixed number of steps since the short-term information is enough to predict the long-term behavior reasonably well. However, this is

not possible in SSP: the hitting time of an arbitrary policy can be arbitrarily large in SSP, and only looking at a fixed number of steps can not always provide accurate information.

A natural solution would be to approximate SSP by other MDP models, and then apply PO in the reduced model. Approximating SSP instances by finite-horizon MDPs (Chen et al., 2021a,d; Cohen et al., 2021) or discounted MDPs (Tarbouriech et al., 2021; Min et al., 2021) is a common practice in the literature, but both have their pros and cons. Finite-horizon approximation shrinks the estimation error exponentially fast and usually leads to optimal regret (Chen et al., 2021d; Cohen et al., 2021). However, it greatly increases the space complexity of the algorithm as it needs to store non-stationary policies with horizon of order $\tilde{\mathcal{O}}(T_{\max})$ or $\tilde{\mathcal{O}}(\frac{B_*}{c_{\min}})$ (as shown in Table 1). Discounted approximation, on the other hand, produces stationary policies, but the estimation error decreases only linearly in the effective horizon $(1 - \gamma)^{-1}$, where γ is the discounted factor. This often leads to sub-optimal regret bounds and large time complexity (Tarbouriech et al., 2021). We include a detailed discussion on limitations of existing approximation schemes in Appendix B.1. These issues greatly limit the practical potential of these methods, and PO methods built on top of them would be less interesting.

To address these issues and achieve optimal regret with small space complexity, we introduce a new approximation scheme called *Stacked Discounted Approximation*, which is a hybrid of finite-horizon and discounted approximations. The key idea is as follows: the finite-horizon approximation requires a horizon of order $\mathcal{O}(T_{\max} \ln K)$, but one can imagine that policies at nearby layers are close to each other and can be approximated by one stationary policy. Thus, we propose to achieve the best of both worlds by dividing the layers into $\mathcal{O}(\ln K)$ parts and performing discounted approximation within each part with an effective horizon $\mathcal{O}(T_{\max})$. Formally, we define the following.

Definition 1 For an SSP instance $\mathcal{M} = (\mathcal{S}, s_{\text{init}}, g, \mathcal{A}, P)$, we define, for number of layers H , discounted factor γ , and terminal cost c_f , another SSP instance $\dot{\mathcal{M}} = (\dot{\mathcal{S}}, \dot{s}_{\text{init}}, g, \mathcal{A}, \dot{P})$ as follows:

1. $\dot{\mathcal{S}} = \mathcal{S} \times [H + 1]$, $\dot{s}_{\text{init}} = (s_{\text{init}}, 1)$, and the goal state g remains the same.
2. Transition from (s, h) to (s', h') is only possible for $h' \in \{h, h + 1\}$: for any $h \leq H$ and $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, we have $\dot{P}_{(s,h),a}(s', h) = \gamma P_{s,a}(s')$ (stay in the same layer with probability γ), $\dot{P}_{(s,h),a}(s', h + 1) = (1 - \gamma) P_{s,a}(s')$ (proceed to the next layer with probability $1 - \gamma$), and $\dot{P}_{(s,h),a}(g) = P_{s,a}(g)$; for $h = H + 1$, we have $\dot{P}_{(s,H+1),a}(g) = 1$ for any (s, a) (immediately reach the goal if at layer $H + 1$). For notational convenience, we also write $\dot{P}_{(s,h),a}(s', h')$ as $P_{(s,h),a}(s', h')$ or $P_{s,a,h}(s', h')$, and $\dot{P}_{(s,h),a}(g)$ as $P_{(s,h),a}(g)$ or $P_{s,a,h}(g)$.
3. For any cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ in \mathcal{M} , we define a cost function \dot{c} for $\dot{\mathcal{M}}$ such that $\dot{c}((s, h), a) = c(s, a)$ for $h \in [H]$ and $\dot{c}((s, H + 1), a) = c_f$ (terminal cost). For notational convenience, we also write $\dot{c}((s, h), a)$ as $c((s, h), a)$ or $c(s, a, h)$.

For any stationary policy π in $\dot{\mathcal{M}}$, we write $\pi(a|(s, h))$ as $\pi(a|s, h)$, and we often abuse the notation $Q^{\pi, P, c}$ and $V^{\pi, P, c}$ to represent the value functions with respect to policy π , transition \dot{P} , and cost function \dot{c} . We also often use (s, a, h) in place of $((s, h), a)$ for function input, that is, we write $f((s, h), a)$ as $f(s, a, h)$.

Define $\hat{\pi}^*$ for $\dot{\mathcal{M}}$ that mimics the behavior of π^* , in the sense that $\hat{\pi}^*(\cdot|s, h) = \pi^*(\cdot|s)$. If we set $\gamma = 1 - \frac{1}{2T_{\max}}$, by the definition of T_{\max} , it can be shown that the probability of $\hat{\pi}^*$ transiting to the next layer before reaching g is upper bounded by $1/2$. If we further set $H = \mathcal{O}(\ln K)$, then

the probability of transiting to the $(H + 1)$ -th layer before reaching g is at most $\frac{1}{2^H} = \tilde{\mathcal{O}}(1/K)$. As a result, the estimation error decreases exponentially in the number of layers while the policy only changes for $\mathcal{O}(\ln K)$ many times. More importantly, due to the discounted factor, the expected hitting time of any policy is of order $\mathcal{O}(\frac{H}{1-\gamma}) = \mathcal{O}(T_{\max} \ln K)$, which controls the cost of exploration and enables the learner to only update its policy at the end of an episode. We summarize the intuition above in the following lemma.

Lemma 2 *For any cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ and terminal cost c_f , we have $V^{\pi, P, c}(s, h) \leq \frac{H-h+1}{1-\gamma} + c_f$ for any $h \in [H]$, $s \in \mathcal{S}$, and policy π in $\dot{\mathcal{M}}$. Moreover, if $\gamma = 1 - \frac{1}{2T_{\max}}$, we further have $Q^{\hat{\pi}^*, P, c}(s, a, h) \leq Q^{\pi^*, P, c}(s, a) + \frac{c_f}{2^{H-h+1}}$ for any $h \in [H]$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$.*

Proof The first statement is because in expectation it takes any policy $\frac{1}{1-\gamma}$ steps to transit from one layer to the next and each step incurs at most 1 cost (except for the terminal cost). For the second statement, note that $V^{\pi, P, c}(s, H+1) = Q^{\pi, P, c}(s, a, H+1) = c_f$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, and for any $h \in [H]$, $V^{\pi, P, c}(s, h) = \sum_{a \in \mathcal{A}} \pi(a|s, h) Q^{\pi, P, c}(s, a, h)$ and

$$Q^{\pi, P, c}(s, a, h) = c(s, a) + \gamma P_{s,a} V^{\pi, P, c}(\cdot, h) + (1 - \gamma) P_{s,a} V^{\pi, P, c}(\cdot, h+1),$$

where we abuse the notation and define $V^{\pi, P, c}(g, h) = 0$ for all $h \in [H+1]$. Now we prove the second statement by induction for $h = H+1, \dots, 1$. The base case $h = H+1$ is clearly true. For $h \leq H$, we bound $Q^{\hat{\pi}^*, P, c}(s, a, h) - Q^{\pi^*, P, c}(s, a)$ as follows:

$$\begin{aligned} & \gamma P_{s,a} V^{\hat{\pi}^*, P, c}(\cdot, h) + (1 - \gamma) P_{s,a} V^{\hat{\pi}^*, P, c}(\cdot, h+1) - P_{s,a} V^{\pi^*, P, c} \\ & \leq \gamma P_{s,a} (V^{\hat{\pi}^*, P, c}(\cdot, h) - V^{\pi^*, P, c}) + (1 - \gamma) \frac{c_f}{2^{H-h}} \\ & \quad (V^{\hat{\pi}^*, P, c}(s, h+1) - V^{\pi^*, P, c}(s)) \leq \frac{c_f}{2^{H-h}} \text{ by induction} \\ & = \gamma \mathbb{E}_{s' \sim P_{s,a}, a' \sim \pi^*(s')} \left[Q^{\hat{\pi}^*, P, c}(s', a', h) - Q^{\pi^*, P, c}(s', a') \right] + (1 - \gamma) \frac{c_f}{2^{H-h}}. \end{aligned}$$

By repeating the arguments above, we arrive at

$$Q^{\hat{\pi}^*, P, c}(s, a, h) - Q^{\pi^*, P, c}(s, a) \leq \mathbb{E} \left[\sum_{t=1}^I \gamma^{t-1} (1 - \gamma) \frac{c_f}{2^{H-h}} \middle| \pi^*, P, s_1 = s, a_1 = a \right],$$

where I is the (random) number of steps it takes for π^* to reach the goal in \mathcal{M} starting from (s, a) . Bounding γ^{t-1} by 1 and $\mathbb{E}[I]$ by T_{\max} , we then obtain the upper bound $\frac{(1-\gamma)T_{\max}c_f}{2^{H-h}} = \frac{c_f}{2^{H-h+1}}$, which finishes the induction. \blacksquare

Remark 3 *Applying the first statement of Lemma 2 with $c(s, a) = 1$ and $c_f = 1$, we have the expected hitting time of any policy in $\dot{\mathcal{M}}$ bounded by $\frac{H}{1-\gamma} + 1$ starting from any state in any layer.*

Now we complete the approximation by showing how to solve the original problem via solving its stacked discounted version. Given a policy π for $\dot{\mathcal{M}}$, define a non-stationary randomized policy $\sigma(\pi)$ for \mathcal{M} as follows: it maintains an internal counter h initialized as 1. In each time step before reaching the goal, it first follows $\pi(\cdot|s, h)$ for one step, where s is the current state. Then, it samples a Bernoulli random variable X with mean γ , and it increases h by 1 if $X = 0$. When $h = H+1$,

Algorithm 1 Template for Policy Optimization with Stacked Discounted Approximation

Initialize: \mathcal{P}_1 , the set of all possible transition functions in $\dot{\mathcal{M}}$ (Eq. (3)); $\eta > 0$, some learning rate.
for $k = 1, \dots, K$ **do**

- Compute $\pi_k(a|s, h) \propto \exp\left(-\eta \sum_{j=1}^{k-1} (\tilde{Q}_j(s, a, h) - B_j(s, a, h))\right)$.
 - Execute $\sigma(\pi_k)$ for one episode (see the paragraph before Lemma 4).
 - Compute some optimistic action-value estimator \tilde{Q}_k and exploration bonus function B_k using \mathcal{P}_k and observations from episode k .
 - Compute transition confidence set \mathcal{P}_{k+1} , as defined in Eq. (4).
-

it executes the fast policy π_f until reaching the goal state. Clearly, the trajectory of $\sigma(\pi)$ indeed follows the same distribution of the trajectory of π in $\dot{\mathcal{M}}$. We show that as long as H is large enough and c_f is of order $\tilde{O}(D)$, this reduction makes sure that the regret between these two problems are similar. The proof is deferred to Appendix B.

Lemma 4 *Let $\gamma = 1 - \frac{1}{2T_{\max}^{\circ}}$, $H = \lceil \log_2(c_f K) \rceil$, $c_f = \lceil 4D \ln \frac{2K}{\delta} \rceil$ for some $\delta \in (0, 1)$, and π_1, \dots, π_K be policies for \mathcal{M} . Then the regret of executing $\sigma(\pi_1), \dots, \sigma(\pi_K)$ in \mathcal{M} satisfies $R_K \leq \hat{R}_K + \tilde{O}(1)$ with probability at least $1 - \delta$, where $\hat{R}_K = \sum_{k=1}^K \left(\sum_{i=1}^{J_k} c_i^k + \check{c}_{J_k+1}^k - V^{\hat{\pi}^*, P, c}(s_1^k, 1) \right)$ for stochastic environments, and $\hat{R}_K = \sum_{k=1}^K \left(\sum_{i=1}^{J_k} c_i^k + \check{c}_{J_k+1}^k - V^{\hat{\pi}^*, P, c_k}(s_1^k, 1) \right)$ for adversarial environments. Here, J_k is the number of time steps in episode k before the learner reaching g or the counter of $\sigma(\pi_k)$ reaching $H + 1$, and $\check{c}_{J_k+1}^k = c_f \mathbb{I}\{s_{J_k+1}^k \neq g\}$.*

Computing Fast Policy and Estimating Diameter For simplicity, we assume knowledge of the diameter and the fast policy above. When these are unknown, one can follow the ideas in (Chen and Luo, 2021) for estimating the fast policy with constant overhead and then adopt their template for learning without knowing the diameter; see (Chen and Luo, 2021, Lemma 1, Appendix E).

Policy Optimization in Stacked Discounted MDPs Now we describe a template of performing policy optimization with the stacked discounted approximation. The pseudocode is shown in Algorithm 1. To handle unknown transition, we maintain standard Bernstein-style transition confidence sets $\{\mathcal{P}_k\}_{k=1}^K$ whose definition is deferred to Appendix A.1. In episode k , the algorithm first computes policy π_k in $\dot{\mathcal{M}}$ following the multiplicative weights update with some learning rate $\eta > 0$, such that $\pi_k(a|s, h) \propto e^{-\eta \sum_{j=1}^{k-1} (\tilde{Q}_j(s, a, h) - B_j(s, a, h))}$ for some optimistic action-value estimator \tilde{Q}_j and exploration bonus function B_j (computed from past observations and confidence sets). Then, it executes $\sigma(\pi_k)$ for this episode. Finally, it computes confidence set \mathcal{P}_{k+1} . All algorithms introduced in this work follow this template and differ from each other in the definition of \tilde{Q}_k and B_k . Ideally, $\tilde{Q}_k - B_k$ should be the action-value function with respect to the true transition, the true cost function, and policy π_k , but since the transition and cost functions are unknown, the key challenge lies in constructing accurate estimators that simultaneously encourage sufficient exploration.

Optimistic Transitions Our algorithms require using some optimistic transitions. Specifically, for a policy π , a confidence set \mathcal{P} , and a cost function c , let $\Gamma(\pi, \mathcal{P}, c)$ be the corresponding optimistic transition such that $\Gamma(\pi, \mathcal{P}, c) \in \operatorname{argmin}_{P \in \mathcal{P}} V^{\pi, P, c}(s, h)$ for all state (s, h) . The existence of such an optimistic transition and how it can be efficiently approximated via Extended Value Iteration (in

at most $\tilde{O}(T_{\max})$ iterations) are deferred to [Appendix A.2](#). We abuse the notation and denote by $V^{\pi, \mathcal{P}, c}$ and $Q^{\pi, \mathcal{P}, c}$ the value function $V^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$ and action-value function $Q^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$.

Occupancy Measure Another important concept for subsequent discussions is *occupancy measure*. Given a policy $\pi : \mathring{S} \rightarrow \Delta_{\mathcal{A}}$ and a transition function $P = \{P_{s,a,h}\}_{(s,h) \in \mathring{S}, a \in \mathcal{A}}$ with $P_{s,a,h} \in \Delta_{\mathring{S}_+}$ and $\mathring{S}_+ = \mathring{S} \cup \{g\}$, define $q_{\pi, P} : \mathring{S} \times \mathcal{A} \times \mathring{S}_+ \rightarrow \mathbb{R}_+$ such that $q_{\pi, P}(\mathring{s}, a, \mathring{s}') = \mathbb{E}[\sum_{i=1}^I \mathbb{I}\{s_i = \mathring{s}, a_i = a, s_{i+1} = \mathring{s}'\} | \pi, P, s_1 = \mathring{s}_{\text{init}}]$ is the expected number of visits to $(\mathring{s}, a, \mathring{s}')$ following policy π in a stacked discounted MDP with transition P . We also let $q_{\pi, P}(s, a, h) = \sum_{\mathring{s}'} q_{\pi, P}((s, h), a, \mathring{s}')$ be the expected number of visits to $((s, h), a)$ and $q_{\pi, P}(s, h) = \sum_a q_{\pi, P}(s, a, h)$ be the number of visits to (s, h) . Note that if a function $q : \mathring{S} \times \mathcal{A} \times \mathring{S}_+ \rightarrow \mathbb{R}_+$ is an occupancy measure, then the corresponding policy π_q satisfies $\pi_q(a|s, h) \propto q(s, a, h)$ and the corresponding transition function P_q satisfies $P_{q,s,a,h}(s', h') \propto q((s, h), a, (s', h'))$. Moreover, $V^{\pi, P, c}(\mathring{s}_{\text{init}}) = \langle q_{\pi, P}, c \rangle$ holds for any policy π , transition function P and cost function c .

Other Notations In the rest of the paper, following [Lemma 4](#) we set $\gamma = 1 - \frac{1}{2T_{\max}}$, $H = \lceil \log_2(c_f K) \rceil$, and $c_f = \lceil 4D \ln \frac{2K}{\delta} \rceil$ for some failure probability $\delta \in (0, 1)$. Define $\chi = 2HT_{\max} + c_f$ as the value function upper bound in \mathcal{M} (according to the first statement of [Lemma 2](#)). Also define $q_k = q_{\pi_k, P}$, $q^* = q_{\pi^*, P}$, and $L = \lceil \frac{8H}{1-\gamma} \ln(2T_{\max}K/\delta) \rceil$.

4. Algorithms and Results for Stochastic Environments

In this section, we consider policy optimization in stochastic environments with three types of feedback introduced in [Section 2](#). We show that a simple policy optimization framework can be used to achieve near-optimal regret for all three settings. In contrast, previous works treat stochastic costs and stochastic adversaries as different problems and solve them via different approaches. Below, we start by describing the algorithm and its guarantees, followed by some explanation behind the algorithm design and then some key ideas and novelty in the analysis.

Algorithm As mentioned, the only elements left to be specified in [Algorithm 1](#) are \tilde{Q}_k and B_k . For stochastic environments, we simply set $B_k(s, a, h) = 0$ for all (s, a, h) since exploration is relatively easier in this case. We now discuss how to construct \tilde{Q}_k .

- **Action-value estimator** \tilde{Q}_k is defined as $Q^{\pi_k, \mathcal{P}_k, \tilde{c}_k}$ for some corrected cost estimator \tilde{c}_k :

$$\tilde{c}_k(s, a, h) = (1 + \lambda \hat{Q}_k(s, a, h)) \hat{c}_k(s, a, h) + e_k(s, a, h), \quad (1)$$

where λ is some parameter, $\hat{Q}_k = Q^{\pi_k, \mathcal{P}_k, \hat{c}_k}$ is another action-value estimator with respect to some optimistic cost estimator \hat{c}_k , and e_k is some correction term (all to be specified below).

- **Optimistic cost estimator** \hat{c}_k is defined as

$$\begin{aligned} \hat{c}_k(s, a, h) &= \bar{c}_k(s, a) \mathbb{I}\{h \leq H\} + c_f \mathbb{I}\{h = H + 1\}, \\ \bar{c}_k(s, a) &= \max \{0, \bar{c}_k(s, a) - 2\sqrt{\bar{c}_k(s, a)\alpha_k(s, a)} - 7\alpha_k(s, a)\}, \end{aligned}$$

where $\bar{c}_k(s, a)$ is the average of all costs that are observed for (s, a) in episode $j = 1, \dots, k-1$ before $\sigma(\pi_j)$ switches to the fast policy, and $\alpha_k(s, a)$ is $\iota = \ln(2SALK/\delta)$ divided by the number of samples used in computing $\bar{c}_k(s, a)$, such that $2\sqrt{\bar{c}_k(s, a)\alpha_k(s, a)} + 7\alpha_k(s, a)$ is a

standard Bernstein-style deviation term (thus making $\widehat{c}_k(s, a)$ an optimistic underestimator). We note that naturally, the way to compute $\widehat{c}_k(s, a)$ is different for different types of feedback — for stochastic costs, we might have multiple samples for (s, a) in one episode, while for stochastic adversaries, we have exactly one sample in each episode in the full-information setting, and one or zero samples in the bandit setting.

- **Correction term** $e_k(s, a, h)$ is defined as 0 for stochastic costs; $(8\ell\sqrt{\widehat{c}_k(s, a, h)/k} + \beta'\widehat{Q}_k(s, a, h))\mathbb{I}\{h \leq H\}$ with $\beta' = \min\{1/T_{\max}, 1/\sqrt{DT_*K}\}$ for stochastic adversary with full information; and $\beta\widehat{Q}_k(s, a, h)\mathbb{I}\{h \leq H\}$ with $\beta = \min\{1/T_{\max}, \sqrt{SA/DT_*K}\}$ for stochastic adversary with bandit feedback.
- **Parameter tuning:** learning rate η (for the multiplicative weights update) is set to $\min\{1/3T_{\max}(8\ell + \times/T_{\max})^2, 1/\sqrt{\lambda T_{\max}^4 K}\}$, and the parameter λ is set to $\min\{1/T_{\max}, \sqrt{S^2 A/\square^2 K}\}$ where \square is B_* for stochastic costs and D for stochastic adversaries.

We now state the regret guarantees of our algorithm for each of the three settings (proofs are deferred to [Appendix C.2.1](#) to [Appendix C.2.3](#)).

Theorem 5 *For stochastic costs, [Algorithm 1](#) with the instantiation above achieves $R_K = \tilde{O}(B_* S \sqrt{AK} + T_{\max}^3 (S^2 AK)^{1/4} + S^4 A^{2.5} T_{\max}^4)$ with probability at least $1 - 32\delta$.*

Ignoring lower-order terms, our bound almost matches the minimax bound $\tilde{O}(B_* \sqrt{SAK})$ of ([Cohen et al., 2021](#)), with a \sqrt{S} factor gap.

Theorem 6 *For stochastic adversary with full information, [Algorithm 1](#) with the instantiation above achieves $R_K = \tilde{O}(\sqrt{DT_*K} + DS\sqrt{AK} + T_{\max}^3 (S^2 A^3 K)^{1/4} + S^4 A^{2.5} T_{\max}^4)$ with probability at least $1 - 50\delta$.*

Theorem 7 *For stochastic adversary with bandit feedback, [Algorithm 1](#) with the instantiation above achieves $R_K = \tilde{O}(\sqrt{SADT_*K} + DS\sqrt{AK} + T_{\max}^3 SA^{5/4} K^{1/4} + S^4 A^{2.5} T_{\max}^4)$ with probability at least $1 - 50\delta$.*

Ignoring lower-order terms again, these bounds for stochastic adversary match the best known results from ([Chen and Luo, 2021](#)), and they all exhibit a \sqrt{S} gap in the term $DS\sqrt{AK}$ compared to the best existing lower bounds ([Chen and Luo, 2021](#)).

We emphasize again that besides the simplicity of PO, one algorithmic advantage of our method compared to those based on finite-horizon approximation is its low space complexity to store policies — the horizon H for our method is only $\mathcal{O}(\ln K)$, while the horizon for other works ([Chen and Luo, 2021](#); [Cohen et al., 2021](#)) is $\tilde{O}(T_{\max})$ when T_{\max} is known or otherwise $\tilde{O}(B_*/c_{\min})$. Note that when $c_{\min} = 0$, a common technique is to perturb the cost and deal with a modified problem with $c_{\min} = 1/\text{poly}(K)$, in which case our space complexity is exponentially better. In fact, even for time complexity, although our method requires calculating optimistic transition and might need $\tilde{O}(T_{\max})$ rounds of Extended Value Iteration, this procedure could terminate much earlier, while the finite-horizon approximation approaches always need at least $\Omega(T_{\max})$ time complexity since that is the horizon of the MDP they are dealing with.

Analysis highlights We start by explaining the design of the corrected cost estimator Eq. (1). Roughly speaking, standard analysis of PO (specifically, by (Chen and Luo, 2021, Lemma 9) and then Lemma 26) leads to a term of order $\lambda \sum_{k=1}^K \langle q_k, c \circ \widehat{Q}_k \rangle$ due to the transition estimation error, which can be prohibitively large (for functions f and g with the same domain, we define $(f \circ g)(x) = f(x)g(x)$). Introducing the correction bias $\lambda \widehat{Q}_k(s, a, h) \widehat{c}_k(s, a, h)$ in Eq. (1), on the other hand, has the effect of transforming this problematic term into its counterpart $\lambda \sum_{k=1}^K \langle q^*, c \circ \widehat{Q}_k \rangle$ in terms of q^* instead of q_k . Bounding the latter term, however, requires a property that PO enjoys, that is, a regret bound for any initial state-action pair: $\sum_{k=1}^K (\widehat{Q}_k - Q^{\hat{\pi}^*, P, c})(s, a, h) = \tilde{O}(\sqrt{K})$ for any (s, a, h) . In contrast, approaches based on occupancy measure (Chen and Luo, 2021) only guarantee a regret bound starting from s_{init} . This makes PO especially compatible with our stacked discounted approximation. Based on this observation, we further have $\lambda \sum_{k=1}^K \langle q^*, c \circ \widehat{Q}_k \rangle \approx \lambda \sum_{k=1}^K \langle q^*, c \circ Q^{\hat{\pi}^*, P, c} \rangle$, where the latter term is only about the behavior of the optimal policy and is thus nicely bounded (see e.g. Lemma 20). To sum up, the correction term $\lambda \widehat{Q}_k(s, a, h) \widehat{c}_k(s, a, h)$ in Eq. (1) together with a favorable property of PO helps us control the transition estimation error in a near-optimal way.

For stochastic adversaries, an extra complication arises due to the cost estimation error $\sum_{k=1}^K \langle q_k, c - \widehat{c}_k \rangle$, which results in the extra $\sqrt{DT_*K}$ or $\sqrt{SADT_*K}$ term in the minimax regret bound (depending on the feedback type). Obtaining this optimal cost estimation error requires us to add yet another correction term e_k in Eq. (1). Specifically, we show that $\sum_{k=1}^K \langle q_k, c - \widehat{c}_k \rangle \approx \sum_{k=1}^K \langle q_k, e_k \rangle$ for e_k defined as in our algorithm description. Then, the role of adding e_k in Eq. (1) is again to turn the term above to its counterpart $\sum_{k=1}^K \langle q^*, e_k \rangle$ in terms of the optimal policy's behavior, which can then be nicely bounded. As a side product, we note that this also provides a much cleaner analysis on bounding the cost estimation error compared to (Chen and Luo, 2021), where they require explicitly forcing the expected hitting time of the learner's policy to be bounded.

Finally, we point out another novelty in our analysis. Compared to other approaches that act according to the exact optimal policy of an estimated MDP, PO incurs an additional cost due to only updating the policy incrementally in each episode. This cost is often of order $\tilde{O}(\sqrt{K})$ and is one of the dominating terms in the regret bound; see e.g. (Shani et al., 2020; Wu et al., 2021) for the finite-horizon case. For SSP, this is undesirable because it also depends on T_* or even T_{\max} . Reducing this cost has been studied from the optimization perspective — for example, an improved $\tilde{O}(1/K)$ convergence rate of PO has been established recently by (Agarwal et al., 2021). However, adopting their analysis to regret minimization requires additional efforts. Specifically, we need to carefully bound the bias from using an action-value estimator in the policy's update, which can be shown to be approximately bounded by $\sum_{k=1}^K (\tilde{Q}_{k+1} - \tilde{Q}_k)(s, a, h)$. In Lemma 25, we show that this term is of lower order by carefully analyzing the drift $(\tilde{Q}_{k+1} - \tilde{Q}_k)(s, a, h)$ in each episode.

Remark 8 We remark that our algorithm can be applied to finite-horizon MDPs with inhomogeneous transition and gives a $\tilde{O}(\sqrt{S^2 AH^3 K})$ regret bound, improving over that of (Shani et al., 2020) by a factor of \sqrt{H} where H is the horizon. We omit the details but only mention that the improvement comes from two sources: first, the aforementioned improved PO analysis turns a $\tilde{O}(H^2 \sqrt{K})$ regret term into a lower order term; second, we use Bernstein-style transition confidence set to obtain an improved $\tilde{O}(\sqrt{S^2 AH^3 K})$ transition estimation error.

5. Algorithms and Results for Adversarial Environments

We move on to consider the more challenging environments with adversarial costs, where the extra exploration bonus function B_k in [Algorithm 1](#) now plays an important role. Even in the finite-horizon setting, developing efficient PO methods in this case can be challenging, and [Luo et al. \(2021\)](#) proposed the so-called ‘‘dilated bonuses’’ to guide better exploration, which we also adopt and extend to SSP. Specifically, for a policy π , a transition confidence set \mathcal{P} , and some bonus function $b : \mathcal{S} \times \mathcal{A} \times [H + 1] \rightarrow \mathbb{R}$, we define the corresponding dilated bonus function $B^{\pi, \mathcal{P}, b} : \mathcal{S} \times \mathcal{A} \times [H + 1] \rightarrow \mathbb{R}$ as: $B^{\pi, \mathcal{P}, b}(s, a, H + 1) = b(s, a, H + 1)$ and for $h \in [H]$,

$$B^{\pi, \mathcal{P}, b}(s, a, h) = b(s, a, h) + \left(1 + \frac{1}{H'}\right) \max_{\hat{P} \in \mathcal{P}} \hat{P}_{s, a, h} \left(\sum_{a'} \pi(a' | \cdot, \cdot) B^{\pi, \mathcal{P}, b}(\cdot, a', \cdot) \right), \quad (2)$$

where $H' = \frac{8(H+1)\ln(2K)}{1-\gamma}$ is the dilated coefficient. Intuitively, $B^{\pi, \mathcal{P}, b}$ is the dilated (by a factor of $1 + 1/H'$) and optimistic (by maximizing over \mathcal{P}) version of the action-value function with respect to π and b . In the finite-horizon setting ([Luo et al., 2021](#)), this can be computed directly via dynamic programming, but how to compute it in a stacked discounted MDP (or even why it exists) is less clear. Fortunately, we show that this can indeed be computed efficiently via a combination of dynamic programming and Extended Value Iteration; see [Appendix D.4](#).

Algorithm (full information) We now describe our algorithm for the adversarial full-information case (where c_k is revealed at the end of episode k). It suffices to specify \tilde{Q}_k and B_k in [Algorithm 1](#).

- **Action-value estimator** \tilde{Q}_k is defined as $\tilde{Q}_k = Q^{\pi_k, \mathcal{P}_k, \tilde{c}_k}$, where $\tilde{c}_k(s, a, h) = (1 + \lambda \tilde{Q}_k(s, a, h))c_k(s, a, h)$ for some parameter λ and $\tilde{Q}_k = Q^{\pi_k, \mathcal{P}_k, c_k}$.
- **Dilated bonus** B_k is defined as $B^{\pi_k, \mathcal{P}_k, b_k}$ with $b_k(s, a, h) = 2\eta \sum_{a \in \mathcal{A}} \pi_k(a | s, h) \tilde{A}_k(s, a, h)^2$, where $\tilde{A}_k(s, a, h) = \tilde{Q}_k(s, a, h) - \tilde{V}_k(s, h)$ (advantage function) and $\tilde{V}_k = V^{\pi_k, \mathcal{P}_k, \tilde{c}_k}$.
- **Parameter tuning:** $\eta = \min\{1/(64\chi^2\sqrt{HH'}), 1/\sqrt{DK}\}$ and $\lambda = \min\{1/\chi, 48\eta + \sqrt{S^2A/DT_*K}\}$.

Our algorithm enjoys the following guarantee (whose proof can be found in [Appendix D.1](#)).

Theorem 9 *For adversarial costs with full information, [Algorithm 1](#) with the instantiation above achieves $R_K = \tilde{\mathcal{O}}(T_*\sqrt{DK} + \sqrt{S^2ADT_*K} + S^4A^2T_{\max}^5)$ with probability at least $1 - 20\delta$.*

The best existing bound is $\tilde{\mathcal{O}}(\sqrt{S^2ADT_*K})$ from ([Chen and Luo, 2021](#)). Ignoring the lower order term, our result matches theirs when $T_* \leq S^2A$ (and is worse by a $\sqrt{T_*/S^2A}$ factor otherwise). Our algorithm enjoys better time and space complexity though, similar to earlier discussions.

Analysis highlights For simplicity we assume that the true transition is known, in which case our bound is only $\tilde{\mathcal{O}}(T_*\sqrt{DK})$ (the other term $\sqrt{S^2ADT_*K}$ is only due to transition estimation error). A naive way to implement PO would lead to a penalty term T_*/η plus a stability term $\eta \sum_{k=1}^K \sum_{s, h} q^*(s, h) \sum_a \pi_k(a | s, h) Q^{\pi_k, \mathcal{P}_k, c_k}(s, a, h)^2$, which eventually leads to a bound of order $\tilde{\mathcal{O}}(T_*T_{\max}\sqrt{K})$ if one bounds $Q^{\pi_k, \mathcal{P}_k, c_k}(s, a, h)$ by $\tilde{\mathcal{O}}(T_{\max})$. Our improvement comes from the following five steps: 1) first, through a careful shifting argument, we show that the stability term can be improved to $\eta \sum_{k=1}^K \sum_{s, h} q^*(s, h) \sum_a \pi_k(a | s, h) A^{\pi_k, \mathcal{P}_k, c_k}(s, a, h)^2$ (recall that A is the advantage function); 2) second, similarly to ([Luo et al., 2021](#)), the dilated bonus B_k helps transform

q^* to q_k in the term above, leading to $\eta \sum_{k=1}^K \langle q_k, (A^{\pi_k, P, c_k})^2 \rangle$; 3) third, in [Lemma 26](#) we show that the previous term is bounded by the variance of the learner's cost, which in turn is at most $\eta \sum_{k=1}^K \langle q_k, c_k \circ Q^{\pi_k, P, c_k} \rangle$; 4) fourth, similarly to [Section 4](#), the correction term $\lambda c_k \circ \hat{Q}_k$ in the definition of \tilde{c}_k helps transform q_k back to q^* , resulting in $\eta \sum_{k=1}^K \langle q^*, c_k \circ Q^{\pi_k, P, c_k} \rangle$; 5) finally, since PO guarantees a regret bound for any initial state (as mentioned in [Section 4](#)), the previous term is close to $\eta \sum_{k=1}^K \langle q^*, c_k \circ Q^{\pi_k^*, P, c_k} \rangle$, which is now only related to the optimal policy and can be shown to be at most $\tilde{O}(\eta DT_* K)$. Combining this with the penalty term T_*/η and picking the best η then results in the claimed $\tilde{O}(T_* \sqrt{DK})$ regret bound.

Algorithm (Bandit Feedback) Finally, we describe our algorithm for the adversarial setting with bandit feedback, starting with the instantiation of B_k followed by that of \tilde{Q}_k .

- Dilated bonus** B_k is again defined as B^{π_k, P_k, b_k} , but with a different b_k function similar to that of [\(Luo et al., 2021\)](#): $b_k(s, a, h) = L' \mathbb{I}\{h \leq H\} \sum_{a'} \pi_k(a'|s, h) \frac{\bar{x}_k(s, a', h) - \underline{x}_k(s, a', h) + 4\theta}{\bar{x}_k(s, a', h) + \theta}$, for some parameters L' and θ . Here, $\bar{x}_k(s, a, h)$ and $\underline{x}_k(s, a, h)$ are respectively the largest and smallest possible probability that $((s, h), a)$ is ever visited in episode k following policy π_k if the transition lies in P_k , and they can be computed efficiently as shown in [Appendix D.5](#).
- Action-value estimator** \tilde{Q}_k is defined as $\tilde{Q}_k(s, a, h) = \frac{G_{k, s, a, h}}{\bar{x}_k(s, a, h) + \theta} \mathbb{I}\{h \leq H\} + c_f \mathbb{I}\{h = H + 1\}$, where $G_{k, s, a, h}$ is the learner's total cost in \mathcal{M} starting from the first visit to $((s, h), a)$ during the first $L + 1$ steps of episode k . Recall the definition of L stated at the end of [Section 3](#), which is a high-probability upper bound on the number of steps any policy in \mathcal{M} takes to reach the last layer (so counting only the first $L + 1$ steps is simply to make sure that $G_{k, s, a, h}$ is always bounded).
- Parameter tuning:** $\eta = \min \left\{ 1/(300HH'T_{\max}L'), \sqrt{1/T_{\max}^2 SAK} \right\}$, $\theta = 2\eta L'$, and $L' = L + c_f$.

We note that this algorithm is in spirit very similar to that of [\(Luo et al., 2021\)](#) for the finite-horizon case. Unfortunately, the correction terms we use throughout other algorithms in this work do not work here for technical reasons, resulting in the following sub-optimal guarantee which still has T_{\max} dependency in the dominating term (see [Appendix D.2](#) for the proof). We remark that the best existing bound is $\tilde{O}(\sqrt{DT_* S^3 A^2 K})$ from [\(Chen and Luo, 2021\)](#).

Theorem 10 *For adversarial costs with bandit feedback, [Algorithm 1](#) with the instantiation above achieves $R_K = \tilde{O}(\sqrt{S^2 AT_{\max}^5 K} + S^{5.5} A^{3.5} T_{\max}^5)$ with probability at least $1 - 28\delta$.*

6. Conclusion

Our work initiates the study of policy optimization for SSP and systematically develops a set of novel algorithms suitable for different settings. Many questions remain open, such as closing the gap between some of our results and the best known results achieved by other types of methods. Moreover, as mentioned, one of the reasons to study PO for SSP is that PO usually works well when combined with function approximation. Our stacked discounted approximation scheme also does not make use of any modeling assumption and should be applicable in more general settings. Although our work is only for the tabular setting, we believe that our results lay a solid foundation for future studies on SSP with function approximation.

Acknowledgments

LC thanks Chen-Yu Wei for many helpful discussions.

References

- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR, 2017.
- Dimitri P Bertsekas and Huizhen Yu. Stochastic shortest path problems under weak conditions. *Lab. for Information and Decision Systems Report LIDS-P-2909, MIT*, 2013.
- Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. In *International Conference on Machine Learning*, pages 1283–1294. PMLR, 2020.
- Liyu Chen and Haipeng Luo. Finding the stochastic shortest path with low regret: The adversarial cost and unknown transition case. In *International Conference on Machine Learning*, 2021.
- Liyu Chen, Mehdi Jafarnia-Jahromi, Rahul Jain, and Haipeng Luo. Implicit finite-horizon approximation and efficient optimal algorithms for stochastic shortest path. *Advances in Neural Information Processing Systems*, 2021a.
- Liyu Chen, Rahul Jain, and Haipeng Luo. Improved no-regret algorithms for stochastic shortest path with linear MDP. *arXiv preprint arXiv:2112.09859*, 2021b.
- Liyu Chen, Haipeng Luo, and Chen-Yu Wei. Impossible tuning made possible: A new expert algorithm and its applications. In *Conference on Learning Theory*, pages 1216–1259. PMLR, 2021c.
- Liyu Chen, Haipeng Luo, and Chen-Yu Wei. Minimax regret for stochastic shortest path with adversarial costs and known transition. In *Conference on Learning Theory*, pages 1180–1215. PMLR, 2021d.
- Alon Cohen, Haim Kaplan, Yishay Mansour, and Aviv Rosenberg. Near-optimal regret bounds for stochastic shortest path. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 8210–8219. PMLR, 2020.
- Alon Cohen, Yonathan Efroni, Yishay Mansour, and Aviv Rosenberg. Minimax regret for stochastic shortest path. *Advances in Neural Information Processing Systems*, 34, 2021.
- Yingjie Fei, Zhuoran Yang, Zhaoran Wang, and Qiaomin Xie. Dynamic regret of policy optimization in non-stationary environments. *Advances in Neural Information Processing Systems*, 33: 6743–6754, 2020.

- Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- Mehdi Jafarnia-Jahromi, Liyu Chen, Rahul Jain, and Haipeng Luo. Online learning for stochastic shortest path model via posterior sampling. *arXiv preprint arXiv:2106.05335*, 2021.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(4), 2010.
- Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- Yeoneung Kim, Insoon Yang, and Kwang-Sung Jun. Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture MDPs. *arXiv preprint arXiv:2111.03289*, 2021.
- Tor Lattimore and Marcus Hutter. Pac bounds for discounted MDPs. In *International Conference on Algorithmic Learning Theory*, pages 320–334. Springer, 2012.
- Boyi Liu, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural trust region/proximal policy optimization attains globally optimal policy. *Advances in neural information processing systems*, 32, 2019.
- Haipeng Luo, Chen-Yu Wei, and Chung-Wei Lee. Policy optimization in adversarial MDPs: Improved exploration via dilated bonuses. *Advances in Neural Information Processing Systems*, 34, 2021.
- Yifei Min, Jiafan He, Tianhao Wang, and Quanquan Gu. Learning stochastic shortest path with linear function approximation. *arXiv preprint arXiv:2110.12727*, 2021.
- Aviv Rosenberg and Yishay Mansour. Stochastic shortest path with adversarially changing costs. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pages 2936–2942. ijcai.org, 2021. doi: 10.24963/ijcai.2021/404. URL <https://doi.org/10.24963/ijcai.2021/404>.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- Lior Shani, Yonathan Efroni, Aviv Rosenberg, and Shie Mannor. Optimistic policy optimization with bandit feedback. In *Proceedings of the 37th International Conference on Machine Learning*, pages 8604–8613, 2020.
- Jacob Steinhardt and Percy Liang. Adaptivity and optimism: An improved exponentiated gradient algorithm. In *International Conference on Machine Learning*, pages 1593–1601. PMLR, 2014.
- Jean Tarbouriech, Evrard Garcelon, Michal Valko, Matteo Pirotta, and Alessandro Lazaric. No-regret exploration in goal-oriented reinforcement learning. In *International Conference on Machine Learning*, pages 9428–9437. PMLR, 2020.

- Jean Tarbouriech, Runlong Zhou, Simon S Du, Matteo Pirota, Michal Valko, and Alessandro Lazaric. Stochastic shortest path: Minimax, parameter-free and towards horizon-free regret. *Advances in Neural Information Processing Systems*, 34, 2021.
- Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. *International Conference on Learning Representations (ICLR)*, 2020.
- Chen-Yu Wei and Haipeng Luo. More adaptive algorithms for adversarial bandits. In *Conference On Learning Theory*, pages 1263–1291. PMLR, 2018.
- Chen-Yu Wei, Mehdi Jafarnia Jahromi, Haipeng Luo, Hiteshi Sharma, and Rahul Jain. Model-free reinforcement learning in infinite-horizon average-reward Markov decision processes. In *International Conference on Machine Learning*, pages 10170–10180. PMLR, 2020.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Tianhao Wu, Yunchang Yang, Han Zhong, Liwei Wang, Simon S Du, and Jiantao Jiao. Nearly optimal policy optimization with stable at any time guarantee. *arXiv preprint arXiv:2112.10935*, 2021.
- Zihan Zhang, Jiaqi Yang, Xiangyang Ji, and Simon S Du. Variance-aware confidence set: Variance-dependent bound for linear bandits and horizon-free bound for linear mixture MDP. *Advances in Neural Information Processing Systems*, 2021.
- Dongruo Zhou, Quanquan Gu, and Csaba Szepesvari. Nearly minimax optimal reinforcement learning for linear mixture Markov decision processes. In *Conference on Learning Theory*, pages 4532–4576. PMLR, 2021.
- Alexander Zimin and Gergely Neu. Online learning in episodic Markovian decision processes by relative entropy policy search. In *Advances in Neural Information Processing Systems*, pages 1583–1591, 2013.

Appendix A. Preliminary for Appendix

Extra Notations Define $\mathring{s}_i^k = (s_i^k, h_i^k)$ as the i -th step in $\mathring{\mathcal{M}}$ in episode k . Define $n_k(s, a, h)$ as the number of visits to $((s, h), a)$ in $\mathring{\mathcal{M}}$ in episode k , and $n_k(s, a) = \sum_{h \leq H} n_k(s, a, h)$ (excluding layer $H + 1$). Define $\bar{J}_k = \min\{L, J_k\}$, $\bar{n}_k(s, a) = \min\{L, n_k(s, a)\}$, and $\bar{n}_k(s, a, h) = \min\{L, n_k(s, a, h)\}$. For any sequence of scalars or functions $\{z_k\}_k$, define $dz_k = z_{k+1} - z_k$. By default we assume $\sum_h = \sum_{h=1}^{H+1}$. For inner product $\langle u, v \rangle$, if $u(s, a)$, $u(s, a, h)$, $v(s, a)$, and $v(s, a, h)$ are all defined, we let $\langle u, v \rangle = \sum_{s,a,h} u(s, a, h)v(s, a, h)$. For functions f and g with the same domain, define function $(f \circ g)(x) = f(x)g(x)$. For any random variable X , define conditional variance $\text{Var}_k[X] = \mathbb{E}_k[(X - \mathbb{E}_k[X])^2]$.

For an occupancy measure q w.r.t policy π and transition P , define $q_{(s,h)}$ as the occupancy measure w.r.t policy π , transition P , and initial state (s, h) , and $q_{(s,a,h)}$ as the occupancy measure w.r.t policy π , transition P , initial state (s, h) , and initial action a . Denote by $x_k(s, a, h)$ the probability that $((s, h), a)$ is ever visited in episode k , $x_k(s, a) = \sum_{h=1}^H x_k(s, a, h)$ the probability that (s, a) is ever visited before layer $H + 1$ in episode k , and $y_k(s, a, h)$ the probability of visiting $((s, h), a)$ again if the agent starts from $((s, h), a)$. For any occupancy measure $q(s, a, h)$, we define $q(s, a) = \sum_{h \leq H} q(s, a, h)$ (excluding layer $H + 1$). Note that $q_k(s, a, h) = \frac{x_k(s, a, h)}{1 - y_k(s, a, h)}$ and $y_k(s, a, h) \leq \gamma$. Thus, we have $q_k(s, a, h) = \mathcal{O}(T_{\max} x_k(s, a, h))$.

Define $\Lambda_{\mathcal{M}}$ as the set of possible transition functions of \mathcal{M} :

$$\Lambda_{\mathcal{M}} = \left\{ P = \{P_{s,a,h}\}_{(s,h) \in \mathring{\mathcal{S}}, a \in \mathcal{A}}, P_{s,a,h} \in \Delta_{\mathring{\mathcal{S}}_+} : P_{s,a,H+1}(g) = 1, \right. \\ \left. \sum_{s' \in \mathcal{S}} P_{s,a,h}(s', h) \leq \gamma, \sum_{s' \in \mathcal{S}} P_{s,a,h}(s', h+1) \leq 1 - \gamma, \right. \\ \left. P_{s,a,h}(s', h') = 0, \forall (s, a) \in \mathcal{S} \times \mathcal{A}, h \in [H], h' \notin \{h, h+1\} \right\}, \quad (3)$$

where $\gamma \cdot \mathcal{X} = \{\gamma x : x \in \mathcal{X}\}$ for some set \mathcal{X} . By definition, the expected hitting time of any stationary policy in an MDP with transition $P \in \Lambda_{\mathcal{M}}$ is upper bounded by $(H + 1)(1 - \gamma)^{-1}$ starting from any state. Therefore, for any occupancy measure q with $P_q \in \Lambda_{\mathcal{M}}$ (for example, q_k and q^*), we have $\sum_{s,a,h} q(s, a, h) \leq (H + 1)(1 - \gamma)^{-1} = \tilde{\mathcal{O}}(T_{\max})$.

Finally define $\mathcal{C}_{\mathcal{M}}$ as the set of possible cost functions of $\mathring{\mathcal{M}}$:

$$\mathcal{C}_{\mathcal{M}} = \left\{ c : \mathring{\mathcal{S}} \rightarrow \mathbb{R}_+ : c(s, a, h) = \tilde{\mathcal{O}}(1), \forall h \leq H, \text{ and } \exists C_0 = \tilde{\mathcal{O}}(T_{\max}), c(s, a, H + 1) = C_0, \forall a \right\}.$$

A.1. Transition Estimation

In this section, we present important lemmas regarding the transition confidence sets $\{\mathcal{P}_k\}_{k=1}^K$. We first prove an auxiliary lemma saying that the number of steps taken by the learner before reaching g or switching to fast policy is well bounded with high probability.

Lemma 11 *With probability at least $1 - \delta$, we have $J_k = \bar{J}_k$ for all $k \in [K]$.*

Proof We want to show that $J_k \leq L = \lceil \frac{8H}{1-\gamma} \ln(2T_{\max}K/\delta) \rceil$ for all $k \in [K]$ with probability at least $1 - \delta$. Let $k \in [K]$, it suffices to show that the expected hitting time of π_k is upper bounded by $\frac{H}{1-\gamma}$ starting from any (s, h) , because then we can apply [Lemma 31](#) and take a union bound over all K episodes.

Note that the expected hitting time (w.r.t J_k) is simply the value function with respect to a cost function that is 1 for all state-action pairs except for 0 cost in the goal state g and layer $H + 1$ (i.e., $c_f = 0$). Thus, by [Lemma 2](#), the expected hitting time starting from (s, h) is bounded by $\frac{H-h+1}{1-\gamma} \leq \frac{H}{1-\gamma}$. \blacksquare

Definition of \mathcal{P}_k We define $\mathcal{P}_k = \bigcap_{s,a,h \leq H} \mathcal{P}_{k,s,a,h}$, where:

$$\begin{aligned} \mathcal{P}_{k,s,a,h} = \{P' \in \Lambda_{\mathcal{M}} : & |\bar{P}_{k,s,a}(s') - P'_{s,a,h}(s', h)/\gamma| \leq \epsilon_k(s, a, s'), \\ & |\bar{P}_{k,s,a}(s') - P'_{s,a,h}(s', h+1)/(1-\gamma)| \leq \epsilon_k(s, a, s'), \\ & |\bar{P}_{k,s,a}(g) - P'_{s,a,h}(g)| \leq \epsilon_k(s, a, g), \forall s' \in \mathcal{S}\}, \end{aligned} \quad (4)$$

where $\epsilon_k(s, a, s') = 4\sqrt{\bar{P}_{k,s,a}(s')\alpha'_k(s, a) + 28\alpha'_k(s, a)}$, $\alpha'_k(s, a) = \frac{L}{N_k^+(s, a)}$, $\bar{P}_{k,s,a}(s') = \frac{N_k(s, a, s')}{N_k^+(s, a)}$ is the empirical transition, $N_k^+(s, a) = \max\{1, N_k(s, a)\}$, $N_k(s, a)$ is the number of visits to (s, a) in episode $j = 1, \dots, k-1$ before $\sigma(\pi_j)$ switches to the fast policy, and $N_k(s, a, s')$ is the number of visits to (s, a, s') in episode $j = 1, \dots, k-1$ before $\sigma(\pi_j)$ switches to the fast policy.

Lemma 12 *Under the event of [Lemma 11](#), we have $\mathring{P} \in \mathcal{P}_k$ for any $k \in [K]$ with probability at least $1 - \delta$.*

Proof Clearly $\mathring{P} \in \Lambda_{\mathcal{M}}$. Moreover, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $s' \in \mathcal{S}_+$ by [Lemma 51](#) and $N_{K+1}(s, a) \leq LK$ under the event of [Lemma 11](#), we have with probability at least $1 - \frac{\delta}{2S^2A}$,

$$|P_{s,a}(s') - \bar{P}_{k,s,a}(s')| \leq \epsilon_k(s, a, s'). \quad (5)$$

By a union bound, we have [Eq. \(5\)](#) holds for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $s' \in \mathcal{S}_+$ with probability at least $1 - \delta$. Then the statement is proved by $\mathring{P}_{s,a,h}(s', h) = \gamma P_{s,a}(s')$, $\mathring{P}_{s,a,h}(s', h+1) = (1-\gamma)P_{s,a}(s')$, and $\mathring{P}_{s,a,h}(g) = P_{s,a}(g)$. \blacksquare

Lemma 13 *Under the event of [Lemma 12](#), for any $P' \in \mathcal{P}_k$, we have for any $\mathring{s}' \in \mathring{\mathcal{S}}_+$:*

$$|P'_{s,a,h}(\mathring{s}') - P_{s,a,h}(\mathring{s}')| \leq 8\sqrt{P_{s,a,h}(\mathring{s}')\alpha'_k(s, a) + 136\alpha'_k(s, a)} \triangleq \epsilon_k^*(s, a, h, \mathring{s}').$$

For simplicity, we also write $\epsilon_k^*(s, a, h, (s', h'))$ as $\epsilon_k^*(s, a, h, s', h')$ for $(s', h') \in \mathring{\mathcal{S}}$.

Proof Under the event of [Lemma 12](#) and by [Eq. \(5\)](#), we have for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, and $s' \in \mathcal{S}_+$:

$$\bar{P}_{k,s,a}(s') \leq P_{s,a}(s') + 4\sqrt{\bar{P}_{k,s,a}(s')\alpha'_k(s, a) + 28\alpha'_k(s, a)}.$$

Applying $x^2 \leq ax + b \implies x \leq a + \sqrt{b}$ with $a = 4\sqrt{\alpha'_k(s, a)}$ and $b = P_{s,a}(s') + 28\alpha'_k(s, a)$, we have

$$\sqrt{\bar{P}_{k,s,a}(s')} \leq 4\sqrt{\alpha'_k(s, a)} + \sqrt{P_{s,a}(s') + 28\alpha'_k(s, a)} \leq \sqrt{P_{s,a}(s')} + 10\sqrt{\alpha'_k(s, a)}.$$

Substituting this back to the definition of ϵ_k , we have

$$\epsilon_k(s, a, s') = 4\sqrt{\bar{P}_{k,s,a}(s')\alpha'_k(s, a)} + 28\alpha'_k(s, a) \leq 4\sqrt{P_{s,a}(s')\alpha'_k(s, a)} + 68\alpha'_k(s, a).$$

Now we start to prove the statement. The statement is clearly true for $\dot{s}' = (s', h')$ with $h' \notin \{h, h+1\}$ since the left-hand side equals to 0. Moreover, by the definition of \mathcal{P}_k , [Lemma 12](#), and $x \leq \sqrt{x}$ for $x \in (0, 1)$,

$$\begin{aligned} |P'_{s,a,h}(s', h) - P_{s,a,h}(s', h)| &\leq |P'_{s,a,h}(s', h) - \gamma\bar{P}_{k,s,a}(s')| + |\gamma\bar{P}_{k,s,a}(s') - P_{s,a,h}(s', h)| \\ &\leq 2\gamma\epsilon_k(s, a, s') \leq \epsilon_k^*(s, a, h, s', h), \end{aligned}$$

$$\begin{aligned} &|P'_{s,a,h}(s', h+1) - P_{s,a,h}(s', h+1)| \\ &\leq |P'_{s,a,h}(s', h+1) - (1-\gamma)\bar{P}_{k,s,a}(s')| + |(1-\gamma)\bar{P}_{k,s,a}(s') - P_{s,a,h}(s', h+1)| \\ &\leq 2(1-\gamma)\epsilon_k(s, a, s') \leq \epsilon_k^*(s, a, h, s', h), \end{aligned}$$

$$\begin{aligned} |P'_{s,a,h}(g) - P_{s,a,h}(g)| &\leq |P'_{s,a,h}(g) - \bar{P}_{k,s,a}(g)| + |\bar{P}_{k,s,a}(g) - P_{s,a,h}(g)| \\ &\leq 2\epsilon_k(s, a, g) \leq 2\epsilon_k^*(s, a, h, g). \end{aligned}$$

This completes the proof. ■

A.2. Approximation of $Q^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$

We show that $Q^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$ can be approximated efficiently by Extended Value Iteration similar to ([Jaksch et al., 2010](#)). Note that finding $\Gamma(\pi, \mathcal{P}, c)$ is equivalent to computing the optimal policy in an augmented MDP $\dot{\mathcal{M}}$ with state space $\dot{\mathcal{S}}$ and extended action space \mathcal{P} , such that for any extended action $P \in \mathcal{P}$, the cost at $((s, h), P)$ is $\sum_a \pi(a|s, h)c(s, a, h)$, and the transition probability to $\dot{s}' \in \dot{\mathcal{S}}_+$ is $\sum_a \pi(a|s, h)P_{s,a,h}(\dot{s}')$. In this work, we have $\mathcal{P} \in \{\mathcal{P}_k\}_{k=1}^K$, and $\mathcal{P}_k = \bigcap_{s,a,h} \mathcal{P}_{k,s,a,h}$, where $\mathcal{P}_{k,s,a,h}$ is a convex set that specifies constraints on $((s, h), a)$. In other words, \mathcal{P}_k is a product of constraints on each $((s, h), a)$ (note that $\Lambda_{\mathcal{M}}$ can also be decomposed into shared constraints on $P_{s,a,H+1}$ and independent constraints on each $s, a, h \leq H$). Thus, any policy in $\dot{\mathcal{M}}$ can be represented by an element $P \in \mathcal{P}$. We can now perform value iteration in $\dot{\mathcal{M}}$ to approximate $Q^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$. The Bellman operator of $\dot{\mathcal{M}}$ is \mathcal{T}_0 defined in [Eq. \(19\)](#) with min operator replaced by max operator. Also note that $\dot{\mathcal{M}}$ is an SSP instance where all policies are proper. Thus, $V^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$ is the unique fixed point of \mathcal{T}_0 ([Bertsekas and Yu, 2013](#)). It is straightforward to show that [Lemma 47](#) still holds with min operator replaced by max operator in [Eq. \(19\)](#) and let $V^0(s, H+1) = \max_a c(s, a, H+1)$. Thus, we can approximate $V^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$ efficiently.

Now suppose after n iterations of modified [Eq. \(19\)](#), we obtain V^n such that $\|V^n - V^{\pi, \Gamma(\pi, \mathcal{P}, c), c}\|_{\infty} \leq \epsilon$. Then we can simply use $Q(s, a, h) = c(s, a, h) + \min_{P \in \mathcal{P}} P_{s,a,h} V^n$ to approximate $Q^{\pi, \Gamma(\pi, \mathcal{P}, c), c}$, since

$$\left| Q(s, a, h) - Q^{\pi, \Gamma(\pi, \mathcal{P}, c), c}(s, a, h) \right| \stackrel{(i)}{=} \left| \min_{P \in \mathcal{P}} P_{s,a,h} V^n - \min_{P \in \mathcal{P}} P_{s,a,h} V^{\pi, \Gamma(\pi, \mathcal{P}, c), c} \right|$$

$$\leq \max_{P \in \mathcal{P}} \left| P_{s,a,h}(V^n - V^{\pi, \Gamma(\pi, \mathcal{P}, c), c}) \right| \leq \left\| V^n - V^{\pi, \Gamma(\pi, \mathcal{P}, c), c} \right\|_{\infty} \leq \epsilon,$$

where (i) is by the definition of $\Gamma(\pi, \mathcal{P}, c)$. In this work, setting $\epsilon = 1/K$ is enough for obtaining the desired regret bounds. [Lemma 47](#) (modified) then implies that $\tilde{\mathcal{O}}(T_{\max})$ iterations of modified [Eq. \(19\)](#) suffices.

Appendix B. Omitted Details for [Section 3](#)

In this section, we provide omitted discussions and proofs for [Section 3](#).

B.1. Limitation of Existing Approximation Schemes

Finite-Horizon Approximation Thanks to [Lemma 31](#), the approximation error under finite-horizon approximation decreases exponentially. Specifically, we only need a horizon of order $\mathcal{O}(T_{\max} \ln K)$ to have approximation error of order $\mathcal{O}(\frac{1}{K})$. This gives optimal regret bound under both adversarial costs ([Chen et al., 2021d](#)) and stochastic costs ([Chen et al., 2021a](#)). However, it also clearly brings an extra $\tilde{\mathcal{O}}(T_{\max})$ dependency in the space complexity since we need to store non-stationary policies changing in different layers. [Chen et al. \(2021a\)](#) proposes an implicit finite-horizon approximation analysis that achieves optimal regret bound without storing non-stationary policies. Unfortunately, their approach does not work for adversarial costs.

Discounted Approximation Approximating an SSP by a discounted MDP clearly produces stationary policies. However, the approximation error scales with $1 - \gamma$ (that is, inversely proportional to the effective horizon $(1 - \gamma)^{-1}$) following similar arguments as in ([Wei et al., 2020](#), Lemma 2), where γ is the discounted factor. This leads to a sub-optimal regret bound when the achieved regret bound in the discounted MDP has polynomial dependency on the horizon even in the lower order term ([Wei et al., 2020](#)). In [Tarbouriech et al. \(2021\)](#), they still achieve minimax optimal regret by deriving a horizon-free regret bound (no polynomial dependency on the horizon even in the lower order term), and approximately set $1 - \gamma = \tilde{\mathcal{O}}(\frac{1}{K})$ to achieve small approximation error. The drawback, however, is that the time complexity of updating the learner’s policy scales linearly w.r.t the effective horizon, which is of order $\tilde{\mathcal{O}}(K)$; see ([Tarbouriech et al., 2021](#), Remark 1).

B.2. Proof of [Lemma 4](#)

Proof We only prove the statement for adversarial environment, and the statement for stochastic environment follows directly from setting $c_1 = \dots = c_K = c$. By [Lemma 2](#), we have $V^{\hat{\pi}^*, P, c_k}(s, 1) \leq V^{\pi^*, P, c_k}(s) + \frac{1}{K}$ for any $k \in [K]$. Now by [Lemma 31](#) and the fact that the expected hitting time of fast policy is upper bounded by D , we have with probability at least $1 - \delta$, the learner reaches the goal within $J_k + c_f$ steps for each episode k . Thus by a union bound, we have with probability at least $1 - \delta$, $\sum_{k=1}^K \sum_{i=1}^{I_k} c_i^k \leq \sum_{k=1}^K \left(\sum_{i=1}^{J_k} c_i^k + c_{J_k+1}^k \right)$. Putting everything together, we get:

$$\begin{aligned} R_K &= \sum_{k=1}^K \left(\sum_{i=1}^{I_k} c_i^k - V^{\pi^*, P, c_k}(s_1^k) \right) \leq \sum_{k=1}^K \left(\sum_{i=1}^{J_k} c_i^k + c_{J_k+1}^k - V^{\hat{\pi}^*, P, c_k}(s_1^k, 1) \right) + \tilde{\mathcal{O}}(1) \\ &= \mathring{R}_K + \tilde{\mathcal{O}}(1). \end{aligned}$$

This completes the proof. ■

Appendix C. Omitted Details for Section 4

In this section, we provide all proofs for Section 4. We first provide omitted details for cost estimation under various feedback types. Then, we establish the main results in Appendix C.2. Finally, we provide proofs of auxiliary lemmas in Appendix C.3.

Extra Notations Define optimistic transitions $\tilde{P}_k = \Gamma(\pi_k, \mathcal{P}_k, \tilde{c}_k)$ and $P_k = \Gamma(\pi_k, \mathcal{P}_k, \hat{c}_k)$, such that $\tilde{Q}_k = Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}$ and $\hat{Q}_k = Q^{\pi_k, P_k, \hat{c}_k}$. Also define $\tilde{q}_k = q_{\pi_k, \tilde{P}_k}$ and $Q_k = Q^{\pi_k, P, \hat{c}_k}$.

C.1. Cost Estimation

We provide more details on the definition of \hat{c}_k for the subsequent analysis. Recall that $\hat{c}_k(s, a) = \max\{0, \bar{c}_k(s, a) - 2\sqrt{\bar{c}_k(s, a)\alpha_k(s, a)} - 7\alpha_k(s, a)\}$. Here, $\bar{c}_k(s, a) = \frac{C_k(s, a)}{\mathfrak{N}_k^+(s, a)}$, where $C_k(s, a)$ is the accumulated costs that are observed at (s, a) in episode $j = 1, \dots, k-1$ before $\sigma(\pi_j)$ switches to the fast policy, $\alpha_k(s, a) = \frac{\iota}{\mathfrak{N}_k^+(s, a)}$ (recall $\iota = \ln(2SALK/\delta)$), $\mathfrak{N}_k^+(s, a) = \max\{1, \mathfrak{N}_k(s, a)\}$, and \mathfrak{N}_k is the number of times the learner observes cost at (s, a) in episode $j = 1, \dots, k-1$ before $\sigma(\pi_j)$ switches to the fast policy. The definition of C_k and \mathfrak{N}_k depends on the type of cost feedback. For stochastic costs, $C_k(s, a) = \sum_{j=1}^{k-1} \sum_{i=1}^{J_k} c_i^j \mathbb{I}\{s_i^j = s, a_i^j = a\}$ and $\mathfrak{N}_k = N_k(s, a)$. For stochastic adversary, $C_k(s, a) = \sum_{j=1}^{k-1} m_j(s, a)c_j(s, a)$, where $m_k(s, a)$ is the indicator of whether $c_k(s, a)$ is observed in episode k before $\sigma(\pi_k)$ switches to the fast policy, and $\mathfrak{N}_k(s, a) = M_k(s, a) \triangleq \sum_{j=1}^{k-1} m_j(s, a)$.

Below we show a lemma quantifying the cost estimation error.

Lemma 14 *Under the event of Lemma 11, we have with probability at least $1 - \delta$,*

$$0 \leq c(s, a) - \hat{c}_k(s, a) \leq 4\sqrt{\bar{c}_k(s, a)\alpha_k(s, a)} + 34\alpha_k(s, a),$$

for all definitions of \hat{c}_k .

Proof Only prove the stochastic cost case and the stochastic adversary case follows similarly. Note that under the event of Lemma 11, $N_{k+1}(s, a) \leq LK$. Applying Lemma 51 with $X_k = c_k(s, a)$ for each $(s, a) \in \mathcal{S} \times \mathcal{A}$ and then by a union bound over all $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have with probability at least $1 - \delta$, for all $k \in [K]$:

$$|\bar{c}_k(s, a) - c(s, a)| \leq 2\sqrt{\alpha_k(s, a)\bar{c}_k(s, a)} + 7\alpha_k(s, a).$$

Hence, $c(s, a) \geq \hat{c}_k(s, a)$ by the definition of \hat{c}_k . Applying $x^2 \leq ax + b \implies x \leq a + \sqrt{b}$ with $x = \sqrt{\bar{c}_k(s, a)}$ to the inequality above (ignoring the absolute value operator), we obtain

$$\sqrt{\bar{c}_k(s, a)} \leq 2\sqrt{\alpha_k(s, a)} + \sqrt{c(s, a) + 7\alpha_k(s, a)} \leq \sqrt{c(s, a)} + 5\sqrt{\alpha_k(s, a)},$$

Therefore, $2\sqrt{\alpha_k(s, a)\bar{c}_k(s, a)} + 7\alpha_k(s, a) \leq 2\sqrt{\alpha_k(s, a)c(s, a)} + 17\alpha_k(s, a)$, and

$$\begin{aligned} c(s, a) - \hat{c}_k(s, a) &= c(s, a) - \bar{c}_k(s, a) + \bar{c}_k(s, a) - \hat{c}_k(s, a) \\ &\leq 2 \cdot (2\sqrt{\alpha_k(s, a)\bar{c}_k(s, a)} + 7\alpha_k(s, a)) \leq 4\sqrt{\alpha_k(s, a)c(s, a)} + 34\alpha_k(s, a). \end{aligned}$$

This completes the proof. \blacksquare

C.2. Main Results for Stochastic Costs and Stochastic Adversary

We first show a general regret bound agnostic to the feedback type ([Theorem 15](#)). Then, we present the proofs of [Theorem 5](#) to [Theorem 7](#) ([Appendix C.2.1](#) to [Appendix C.2.3](#)) using the general regret bound.

Theorem 15 *Assuming that there exists a constant G such that for any s, h :*

$$\sum_{k=1}^{K-1} \left\langle \pi_{k+1}(\cdot | s, h), d\tilde{Q}_k(s, \cdot, h) \right\rangle \leq G.$$

Then, [Algorithm 1](#) in stochastic environments with $\lambda \leq \min\{1/T_{\max}, \sqrt{S^2 A/K}\}$ ensures with probability at least $1 - 22\delta$,

$$\begin{aligned} \mathring{R}_K &= \tilde{O} \left(\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \hat{c}_k(\hat{s}_i^k, a_i^k)) - \sum_{k=1}^K \langle q_k, e_k \rangle + \sum_{k=1}^K \langle q^*, e_k \rangle + \frac{S^2 A}{\lambda} \right) \\ &\quad + \tilde{O} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, e_k \circ Q^{\pi_k, P, e_k} \rangle} + S^4 A^{2.5} T_{\max}^3 + \lambda \sum_{k=1}^K \left\langle q^*, c \circ Q^{\hat{\pi}^*, P, \hat{c}_k} \right\rangle \right) \\ &\quad + \tilde{O} \left(\frac{T_{\max}}{\eta} + T_{\max} G + \lambda \sum_{k=1}^K \left\langle q^*, Q^{\hat{\pi}^*, P, e_k} \right\rangle \right). \end{aligned}$$

Proof For notational convenience, define $\omega = S^4 A^{2.5} T_{\max}^3$. By $\langle q^*, \hat{c}_k \rangle \leq \langle q^*, c \rangle$ ([Lemma 14](#)) and [Lemma 11](#) (under which $n_k = \bar{n}_k$), we have with probability at least $1 - 2\delta$,

$$\begin{aligned} \mathring{R}_K &= \sum_{k=1}^K \left(\sum_{i=1}^{\bar{J}_k} c_i^k + \hat{c}_{J_k+1}^k - V^{\hat{\pi}^*, P, c}(\hat{s}_1^k) \right) \leq \sum_{k=1}^K \left(\sum_{i=1}^{\bar{J}_k} c_i^k + \hat{c}_{J_k+1}^k - V^{\hat{\pi}^*, P, \hat{c}_k}(\hat{s}_1^k) \right) \\ &\leq \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \hat{c}_k(\hat{s}_i^k, a_i^k)) + \sum_{k=1}^K \langle \bar{n}_k - q^*, \hat{c}_k \rangle. \end{aligned}$$

For the second term, by the definition of \tilde{c}_k ,

$$\begin{aligned} \sum_{k=1}^K \langle \bar{n}_k - q^*, \hat{c}_k \rangle &= \sum_{k=1}^K \langle \bar{n}_k - q_k, \hat{c}_k \rangle + \sum_{k=1}^K \langle q_k - \tilde{q}_k, \tilde{c}_k \rangle + \sum_{k=1}^K \langle \tilde{q}_k - q^*, \tilde{c}_k \rangle \\ &\quad - \sum_{k=1}^K \langle q_k, e_k \rangle + \sum_{k=1}^K \langle q^*, e_k \rangle - \lambda \sum_{k=1}^K \left\langle q_k, \hat{c}_k \circ \hat{Q}_k \right\rangle + \lambda \sum_{k=1}^K \left\langle q^*, \hat{c}_k \circ \hat{Q}_k \right\rangle \\ &\leq \underbrace{\sum_{k=1}^K \langle \bar{n}_k - q_k, \hat{c}_k \rangle + \sum_{k=1}^K \langle q_k - \tilde{q}_k, \tilde{c}_k \rangle - \lambda \sum_{k=1}^K \langle q_k, \hat{c}_k \circ \hat{Q}_k \rangle}_{\xi_1} \\ &\quad - \sum_{k=1}^K \langle q_k, e_k \rangle + \sum_{k=1}^K \langle q^*, e_k \rangle + \underbrace{\lambda \sum_{k=1}^K \left\langle q_k, \hat{c}_k \circ (Q_k - \hat{Q}_k) \right\rangle}_{\xi_2} \end{aligned}$$

$$\begin{aligned}
 & + \lambda \sum_{k=1}^K \left\langle q^*, c \circ Q^{\hat{\pi}^*, P, \hat{c}_k} \right\rangle + \underbrace{\sum_{k=1}^K \left\langle \tilde{q}_k - q^*, \tilde{c}_k \right\rangle + \lambda \sum_{k=1}^K \left\langle q^*, c \circ (\hat{Q}_k - Q^{\hat{\pi}^*, P, \hat{c}_k}) \right\rangle}_{\xi_3} \\
 & \hspace{15em} (\hat{c}_k(s, a, h) \leq c(s, a, h))
 \end{aligned}$$

For ξ_1 , with probability at least $1 - 17\delta$:

$$\begin{aligned}
 & \sum_{k=1}^K \langle \bar{n}_k - q_k, \hat{c}_k \rangle + \sum_{k=1}^K \langle q_k - \tilde{q}_k, \tilde{c}_k \rangle - \lambda \sum_{k=1}^K \langle q_k, \hat{c}_k \circ Q_k \rangle \leq \tilde{O} \left(\sqrt{\sum_{k=1}^K \langle q_k, \hat{c}_k \circ Q_k \rangle} + SAT_{\max} \right) \\
 & + \sum_{k=1}^K \left\langle q_k - \tilde{q}_k, (1 + \lambda \hat{Q}_k) \circ \hat{c}_k \right\rangle + \sum_{k=1}^K \langle q_k - \tilde{q}_k, e_k \rangle - \lambda \sum_{k=1}^K \langle q_k, \hat{c}_k \circ Q_k \rangle \\
 & \hspace{10em} (\mathbb{E}_k[\bar{n}_k(s, a, h)] \leq q_k(s, a, h), \text{ Lemma 50, Lemma 26, and } \bar{n}_k(s, a, h) \leq L = \tilde{O}(T_{\max})) \\
 & = \tilde{O} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, \hat{c}_k \circ Q_k \rangle} + \sqrt{S^2 A \sum_{k=1}^K \langle q_k, e_k \circ Q^{\pi_k, P, e_k} \rangle} + \omega \right) - \lambda \sum_{k=1}^K \langle q_k, \hat{c}_k \circ Q_k \rangle \\
 & \hspace{10em} (\text{Lemma 28 and } (1 + \lambda \hat{Q}_k(s, a, h)) \hat{c}_k(s, a, h) = \tilde{O}(\hat{c}_k(s, a, h))) \\
 & = \tilde{O} \left(\frac{S^2 A}{\lambda} + \sqrt{S^2 A \sum_{k=1}^K \langle q_k, e_k \circ Q^{\pi_k, P, e_k} \rangle} + \omega \right). \hspace{10em} (\text{AM-GM inequality})
 \end{aligned}$$

For ξ_2 , by Lemma 30 and Lemma 13, with probability at least $1 - 2\delta$,

$$\begin{aligned}
 & Q_k(s, a, h) - \hat{Q}_k(s, a, h) = \sum_{s', a', h'} q_{k, (s, a, h)}(s', a', h') (P_{s', a', h'} - P_{k, s', a', h'}) V^{\pi_k, P_k, \hat{c}_k} \\
 & = \tilde{O} \left(\sum_{s', a'} q_{k, (s, a, h)}(s', a') \left(\frac{\sqrt{ST_{\max}}}{\sqrt{N_k^+(s', a')}} + \frac{ST_{\max}}{N_k^+(s', a')} \right) \right). \hspace{5em} (6)
 \end{aligned}$$

By $q_k(s, a, h) = \frac{x_k(s, a, h)}{1 - y_k(s, a, h)}$ and $y_k(s, a, h) \leq \gamma = 1 - \frac{1}{2T_{\max}}$, we have

$$\begin{aligned}
 \sum_{s, a, h \leq H} q_k(s, a, h) q_{k, (s, a, h)}(s', a') & \leq 2T_{\max} \sum_{s, a, h \leq H} x_k(s, a, h) q_{k, (s, a, h)}(s', a') \\
 & \leq 2T_{\max} \sum_{s, a, h \leq H} q_k(s', a') = 2T_{\max} SAH q_k(s', a'). \hspace{5em} (7)
 \end{aligned}$$

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

$$\begin{aligned}
 \xi_2 & = \lambda \sum_{k=1}^K \left\langle q_k, \hat{c}_k \circ (Q_k - \hat{Q}_k) \right\rangle \\
 & = \tilde{O} \left(\lambda \sum_{k=1}^K \sum_{s, a, h} q_k(s, a, h) \sum_{s', a'} q_{k, (s, a, h)}(s', a') \left(\frac{\sqrt{ST_{\max}}}{\sqrt{N_k^+(s', a')}} + \frac{ST_{\max}}{N_k^+(s', a')} \right) \right)
 \end{aligned}$$

$$\begin{aligned}
 &= \tilde{\mathcal{O}} \left(\lambda T_{\max}^2 S^{3/2} A \sum_{s', a'} \sum_{k=1}^K \frac{q_k(s', a')}{\sqrt{N_k^+(s', a')}} + \lambda T_{\max}^2 S^2 A \sum_{s', a'} \sum_{k=1}^K \frac{q_k(s', a')}{N_k^+(s', a')} \right) \quad (\text{Eq. (7)}) \\
 &= \tilde{\mathcal{O}} \left(\lambda T_{\max}^2 S^{3/2} A \sqrt{S A T_{\max} K} + \lambda S^3 A^2 T_{\max}^3 \right) = \tilde{\mathcal{O}}(\omega). \\
 &\quad (\text{Lemma 32 and } \sum_{s, a} q_k(s, a) = \tilde{\mathcal{O}}(T_{\max}))
 \end{aligned}$$

For ξ_3 , first note that $\|\tilde{Q}_1\|_{\infty} = \tilde{\mathcal{O}}(T_{\max})$ under all definitions of \tilde{c}_k , and by Lemma 30:

$$\begin{aligned}
 \sum_{k=1}^K \langle \tilde{q}_k - q^*, \tilde{c}_k \rangle &= \sum_{k=1}^K \sum_{s, h} q^*(s, h) \sum_a (\pi_k(a|s, h) - \pi^*(a|s, h)) \tilde{Q}_k(s, a, h) \\
 &\quad + \sum_{k=1}^K \sum_{s, a, h} q^*(s, a, h) \left(\tilde{Q}_k(s, a, h) - \tilde{c}_k(s, a, h) - P_{s, a, h} V^{\pi_k, \tilde{P}_k, \tilde{c}_k} \right) \\
 &= \tilde{\mathcal{O}} \left(\frac{T_{\star}}{\eta} + T_{\star} G + T_{\max}^2 \right). \\
 &\quad (\text{Lemma 24, the definition of } \tilde{P}_k \text{ and } \sum_{s, a, h} q^*(s, a, h) = \tilde{\mathcal{O}}(T_{\star}) \text{ by Lemma 2})
 \end{aligned}$$

Next, note that

$$\begin{aligned}
 \sum_{k=1}^K (\hat{Q}_k(s, a, h) - Q^{\hat{\pi}^*, P, \hat{c}_k}(s, a, h)) &\leq \sum_{k=1}^K (Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\hat{\pi}^*, P, \hat{c}_k}(s, a, h)) \\
 &\quad (P_k, \tilde{P}_k \in \mathcal{P}_k) \\
 &\leq \sum_{k=1}^K \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\hat{\pi}^*, P, \tilde{c}_k}(s, a, h) \right) + \sum_{k=1}^K Q^{\hat{\pi}^*, P, \lambda \hat{Q}_k + e_k}(s, a, h) \quad (\text{definition of } \tilde{c}_k)
 \end{aligned}$$

Also note that $\lambda \sum_{k=1}^K \langle q^*, Q^{\hat{\pi}^*, P, \lambda \hat{Q}_k} \rangle = \tilde{\mathcal{O}}(\lambda^2 T_{\max}^3 K) = \tilde{\mathcal{O}}(S^2 A T_{\max}^3)$ by $\lambda \leq \sqrt{S^2 A / K}$. Thus,

$$\begin{aligned}
 &\lambda \sum_{k=1}^K \langle q^*, c \circ (\hat{Q}_k - Q^{\hat{\pi}^*, P, \hat{c}_k}) \rangle \\
 &= \tilde{\mathcal{O}} \left(\lambda \sum_{k=1}^K \langle q^*, c \circ (\tilde{Q}_k - Q^{\hat{\pi}^*, P, \tilde{c}_k}) \rangle + \lambda \sum_{k=1}^K \langle q^*, Q^{\hat{\pi}^*, P, e_k} \rangle + S^2 A T_{\max}^3 \right).
 \end{aligned}$$

Now by Lemma 30 and the definition of \tilde{P}_k :

$$\begin{aligned}
 &\sum_{k=1}^K (\tilde{Q}_k(s, a, h) - Q^{\hat{\pi}^*, P, \tilde{c}_k}(s, a, h)) \quad (8) \\
 &\leq \sum_{k=1}^K \sum_{s'', h''} P_{s, a, h}(s'', h'') \sum_{s', a', h'} q_{(s'', h'')}^*(s', h') (\pi_k(a'|s', h') - \hat{\pi}^*(a'|s', h')) \tilde{Q}_k(s', a', h')
 \end{aligned}$$

$$= \tilde{\mathcal{O}} \left(\frac{T_{\max}}{\eta} + T_{\max}G + T_{\max}^2 \right). \quad (\text{Lemma 24})$$

Thus, by $\lambda T_{\max} \leq 1$, we have $\lambda \sum_{k=1}^K \langle q^*, c \circ (\tilde{Q}_k - Q^{\hat{\pi}^*, P, \tilde{c}_k}) \rangle = \tilde{\mathcal{O}}(\frac{T_{\max}}{\eta} + T_{\max}G + T_{\max}^2)$. Putting everything together completes the proof. \blacksquare

C.2.1. PROOF OF THEOREM 5

Proof By Lemma 25 with $n_k = n_k$, $\mathfrak{N}_k = N_k$, and $e_k(s, a, h) = 0$, with probability at least $1 - 2\delta$:

$$\begin{aligned} \sum_{k=1}^{K-1} \langle \pi_{k+1}(\cdot | s, h), d\tilde{Q}_k(s, \cdot, h) \rangle &= \tilde{\mathcal{O}} \left(T_{\max}^2 \sum_{k=1}^K \sum_{s', a'} \frac{S n_k(s', a')}{N_k^+(s', a')} + \lambda \eta T_{\max}^4 K \right) \\ &= \tilde{\mathcal{O}} \left(S^2 A T_{\max}^3 + T_{\max}^2 (S^2 A K)^{1/4} \right). \end{aligned}$$

(definition of λ and η , $n_k(s, a) = \bar{n}_k(s, a)$ under the event of Lemma 11, and Lemma 32)

Thus, by Theorem 15, Lemma 16, definition of λ , and replacing G by the bound above, we have with probability at least $1 - 28\delta$,

$$\begin{aligned} \dot{R}_K &= \tilde{\mathcal{O}} \left(\sqrt{SA \sum_{k=1}^K \sum_{i=1}^{J_k} c_i^k} + \frac{S^2 A}{\lambda} + T_{\max}^3 (S^2 A K)^{1/4} + S^4 A^{2.5} T_{\max}^4 + \lambda \sum_{k=1}^K \langle q^*, c \circ Q^{\hat{\pi}^*, P, \hat{c}_k} \rangle \right) \\ &= \tilde{\mathcal{O}} \left(\sqrt{SA \sum_{k=1}^K \sum_{i=1}^{J_k} c_i^k} + B_* S \sqrt{AK} + T_{\max}^3 (S^2 A K)^{1/4} + S^4 A^{2.5} T_{\max}^4 \right). \quad (\text{Lemma 20}) \end{aligned}$$

Now by $\dot{R}_k = \sum_{k=1}^K \sum_{i=1}^{J_k} c_i^k - K \cdot V^{\hat{\pi}^*, P, c}(\hat{s}_1^k)$ and Lemma 48, we have $\sum_{k=1}^K \sum_{i=1}^{J_k} c_i^k = \tilde{\mathcal{O}}(B_* K)$. Plugging this back, we get $\dot{R}_K = \tilde{\mathcal{O}}(B_* S \sqrt{AK} + T_{\max}^3 (S^2 A K)^{1/4} + S^4 A^{2.5} T_{\max}^4)$. Applying Lemma 4 then completes the proof. \blacksquare

C.2.2. PROOF OF THEOREM 6

Proof First note that with probability at least $1 - 3\delta$,

$$\begin{aligned} \sum_{k=1}^K \|de_k\|_1 &\leq \sum_{k=1}^K \sum_{s, a, h \leq H} \left| \sqrt{\frac{\hat{c}_k(s, a, h)}{k}} - \sqrt{\frac{\hat{c}_{k+1}(s, a, h)}{k+1}} \right| + \beta' \sum_{k=1}^K \sum_{s, a, h \leq H} |d\hat{Q}_k(s, a, h)| \\ &= \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^2 + S^{1/2} A^{3/4} T_{\max} K^{1/4} \right), \end{aligned}$$

where in the last inequality we apply

$$\sum_{k=1}^K \sum_{s, a, h \leq H} \left| \sqrt{\frac{\hat{c}_k(s, a, h)}{k}} - \sqrt{\frac{\hat{c}_{k+1}(s, a, h)}{k+1}} \right|$$

$$\begin{aligned}
 &\leq \sum_{k=1}^K \sum_{s,a,h \leq H} \left(\frac{1}{\sqrt{k}} - \frac{1}{\sqrt{k+1}} \right) + \sum_{k=1}^K \sum_{s,a,h \leq H} \frac{\sqrt{|\widehat{c}_k(s,a,h) - \widehat{c}_{k+1}(s,a,h)|}}{\sqrt{k+1}} \\
 &\quad \text{(add and subtract } \sqrt{\widehat{c}_k(s,a,h)/(k+1)}, \text{ and } |\sqrt{a} - \sqrt{b}| \leq \sqrt{|a-b|}) \\
 &= \tilde{\mathcal{O}} \left(SA + \sqrt{\sum_{k=1}^K \sum_{s,a,h \leq H} \frac{1}{k+1}} \sqrt{\sum_{k=1}^K \sum_{s,a,h \leq H} \frac{m_k(s,a)}{M_k^+(s,a)}} \right) = \tilde{\mathcal{O}}(SA), \\
 &\quad \text{(Cauchy-Schwarz inequality, Lemma 25, and Lemma 32)}
 \end{aligned}$$

and by Lemma 11,

$$\begin{aligned}
 &\beta' \sum_{k=1}^K \sum_{s,a,h \leq H} |d\widehat{Q}_k(s,a,h)| \\
 &= \tilde{\mathcal{O}} \left(\beta' \sum_{k=1}^K \sum_{s',a' \leq H} \left(T_{\max}^2 \sum_{s',a'} \frac{S n_k(s',a')}{N_k^+(s',a')} + T_{\max} \sum_{s',a'} \frac{m_k(s',a')}{M_k^+(s',a')} + \eta T_{\max}^3 \right) \right) \\
 &\quad \text{(Lemma 25)} \\
 &= \tilde{\mathcal{O}} \left(\beta' S^3 A^2 T_{\max}^3 + \beta' \eta S A T_{\max}^3 K \right) = \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^2 + S^{1/2} A^{3/4} T_{\max} K^{1/4} \right). \quad \text{(Lemma 32)}
 \end{aligned}$$

Moreover, by Lemma 25 with $n_k = m_k$, $\mathfrak{N}_k = M_k$, and $\lambda \leq \frac{1}{T_{\max}}$, we have with probability at least $1 - \delta$:

$$\begin{aligned}
 &\sum_{k=1}^{K-1} \left\langle \pi_{k+1}(\cdot|s,h), d\tilde{Q}_k(s,\cdot,h) \right\rangle \\
 &= \tilde{\mathcal{O}} \left(T_{\max}^2 \sum_{k=1}^K \sum_{s',a'} \frac{S n_k(s',a')}{N_k^+(s',a')} + \lambda T_{\max}^2 \sum_{k=1}^K \sum_{s',a'} \frac{m_k(s',a')}{M_k^+(s',a')} + \lambda \eta T_{\max}^4 K + T_{\max} \sum_{k=1}^K \|de_k\|_1 \right) \\
 &= \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^3 + T_{\max}^2 (S^2 A^3 K)^{1/4} \right), \quad (9)
 \end{aligned}$$

where the last step is by Lemma 32, the definition of η and λ , and the bound on $\sum_{k=1}^K \|de_k\|_1$. Moreover, by Lemma 17 and definition of e_k , we have with probability at least $1 - 16\delta$:

$$\begin{aligned}
 &\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} \left(c_i^k - \widehat{c}_k(\bar{s}_i^k, a_i^k) \right) - \sum_{k=1}^K \langle q_k, e_k \rangle = \tilde{\mathcal{O}} \left(\beta' \sum_{k=1}^K \langle q_k, Q_k - \widehat{Q}_k \rangle + \frac{1}{\beta'} + \sqrt{S^3 A^3 T_{\max}^3} \right) \\
 &= \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^3 + \sqrt{DT_{\star}K} \right), \\
 &\quad \text{(Eq. (6), Eq. (7) similar to bounding } \xi_2, \text{ and the definition of } \beta')
 \end{aligned}$$

$$\begin{aligned}
 &\sum_{k=1}^K \langle q^*, e_k \rangle = \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{s,a} q^*(s,a) \sqrt{\frac{c(s,a)}{k}} + \beta' \sum_{k=1}^K \langle q^*, \widehat{Q}_k \rangle \right) \quad \text{(Lemma 14)} \\
 &\stackrel{(i)}{=} \tilde{\mathcal{O}} \left(\sqrt{DT_{\star}K} + S^3 A^2 T_{\max}^4 \right), \\
 &\sqrt{S^2 A \sum_{k=1}^K \langle q_k, e_k \circ Q^{\pi_k, P, e_k} \rangle} = \tilde{\mathcal{O}} \left(\sqrt{S^2 A T_{\max}^4} \right),
 \end{aligned}$$

$$\lambda \sum_{k=1}^K \langle q^*, Q^{\hat{\pi}^*, P, e_k} \rangle = \tilde{\mathcal{O}} \left(\lambda T_{\max}^2 \sqrt{K} + \lambda \beta' T_{\max}^3 K \right) = \tilde{\mathcal{O}} \left(\sqrt{S^2 A T_{\max}^3} \right),$$

where (i) is by

$$\sum_{k=1}^K \sum_{s,a} q^*(s,a) \sqrt{\frac{c(s,a)}{k}} = \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \sum_{s,a} q^*(s,a) c(s,a)} \sqrt{\sum_{k=1}^K \sum_{s,a} \frac{q^*(s,a)}{k}} \right) = \tilde{\mathcal{O}} \left(\sqrt{DT_* K} \right),$$

(Cauchy-Schwarz inequality)

definition of β' , and

$$\begin{aligned} \beta' \sum_{k=1}^K \langle q^*, \hat{Q}_k \rangle &= \beta' \sum_{k=1}^K \langle q^*, \hat{Q}_k - Q^{\hat{\pi}^*, P, \hat{c}_k} \rangle + \beta' \sum_{k=1}^K \langle q^*, Q^{\hat{\pi}^*, P, \hat{c}_k} \rangle \\ &= \tilde{\mathcal{O}} \left(\beta' \sum_{k=1}^K \sum_{s,a,h} q^*(s,a,h) \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s,a,h) - Q^{\hat{\pi}^*, P, \tilde{c}_k}(s,a,h) \right) \right) \\ &\quad + \tilde{\mathcal{O}} \left(\beta' \sum_{s,a,h} \sum_{k=1}^K q^*(s,a,h) Q^{\hat{\pi}^*, P, \lambda \hat{Q}_k + e_k}(s,a,h) + \sqrt{DT_* K} \right), \\ &\quad \left(\sum_{s,a,h} q^*(s,a,h) = \mathcal{O}(T_*) \text{, and } \|Q^{\hat{\pi}^*, P, \hat{c}_k}\|_{\infty} = \mathcal{O}(D) \right) \\ &= \tilde{\mathcal{O}} \left(\frac{\beta' T_{\max}^2}{\eta} + \beta' T_{\max}^2 G + \beta' T_{\max}^3 + (\lambda \beta' + \beta'^2) T_{\max}^3 K + \beta' T_{\max}^2 \sqrt{K} + \sqrt{DT_* K} \right), \\ &\quad \text{(Eq. (8))} \\ &= \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^4 + \sqrt{DT_* K} \right). \quad \text{(replace } G \text{ by Eq. (9))} \end{aligned}$$

Thus, by [Theorem 15](#), [Lemma 21](#), and definition of η , λ , we have with probability at least $1 - 22\delta$,

$$\dot{R}_K = \tilde{\mathcal{O}} \left(\sqrt{DT_* K} + DS\sqrt{AK} + T_{\max}^3 (S^2 A^3 K)^{1/4} + S^4 A^{2.5} T_{\max}^4 \right).$$

Applying [Lemma 4](#) completes the proof. ■

C.2.3. PROOF OF [THEOREM 7](#)

Proof By [Lemma 25](#) with $\mathbf{n}_k = m_k$, $\mathfrak{N}_k = M_k$, and $\lambda \leq \frac{1}{T_{\max}}$, we have with probability at least $1 - 2\delta$:

$$\begin{aligned} &\sum_{k=1}^{K-1} \left\langle \pi_{k+1}(\cdot | s, h), d\tilde{Q}_k(s, \cdot, h) \right\rangle \quad (10) \\ &= \tilde{\mathcal{O}} \left(T_{\max}^2 \sum_{k=1}^K \sum_{s', a'} \frac{S n_k(s', a')}{N_k^+(s', a')} + \lambda T_{\max}^2 \sum_{k=1}^K \sum_{s', a'} \frac{m_k(s', a')}{M_k^+(s', a')} + \lambda \eta T_{\max}^4 K + T_{\max} \sum_{k=1}^K \|de_k\|_1 \right) \\ &= \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^3 + T_{\max}^2 S A^{5/4} K^{1/4} \right), \quad \text{(definition of } \eta \text{ and } \text{Lemma 32)} \end{aligned}$$

where in the last step we apply

$$\begin{aligned}
 \sum_{k=1}^K \|de_k\|_1 &= \beta \sum_{k=1}^K \sum_{s,a,h \leq H} \left| d\widehat{Q}_k(s, a, h) \right| \\
 &= \tilde{\mathcal{O}} \left(\beta \sum_{k=1}^K \sum_{s,a,h \leq H} \left(T_{\max}^2 \sum_{s',a'} \frac{Sn_k(s', a')}{N_k^+(s', a')} + SAT_{\max} \frac{m_k(s, a)}{M_k^+(s, a)} + \eta T_{\max}^3 \right) \right) \quad (\text{Lemma 25}) \\
 &= \tilde{\mathcal{O}} \left(\beta S^3 A^2 T_{\max}^3 + \beta \eta SAT_{\max}^3 K \right) = \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^2 + SA^{5/4} T_{\max} K^{1/4} \right). \quad (\text{Lemma 32})
 \end{aligned}$$

By Lemma 18 and the definition of e_k , we have with probability at least $1 - 11\delta$:

$$\begin{aligned}
 \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} \left(c_i^k - \widehat{c}_k(s_i^k, a_i^k) \right) - \sum_{k=1}^K \langle q_k, e_k \rangle &= \tilde{\mathcal{O}} \left(\beta \sum_{k=1}^K \langle q_k, Q_k - \widehat{Q}_k \rangle + \frac{SA}{\beta} + \sqrt{S^3 A^3 T_{\max}^3} \right) \\
 &= \tilde{\mathcal{O}} \left(\sqrt{SADT_{\star} K} + S^{2.5} A^2 T_{\max}^3 \right), \\
 &\quad (\text{Eq. (6), Eq. (7) similar to bounding } \xi_2, \text{ and the definition of } \beta)
 \end{aligned}$$

$$\begin{aligned}
 \sum_{k=1}^K \langle q^*, e_k \rangle &\leq \beta \sum_{k=1}^K \langle q^*, \widehat{Q}_k - Q^{\hat{\pi}^*, P, \widehat{c}_k} \rangle + \beta \sum_{k=1}^K \langle q^*, Q^{\hat{\pi}^*, P, \widehat{c}_k} \rangle \\
 &= \tilde{\mathcal{O}} \left(\beta \sum_{k=1}^K \sum_{s,a,h} q^*(s, a, h) \left(Q^{\pi_k, \widehat{P}_k, \widehat{c}_k}(s, a, h) - Q^{\hat{\pi}^*, P, \widehat{c}_k}(s, a, h) \right) \right) \\
 &\quad + \tilde{\mathcal{O}} \left(\beta \sum_{s,a,h} \sum_{k=1}^K q^*(s, a, h) Q^{\hat{\pi}^*, P, \lambda \widehat{Q}_k + e_k}(s, a, h) + \sqrt{SADT_{\star} K} \right), \\
 &\quad (\sum_{s,a,h} q^*(s, a, h) = \mathcal{O}(T_{\star}), \text{ and } \|Q^{\hat{\pi}^*, P, \widehat{c}_k}\|_{\infty} = \mathcal{O}(D)) \\
 &= \tilde{\mathcal{O}} \left(\frac{\beta T_{\max}^2}{\eta} + \beta T_{\max}^2 G + \beta T_{\max}^3 + (\lambda\beta + \beta^2) T_{\max}^3 K + \sqrt{SADT_{\star} K} \right), \quad (\text{Eq. (8)}) \\
 &= \tilde{\mathcal{O}} \left(S^3 A^2 T_{\max}^4 + \sqrt{SADT_{\star} K} \right), \quad (\text{replace } G \text{ by Eq. (10)})
 \end{aligned}$$

$$\begin{aligned}
 \sqrt{S^2 A \sum_{k=1}^K \langle q_k, e_k \circ Q^{\pi_k, P, e_k} \rangle} &= \tilde{\mathcal{O}} \left(\sqrt{\beta^2 S^2 A T_{\max}^4 K} \right) = \tilde{\mathcal{O}} \left(\sqrt{S^3 A^2 T_{\max}^4} \right), \\
 \lambda \sum_{k=1}^K \langle q^*, Q^{\hat{\pi}^*, P, e_k} \rangle &= \tilde{\mathcal{O}} \left(\lambda \beta T_{\max}^3 K \right) = \tilde{\mathcal{O}} \left(S^{3/2} A T_{\max}^3 \right).
 \end{aligned}$$

Thus, by Theorem 15, definition of η , λ , and β , and Lemma 21, with probability at least $1 - 22\delta$,

$$\dot{R}_K = \tilde{\mathcal{O}} \left(\sqrt{SADT_{\star} K} + DS\sqrt{AK} + T_{\max}^3 SA^{5/4} K^{1/4} + S^4 A^{2.5} T_{\max}^4 \right).$$

Applying Lemma 4 completes the proof. \blacksquare

C.3. Extra Lemmas for Section 4

We give an outline of this section: [Lemma 16](#) to [Lemma 18](#) bound the term $\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \widehat{c}_k(s_i^k, a_i^k))$ under various feedback types. [Lemma 20](#) and [Lemma 21](#) bound the term $\sum_{k=1}^K \langle q^*, c \circ Q^{\tilde{\pi}^*, P, c} \rangle$ under stochastic costs and stochastic adversaries respectively. [Lemma 23](#) establishes stability of PO updates. [Lemma 24](#) provide a refined analysis of PO. [Lemma 25](#) bounds the drift of various quantities (such as $d\widehat{c}_k$ and $d\widetilde{Q}_k$) across episodes. [Lemma 26](#) provide bounds on variance of learner's costs. [Lemma 28](#) gives a bound on the estimation error of value functions due to transition estimation.

Lemma 16 *Under stochastic costs, we have with probability at least $1 - 6\delta$:*

$$\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \widehat{c}_k(s_i^k, a_i^k)) = \tilde{O} \left(\sqrt{SA \sum_{k=1}^K \sum_{i=1}^{J_k} c_i^k + SAT_{\max}} \right).$$

Proof First note that:

$$\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \widehat{c}_k(s_i^k, a_i^k)) = \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - c(s_i^k, a_i^k)) + \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c(s_i^k, a_i^k) - \widehat{c}_k(s_i^k, a_i^k)).$$

For the first term, by [Lemma 50](#) and [Lemma 52](#), we have with probability at least $1 - 2\delta$,

$$\begin{aligned} \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - c(s_i^k, a_i^k)) &= \tilde{O} \left(\sqrt{\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} \mathbb{E}[(c_i^k)^2 | s_i^k, a_i^k]} \right) = \tilde{O} \left(\sqrt{\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} c(s_i^k, a_i^k)} \right) \\ &= \tilde{O} \left(\sqrt{\sum_{k=1}^K \sum_{i=1}^{J_k} c_i^k} \right). \end{aligned}$$

For the second term, with probability at least $1 - 4\delta$,

$$\begin{aligned} \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c(s_i^k, a_i^k) - \widehat{c}_k(s_i^k, a_i^k)) &= \tilde{O} \left(\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} \left(\sqrt{\frac{c(s_i^k, a_i^k)}{N_k^+(s_i^k, a_i^k)}} + \frac{1}{N_k^+(s_i^k, a_i^k)} \right) \right) \\ &\quad \text{(Lemma 14 and } \widehat{c}_k(s, a) \leq c(s, a)) \\ &= \tilde{O} \left(\sum_{s,a} \sum_{k=1}^K \left(\bar{n}_k(s, a) \sqrt{\frac{c(s, a)}{N_k^+(s, a)}} + \frac{\bar{n}_k(s, a)}{N_k^+(s, a)} \right) \right) \\ &= \tilde{O} \left(\sqrt{SA \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} c(s_i^k, a_i^k) + SAT_{\max}} \right) \quad \text{(Lemma 32 and } J_k = \bar{J}_k) \\ &= \tilde{O} \left(\sqrt{SA \sum_{k=1}^K \sum_{i=1}^{J_k} c_i^k + SAT_{\max}} \right). \quad \text{(Lemma 52)} \end{aligned}$$

This completes the proof. ■

Lemma 17 *Under stochastic adversary with full information, with probability at least $1 - 8\delta$,*

$$\begin{aligned} \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \widehat{c}_k(s_i^k, a_i^k)) &= 8\iota \cdot \sum_{k=1}^K \sum_{s,a} q_k(s, a) \sqrt{\widehat{c}_k(s, a)/k} \\ &+ \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \sum_{s,a,h \leq H} q_k(s, a, h) Q_k(s, a, h)} + \sqrt{S^3 A^3 T_{\max}^3} \right). \end{aligned}$$

Proof First note that by $c_i^k = c_k(s_i^k, a_i^k)$:

$$\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \widehat{c}_k(s_i^k, a_i^k)) = \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_k(s_i^k, a_i^k) - c(s_i^k, a_i^k)) + \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c(s_i^k, a_i^k) - \widehat{c}_k(s_i^k, a_i^k)).$$

For the first term, with probability at least $1 - \delta$,

$$\begin{aligned} \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_k(s_i^k, a_i^k) - c(s_i^k, a_i^k)) &= \sum_{k=1}^K \sum_{s,a} \bar{n}_k(s, a) (c_k(s, a) - c(s, a)) \\ &\stackrel{(i)}{=} \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \mathbb{E}_k \left[\left(\sum_{s,a} \bar{n}_k(s, a) c_k(s, a) \right)^2 \right] + T_{\max}} \right) \\ &= \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \mathbb{E}_{c_k} \left[\sum_{s,a,h \leq H} q_k(s, a, h) Q_k^{\pi_k, P, c_k}(s, a, h) \right] + T_{\max}} \right) \\ &\quad (\mathbb{E}_k[\cdot] = \mathbb{E}_{c_k, \bar{n}_k}[\cdot], \text{ Lemma 26 and } c_k(s, a) \leq 1) \\ &= \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \sum_{s,a,h \leq H} q_k(s, a, h) Q_k(s, a, h) + \sum_{k=1}^K \langle q_k, Q_k^{\pi_k, P, c} - Q_k \rangle + T_{\max}} \right), \\ &\quad (\widehat{c}_k(s, a, H+1) = c(s, a, H+1)) \end{aligned}$$

where in (i) we apply Lemma 50, $\mathbb{E}_k[\cdot] = \mathbb{E}_{c_k, \bar{n}_k}[\cdot]$, and

$$\mathbb{E}_{c_k} \left[\left(\sum_{s,a} \bar{n}_k(s, a) (c_k(s, a) - c(s, a)) \right)^2 \middle| \bar{n}_k \right] \leq \mathbb{E}_{c_k} \left[\left(\sum_{s,a} \bar{n}_k(s, a) c_k(s, a) \right)^2 \middle| \bar{n}_k \right].$$

Now note that for $h \leq H$, by Lemma 30, Lemma 14, and $\widehat{c}_k(s, a, H+1) = c(s, a, H+1)$, we have with probability at least $1 - 2\delta$:

$$\begin{aligned} Q_k^{\pi_k, P, c}(s, a, h) - Q_k(s, a, h) &= \sum_{s', a', h' \leq H} q_{k,(s,a,h)}(s', a', h') (c(s', a', h') - \widehat{c}_k(s', a', h')) \\ &= \tilde{\mathcal{O}} \left(\sum_{s', a'} q_{k,(s,a,h)}(s', a') \left(\sqrt{\frac{\widehat{c}_k(s', a')}{M_k^+(s', a')}} + \frac{1}{M_k^+(s', a')} \right) \right). \end{aligned}$$

Note that $q_k(s, a, h)q_{k,(s,a,h)}(s', a') = \tilde{\mathcal{O}}(T_{\max}x_k(s, a, h)q_{k,(s,a,h)}(s', a')) = \tilde{\mathcal{O}}(T_{\max}q_k(s', a'))$. Therefore, we have with probability at least $1 - \delta$:

$$\begin{aligned}
 \sum_{k=1}^K \langle q_k, Q^{\pi_k, P, c} - Q_k \rangle &= \tilde{\mathcal{O}} \left(T_{\max} \sum_{k=1}^K \sum_{s,a,h \leq H} \sum_{s',a'} q_k(s', a') \left(\sqrt{\frac{\hat{c}_k(s', a')}{M_k^+(s', a')}} + \frac{1}{M_k^+(s', a')} \right) \right) \\
 &= \tilde{\mathcal{O}} \left(SAT_{\max} \sum_{k=1}^K \sum_{s',a'} q_k(s', a') \left(\sqrt{\frac{\hat{c}_k(s', a')}{M_k^+(s', a')}} + \frac{1}{M_k^+(s', a')} \right) \right) \\
 &\hspace{15em} (q_k(s', a', h') = \mathcal{O}(T_{\max}x_k(s', a', h'))) \\
 &= \tilde{\mathcal{O}} \left(SAT_{\max} \left(\sqrt{\sum_{k=1}^K \sum_{s',a'} \frac{q_k(s', a')}{M_k^+(s', a')}} \sqrt{\sum_{k=1}^K \sum_{s',a'} q_k(s', a') \hat{c}_k(s', a')} + \sum_{k=1}^K \sum_{s',a'} \frac{q_k(s', a')}{M_k^+(s', a')} \right) \right) \\
 &\hspace{15em} (\text{Cauchy-Schwarz inequality}) \\
 &= \tilde{\mathcal{O}} \left(\sqrt{S^3 A^3 T_{\max}^3 \sum_{k=1}^K \sum_{s',a'} q_k(s', a') \hat{c}_k(s', a')} + S^2 A^2 T_{\max}^2 \right) \\
 &\hspace{15em} (q_k(s', a') \leq T_{\max}x_k(s', a') \text{ and Lemma 32}) \\
 &= \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{s',a'} q_k(s', a') \hat{c}_k(s', a') + S^3 A^3 T_{\max}^3 \right). \hspace{5em} (\text{AM-GM inequality})
 \end{aligned}$$

Substituting these back, we have

$$\begin{aligned}
 &\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_k(s_i^k, a_i^k) - c(s_i^k, a_i^k)) \\
 &= \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \sum_{s,a,h \leq H} q_k(s, a, h) Q_k(s, a, h)} + \sqrt{S^3 A^3 T_{\max}^3} \right). \hspace{5em} (11)
 \end{aligned}$$

For the second term, with probability at least $1 - 4\delta$,

$$\begin{aligned}
 \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c(s_i^k, a_i^k) - \hat{c}_k(s_i^k, a_i^k)) &\leq \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} \left(4 \sqrt{\frac{\hat{c}_k(s_i^k, a_i^k) \iota}{M_k^+(s_i^k, a_i^k)}} + \frac{34\iota}{M_k^+(s_i^k, a_i^k)} \right) \hspace{2em} (\text{Lemma 14}) \\
 &\leq \sum_{k=1}^K \sum_{s,a} 8 \cdot q_k(s, a) \iota \sqrt{\frac{\hat{c}_k(s, a)}{k}} + \sum_{k=1}^K \sum_{s,a} \frac{68q_k(s, a) \iota}{k} + \tilde{\mathcal{O}}(T_{\max}) \hspace{2em} (\text{Lemma 52}) \\
 &\leq 8\iota \cdot \sum_{k=1}^K \sum_{s,a} q_k(s, a) \sqrt{\hat{c}_k(s, a)/k} + \tilde{\mathcal{O}}(T_{\max}).
 \end{aligned}$$

Putting everything together completes the proof. ■

Lemma 18 *Under stochastic adversary with bandit feedback, with probability at least $1 - 8\delta$,*

$$\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \widehat{c}_k(s_i^k, a_i^k)) = \tilde{\mathcal{O}} \left(\sqrt{SA \sum_{k=1}^K \sum_{s,a} \sum_{h=1}^H q_k(s, a, h) Q_k(s, a, h)} + \sqrt{S^3 A^3 T_{\max}^3} \right).$$

Proof First note that by $c_i^k = c_k(s_i^k, a_i^k)$:

$$\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_i^k - \widehat{c}_k(s_i^k, a_i^k)) = \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c_k(s_i^k, a_i^k) - c(s_i^k, a_i^k)) + \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c(s_i^k, a_i^k) - \widehat{c}_k(s_i^k, a_i^k)).$$

For the first term, [Eq. \(11\)](#) holds by the same arguments as in [Lemma 17](#) with probability at least $1 - 4\delta$. For the second term, we have with probability at least $1 - 4\delta$,

$$\begin{aligned} \sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} (c(s_i^k, a_i^k) - \widehat{c}_k(s_i^k, a_i^k)) &= \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{i=1}^{\bar{J}_k} \left(\sqrt{\frac{\widehat{c}_k(s_i^k, a_i^k)}{M_k^+(s_i^k, a_i^k)}} + \frac{1}{M_k^+(s_i^k, a_i^k)} \right) \right) \\ &\hspace{15em} \text{(Lemma 14)} \\ &= \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{s,a} \sum_{h=1}^H q_k(s, a, h) \sqrt{\frac{\widehat{c}_k(s, a)}{M_k^+(s, a)}} + \sum_{k=1}^K \sum_{s,a} \frac{q_k(s, a)}{M_k^+(s, a)} + T_{\max} \right) \\ &\hspace{15em} \text{(Lemma 52)} \\ &= \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \sum_{s,a} \sum_{h=1}^H \frac{q_k^2(s, a, h)}{x_k(s, a, h)} \widehat{c}_k(s, a)} \sqrt{\sum_{k=1}^K \sum_{s,a} \sum_{h=1}^H \frac{x_k(s, a, h)}{M_k^+(s, a)} + SAT_{\max}} \right) \\ &\hspace{5em} \text{(Cauchy-Schwarz inequality, Lemma 32, and } q_k(s, a) = \tilde{\mathcal{O}}(T_{\max} x_k(s, a))\text{)} \\ &= \tilde{\mathcal{O}} \left(\sqrt{SA \sum_{k=1}^K \sum_{s,a} \sum_{h=1}^H q_k(s, a, h) Q_k(s, a, h)} + SAT_{\max} \right) \\ &\hspace{10em} \text{(Lemma 32 and } \frac{q_k(s, a, h)}{x_k(s, a, h)} \widehat{c}_k(s, a) \leq Q_k(s, a, h)\text{)} \end{aligned}$$

■

Lemma 19 *For $h \in [H + 1]$, we have $\sum_{s,a} q^*(s, a, h) \leq (\frac{1}{2})^{h-1} T_{\max}$.*

Proof Denote by $p(s)$ the probability that the learner starts at state s in layer h and eventually reaches layer $h + 1$ following $\hat{\pi}^*$. Clearly, $p(g) = 0$, and

$$p(s) \leq 1 - \gamma + \gamma P_{s, \pi^*(s)} p \stackrel{(i)}{\leq} \mathbb{E} \left[\sum_{t=1}^I (1 - \gamma) \gamma^{t-1} \mathbb{1}_{\pi^*, P, s_1 = s} \right] \leq \frac{1}{2},$$

where (i) is by repeatedly applying the first inequality. By a recursive argument, we have the probability of reaching layer h is upper bounded by $(\frac{1}{2})^{h-1}$. Then by $\sum_{s,a} q_{(s', h)}^*(s, a, h) \leq T_{\max}$ for any s' , we have $\sum_{s,a} q^*(s, a, h) \leq (\frac{1}{2})^{h-1} T_{\max}$. ■

Lemma 20 Under stochastic costs, $\langle q^*, c \circ Q^{\hat{\pi}^*, P, c} \rangle \leq 2B_*^2 + \frac{(H+1)T_{\max}}{K}$.

Proof By Lemma 2 and Lemma 19, we have:

$$\begin{aligned}
 \langle q^*, c \circ Q^{\hat{\pi}^*, P, c} \rangle &= \sum_{h=1}^H \sum_{s,a} q^*(s, a, h) c(s, a) Q^{\hat{\pi}^*, P, c}(s, a, h) + \sum_s q^*(s, H+1) c_f \\
 &\leq \sum_{h=1}^H \sum_{s,a} q^*(s, a, h) c(s, a) \left(Q^{\pi^*, P, c}(s, a) + \frac{c_f}{2^{H-h+1}} \right) + \frac{c_f T_{\max}}{2^H} \\
 &\leq 2B_*^2 + \sum_{h=1}^H \frac{T_{\max}}{2^{h-1}} \frac{c_f}{2^{H-h+1}} + \frac{c_f T_{\max}}{2^H} \\
 &\quad \left(\sum_{h=1}^H q^*(s, a, h) c(s, a) \leq B_* \text{ and } Q^{\pi^*, P, c}(s, a) \leq 1 + B_* \right) \\
 &\leq 2B_*^2 + (H+1) \frac{c_f T_{\max}}{2^H} \leq 2B_*^2 + \frac{(H+1)T_{\max}}{K}.
 \end{aligned}$$

■

Lemma 21 For stochastic adversary, we have $\sum_{k=1}^K \langle q^*, c \circ Q^{\hat{\pi}^*, P, c} \rangle = \tilde{\mathcal{O}}(D^2 K)$.

Proof $\sum_{k=1}^K \langle q^*, c \circ Q^{\hat{\pi}^*, P, c} \rangle = \tilde{\mathcal{O}}(DK \langle q^*, c \rangle) = \tilde{\mathcal{O}}(D^2 K)$.

■

Lemma 22 $\eta \left\| \tilde{Q}_k \right\|_{\infty} \leq 1$ under all definitions of \tilde{c}_k .

Proof It suffices to bound $\left\| \tilde{Q}_k \right\|_{\infty}$. By Lemma 2, $\hat{Q}_k(s, a, h) \leq \frac{H}{1-\gamma} + c_f = \chi$. Therefore, $e_k(s, a, h) \leq 8\iota + \chi/T_{\max}$ under all feedback types. This gives $\tilde{c}_k(s, a, h) \leq (1 + \lambda \hat{Q}_k(s, a, h)) + e_k(s, a, h) \leq 3(8\iota + \chi/T_{\max})$ for $h \leq H$ and $\tilde{c}_k(s, a, H+1) \leq (1 + \lambda \hat{Q}_k(s, a, H+1))c_f \leq 3c_f \chi/T_{\max}$. Lemma 2 then gives $\tilde{Q}_k(s, a, h) \leq \frac{H}{1-\gamma} \cdot 3(8\iota + \chi/T_{\max}) + 3c_f \chi/T_{\max} \leq 3T_{\max}(8\iota + \chi/T_{\max})^2$, and the statement is proved by the definition of η .

■

Lemma 23 Under all definitions of \tilde{c}_k , we have $|d\pi_k(a|s, h)| = \tilde{\mathcal{O}}(\eta T_{\max} \pi_k(a|s, h))$ and $\left\| dQ^{\pi_k, P', c'} \right\|_{\infty} = \tilde{\mathcal{O}}(\eta T_{\max}^3)$ for $P' \in \Lambda_{\mathcal{M}}$ and $c' \in \mathcal{C}_{\mathcal{M}}$.

Proof Note that:

$$\begin{aligned}
 \pi_{k+1}(a|s, h) - \pi_k(a|s, h) &= \frac{\pi_k(a|s, h) \exp(-\eta \tilde{Q}_k(s, a, h))}{\sum_{a'} \pi_k(a'|s, h) \exp(-\eta \tilde{Q}_k(s, a', h))} - \pi_k(a|s, h) \\
 &\leq \frac{\pi_k(a|s, h)}{\sum_{a'} \pi_k(a'|s, h)} \exp(\max_{a'} |\eta \tilde{Q}_k(s, a', h)|) - \pi_k(a|s, h) = \tilde{\mathcal{O}}(\eta T_{\max} \pi_k(a|s, h)).
 \end{aligned}$$

(Lemma 22 and $|e^x - 1| \leq 2|x|$ for $x \in [-1, 1]$)

The other direction can be proved similarly. Then by [Lemma 30](#),

$$\begin{aligned}
 & \left| Q^{\pi_{k+1}, P', c'}(s, a, h) - Q^{\pi_k, P', c'}(s, a, h) \right| \\
 &= \left| \sum_{s'', h''} P_{s, a, h}(s'', h'') \sum_{s', a', h'} q_{\pi_k, P', (s'', h'')}(s', h') (d\pi_k(a'|s', h')) Q^{\pi_{k+1}, P', c'}(s', a', h') \right| \\
 &= \tilde{\mathcal{O}}(\eta T_{\max}^3).
 \end{aligned}$$

This completes the proof. ■

Lemma 24 *Suppose $\pi_k(a|s, h) \propto \exp(\sum_{j < k} \tilde{Q}_j(s, a, h))$. Then,*

$$\begin{aligned}
 & \sum_{k=1}^K \sum_{a \in \mathcal{A}} (\pi_k(a|s, h) - \pi^*(a|s, h)) \tilde{Q}_k(s, a, h) \\
 & \leq \frac{\ln A}{\eta} + \left\langle \pi_1(\cdot|s, h), \tilde{Q}_1(s, \cdot, h) \right\rangle + \sum_{k=1}^{K-1} \left\langle \pi_{k+1}(\cdot|s, h), \tilde{Q}_{k+1}(s, \cdot, h) - \tilde{Q}_k(s, \cdot, h) \right\rangle.
 \end{aligned}$$

Proof First note that:

$$\pi_{k+1}(\cdot|s, h) = \underset{\pi(\cdot|s, h) \in \Delta(A)}{\operatorname{argmin}} \eta \left\langle \pi(\cdot|s, h), \tilde{Q}_k(s, \cdot, h) \right\rangle + \operatorname{KL}(\pi(\cdot|s, h), \pi_k(\cdot|s, h)), \quad (12)$$

where $\operatorname{KL}(p, q) = \sum_a (p(a) \ln \frac{p(a)}{q(a)} - p(a) + q(a))$, and

$$\pi_{k+1}(\cdot|s, h) \propto \pi'_{k+1}(a|s, h) \triangleq \pi_k(a|s, h) \exp(-\eta \tilde{Q}_k(s, a, h)),$$

where π'_{k+1} is the solution of the unconstrained variant of [Eq. \(12\)](#) (that is, replacing $\operatorname{argmin}_{\pi(\cdot|s, h) \in \Delta(A)}$ by $\operatorname{argmin}_{\pi(\cdot|s, h) \in \mathbb{R}^A}$). It is easy to verify that:

$$\begin{aligned}
 & \operatorname{KL}(\pi_k(\cdot|s, h), \pi_{k+1}(\cdot|s, h)) + \operatorname{KL}(\pi_{k+1}(\cdot|s, h), \pi_k(\cdot|s, h)) \\
 &= \left\langle \pi_k(\cdot|s, h), \ln \frac{\pi_k(\cdot|s, h)}{\pi_{k+1}(\cdot|s, h)} \right\rangle + \left\langle \pi_{k+1}(\cdot|s, h), \ln \frac{\pi_{k+1}(\cdot|s, h)}{\pi_k(\cdot|s, h)} \right\rangle \\
 &= \left\langle \pi_k(\cdot|s, h) - \pi_{k+1}(\cdot|s, h), \ln \frac{\pi_k(\cdot|s, h)}{\pi'_{k+1}(\cdot|s, h)} \right\rangle \quad (\pi_{k+1}(\cdot|s, h) \propto \pi'_{k+1}(\cdot|s, h)) \\
 &= \left\langle \pi_k(\cdot|s, h) - \pi_{k+1}(\cdot|s, h), \eta \tilde{Q}_k(s, \cdot, h) \right\rangle \geq 0. \quad (13)
 \end{aligned}$$

By the standard OMD analysis ([Hazan et al., 2016](#)) (note that KL is the Bregman divergence w.r.t the negative entropy regularizer),

$$\sum_{k=1}^K \left\langle \pi_k(\cdot|s, h) - \pi^*(\cdot|s, h), \tilde{Q}_k(s, \cdot, h) \right\rangle$$

$$\begin{aligned}
 &= \frac{1}{\eta} \sum_{k=1}^K (\text{KL}(\pi^*(\cdot|s, h), \pi_k(\cdot|s, h)) - \text{KL}(\pi^*(\cdot|s, h), \pi'_{k+1}(\cdot|s, h)) + \text{KL}(\pi_k(\cdot|s, h), \pi'_{k+1}(\cdot|s, h))) \\
 &= \frac{1}{\eta} \sum_{k=1}^K (\text{KL}(\pi^*(\cdot|s, h), \pi_k(\cdot|s, h)) - \text{KL}(\pi^*(\cdot|s, h), \pi_{k+1}(\cdot|s, h)) + \text{KL}(\pi_k(\cdot|s, h), \pi_{k+1}(\cdot|s, h))) \\
 &\leq \frac{\text{KL}(\pi^*(\cdot|s, h), \pi_1(\cdot|s, h))}{\eta} + \sum_{k=1}^K \left\langle \pi_k(\cdot|s, h) - \pi_{k+1}(\cdot|s, h), \tilde{Q}_k(s, \cdot, h) \right\rangle \quad (\text{Eq. (13)}) \\
 &\leq \frac{\ln A}{\eta} + \sum_{k=1}^{K-1} \left\langle \pi_{k+1}(\cdot|s, h), \tilde{Q}_{k+1}(s, \cdot, h) - \tilde{Q}_k(s, \cdot, h) \right\rangle \\
 &\quad + \left\langle \pi_1(\cdot|s, h), \tilde{Q}_1(s, \cdot, h) \right\rangle - \left\langle \pi_{K+1}(\cdot|s, h), \tilde{Q}_K(s, \cdot, h) \right\rangle.
 \end{aligned}$$

This completes the proof. \blacksquare

Lemma 25 Define $\mathbf{n}_k(s, a) = \mathfrak{N}_{k+1}(s, a) - \mathfrak{N}_k(s, a)$. We have:

$$\begin{aligned}
 |d\hat{c}_k(s, a)| &= \mathcal{O} \left(\frac{\mathbf{n}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)} \right), \\
 |d\hat{Q}_k(s, a, h)| &= \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{S n_k(s', a')\iota}{N_k^+(s', a')} + T_{\max} \sum_{s', a'} \frac{\mathbf{n}_k(s', a')\iota}{\mathfrak{N}_k^+(s', a')} + \eta T_{\max}^3 \right), \\
 |d\tilde{c}_k(s, a)| &= \mathcal{O} \left(\frac{\mathbf{n}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)} + \lambda T_{\max}^2 \sum_{s', a'} \frac{S n_k(s', a')\iota}{N_k^+(s', a')} + \lambda T_{\max} \sum_{s', a'} \frac{\mathbf{n}_k(s', a')}{\mathfrak{N}_k^+(s', a')} + \lambda \eta T_{\max}^3 + |de_k(s, a, h)| \right), \\
 d\tilde{Q}_k(s, a, h) &= \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{S n_k(s', a')\iota}{N_k^+(s', a')} + \lambda T_{\max}^2 \sum_{s', a'} \frac{\mathbf{n}_k(s', a')}{\mathfrak{N}_k^+(s', a')} + \lambda \eta T_{\max}^4 + T_{\max} \|de_k\|_1 \right).
 \end{aligned}$$

Proof First statement: Note that for all definitions of \hat{c}_k used in this paper, we have $\|\hat{c}_k\|_{\infty} \leq 1$. Then by the definition of \hat{c}_k and $|\max\{0, a\} - \max\{0, b\}| \leq |a - b|$:

$$\begin{aligned}
 &|\hat{c}_{k+1}(s, a) - \hat{c}_k(s, a)| \\
 &= \mathcal{O} \left(|\bar{c}_{k+1}(s, a) - \bar{c}_k(s, a)| + \left| \sqrt{\frac{\bar{c}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)}} - \sqrt{\frac{\bar{c}_{k+1}(s, a)\iota}{\mathfrak{N}_{k+1}^+(s, a)}} \right| + \frac{\iota}{\mathfrak{N}_k^+(s, a)} - \frac{\iota}{\mathfrak{N}_{k+1}^+(s, a)} \right).
 \end{aligned}$$

Note that:

$$\begin{aligned}
 |\bar{c}_{k+1}(s, a) - \bar{c}_k(s, a)| &= \left| \frac{C_{k+1}(s, a)}{\mathfrak{N}_{k+1}^+(s, a)} - \frac{C_k(s, a)}{\mathfrak{N}_k^+(s, a)} \right| \\
 &\leq \left| \frac{C_{k+1}(s, a) - C_k(s, a)}{\mathfrak{N}_{k+1}^+(s, a)} \right| + \mathfrak{N}_k(s, a) \left| \frac{1}{\mathfrak{N}_k^+(s, a)} - \frac{1}{\mathfrak{N}_{k+1}^+(s, a)} \right| \quad (C_k(s, a) \leq \mathfrak{N}_k(s, a))
 \end{aligned}$$

$$\leq \frac{\mathbf{n}_k(s, a)}{\mathfrak{N}_{k+1}^+(s, a)} + \frac{\mathfrak{N}_k(s, a)\mathbf{n}_k(s, a)}{\mathfrak{N}_k^+(s, a)\mathfrak{N}_{k+1}^+(s, a)} \leq \frac{2\mathbf{n}_k(s, a)}{\mathfrak{N}_k^+(s, a)},$$

and by $|\sqrt{a} - \sqrt{b}| \leq \sqrt{|a - b|}$, $\mathbf{n}_k(s, a) \in \mathbb{N}$:

$$\begin{aligned} & \left| \sqrt{\frac{\bar{c}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)}} - \sqrt{\frac{\bar{c}_{k+1}(s, a)\iota}{\mathfrak{N}_{k+1}^+(s, a)}} \right| \\ & \leq \sqrt{\frac{|\bar{c}_k(s, a) - \bar{c}_{k+1}(s, a)|\iota}{\mathfrak{N}_k^+(s, a)}} + \sqrt{\bar{c}_{k+1}(s, a)\iota} \left(\frac{1}{\sqrt{\mathfrak{N}_k^+(s, a)}} - \frac{1}{\sqrt{\mathfrak{N}_{k+1}^+(s, a)}} \right) \\ & \leq \frac{2\mathbf{n}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)} + \left(\sqrt{\frac{\iota}{\mathfrak{N}_k^+(s, a)}} - \sqrt{\frac{\iota}{\mathfrak{N}_{k+1}^+(s, a)}} \right) = \mathcal{O}\left(\frac{\mathbf{n}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)}\right), \end{aligned}$$

where in the last inequality we apply

$$\begin{aligned} \frac{1}{\sqrt{\mathfrak{N}_k^+(s, a)}} - \frac{1}{\sqrt{\mathfrak{N}_{k+1}^+(s, a)}} &= \left(\frac{1}{\mathfrak{N}_k^+(s, a)} - \frac{1}{\mathfrak{N}_{k+1}^+(s, a)} \right) / \left(\frac{1}{\sqrt{\mathfrak{N}_k^+(s, a)}} + \frac{1}{\sqrt{\mathfrak{N}_{k+1}^+(s, a)}} \right) \\ &\leq \sqrt{\mathfrak{N}_{k+1}^+(s, a)} \cdot \frac{\mathbf{n}_k(s, a)}{\mathfrak{N}_k^+(s, a)\mathfrak{N}_{k+1}^+(s, a)} \leq \frac{\mathbf{n}_k(s, a)}{\mathfrak{N}_k^+(s, a)}. \end{aligned} \quad (14)$$

Thus, $|d\hat{c}_k(s, a)| = \mathcal{O}\left(\frac{\mathbf{n}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)}\right)$.

Second statement: Define $\Pi_k(P') = \operatorname{argmin}_{P'' \in \mathcal{P}_{k+1}} \sum_{s, a, h} \|P''_{s, a, h} - P'_{s, a, h}\|_1$ for any $P' \in \mathcal{P}_k$. By the definition of \mathcal{P}_k , we have (note that $P'_{s, a, h}(s', h') = 0$ for $h' \notin \{h, h + 1\}$):

$$\|\Pi_k(P')_{s, a, h} - P'_{s, a, h}\|_1 \leq 2 \sum_{s'} |\bar{P}_{k, s, a}(s') - \bar{P}_{k+1, s, a}(s')| + 2 \sum_{s'} |\epsilon_{k+1}(s, a, s') - \epsilon_k(s, a, s')|.$$

Denote by $n_k(s, a, s')$ the number of visits to (s, a, s') (before policy switch or goal state is reached) in episode k . Note that:

$$\begin{aligned} |\bar{P}_{k, s, a}(s') - \bar{P}_{k+1, s, a}(s')| &= \left| \frac{N_k(s, a, s') + n_k(s, a, s')}{N_{k+1}^+(s, a)} - \frac{N_k(s, a, s')}{N_k^+(s, a)} \right| \\ &\leq N_k(s, a, s') \left(\frac{1}{N_k^+(s, a)} - \frac{1}{N_{k+1}^+(s, a)} \right) + \frac{n_k(s, a, s')}{N_{k+1}^+(s, a)} \leq \frac{2n_k(s, a)}{N_k^+(s, a)}. \end{aligned}$$

and by $|\sqrt{a} - \sqrt{b}| \leq \sqrt{|a - b|}$,

$$\begin{aligned} |\epsilon_k(s, a, s') - \epsilon_{k+1}(s, a, s')| &= \mathcal{O}\left(\left| \sqrt{\frac{\bar{P}_{k, s, a}(s')\iota}{N_k^+(s, a)}} - \sqrt{\frac{\bar{P}_{k+1, s, a}(s')\iota}{N_{k+1}^+(s, a)}} \right| + d\left(\frac{-\iota}{N_k^+(s, a)}\right) \right) \\ &= \mathcal{O}\left(\sqrt{\frac{|\bar{P}_{k, s, a}(s') - \bar{P}_{k+1, s, a}(s')|\iota}{N_k^+(s, a)}} + \sqrt{\bar{P}_{k+1, s, a}(s')\iota} d\left(\frac{-1}{\sqrt{N_k^+(s, a)}}\right) + d\left(\frac{-\iota}{N_k^+(s, a)}\right) \right) \end{aligned}$$

$$= \mathcal{O} \left(\frac{n_k(s, a)\iota}{N_k^+(s, a)} + \sqrt{\bar{P}_{k+1, s, a}(s')} d \left(\frac{-\sqrt{\iota}}{\sqrt{N_k^+(s, a)}} \right) \right).$$

Plugging these back, and by Cauchy-Schwarz inequality and Eq. (14) with $\mathfrak{N}_k = N_k$, we have

$$\|\Pi_k(P')_{s, a, h} - P'_{s, a, h}\|_1 = \mathcal{O} \left(\frac{Sn_k(s, a)\iota}{N_k^+(s, a)} + d \left(\frac{-\sqrt{S\iota}}{\sqrt{N_k^+(s, a)}} \right) \right) = \mathcal{O} \left(\frac{Sn_k(s, a)\iota}{N_k^+(s, a)} \right). \quad (15)$$

Thus, for any policy π' and cost function $c' \in \mathcal{C}_{\mathcal{M}}$ with $c'(s, a, h) \in [0, 1]$ for $h \leq H$, by Lemma 30 and Eq. (15),

$$\begin{aligned} & \left| Q^{\pi', \Pi_k(P'), c'}(s, a, h) - Q^{\pi', P', c'}(s, a, h) \right| \\ &= \left| \sum_{s', a', h'} q_{\pi', P', (s, a, h)}(s', a', h') (\Pi_k(P')_{s', a', h'} - P'_{s', a', h'}) V^{\pi', \Pi_k(P'), c'} \right| \\ &= \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{Sn_k(s', a')\iota}{N_k^+(s', a')} \right). \end{aligned} \quad (16)$$

Now define $P'_k = \Pi_k(P_k)$. We have

$$\begin{aligned} & \widehat{Q}_{k+1}(s, a, h) - \widehat{Q}_k(s, a, h) = Q^{\pi_{k+1}, P_{k+1}, \widehat{c}_{k+1}}(s, a, h) - Q^{\pi_k, P_k, \widehat{c}_k}(s, a, h) \\ & \leq Q^{\pi_{k+1}, P'_k, \widehat{c}_{k+1}}(s, a, h) - Q^{\pi_{k+1}, P_k, \widehat{c}_{k+1}}(s, a, h) + Q^{\pi_{k+1}, P_k, \widehat{c}_{k+1}}(s, a, h) - Q^{\pi_k, P_k, \widehat{c}_k}(s, a, h) \\ & = \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{Sn_k(s', a')\iota}{N_k^+(s', a')} \right) + (Q^{\pi_{k+1}, P_k, \widehat{c}_{k+1}}(s, a, h) - Q^{\pi_{k+1}, P_k, \widehat{c}_k}(s, a, h)) \quad (\text{Eq. (16)}) \\ & \quad + (Q^{\pi_{k+1}, P_k, \widehat{c}_k}(s, a, h) - Q^{\pi_k, P_k, \widehat{c}_k}(s, a, h)) \\ & = \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{Sn_k(s', a')\iota}{N_k^+(s', a')} + T_{\max} \sum_{s', a'} |\widehat{c}_{k+1}(s', a') - \widehat{c}_k(s', a')| + \eta T_{\max}^3 \right) \\ & \hspace{15em} (\text{Lemma 30 and Lemma 23}) \\ & = \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{Sn_k(s', a')\iota}{N_k^+(s', a')} + T_{\max} \sum_{s', a'} \frac{\mathfrak{n}_k(s', a')\iota}{\mathfrak{N}_k^+(s', a')} + \eta T_{\max}^3 \right). \end{aligned}$$

The other direction can be proved similarly.

Third statement: Note that $|d\tilde{c}_k(s, a, H+1)| = 0$, and for $h \leq H$,

$$\begin{aligned} & |\tilde{c}_{k+1}(s, a, h) - \tilde{c}_k(s, a, h)| \\ & \leq |d\widehat{c}_k(s, a)| + \lambda \left| \widehat{c}_{k+1}(s, a) \widehat{Q}_{k+1}(s, a, h) - \widehat{Q}_k(s, a, h) \widehat{c}_k(s, a) \right| + |de_k(s, a, h)| \\ & \leq |d\widehat{c}_k(s, a)| + \lambda \widehat{Q}_{k+1}(s, a, h) |d\widehat{c}_k(s, a)| + \lambda \widehat{c}_k(s, a) \left| d\widehat{Q}_k(s, a, h) \right| + |de_k(s, a, h)| \end{aligned}$$

$$= \mathcal{O} \left(\frac{\mathbf{n}_k(s, a)\iota}{\mathfrak{N}_k^+(s, a)} + \lambda T_{\max}^2 \sum_{s', a'} \frac{S n_k(s', a')\iota}{N_k^+(s', a')} + \lambda T_{\max} \sum_{s', a'} \frac{\mathbf{n}_k(s', a')}{\mathfrak{N}_k^+(s', a')} + \lambda \eta T_{\max}^3 + |de_k(s, a, h)| \right).$$

Fourth statement: Define $\tilde{P}'_k = \Pi_k(\tilde{P}_k)$. By Eq. (15), $\|\tilde{P}'_{k,s,a,h} - \tilde{P}_{k,s,a,h}\|_1 = \mathcal{O}\left(\frac{S n_k(s,a)\iota}{N_k^+(s,a)}\right)$, and $\lambda \leq 1/T_{\max}$, we have

$$\begin{aligned} \tilde{Q}_{k+1}(s, a, h) - \tilde{Q}_k(s, a, h) &\leq Q^{\pi_{k+1}, \tilde{P}'_k, \tilde{c}_{k+1}}(s, a, h) - Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) \\ &= \left(Q^{\pi_{k+1}, \tilde{P}'_k, \tilde{c}_{k+1}}(s, a, h) - Q^{\pi_{k+1}, \tilde{P}_k, \tilde{c}_{k+1}}(s, a, h) \right) \\ &\quad + \left(Q^{\pi_{k+1}, \tilde{P}_k, \tilde{c}_{k+1}}(s, a, h) - Q^{\pi_{k+1}, \tilde{P}_k, \tilde{c}_k}(s, a, h) \right) + \left(Q^{\pi_{k+1}, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) \right) \\ &\stackrel{(i)}{\leq} \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{S n_k(s', a')\iota}{N_k^+(s', a')} \right) + \sum_{s', a', h'} q_{\pi_{k+1}, \tilde{P}_k, (s', a', h)}(s', a', h') |\tilde{c}_{k+1}(s', a', h') - \tilde{c}_k(s', a', h')| \\ &= \mathcal{O} \left(T_{\max}^2 \sum_{s', a'} \frac{S n_k(s', a')\iota}{N_k^+(s', a')} + \lambda T_{\max}^2 \sum_{s', a'} \frac{\mathbf{n}_k(s', a')}{\mathfrak{N}_k^+(s', a')} + \lambda \eta T_{\max}^4 + T_{\max} \|de_k\|_1 \right), \end{aligned}$$

where in (i) we apply Eq. (16), Lemma 30, and

$$\begin{aligned} &Q^{\pi_{k+1}, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) \\ &= \sum_{s'', h''} \tilde{P}_{k,s,a,h}(s'', h'') \sum_{s', a', h'} q_{\pi_{k+1}, \tilde{P}_k, (s', a', h)}(s', a', h') (d\pi_k(a'|s', h')) Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s', a', h') \\ &\leq 0. \end{aligned} \tag{Lemma 30} \tag{Eq. (13)}$$

This completes the proof. \blacksquare

Lemma 26 For any cost function c in \mathcal{M} such that $c((s, h), a) \geq 0$, we have:

$$\begin{aligned} \text{Var}_k[\langle n_k, c \rangle] &= \sum_{s, a, h} q_k(s, a, h) (A^{\pi_k, P, c}(s, a, h)^2 + \mathbb{V}(P_{s, a, h}, V^{\pi_k, P, c})) \\ &\leq \mathbb{E}_k[\langle n_k, c \rangle^2] \leq 2 \langle q_k, c \circ Q^{\pi_k, P, c} \rangle. \end{aligned}$$

Proof Let $Q = Q^{\pi_k, P, c}$, $V = V^{\pi_k, P, c}$, $A = A^{\pi_k, P, c}$ and define $c(g, a) = 0$. Then,

$$\begin{aligned} \text{Var}_k[\langle n_k, c \rangle] &= \mathbb{E}_k \left[\left(\sum_{i=1}^{J_k+1} c(\hat{s}_i^k, a_i^k) - V(\hat{s}_1^k) \right)^2 \right] \\ &= \mathbb{E}_k \left[\left(\sum_{i=2}^{J_k+1} c(\hat{s}_i^k, a_i^k) + Q(\hat{s}_1^k, a_1^k) - P_{\hat{s}_1^k, a_1^k} V - V(\hat{s}_1^k) \right)^2 \right] \quad (Q(\hat{s}, a) = c(\hat{s}, a) + P_{\hat{s}, a} V) \end{aligned}$$

$$\begin{aligned}
 &\stackrel{(i)}{=} \mathbb{E}_k \left[\left(Q(\hat{s}_1^k, a_1^k) - V(\hat{s}_1^k) \right)^2 \right] + \mathbb{E}_k \left[\left(\sum_{i=2}^{J_k+1} c(\hat{s}_i^k, a_i^k) - P_{\hat{s}_1^k, a_1^k} V \right)^2 \right] \\
 &\stackrel{(ii)}{=} \mathbb{E}_k \left[\left(Q(\hat{s}_1^k, a_1^k) - V(\hat{s}_1^k) \right)^2 \right] + \mathbb{E}_k \left[\left(\sum_{i=2}^{J_k+1} c(\hat{s}_i^k, a_i^k) - V(\hat{s}_2^k) \right)^2 \right] + \mathbb{E}_k \left[\left(V(\hat{s}_2^k) - P_{\hat{s}_1^k, a_1^k} V \right)^2 \right] \\
 &= \mathbb{E}_k \left[\sum_{i=1}^{J_k+1} \left[\left(Q(\hat{s}_i^k, a_i^k) - V(\hat{s}_i^k) \right)^2 + \left(V(\hat{s}_{i+1}^k) - P_{\hat{s}_i^k, a_i^k} V \right)^2 \right] \right] \quad (\text{recursive argument}) \\
 &= \sum_{s,a,h} q_k(s, a, h) \left(A^2(s, a, h) + \mathbb{V}(P_{s,a,h}, V) \right),
 \end{aligned}$$

where (i) is by $Q(\hat{s}_1^k, a_1^k) - V(\hat{s}_1^k) \in \sigma(\hat{s}_1^k, a_1^k)$ (the σ -algebra of events defined on (\hat{s}_1^k, a_1^k)) and

$$\mathbb{E}_k \left[\sum_{i=2}^{J_k+1} c(\hat{s}_i^k, a_i^k) - P_{\hat{s}_1^k, a_1^k} V \middle| \hat{s}_1^k, a_1^k \right] = 0;$$

(ii) is by $V(\hat{s}_2^k) - P_{\hat{s}_1^k, a_1^k} V \in \sigma(\hat{s}_1^k, a_1^k, \hat{s}_2^k)$ and

$$\mathbb{E}_k \left[\sum_{i=2}^{J_k+1} c(\hat{s}_i^k, a_i^k) - V(\hat{s}_2^k) \middle| \hat{s}_1^k, a_1^k, \hat{s}_2^k \right] = 0.$$

Moreover, by $(\sum_{i=1}^n a_i)^2 \leq 2a_i(\sum_{i'=i}^n a_{i'})$ for any $n \geq 1$ and $P(J_k = \infty) = 0$,

$$\begin{aligned}
 \text{Var}_k[\langle n_k, c \rangle] &\leq \mathbb{E}_k[\langle n_k, c \rangle^2] = \mathbb{E}_k \left[\left(\sum_{i=1}^{J_k+1} c(\hat{s}_i^k, a_i^k) \right)^2 \right] \leq 2\mathbb{E}_k \left[\sum_{i=1}^{J_k+1} c(\hat{s}_i^k, a_i^k) \sum_{i'=i}^{J_k+1} c(\hat{s}_{i'}^k, a_{i'}^k) \right] \\
 &= 2\mathbb{E}_k \left[\sum_{i=1}^{\infty} \mathbb{I}\{J_k + 1 \geq i\} c(\hat{s}_i^k, a_i^k) \sum_{i'=i}^{J_k+1} c(\hat{s}_{i'}^k, a_{i'}^k) \right] \\
 &\stackrel{(i)}{=} 2\mathbb{E}_k \left[\sum_{i=1}^{J_k+1} c(\hat{s}_i^k, a_i^k) Q(\hat{s}_i^k, a_i^k) \right] = 2 \langle q_k, c \circ Q \rangle,
 \end{aligned}$$

where in (i) we apply $Q(\hat{s}_i^k, a_i^k) = \mathbb{E}[\sum_{i'=i}^{J_k+1} c(\hat{s}_{i'}^k, a_{i'}^k) | \hat{s}_1^k, a_1^k, \dots, \hat{s}_i^k, a_i^k]$ and $\{J_k + 1 \geq i\} \in \sigma(\hat{s}_1^k, a_1^k, \dots, \hat{s}_i^k, a_i^k)$. \blacksquare

Lemma 27 For every $k \in [K]$ it holds that $q_k(s, a, h) \leq \mathbb{E}_k[\bar{n}_k(s, a, h)] + \tilde{O}(1/K)$.

Proof By definition of $n_k(s, a, h)$, $x_k(s, a, h)$, and $y_k(s, a, h)$ we have:

$$\begin{aligned}
 \Pr(n_k(s, a, h) > n) &= \Pr(n_k(s, a, h) > n \mid n_k(s, a, h) > n - 1) \Pr(n_k(s, a, h) > n - 1) \\
 &= \Pr(\text{return to } (s, a, h)) \Pr(n_k(s, a, h) > n - 1) \\
 &= y_k(s, a, h) \Pr(n_k(s, a, h) > n - 1)
 \end{aligned}$$

$$= \dots = y_k^n(s, a, h) \Pr(n_k(s, a, h) > 0) = y_k^n(s, a, h)x_k(s, a, h).$$

Now, since $q_k(s, a, h)$ is the expected number of visits to (s, a, h) ,

$$\begin{aligned} q_k(s, a, h) &= \mathbb{E}_k[n_k(s, a, h)] = \sum_{n=0}^{\infty} \Pr(n_k(s, a, h) > n) = x_k(s, a, h) \sum_{n=0}^{\infty} y_k^n(s, a, h) \\ &= x_k(s, a, h) \sum_{n=0}^{L-1} y_k^n(s, a, h) + x_k(s, a, h) \sum_{n=L}^{\infty} y_k^n(s, a, h). \end{aligned}$$

To finish we bound each of the sums separately. By definition of $\bar{n}_k(s, a, h)$:

$$\begin{aligned} x_k(s, a, h) \sum_{n=0}^{L-1} y_k^n(s, a, h) &= \sum_{n=0}^{L-1} \Pr(n_k(s, a, h) > n) = \sum_{n=0}^{L-1} \Pr(\min\{L, n_k(s, a, h)\} > n) \\ &\leq \sum_{n=0}^{\infty} \Pr(\min\{L, n_k(s, a, h)\} > n) \\ &= \sum_{n=0}^{\infty} \Pr(\bar{n}_k(s, a, h) > n) = \mathbb{E}_k[\bar{n}_k(s, a, h)]. \end{aligned}$$

In each step there's a probability of at most γ to stay in layer h . So $y_k(s, a, h) \leq \gamma$, which implies:

$$\begin{aligned} x_k(s, a, h) \sum_{n=L}^{\infty} y_k^n(s, a, h) &\leq \sum_{n=L}^{\infty} \gamma^n = \frac{\gamma^L}{1-\gamma} \leq \frac{\gamma^{\frac{8H}{1-\gamma} \ln(2T_{\max}K/\delta)}}{1-\gamma} \leq \frac{e^{-8H \ln(2T_{\max}K/\delta)}}{1-\gamma} \\ &\leq 2T_{\max} \left(\frac{\delta}{2T_{\max}K} \right)^{8 \log_2(c_f K)} = \tilde{\mathcal{O}}(1/K), \end{aligned}$$

where the second inequality uses $\gamma^{\frac{1}{1-\gamma}} \leq e^{-1}$. ■

Lemma 28 Consider a sequence of cost functions $\{c_k\}_{k=1}^K$ and transition functions $\{P_k\}_{k=1}^K$ such that $c_k \in \mathcal{C}_{\mathcal{M}}$ and $P_k \in \mathcal{P}_k$. Also define $\hat{q}_k = q_{\pi_k, P_k}$. Then with probability at least $1 - 8\delta$,

$$\sum_{k=1}^K |\langle q_k - \hat{q}_k, c_k \rangle| = \tilde{\mathcal{O}} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, c_k \circ Q^{\pi_k, P, c_k} \rangle} + S^{2.5} A^{1.5} T_{\max}^3 \right).$$

Proof Define $v_{k,s,a,h}(\hat{s}') = V^{\pi_k, P, c_k}(\hat{s}') - P_{s,a,h} V^{\pi_k, P, c_k}$ for $\hat{s}' \in \hat{\mathcal{S}}_+$. Note that with probability at least $1 - 4\delta$:

$$\begin{aligned} \sum_{k=1}^K |\langle q_k - \hat{q}_k, c_k \rangle| &= \sum_{k=1}^K \left| \sum_{s,a,h} q_k(s, a, h) (P_{s,a,h} - P_{k,s,a,h}) V^{\pi_k, P, c_k} \right| \quad (\text{Lemma 30}) \\ &= \sum_{k=1}^K \left| \sum_{s,a,h} q_k(s, a, h) (P_{s,a,h} - P_{k,s,a,h}) V^{\pi_k, P, c_k} \right| + \tilde{\mathcal{O}}(S^{2.5} A^{1.5} T_{\max}^3). \end{aligned}$$

(Lemma 13 and Lemma 29)

Below we bound the first term. We continue with:

$$\begin{aligned}
 &= \sum_{k=1}^K \left| \sum_{s,a} \sum_{h=1}^H q_k(s, a, h) (P_{s,a,h} - P_{k,s,a,h}) v_{k,s,a,h} \right| && (P_{s,a,H+1} = P_{k,s,a,H+1}) \\
 &= \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{s,a,h \leq H, \hat{s}'} q_k(s, a, h) \sqrt{\frac{P_{s,a,h}(\hat{s}') v_{k,s,a,h}^2(\hat{s}')}{N_k^+(s, a)}} + ST_{\max} \sum_{k=1}^K \sum_{s,a} \frac{q_k(s, a)}{N_k^+(s, a)} \right). && \text{(Lemma 13)}
 \end{aligned}$$

By Lemma 27, we have $q_k(s, a, h) \leq \mathbb{E}_k[\bar{n}_k(s, a, h)] + \tilde{\mathcal{O}}(1/K)$. Therefore, we continue with

$$= \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{s,a,h \leq H, \hat{s}'} \mathbb{E}_k[\bar{n}_k(s, a, h)] \sqrt{\frac{P_{s,a,h}(\hat{s}') v_{k,s,a,h}^2(\hat{s}')}{N_k^+(s, a)}} + S^2 AT_{\max} \right) \quad \text{(Lemma 32)}$$

$$= \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{s,a,h \leq H, \hat{s}'} \bar{n}_k(s, a, h) \sqrt{\frac{P_{s,a,h}(\hat{s}') v_{k,s,a,h}^2(\hat{s}')}{N_k^+(s, a)}} + S^2 AT_{\max}^2 \right) \quad \text{(Lemma 52)}$$

$$\begin{aligned}
 &= \tilde{\mathcal{O}} \left(\sum_{k=1}^K \sum_{s,a,h \leq H, \hat{s}'} \bar{n}_k(s, a, h) \sqrt{\frac{P_{s,a,h}(\hat{s}') v_{k,s,a,h}^2(\hat{s}')}{N_{k+1}^+(s, a)}} \right. \\
 &\quad \left. + \tilde{\mathcal{O}} \left(ST_{\max}^2 \sum_{s,a} \sum_{k=1}^K \left(\frac{1}{\sqrt{N_k^+(s, a)}} - \frac{1}{\sqrt{N_{k+1}^+(s, a)}} \right) + S^2 AT_{\max}^2 \right) \right) \\
 &= \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \sum_{s,a, \hat{s}'} \frac{\bar{n}_k(s, a)}{N_{k+1}^+(s, a)}} \sqrt{\sum_{k=1}^K \sum_{s,a,h \leq H, \hat{s}'} \bar{n}_k(s, a, h) P_{s,a,h}(\hat{s}') v_{k,s,a,h}^2(\hat{s}') + S^2 AT_{\max}^2} \right) \\
 &\hspace{15em} \text{(Cauchy-Schwarz inequality)} \\
 &= \tilde{\mathcal{O}} \left(\sqrt{S^2 A} \sqrt{\sum_{k=1}^K \sum_{s,a,h \leq H, \hat{s}'} q_k(s, a, h) P_{s,a,h}(\hat{s}') v_{k,s,a,h}^2(\hat{s}') + SAT_{\max}^3 + S^2 AT_{\max}^2} \right) \\
 &\hspace{15em} \text{(Lemma 52)}
 \end{aligned}$$

$$\begin{aligned}
 &= \tilde{\mathcal{O}} \left(\sqrt{S^2 A \sum_{k=1}^K \text{Var}_k[\langle n_k, c_k \rangle] + S^2 AT_{\max}^2} \right) \\
 &\hspace{15em} (\sum_{\hat{s}'} P_{s,a,h}(\hat{s}') v_{k,s,a,h}^2(\hat{s}') = \mathbb{V}(P_{s,a,h}, V^{\pi_k, P, c_k}) \text{ and Lemma 26}) \\
 &= \tilde{\mathcal{O}} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, c_k \circ Q^{\pi_k, P, c_k} \rangle + S^2 AT_{\max}^2} \right). \quad \text{(Lemma 26)}
 \end{aligned}$$

Substituting these back completes the proof. \blacksquare

Lemma 29 Consider a sequence of cost functions $\{c_k\}_{k=1}^K$ and transition functions $\{P_k\}_{k=1}^K$ such that $c_k \in \mathcal{C}_{\mathcal{M}}$ and $P_k \in \mathcal{P}_k$. Then, we have with probability at least $1 - 4\delta$:

$$\sum_{k=1}^K \sum_{s,a,h,s',h'} q_k(s,a,h) \epsilon_k^*(s,a,h,s',h') |V^{\pi_k, P, c_k}(s',h') - V^{\pi_k, P_k, c_k}(s',h')| = \tilde{\mathcal{O}}(S^{2.5} A^{1.5} T_{\max}^3).$$

Proof Below \lesssim is equivalent to $\tilde{\mathcal{O}}(\cdot)$. Also denote $z = (s, a, h, s', h')$ and $\tilde{z} = (\tilde{s}, \tilde{a}, \tilde{h}, \tilde{s}', \tilde{h}')$. By Lemma 30 we have with probability at least $1 - 2\delta$:

$$\begin{aligned} |V^{\pi_k, P, c_k}(s',h') - V^{\pi_k, P_k, c_k}(s',h')| &\lesssim \sum_{\tilde{s}, \tilde{a}, \tilde{h}} q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h}) \left| P_{\tilde{s}, \tilde{a}, \tilde{h}} V^{\pi_k, P, c_k} - P_{k, \tilde{s}, \tilde{a}, \tilde{h}} V^{\pi_k, P, c_k} \right| \\ &\lesssim T_{\max} \sum_{\tilde{s}, \tilde{a}, \tilde{h}} q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h}) \left\| P_{\tilde{s}, \tilde{a}, \tilde{h}} - P_{k, \tilde{s}, \tilde{a}, \tilde{h}} \right\|_1 \\ &\lesssim T_{\max} \sum_{\tilde{s}, \tilde{a}, \tilde{h}, \tilde{s}', \tilde{h}'} q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h}) \epsilon_k^*(\tilde{s}, \tilde{a}, \tilde{h}, \tilde{s}', \tilde{h}'), \end{aligned}$$

where the second inequality is by Lemma 2, and the third is by Lemma 13. Thus, using Lemma 13 and the Cauchy-Schwarz inequality, we get:

$$\begin{aligned} &\sum_{k=1}^K \sum_{s,a,h,s',h'} q_k(s,a,h) \epsilon_k^*(s,a,h,s',h') |V^{\pi_k, P, c_k}(s',h') - V^{\pi_k, P_k, c_k}(s',h')| \\ &\lesssim T_{\max} \sum_{k=1}^K \sum_z q_k(s,a,h) \epsilon_k^*(s,a,h,s',h') \sum_{\tilde{z}} q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h}) \epsilon_k^*(\tilde{s}, \tilde{a}, \tilde{h}, \tilde{s}', \tilde{h}') \\ &\lesssim T_{\max} \sum_{k=1}^K \sum_z q_k(s,a,h) \sqrt{\frac{P_{s,a,h}(s',h')}{N_k^+(s,a)}} \sum_{\tilde{z}} q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h}) \sqrt{\frac{P_{\tilde{s}, \tilde{a}, \tilde{h}}(\tilde{s}', \tilde{h}')}{N_k^+(\tilde{s}, \tilde{a})}} \\ &\lesssim T_{\max} \sqrt{\sum_{k,z,\tilde{z}} \frac{q_k(s,a,h) P_{\tilde{s}, \tilde{a}, \tilde{h}}(\tilde{s}', \tilde{h}') q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h})}{N_k^+(s,a)}} \sqrt{\sum_{k,z,\tilde{z}} \frac{q_k(s,a,h) P_{s,a,h}(s',h') q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h})}{N_k^+(\tilde{s}, \tilde{a})}}. \end{aligned}$$

Note that we ignore some lower order terms in the calculation above. To finish the proof we bound each of the terms separately. For the first term we have with probability at least $1 - \delta$:

$$\begin{aligned} &\sum_{k,z,\tilde{z}} \frac{q_k(s,a,h) P_{\tilde{s}, \tilde{a}, \tilde{h}}(\tilde{s}', \tilde{h}') q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h})}{N_k^+(s,a)} \\ &= \sum_{k,s,a} \frac{(\sum_h q_k(s,a,h)) \sum_{s',h',\tilde{s},\tilde{a},\tilde{h}} q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h}) \sum_{\tilde{s}',\tilde{h}'} P_{\tilde{s}, \tilde{a}, \tilde{h}}(\tilde{s}', \tilde{h}')}{N_k^+(s,a)} \\ &\lesssim T_{\max} S \sum_{k,s,a} \frac{q_k(s,a)}{N_k^+(s,a)} \lesssim T_{\max}^2 S^2 A, \end{aligned}$$

where the last inequality is by [Lemma 32](#). For the second term we have with probability at least $1 - \delta$:

$$\begin{aligned}
 & \sum_{k,z,\tilde{z}} \frac{q_k(s, a, h) P_{s,a,h}(s', h') q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h})}{N_k^+(\tilde{s}, \tilde{a})} \\
 & \lesssim S \sum_{k,s,a,h,\tilde{s},\tilde{a},\tilde{h}} \frac{q_k(s, a, h) \sum_{s',h'} P_{s,a,h}(s', h') q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h})}{N_k^+(\tilde{s}, \tilde{a})} \\
 & \lesssim T_{\max} S \sum_{k,s,a,h,\tilde{s},\tilde{a},\tilde{h}} \frac{x_k(s, a, h) \sum_{s',h'} P_{s,a,h}(s', h') q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h})}{N_k^+(\tilde{s}, \tilde{a})} \\
 & \lesssim T_{\max} S \sum_{k,s,a,h,\tilde{s},\tilde{a},\tilde{h}} \frac{q_k(\tilde{s}, \tilde{a}, \tilde{h})}{N_k^+(\tilde{s}, \tilde{a})} \lesssim T_{\max} S^2 A \sum_{k,\tilde{s},\tilde{a}} \frac{q_k(\tilde{s}, \tilde{a})}{N_k^+(\tilde{s}, \tilde{a})} \lesssim S^3 A^2 T_{\max}^2,
 \end{aligned}$$

where the second inequality follows by $q_k(s, a, h) \lesssim T_{\max} x_k(s, a, h)$, the third by $x_k(s, a, h) \sum_{s',h'} P_{s,a,h}(s', h') q_{k,(s',h')}(\tilde{s}, \tilde{a}, \tilde{h}) \leq q_k(\tilde{s}, \tilde{a}, \tilde{h})$, and the last one by [Lemma 32](#). \blacksquare

Lemma 30 (Extended Value Difference) *For any policies π, π' , transitions P, P' , and cost functions c, c' in \mathcal{M} , we have:*

$$\begin{aligned}
 & Q^{\pi, P, c}(s, a, h) - Q^{\pi', P', c'}(s, a, h) \\
 & = \sum_{s'', h''} P'_{s,a,h}(s'', h'') \sum_{s', h'} q_{\pi', P', (s'', h'')} (s', h') \sum_{a'} (\pi(a'|s', h') - \pi'(a'|s', h')) Q^{\pi, P, c}(s', a', h') \\
 & \quad + \sum_{s', a', h'} q_{\pi', P', (s,a,h)}(s', a', h') (Q^{\pi, P, c}(s', a', h') - c'(s', a', h') - P'_{s', a', h'} V^{\pi, P, c}).
 \end{aligned}$$

and

$$\begin{aligned}
 & V^{\pi, P, c}(s, h) - V^{\pi', P', c'}(s, h) \\
 & = \sum_{s', h'} q_{\pi', P', (s,h)}(s', h') \sum_{a'} (\pi(a'|s', h') - \pi'(a'|s', h')) Q^{\pi, P, c}(s', a', h') \\
 & \quad + \sum_{s', a', h'} q_{\pi', P', (s,h)}(s', a', h') (Q^{\pi, P, c}(s', a', h') - c'(s', a', h') - P'_{s', a', h'} V^{\pi, P, c}).
 \end{aligned}$$

Proof We first prove the second statement, note that:

$$\begin{aligned}
 & V^{\pi, P, c}(s, h) - V^{\pi', P', c'}(s, h) = \sum_{a'} (\pi(a'|s, h) - \pi'(a'|s, h)) Q^{\pi, P, c}(s, a', h) \\
 & \quad + \sum_{a'} \pi'(a'|s, h) (Q^{\pi, P, c}(s, a', h) - Q^{\pi', P', c'}(s, a', h)) \\
 & = \sum_{a'} (\pi(a'|s, h) - \pi'(a'|s, h)) Q^{\pi, P, c}(s, a', h)
 \end{aligned}$$

$$\begin{aligned}
 & + \sum_{a'} \pi'(a'|s, h) (Q^{\pi, P, c}(s, a', h) - c'(s, a', h) - P'_{s, a', h} V^{\pi, P, c}) \\
 & + \sum_{a'} \pi'(a'|s, h) P'_{s, a', h} (V^{\pi, P, c} - V^{\pi', P', c'}).
 \end{aligned}$$

Applying the equality above recursively and by the definition of $q_{\pi', P', (s, h)}$, we prove the second statement. For the first statement, note that:

$$\begin{aligned}
 & Q^{\pi, P, c}(s, a, h) - Q^{\pi', P', c'}(s, a, h) \\
 & = (Q^{\pi, P, c}(s, a, h) - c'(s, a, h) - P'_{s, a, h} V^{\pi, P, c}) + P'_{s, a, h} (V^{\pi, P, c} - V^{\pi', P', c'}).
 \end{aligned}$$

Applying the second statement and the definition of $q_{\pi', P', (s, a, h)}$ completes the proof. \blacksquare

Lemma 31 (*Rosenberg and Mansour, 2021, Lemma 6*) *Let π be a policy with expected hitting time at most τ starting from any state. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, π takes no more than $4\tau \ln \frac{2}{\delta}$ steps to reach the goal state.*

Lemma 32 *For any $z_k : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, with probability at least $1 - \delta$,*

$$\begin{aligned}
 \sum_{k=1}^K \sum_{s, a} \frac{\bar{n}_k(s, a) \sqrt{z_k(s, a)}}{\sqrt{N_k^+(s, a)}} &= \tilde{\mathcal{O}} \left(SAT_{\max} + \sqrt{SA \sum_k \sum_{s, a} \bar{n}_k(s, a) z_k(s, a)} \right) \\
 &= \tilde{\mathcal{O}} \left(SAT_{\max} + \sqrt{SA \sum_k \sum_{s, a} q_k(s, a) z_k(s, a)} \right), \\
 \sum_{k=1}^K \sum_{s, a} \frac{q_k(s, a) \sqrt{z_k(s, a)}}{\sqrt{N_k^+(s, a)}} &= \tilde{\mathcal{O}} \left(SAT_{\max} + \sqrt{SA \sum_k \sum_{s, a} \bar{n}_k(s, a) z_k(s, a)} \right) \\
 &= \tilde{\mathcal{O}} \left(SAT_{\max} + \sqrt{SA \sum_k \sum_{s, a} q_k(s, a) z_k(s, a)} \right), \\
 \sum_{k=1}^K \sum_{s, a} \frac{\bar{n}_k(s, a)}{N_k^+(s, a)} &= \tilde{\mathcal{O}}(SAT_{\max}), \quad \sum_{k=1}^K \sum_{s, a} \frac{q_k(s, a)}{N_k^+(s, a)} = \tilde{\mathcal{O}}(SAT_{\max}), \\
 \sum_{k=1}^K \sum_{s, a} \frac{m_k(s, a)}{M_k^+(s, a)} &= \tilde{\mathcal{O}}(SA), \quad \sum_{k=1}^K \sum_{s, a} \frac{x_k(s, a)}{M_k^+(s, a)} = \tilde{\mathcal{O}}(SA).
 \end{aligned}$$

Proof First statement: Since $z_k(s, a) \leq 1$ and $\bar{n}_k(s, a) \leq L = \tilde{\mathcal{O}}(T_{\max})$ we have:

$$\begin{aligned}
 \sum_{k=1}^K \frac{\bar{n}_k(s, a) \sqrt{z_k(s, a)}}{\sqrt{N_k^+(s, a)}} &\leq \sum_{k=1}^K \frac{\bar{n}_k(s, a) \sqrt{z_k(s, a)}}{\sqrt{N_{k+1}^+(s, a)}} + \sum_{k=1}^K L \left(\frac{1}{\sqrt{N_k^+(s, a)}} - \frac{1}{\sqrt{N_{k+1}^+(s, a)}} \right) \\
 &\leq \sum_{k=1}^K \frac{\bar{n}_k(s, a) \sqrt{z_k(s, a)}}{\sqrt{N_{k+1}^+(s, a)}} + \tilde{\mathcal{O}}(T_{\max}).
 \end{aligned}$$

By Cauchy-Schwarz inequality this implies:

$$\begin{aligned} \sum_{s,a} \sum_{k=1}^K \frac{\bar{n}_k(s,a) \sqrt{z_k(s,a)}}{\sqrt{N_k^+(s,a)}} &= \tilde{O} \left(\sqrt{\sum_{k=1}^K \sum_{s,a} \frac{\bar{n}_k(s,a)}{N_{k+1}^+(s,a)}} \sqrt{\sum_{k=1}^K \sum_{s,a} \bar{n}_k(s,a) z_k(s,a) + SAT_{\max}} \right) \\ &= \tilde{O} \left(\sqrt{SA \sum_k \sum_{s,a} \bar{n}_k(s,a) z_k(s,a) + SAT_{\max}} \right). \end{aligned}$$

Finally, $\sum_{s,a} \sum_k \bar{n}_k(s,a) z_k(s,a) = \tilde{O} \left(\sum_{s,a} \sum_k q_k(s,a) z_k(s,a) + SAT_{\max} \right)$ with high probability by [Lemma 52](#).

Second statement: By [Lemma 27](#) we have:

$$\begin{aligned} \sum_{k=1}^K \sum_{s,a} \frac{q_k(s,a) \sqrt{z_k(s,a)}}{\sqrt{N_k^+(s,a)}} &\leq \sum_{k=1}^K \sum_{s,a} \frac{\mathbb{E}_k[\bar{n}_k(s,a)] \sqrt{z_k(s,a)}}{\sqrt{N_k^+(s,a)}} + \tilde{O}(1/K) \sum_{k=1}^K \sum_{s,a} \frac{\sqrt{z_k(s,a)}}{\sqrt{N_k^+(s,a)}} \\ &\leq \sum_{k=1}^K \sum_{s,a} \frac{\mathbb{E}_k[\bar{n}_k(s,a)] \sqrt{z_k(s,a)}}{\sqrt{N_k^+(s,a)}} + \tilde{O}(SA) \\ &= \tilde{O} \left(\sum_{k=1}^K \sum_{s,a} \frac{\bar{n}_k(s,a) \sqrt{z_k(s,a)}}{\sqrt{N_k^+(s,a)}} + T_{\max} SA \right), \end{aligned}$$

where the last relation holds with high probability by [Lemma 52](#). Now the statement follows by the first statement.

Third and forth statements: Similarly to the first statement,

$$\begin{aligned} \sum_{k=1}^K \frac{\bar{n}_k(s,a)}{N_k^+(s,a)} &\leq \sum_{k=1}^K \frac{\bar{n}_k(s,a)}{N_{k+1}^+(s,a)} + \sum_{k=1}^K L \left(\frac{1}{N_k^+(s,a)} - \frac{1}{N_{k+1}^+(s,a)} \right) \\ &\leq \sum_{k=1}^K \frac{\bar{n}_k(s,a)}{\max\{1, \sum_{i \leq k} \bar{n}_i(s,a)\}} + \tilde{O}(T_{\max}) = \tilde{O}(T_{\max}). \end{aligned}$$

Summing over (s, a) proves the third statement. The forth statement is then proved similarly to the second statement.

Fifth and sixth statements: Similarly to the third statement,

$$\begin{aligned} \sum_{k=1}^K \frac{m_k(s,a)}{M_k^+(s,a)} &\leq \sum_{k=1}^K \frac{m_k(s,a)}{M_{k+1}^+(s,a)} + \sum_{k=1}^K \left(\frac{1}{M_k^+(s,a)} - \frac{1}{M_{k+1}^+(s,a)} \right) \\ &\leq \sum_{k=1}^K \frac{m_k(s,a)}{\max\{1, \sum_{i \leq k} m_i(s,a)\}} + 1 = \tilde{O}(1). \end{aligned}$$

Summing over (s, a) proves the fifth statement. The sixth statement is again obtained with high probability by [Lemma 52](#). ■

Appendix D. Omitted Details for Section 5

Extra Notations Define $\tilde{q}_k = q_{\pi_k, \tilde{P}_k}$, $Q_k = Q^{\pi_k, P, c_k}$, $V_k = V^{\pi_k, P, c_k}$, and $A_k = A^{\pi_k, P, c_k}$.

D.1. Proof of Theorem 9

In this part, define $\tilde{P}_k = \Gamma(\pi_k, \mathcal{P}_k, \tilde{c}_k)$ and $P_k = \Gamma(\pi_k, \mathcal{P}_k, c_k)$, such that $\tilde{Q}_k = Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}$, $\tilde{V}_k = V^{\pi_k, \tilde{P}_k, \tilde{c}_k}$, and $\hat{Q}_k = Q^{\pi_k, P, c_k}$. We first provide bounds on some important quantities.

Lemma 33 $\tilde{c}_k \in \mathcal{C}_M$, $\eta \left\| \tilde{A}_k - B_k \right\|_\infty \leq 1$, and $\eta \|B_k\|_\infty \leq \frac{1}{2H'}$.

Proof For the first statement, by $P_k \in \Lambda_M$, we have $\hat{Q}_k(s, a, h) \leq \frac{H}{1-\gamma} + c_f = \chi$. Therefore, $\lambda \hat{Q}_k(s, a, h) \leq 1$ and $\tilde{c}_k \in \mathcal{C}_M$. For the second statement, by $\tilde{c}_k(s, a, h) \leq 2$ for $h \leq H$, we have $|\tilde{A}_k(s, a, h)| \leq |\hat{Q}_k(s, a, h)| + |\tilde{V}_k(s, h)| \leq 4(\frac{H}{1-\gamma} + c_f) = 4\chi$ for $h \leq H$. Therefore, $\|b_k\|_\infty \leq 32\eta\chi^2$, and by Lemma 45, we have $\|B_k\|_\infty \leq \frac{15H\|b_k\|_\infty}{1-\gamma} \leq 960\eta HT_{\max}\chi^2$. Thus by the definition of η , we have $\eta \|B_k\|_\infty \leq \frac{1}{2H'}$ and $\eta \left\| \tilde{A}_k - B_k \right\|_\infty \leq \eta(\left\| \tilde{A}_k \right\|_\infty + \|B_k\|_\infty) \leq 1$. ■

We are now ready to prove Theorem 9. The proof decomposes the regret into several terms, each of which is bounded by a lemma included after the proof.

Proof [Proof of Theorem 9] With probability at least $1 - 10\delta$, we decompose the regret as follows:

$$\begin{aligned} \dot{R}_K &= \sum_{k=1}^K \langle n_k - q_k, c_k \rangle + \langle q_k - q^*, c_k \rangle \stackrel{(i)}{\leq} \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle} + SAT_{\max} \right) \\ &\quad + \sum_{k=1}^K \langle q_k - \tilde{q}_k, \tilde{c}_k \rangle + \sum_{k=1}^K \langle \tilde{q}_k - q^*, \tilde{c}_k \rangle - \lambda \sum_{k=1}^K \langle q_k, c_k \circ \hat{Q}_k \rangle + \lambda \sum_{k=1}^K \langle q^*, c_k \circ \hat{Q}_k \rangle \\ &\stackrel{(ii)}{=} \tilde{\mathcal{O}} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle} + S^{2.5} A^{1.5} T_{\max}^3 \right) + \sum_{k=1}^K \langle \tilde{q}_k - q^*, \tilde{c}_k \rangle - \lambda \sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle \\ &\quad + \lambda \sum_{k=1}^K \langle q_k, c_k \circ (Q_k - \hat{Q}_k) \rangle + \lambda \sum_{k=1}^K \langle q^*, c_k \circ Q^{\pi_k, P, c_k} \rangle + \lambda \sum_{k=1}^K \langle q^*, c_k \circ (\hat{Q}_k - Q^{\pi_k, P, c_k}) \rangle, \end{aligned}$$

where in (i) we apply Lemma 11 and Lemma 50 to have

$$\begin{aligned} \sum_{k=1}^K \langle n_k - q_k, c_k \rangle &= \sum_{k=1}^K \langle \bar{n}_k - q_k, c_k \rangle = \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \mathbb{E}_k[\langle \bar{n}_k, c_k \rangle^2]} + SAT_{\max} \right) \\ &= \tilde{\mathcal{O}} \left(\sqrt{\sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle} + SAT_{\max} \right), \end{aligned} \tag{Lemma 26}$$

and in (ii) we apply Lemma 28 and $\tilde{c}_k \in \mathcal{C}_M$ for $h \leq H$ to have

$$\sum_{k=1}^K \langle q_k - \tilde{q}_k, \tilde{c}_k \rangle = \tilde{\mathcal{O}} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, \tilde{c}_k \circ Q^{\pi_k, P, \tilde{c}_k} \rangle} + S^{2.5} A^{1.5} T_{\max}^3 \right)$$

$$= \tilde{\mathcal{O}} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, c_k \circ Q^{\pi_k, P, c_k} \rangle + S^{2.5} A^{1.5} T_{\max}^3} \right). \\ (\tilde{c}_k(s, a, h) \leq 2c_k(s, a, h))$$

Define $\lambda' = \sqrt{\frac{S^2 A}{DT_* K}}$. By [Lemma 34](#), [Lemma 35](#), [Lemma 36](#), and definition of λ, η , with probability at least $1 - 9\delta$:

$$\begin{aligned} \hat{R}_K &\leq \tilde{\mathcal{O}} \left(\sqrt{S^2 A \sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle + \frac{T_*}{\eta} + S^4 A^2 T_{\max}^5} \right) \\ &\quad + 24\eta \sum_{k=1}^K \langle q_k, A_k^2 \rangle - \lambda \sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle + \lambda \sum_{k=1}^K \langle q^*, c_k \circ Q^{\hat{\pi}^*, P, c_k} \rangle \\ &= \tilde{\mathcal{O}} \left(\frac{S^2 A}{\lambda'} \right) + \lambda' \sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle + \tilde{\mathcal{O}} \left(\frac{T_*}{\eta} + S^4 A^2 T_{\max}^5 \right) \\ &\quad + 48\eta \sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle - \lambda \sum_{k=1}^K \langle q_k, c_k \circ Q_k \rangle + \mathcal{O}(\lambda D T_* K) \\ &\hspace{15em} \text{(AM-GM inequality, [Lemma 26](#), and [Lemma 37](#))} \\ &= \tilde{\mathcal{O}} \left(T_* \sqrt{DK} + \sqrt{S^2 A D T_* K} + S^4 A^2 T_{\max}^5 \right). \quad (K = \tilde{\mathcal{O}}(S^2 A T_{\max}^2) \text{ when } \lambda < 48\eta + \lambda') \end{aligned}$$

Applying [Lemma 4](#) completes the proof. \blacksquare

Lemma 34 *With probability at least $1 - 6\delta$,*

$$\sum_{k=1}^K \langle \tilde{q}_k - q^*, \tilde{c}_k \rangle = 24\eta \langle q_k, A_k^2 \rangle + \tilde{\mathcal{O}} \left(\frac{T_*}{\eta} + S^4 A^2 T_{\max}^{3.5} \right).$$

Proof Note that by [Lemma 49](#) and [Lemma 33](#):

$$\begin{aligned} &\sum_{k=1}^K \sum_{s, h} q^*(s, h) \sum_{a \in \mathcal{A}} (\pi_k(a|s, h) - \hat{\pi}^*(a|s, h)) \left(\tilde{A}_k(s, a, h) - B_k(s, a, h) \right) \\ &\leq \sum_{s, h} q^*(s, h) \left(\frac{\ln A}{\eta} + \eta \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s, h) \left(\tilde{A}_k(s, a, h) - B_k(s, a, h) \right)^2 \right) \\ &\leq \tilde{\mathcal{O}} \left(\frac{T_*}{\eta} \right) + 2\eta \sum_{s, h} q^*(s, h) \left(\sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s, h) \tilde{A}_k(s, a, h)^2 + \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s, h) B_k(s, a, h)^2 \right) \\ &= \tilde{\mathcal{O}} \left(\frac{T_*}{\eta} \right) + \sum_{k=1}^K \langle q^*, b_k \rangle + \frac{1}{H'} \sum_{s, h} q^*(s, h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s, h) B_k(s, a, h). \quad \text{(Lemma 33)} \end{aligned}$$

Define $\tilde{q}'_k = q_{\pi_k, P'_k}$, where P'_k is the optimistic transition defined in B_k . We have

$$\begin{aligned} \sum_{k=1}^K \langle \tilde{q}'_k - q^*, \tilde{c}_k \rangle &= \sum_{k=1}^K \sum_{s,h} q^*(s, h) \sum_{a \in \mathcal{A}} (\pi_k(a|s, h) - \hat{\pi}^*(a|s, h)) \left(\tilde{A}_k(s, a, h) - B_k(s, a, h) \right) \\ &\quad + \sum_{k=1}^K \sum_{s,a,h} q^*(s, a, h) \left(\tilde{Q}_k(s, a, h) - \tilde{c}_k(s, a, h) - P_{s,a,h} \tilde{V}_k \right) \\ &\quad + \sum_{k=1}^K \sum_{s,h} q^*(s, h) \sum_{a \in \mathcal{A}} (\pi_k(a|s, h) - \hat{\pi}^*(a|s, h)) B_k(s, a, h) \end{aligned}$$

(shifting argument and [Lemma 30](#))

$$\begin{aligned} &\stackrel{(i)}{\leq} \tilde{\mathcal{O}} \left(\frac{T^*}{\eta} \right) + 3 \sum_{k=1}^K \langle \tilde{q}'_k, b_k \rangle + \tilde{\mathcal{O}}(T_{\max}) \\ &= \tilde{\mathcal{O}} \left(\frac{T^*}{\eta} \right) + 6\eta \sum_{k=1}^K \langle q_k, \tilde{A}_k^2 \rangle + 3 \sum_{k=1}^K \langle \tilde{q}'_k - q_k, b_k \rangle, \end{aligned}$$

where in (i) we apply [Lemma 46](#), $b_k(s, a, h) = \tilde{\mathcal{O}}(1)$, and the definition of \tilde{P}_k so that

$$\sum_{k=1}^K \sum_{s,a,h} q^*(s, a, h) \left(\tilde{Q}_k(s, a, h) - \tilde{c}_k(s, a, h) - P_{s,a,h} \tilde{V}_k \right) \leq 0.$$

For the second term, by $(a + b + c)^2 \leq 2a^2 + 2(b + c)^2 \leq 2a^2 + 4b^2 + 4c^2$,

$$\begin{aligned} \eta \sum_{k=1}^K \sum_{s,a,h} q_k(s, a, h) \tilde{A}_k(s, a, h)^2 &\leq 2\eta \sum_{k=1}^K \sum_{s,a,h} q_k(s, a, h) A^{\pi_k, P, \tilde{c}_k}(s, a, h)^2 \\ &\quad + 4\eta \sum_{k=1}^K \sum_{s,a,h} q_k(s, a, h) \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, P, \tilde{c}_k}(s, a, h) \right)^2 \\ &\quad + 4\eta \sum_{k=1}^K \sum_{s,h} q_k(s, h) \left(V^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, h) - V^{\pi_k, P, \tilde{c}_k}(s, h) \right)^2 \\ &\leq 2\eta \sum_{k=1}^K \sum_{s,a,h} q_k(s, a, h) A^{\pi_k, P, \tilde{c}_k}(s, a, h)^2 \\ &\quad + 8\eta \sum_{k=1}^K \sum_{s,a,h} q_k(s, a, h) \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, P, \tilde{c}_k}(s, a, h) \right)^2, \end{aligned}$$

where in the last step we apply Cauchy-Schwarz inequality to obtain

$$\left(V^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, h) - V^{\pi_k, P, \tilde{c}_k}(s, h) \right)^2 = \left(\sum_a \pi_k(a|s, h) (Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, P, \tilde{c}_k}(s, a, h)) \right)^2$$

$$\leq \sum_a \pi_k(a|s, h) \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, P, \tilde{c}_k}(s, a, h) \right)^2.$$

Note that with probability at least $1 - 2\delta$,

$$\begin{aligned} & \left| Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, P, \tilde{c}_k}(s, a, h) \right| \tag{17} \\ &= \tilde{O} \left(T_{\max} \sum_{s', a', h' \leq H} q_{\pi_k, P, (s, a, h)}(s', a', h') \left\| P_{s', a', h'} - \tilde{P}_{k, s', a', h'} \right\|_1 \right) \\ & \quad \text{(Lemma 30, Hölder's inequality, and } V^{\pi_k, \tilde{P}_k, \tilde{c}_k} = \tilde{O}(T_{\max})\text{)} \\ &= \tilde{O} \left(T_{\max} S \sum_{s', a'} \frac{q_{\pi_k, P, (s, a, h)}(s', a')}{\sqrt{N_k^+(s', a')}} \right). \tag{Lemma 13} \end{aligned}$$

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

$$\begin{aligned} & \eta \sum_{k=1}^K \sum_{s, a, h} q_k(s, a, h) \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\pi_k, P, \tilde{c}_k}(s, a, h) \right)^2 \\ & \leq \eta \sum_{k=1}^K \sum_{s, a, h \leq H} q_k(s, a, h) T_{\max}^3 S^2 \sum_{s', a'} \frac{q_{\pi_k, P, (s, a, h)}(s', a')}{N_k^+(s', a')} \tag{Cauchy-Schwarz inequality} \\ & \stackrel{(i)}{=} \tilde{O} \left(\eta T_{\max}^4 S^3 A \sum_{k=1}^K \sum_{s', a'} \frac{q_{\pi_k, P}(s', a')}{N_k^+(s', a')} \right) \stackrel{(ii)}{=} \tilde{O}(\eta T_{\max}^5 S^4 A^2), \end{aligned}$$

where in (i) we apply $q_k(s, a, h) \leq 2T_{\max} x_k(s, a, h)$ and $x_k(s, a, h) q_{\pi_k, P, (s, a, h)}(s', a') \leq q_{\pi_k, P}(s', a')$, and in (ii) we apply Lemma 32. Plugging these back, we get:

$$\begin{aligned} & \eta \sum_{k=1}^K \sum_{s, a, h} q_k(s, a, h) \tilde{A}_k(s, a, h)^2 \leq 2\eta \sum_{k=1}^K \left\langle q_k, (A^{\pi_k, P, \tilde{c}_k})^2 \right\rangle + \tilde{O}(\eta T_{\max}^5 S^4 A^2) \\ & \leq 4\eta \sum_{k=1}^K \left\langle q_k, A_k^2 \right\rangle + 4\eta \lambda^2 \sum_{k=1}^K \left\langle q_k, (A^{\pi_k, P, \hat{Q}_k})^2 \right\rangle + \tilde{O}(\eta T_{\max}^5 S^4 A^2) \quad ((a+b)^2 \leq 2a^2 + 2b^2) \\ & \leq 4\eta \left\langle q_k, A_k^2 \right\rangle + \tilde{O}(\eta T_{\max}^5 S^4 A^2 + \eta \lambda^2 T_{\max}^5 K) = 4\eta \left\langle q_k, A_k^2 \right\rangle + \tilde{O}(S^4 A^2 T_{\max}^3). \end{aligned}$$

For the third term, with probability at least $1 - 3\delta$,

$$\begin{aligned} & \sum_{k=1}^K \left\langle \hat{q}'_k - q_k, b_k \right\rangle \leq \sum_{k=1}^K \sum_{s, a, h \leq H} q_k(s, a, h) \left\| \hat{P}'_{k, s, a, h} - P_{s, a, h} \right\|_1 \left\| V^{\pi_k, \hat{P}'_k, b_k} \right\|_{\infty} \\ & \quad \text{(Lemma 30 and Hölder's inequality)} \\ &= \tilde{O} \left(\eta S T_{\max}^3 \sum_{k=1}^K \sum_{s, a, h \leq H} \frac{q_k(s, a, h)}{\sqrt{N_k^+(s, a)}} \right) \quad (b_k(s, a, h) = \tilde{O}(\eta T_{\max}^2) \text{ and Lemma 13}) \end{aligned}$$

$$= \tilde{\mathcal{O}} \left(\eta S T_{\max}^3 \sqrt{S A T_{\max} K} + \eta S^2 A T_{\max}^4 \right) = \tilde{\mathcal{O}} \left(S^2 A T_{\max}^{3.5} \right). \quad (\text{Lemma 32})$$

Putting everything together completes the proof. \blacksquare

$$\textbf{Lemma 35} \quad \lambda \sum_{k=1}^K \left\langle q^*, c_k \circ (\widehat{Q}_k - Q^{\hat{\pi}^*, P, c_k}) \right\rangle = \tilde{\mathcal{O}} \left(S^2 A T_{\max}^5 \right).$$

Proof Define $q'_{s,a,h}(s', h') = \sum_{s'', h''} P_{s,a,h}(s'', h'') q^*_{(s'', h'')}(s', h')$. We have for $h \leq H$:

$$\begin{aligned} & \lambda \sum_{k=1}^K (\widehat{Q}_k(s, a, h) - Q^{\hat{\pi}^*, P, c_k}(s, a, h)) \leq \lambda \sum_{k=1}^K (Q^{\pi_k, \tilde{P}_k, c_k}(s, a, h) - Q^{\hat{\pi}^*, P, c_k}(s, a, h)) \\ & \leq \lambda \sum_{k=1}^K \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\hat{\pi}^*, P, \tilde{c}_k}(s, a, h) \right) + \lambda^2 \sum_{k=1}^K Q^{\hat{\pi}^*, P, \widehat{Q}_k}(s, a, h) \\ & \leq \lambda \sum_{k=1}^K \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s, a, h) - Q^{\hat{\pi}^*, P, \tilde{c}_k}(s, a, h) \right) + \tilde{\mathcal{O}} \left(\lambda^2 T_{\max}^2 K \right) \\ & \leq \lambda \sum_{s', h'} q'_{s,a,h}(s', h') \sum_{k=1}^K \sum_{a'} (\pi_k(a'|s', h') - \hat{\pi}^*(a'|s', h')) \left(\tilde{A}_k(s', a', h') - B_k(s', a', h') \right) \\ & \quad + \lambda \sum_{s', a', h'} q^*_{(s,a,h)}(s', a', h') \sum_{k=1}^K \left(Q^{\pi_k, \tilde{P}_k, \tilde{c}_k}(s', a', h') - \tilde{c}_k(s', a', h') - P_{s', a', h'} V^{\pi_k, \tilde{P}_k, \tilde{c}_k} \right) \\ & \quad + \lambda \sum_{s', h'} q'_{s,a,h}(s', h') \sum_{k=1}^K \sum_{a'} (\pi_k(a'|s', h') - \hat{\pi}^*(a'|s', h')) B_k(s', a', h') + \tilde{\mathcal{O}} \left(\lambda^2 T_{\max}^2 K \right) \end{aligned}$$

(Lemma 30)

$$\begin{aligned} & = \tilde{\mathcal{O}} \left(\lambda \sum_{s', h'} q'_{s,a,h}(s', h') \left(\frac{T_{\max}^*}{\eta} + \eta \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s, h) T_{\max}^2 \right) + \lambda \eta T_{\max}^4 K + \lambda^2 T_{\max}^2 K \right) \\ & \quad \left(\|\tilde{A}_k\|_{\infty} = \tilde{\mathcal{O}}(T_{\max}), \text{ definition of } \tilde{P}_k, \text{ and } B_k(s, a, h) = \tilde{\mathcal{O}}(\eta T_{\max}^3) \right) \\ & = \tilde{\mathcal{O}} \left(\frac{\lambda T_{\max}^2}{\eta} + \lambda \eta T_{\max}^4 K + \lambda^2 T_{\max}^2 K \right). \end{aligned}$$

Plugging this back and by the definition of λ, η :

$$\lambda \sum_{k=1}^K \left\langle q^*, \widehat{Q}_k - Q^{\hat{\pi}^*, P, c_k} \right\rangle = \tilde{\mathcal{O}} \left(\frac{\lambda T_{\max}^3}{\eta} + \lambda \eta T_{\max}^5 K + \lambda^2 T_{\max}^3 K \right) = \tilde{\mathcal{O}} \left(S^2 A T_{\max}^5 \right).$$

This completes the proof. \blacksquare

$$\textbf{Lemma 36} \quad \text{With probability at least } 1 - 3\delta, \lambda \sum_{k=1}^K \left\langle q_k, c_k \circ (Q_k - \widehat{Q}_k) \right\rangle = \tilde{\mathcal{O}} \left(S^{3.5} A^2 T_{\max}^3 \right).$$

Proof By similar arguments as in Eq. (17) with \tilde{P}_k replaced by P_k and \tilde{c}_k replaced by c_k , with probability at least $1 - 3\delta$:

$$\begin{aligned}
 \lambda \sum_{k=1}^K \langle q_k, Q_k - \hat{Q}_k \rangle &= \lambda T_{\max} S \sum_{k=1}^K \sum_{s,a,h \leq H} q_k(s, a, h) \sum_{s',a'} \frac{q_{\pi_k, P, (s,a,h)}(s', a')}{\sqrt{N_k^+(s', a')}} \\
 &= \tilde{\mathcal{O}} \left(\lambda S^2 A T_{\max}^2 \sum_{k=1}^K \sum_{s',a'} \frac{q_k(s', a')}{\sqrt{N_k^+(s', a')}} \right) \\
 &\quad (q_k(s, a, h) = \mathcal{O}(T_{\max} x_k(s, a, h)) \text{ and } x_k(s, a, h) q_{\pi_k, P, (s,a,h)}(s', a') \leq q_k(s', a')) \\
 &= \tilde{\mathcal{O}} \left(\lambda S^2 A T_{\max}^2 \sqrt{S A T_{\max} K} + \lambda S^3 A^2 T_{\max}^3 \right) = \tilde{\mathcal{O}} \left(S^{3.5} A^2 T_{\max}^3 \right). \quad (\text{Lemma 32})
 \end{aligned}$$

■

Lemma 37 $\sum_{k=1}^K \langle q^*, c_k \circ Q^{\hat{\pi}^*, P, c_k} \rangle = \tilde{\mathcal{O}}(DT_* K) + T_{\max}$.

Proof By Lemma 2, for $h \leq H$, $\sum_{k=1}^K c_k(s, a, h) Q^{\hat{\pi}^*, P, c_k}(s, a, h) \leq \sum_{k=1}^K (Q^{\pi^*, P, c_k}(s, a) + c_f) = \tilde{\mathcal{O}}(DK)$. Therefore,

$$\begin{aligned}
 &\sum_{k=1}^K \langle q^*, c_k \circ Q^{\hat{\pi}^*, P, c_k} \rangle \\
 &= \sum_{k=1}^K \sum_{s,a,h \leq H} q^*(s, a, h) c_k(s, a, h) Q^{\hat{\pi}^*, P, c_k}(s, a, h) + \sum_{k=1}^K \sum_{s,a} q^*(s, a, H+1) c_f \\
 &\leq \tilde{\mathcal{O}}(DT_* K) + T_{\max},
 \end{aligned}$$

where the last step is by $\sum_{s,a,h \leq H} q^*(s, a, h) = \mathcal{O}(T_*)$, $\sum_{k=1}^K Q^{\hat{\pi}^*, P, c_k}(s, a, h) = \mathcal{O}(DK)$, and Lemma 19. ■

D.2. Proof of Theorem 10

Here we denote by \hat{q}_k^t the occupancy measure w.r.t policy π_k and the optimistic transition defined in B_k . Also define $\bar{Q}_k(s, a, h) = \mathbb{E}_k[\sum_{i=1}^{\min\{J_k, L\}+1} c(s_i^k, a_i^k, h_i^k) | \pi_k, P, s_1^k = s, a_1^k = a, h_1^k = h]$, such that $\mathbb{E}_k[G_{k,s,a,h}] = x_k(s, a, h) \bar{Q}_k(s, a, h)$. We again decompose the regret into several terms, each of which is bounded by a lemma included after the proof.

Proof With probability at least $1 - 2\delta$, we decompose the regret as follows:

$$\begin{aligned}
 \hat{R}_K &= \sum_{k=1}^K \langle n_k - q^*, c_k \rangle = \sum_{k=1}^K \langle \bar{n}_k - q_k, c_k \rangle + \sum_{k=1}^K \langle q_k - q^*, c_k \rangle \quad (\text{Lemma 11}) \\
 &= \tilde{\mathcal{O}} \left(T_{\max} \sqrt{K} + S A T_{\max} \right) + \sum_{k=1}^K \sum_{s,h} q^*(s, h) \sum_{a \in \mathcal{A}} (\pi_k(a|s, h) - \hat{\pi}^*(a|s, h)) \bar{Q}_k(s, a, h) \\
 &\quad (\text{Lemma 50 and Lemma 30})
 \end{aligned}$$

$$\begin{aligned}
 &= \tilde{\mathcal{O}} \left(T_{\max} \sqrt{K} + SAT_{\max} \right) + \underbrace{\sum_{k=1}^K \sum_{s,h} q^*(s,h) \sum_{a \in \mathcal{A}} (\pi_k(a|s,h) - \hat{\pi}^*(a|s,h)) \tilde{Q}_k(s,a,h)}_{\text{REG}} \\
 &\quad + \underbrace{\sum_{k=1}^K \sum_{s,h} q^*(s,h) \sum_{a \in \mathcal{A}} \pi_k(a|s,h) (Q_k(s,a,h) - \tilde{Q}_k(s,a,h))}_{\text{BIAS}_1} \\
 &\quad + \underbrace{\sum_{k=1}^K \sum_{s,h} q^*(s,h) \sum_{a \in \mathcal{A}} \hat{\pi}^*(a|s,h) (\tilde{Q}_k(s,a,h) - Q_k(s,a,h))}_{\text{BIAS}_2}.
 \end{aligned}$$

Therefore, by [Lemma 40](#), we have $\hat{R}_K = \tilde{\mathcal{O}}(\sqrt{S^2 AT_{\max}^5 K} + S^{5.5} A^{3.5} T_{\max}^5)$ with probability at least $1 - 25\delta$. Applying [Lemma 4](#) completes the proof. \blacksquare

Lemma 38 $\left\| \tilde{Q}_k \right\|_{\infty} \leq L'/\theta$, $\|b_k\|_{\infty} \leq 5L'$, $\|B_k\|_{\infty} \leq 150HT_{\max}L'$, and $\eta \left\| \tilde{Q}_k - B_k \right\|_{\infty} \leq 1$.

Proof The first statement is by the definition of \tilde{Q}_k and $G_{k,s,a,h} \leq L'$. For the second statement, $b_k \leq 5L'$ by definition. For the third statement, by [Lemma 45](#), we have $\|B_k\|_{\infty} \leq \frac{15H\|b_k\|_{\infty}}{1-\gamma} \leq 150HT_{\max}L'$. For the fourth statement, we have $\eta \left\| \tilde{Q}_k - B_k \right\|_{\infty} \leq \eta(\left\| \tilde{Q}_k \right\|_{\infty} + \|B_k\|_{\infty}) \leq 1/2 + \eta 150HT_{\max}L' \leq 1$. \blacksquare

Lemma 39 $Q_k(s,a,h) - \bar{Q}_k(s,a,h) = \tilde{\mathcal{O}}(1/K)$.

Proof Note that:

$$\begin{aligned}
 Q_k(s,a,h) - \bar{Q}_k(s,a,h) &= \mathbb{E}_k \left[\sum_{i=\bar{J}_k+2}^{J_k+1} c(s_i^k, a_i^k, h_i^k) \middle| \pi_k, P, s_1^k = s, a_1^k = a, h_1^k = h \right] \\
 &= \tilde{\mathcal{O}} \left(\frac{T_{\max}}{T_{\max}K} \right) = \tilde{\mathcal{O}}(1/K). \tag{Lemma 31}
 \end{aligned}$$

Lemma 40 With probability at least $1 - 25\delta$,

$$\text{REG} + \text{BIAS}_1 + \text{BIAS}_2 = \tilde{\mathcal{O}} \left(\sqrt{S^2 AT_{\max}^5 K} + S^{5.5} A^{3.5} T_{\max}^5 \right).$$

Proof Define $\xi_B = \sum_{k=1}^K \sum_{s,h \leq H} q^*(s,h) \sum_{a \in \mathcal{A}} (\pi_k(a|s,h) - \hat{\pi}^*(a|s,h)) B_k(s,a,h)$. By [Lemma 38](#) and [Lemma 49](#), with probability at least $1 - \delta$,

$$\text{REG} = \sum_{k=1}^K \sum_{s,h \leq H} q^*(s,h) \sum_{a \in \mathcal{A}} (\pi_k(a|s,h) - \hat{\pi}^*(a|s,h)) \left(\tilde{Q}_k(s,a,h) - B_k(s,a,h) \right) + \xi_B$$

$$\begin{aligned}
 &\leq \frac{T_\star \ln A}{\eta} + \eta \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \left(\tilde{Q}_k(s,a,h) - B_k(s,a,h) \right)^2 + \xi_B \\
 &\leq \frac{T_\star \ln A}{\eta} + 2\eta \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \tilde{Q}_k^2(s,a,h) \\
 &\quad + \frac{1}{H'} \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) B_k(s,a,h) + \xi_B \\
 &\hspace{15em} ((a+b)^2 \leq 2a^2 + 2b^2 \text{ and } \eta \|B_k\|_\infty \leq \frac{1}{H'}) \\
 &\leq \tilde{\mathcal{O}} \left(\frac{T_\star}{\eta} \right) + \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{2\theta L'}{\bar{x}_k(s,a,h) + \theta} \\
 &\quad + \frac{1}{H'} \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) B_k(s,a,h) + \xi_B,
 \end{aligned}$$

where in the last inequality we apply:

$$\begin{aligned}
 &2\eta \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \tilde{Q}_k^2(s,a,h) \\
 &= 2\eta \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{G_{k,s,a,h}^2}{(\bar{x}_k(s,a,h) + \theta)^2} \\
 &\leq L' \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \frac{\theta \pi_k(a|s,h)}{\bar{x}_k(s,a,h) + \theta} \frac{m_k(s,a,h)}{\bar{x}_k(s,a,h) + \theta} \\
 &\hspace{15em} (2\eta = \theta/L' \text{ and } G_{k,s,a,h} \leq L' m_k(s,a,h)) \\
 &\leq L' \sum_{s,h \leq H} q^\star(s,h) \left(2 \sum_{k=1}^K \sum_{a \in \mathcal{A}} \frac{\theta \pi_k(a|s,h)}{\bar{x}_k(s,a,h) + \theta} + \tilde{\mathcal{O}} \left(\frac{1}{\theta} \right) \right) \quad (\text{Lemma 52 and } \frac{x_k(s,a,h)}{\bar{x}_k(s,a,h) + \theta} \leq 1) \\
 &\leq \sum_{s,h \leq H} q^\star(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \frac{2\theta L' \pi_k(a|s,h)}{\bar{x}_k(s,a,h) + \theta} + \tilde{\mathcal{O}} \left(\frac{T_\star L'}{\theta} \right).
 \end{aligned}$$

Therefore, by [Lemma 42](#), [Lemma 43](#), and [Lemma 46](#), with probability at least $1 - 24\delta$,

$$\begin{aligned}
 \text{REG} + \text{BIAS}_1 + \text{BIAS}_2 &\leq \tilde{\mathcal{O}} \left(\frac{T_\star}{\eta} \right) + 3 \sum_{k=1}^K \langle \hat{q}'_k, b_k \rangle + \tilde{\mathcal{O}}(T_{\max}^2) \\
 &\hspace{15em} (\theta = 2\eta L' \text{ and } \|b_k\|_\infty = \tilde{\mathcal{O}}(T_{\max}) \text{ by } \text{Lemma 38}) \\
 &\leq \tilde{\mathcal{O}} \left(\frac{T_\star}{\eta} + \sum_{k=1}^K \sum_{s,a,h \leq H} \hat{q}'_k(s,a,h) \frac{L'(\bar{x}_k(s,a,h) - \underline{x}_k(s,a,h)) + \theta L'}{\bar{x}_k(s,a,h) + \theta} + T_{\max}^2 \right)
 \end{aligned}$$

$$\begin{aligned}
 &\leq \tilde{\mathcal{O}} \left(\frac{T_\star}{\eta} + L' \sum_{k=1}^K \sum_{s,a,h \leq H} (\bar{x}_k(s,a,h) - \underline{x}_k(s,a,h)) \tilde{q}'_{k,(s,a,h)}(s,a,h) + \theta T_{\max} L' SAK + T_{\max}^2 \right). \\
 &\quad (\tilde{q}'_k(s,a,h) \leq \bar{x}_k(s,a,h) \tilde{q}'_{k,(s,a,h)}(s,a,h) \text{ and } \tilde{q}'_k(s,a,h) = \tilde{\mathcal{O}}(T_{\max} \bar{x}_k(s,a,h))) \\
 &= \tilde{\mathcal{O}} \left(\sqrt{S^2 A T_{\max}^5 K} + S^{5.5} A^{3.5} T_{\max}^5 \right). \quad (\text{Lemma 41 and definition of } \eta, \theta)
 \end{aligned}$$

This completes the proof. \blacksquare

Lemma 41 *With probability at least $1 - 22\delta$,*

$$\sum_{k=1}^K \sum_{s,a,h \leq H} (\bar{x}_k(s,a,h) - \underline{x}_k(s,a,h)) \tilde{q}'_{k,(s,a,h)}(s,a,h) = \tilde{\mathcal{O}}(\sqrt{S^2 A T_{\max}^3 K} + S^{5.5} A^{3.5} T_{\max}^4).$$

Proof For any $z \in \mathcal{S} \times \mathcal{A} \times [H]$, denote by \bar{q}_k^z / q_k^z the occupancy measure w.r.t the policy and transition defined in $\bar{x}_k(z) / \underline{x}_k(z)$ (transition at (s,a,h) can be randomly pick as long as $P_{\bar{q}_k^z}, P_{q_k^z} \in \mathcal{P}_k$). For a fixed tuple $z = (s,a,h)$,

$$\begin{aligned}
 &\bar{x}_k(s,a,h) \tilde{q}'_{k,(s,a,h)}(s,a,h) = \bar{q}_k^z(s,a,h) + \bar{x}_k(s,a,h) (\tilde{q}'_{k,(s,a,h)}(s,a,h) - \bar{q}_k^z(s,a,h)(s,a,h)) \\
 &\leq \bar{q}_k^z(s,a,h) + 2\bar{x}_k(s,a,h) \sum_{s',a',h'} \bar{q}_{k,(s,a,h)}^z(s',a',h') \sum_{s'',h''} \epsilon_k^*(s',a',h',s'',h'') \tilde{q}'_{k,(s'',h'')}(s,a,h) \\
 &\quad (\text{Lemma 30 and Lemma 13}) \\
 &\leq \bar{q}_k^z(s,a,h) + 2 \sum_{s',a',h'} \bar{q}_k^z(s',a',h') \sum_{s'',h''} \epsilon_k^*(s',a',h',s'',h'') \tilde{q}'_{k,(s'',h'')}(s,a,h) \\
 &\quad (\bar{x}_k(s,a,h) \bar{q}_{k,(s,a,h)}^z(s',a',h') \leq \bar{q}_k^z(s',a',h')) \\
 &= \bar{q}_k^z(s,a,h) + 2 \sum_{s',a',h'} q_k(s',a',h') \sum_{s'',h''} \epsilon_k^*(s',a',h',s'',h'') \tilde{q}'_{k,(s'',h'')}(s,a,h) \\
 &\quad + 2 \sum_{s',a',h'} (\bar{q}_k^z(s',a',h') - q_k(s',a',h')) \sum_{s'',h''} \epsilon_k^*(s',a',h',s'',h'') \tilde{q}'_{k,(s'',h'')}(s,a,h).
 \end{aligned}$$

Therefore, with probability at least $1 - 7\delta$,

$$\begin{aligned}
 &\sum_{k=1}^K \sum_{s,a,h \leq H} \bar{x}_k(s,a,h) \tilde{q}'_{k,(s,a,h)}(s,a,h) \stackrel{(i)}{\leq} \sum_{k=1}^K \sum_{s,a,h \leq H} \bar{q}_k^z(s,a,h) \\
 &\quad + \tilde{\mathcal{O}} \left(T_{\max} \sum_{k=1}^K \sum_{s',a',h'} q_k(s',a',h') \sum_{s'',h''} \epsilon_k^*(s',a',h',s'',h'') + S^{5.5} A^{3.5} T_{\max}^4 \right) \\
 &\stackrel{(ii)}{\leq} \sum_{k=1}^K \sum_{s,a,h \leq H} \bar{q}_k^z(s,a,h) + \tilde{\mathcal{O}} \left(\sqrt{S^2 A T_{\max}^3 K} + S^{5.5} A^{3.5} T_{\max}^4 \right),
 \end{aligned}$$

where in (i) we apply $\sum_{s,a,h \leq H} \tilde{q}'_{k,(s'',h'')}(s,a,h) = \tilde{\mathcal{O}}(T_{\max})$ and $(z = (s,a,h))$ iterates over $\mathcal{S} \times \mathcal{A} \times [H]$:

$$\sum_{k,z} \sum_{s',a',h'} (\bar{q}_k^z(s',a',h') - q_k(s',a',h')) \sum_{s'',h''} \epsilon_k^*(s',a',h',s'',h'') \tilde{q}'_{k,(s'',h'')}(z)$$

$$\begin{aligned}
 &\leq \sum_{k,z} \sum_{\substack{\tilde{s}, \tilde{a}, \tilde{h} \\ \tilde{s}', \tilde{h}'}} \sum_{\substack{s', a', h' \\ s'', h''}} q_k(\tilde{s}, \tilde{a}, \tilde{h}) \epsilon_k^*(\tilde{s}, \tilde{a}, \tilde{h}, \tilde{s}', \tilde{h}') \bar{q}_{k,(\tilde{s}', \tilde{h}')}^z(s', a', h') \epsilon_k^*(s', a', h', s'', h'') \hat{q}'_{k,(s'', h'')}(z) \\
 &\hspace{25em} \text{(Lemma 30 and Lemma 13)} \\
 &\leq \sum_{k,z} \sum_{\substack{s', a', h' \\ s'', h''}} \sum_{\substack{\tilde{s}, \tilde{a}, \tilde{h} \\ \tilde{s}', \tilde{h}'}} q_k(\tilde{s}, \tilde{a}, \tilde{h}) \epsilon_k^*(\tilde{s}, \tilde{a}, \tilde{h}, \tilde{s}', \tilde{h}') q_{k,(\tilde{s}', \tilde{h}')}^z(s', a', h') \epsilon_k^*(s', a', h', s'', h'') \hat{q}'_{k,(s'', h'')}(z) \\
 &\quad + \tilde{O} \left(S^{2.5} A^{1.5} T_{\max}^3 \sum_z \sum_{\substack{s', a', h' \\ s'', h''}} T_{\max} \right) \\
 &\hspace{15em} (\epsilon_k^*(s', a', h', s'', h'') \hat{q}'_{k,(s'', h'')}(z) = \tilde{O}(T_{\max}) \text{ and Lemma 29}) \\
 &= \tilde{O} (S^{5.5} A^{3.5} T_{\max}^4).
 \end{aligned}$$

and in (ii) we apply:

$$\begin{aligned}
 &T_{\max} \sum_{k=1}^K \sum_{s', a', h'} q_k(s', a', h') \sum_{s'', h''} \epsilon_k^*(s', a', h', s'', h'') \\
 &= \tilde{O} \left(T_{\max} \sum_{k=1}^K \sum_{s', a', h' \leq H} q_k(s', a', h') \sum_{s'', h''} \left(\sqrt{\frac{P_{s', a', h'}(s'', h'')}{N_k^+(s', a')}} + \frac{1}{N_k^+(s', a')} \right) \right) \\
 &\hspace{25em} \text{(definition of } \epsilon_k^*) \\
 &= \tilde{O} \left(T_{\max} \sqrt{S} \sum_{k=1}^K \sum_{s', a'} \frac{q_k(s', a')}{\sqrt{N_k^+(s', a')}} + T_{\max} S \sum_{k=1}^K \sum_{s', a'} \frac{q_k(s', a')}{N_k^+(s', a')} \right) \\
 &= \tilde{O} \left(\sqrt{S^2 A T_{\max}^3 K} + S^2 A T_{\max}^2 \right). \hspace{10em} \text{(Lemma 32)}
 \end{aligned}$$

By similar arguments, we also have with probability at least $1 - 7\delta$,

$$\begin{aligned}
 &-\sum_{k=1}^K \sum_{s, a, h \leq H} \underline{x}_k(s, a, h) \hat{q}'_{k,(s, a, h)}(s, a, h) \\
 &\leq -\sum_{k=1}^K \sum_{s, a, h \leq H} \underline{q}_k^z(s, a, h) + \tilde{O} \left(\sqrt{S^2 A T_{\max}^3 K} + S^{5.5} A^{3.5} T_{\max}^4 \right).
 \end{aligned}$$

Therefore, with probability at least $1 - 7\delta$,

$$\begin{aligned}
 &\sum_{k=1}^K \sum_{s, a, h \leq H} (\bar{x}_k(s, a, h) - \underline{x}_k(s, a, h)) \hat{q}'_{k,(s, a, h)}(s, a, h) \\
 &= \tilde{O} \left(\sum_{k=1}^K \sum_{s, a, h \leq H} (\bar{q}_k^{(s, a, h)}(s, a, h) - \underline{q}_k^{(s, a, h)}(s, a, h)) + \sqrt{S^2 A T_{\max}^3 K} + S^{5.5} A^{3.5} T_{\max}^4 \right)
 \end{aligned}$$

$$= \tilde{\mathcal{O}} \left(\sqrt{S^2 A T_{\max}^3 K} + S^{5.5} A^{3.5} T_{\max}^4 \right),$$

where in the last inequality we apply (similarly for $\sum_{k=1}^K (q_k(s, a, h) - \underline{q}_k^{(s, a, h)}(s, a, h))$):

$$\begin{aligned} & \sum_{k=1}^K \sum_{s, a, h \leq H} (\bar{q}_k^{(s, a, h)}(s, a, h) - q_k(s, a, h)) \\ & \leq \sum_{k, s, a, h \leq H} \sum_{s', a', h'} q_k(s', a', h') \sum_{s'', h''} \epsilon_k^*(s', a', h', s'', h'') \bar{q}_{k, (s'', h'')}^{(s, a, h)}(s, a, h) \\ & \hspace{25em} \text{(Lemma 30 and Lemma 13)} \\ & \leq \sum_{k, s, a, h \leq H} \sum_{s', a', h'} q_k(s', a', h') \sum_{s'', h''} \epsilon_k^*(s', a', h', s'', h'') q_{k, (s'', h'')} (s, a, h) + \tilde{\mathcal{O}} \left(\sum_{s, a, h} S^{2.5} A^{1.5} T_{\max}^3 \right) \\ & \hspace{25em} \text{(Lemma 29)} \\ & = \tilde{\mathcal{O}} \left(T_{\max} \sum_{k, s', a', h' \leq H} q_k(s', a', h') \sum_{s'', h''} \left(\sqrt{\frac{P_{s', a', h'}(s'', h'')}{N_k^+(s', a')}} + \frac{1}{N_k^+(s', a')} \right) + S^{3.5} A^{2.5} T_{\max}^3 \right) \\ & \hspace{25em} \text{(definition of } \epsilon_k^* \text{)} \\ & = \tilde{\mathcal{O}} \left(\sqrt{S^2 A T_{\max}^3 K} + S^{3.5} A^{2.5} T_{\max}^3 \right). \hspace{5em} \text{(Cauchy-Schwarz inequality and Lemma 32)} \end{aligned}$$

This completes the proof. ■

Lemma 42 *With probability at least $1 - \delta$,*

$$\text{BIAS}_1 \leq \sum_{k=1}^K \sum_{s, h \leq H} q^*(s, h) \sum_{a \in \mathcal{A}} \pi_k(a|s, h) \frac{L'(\bar{x}_k(s, a, h) - x_k(s, a, h)) + 2\theta L'}{\bar{x}_k(s, a, h) + \theta} + \tilde{\mathcal{O}} \left(\frac{T_* L'}{\theta} \right).$$

Proof Note that:

$$\begin{aligned} \text{BIAS}_1 &= \sum_{k=1}^K \sum_{s, h \leq H} q^*(s, h) \sum_{a \in \mathcal{A}} \pi_k(a|s, h) (Q_k(s, a, h) - \mathbb{E}_k[\tilde{Q}_k(s, a, h)]) \\ & \quad + \sum_{k=1}^K \sum_{s, h \leq H} q^*(s, h) \sum_{a \in \mathcal{A}} \pi_k(a|s, h) (\mathbb{E}_k[\tilde{Q}_k(s, a, h)] - \tilde{Q}_k(s, a, h)). \end{aligned}$$

For the first term,

$$\begin{aligned} & \sum_{k=1}^K \sum_{s, h \leq H} q^*(s, h) \sum_{a \in \mathcal{A}} \pi_k(a|s, h) (Q_k(s, a, h) - \mathbb{E}_k[\tilde{Q}_k(s, a, h)]) \\ & \leq \sum_{k=1}^K \sum_{s, h \leq H} q^*(s, h) \sum_{a \in \mathcal{A}} \pi_k(a|s, h) \bar{Q}_k(s, a, h) \left(1 - \frac{x_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta} \right) + \tilde{\mathcal{O}}(T_*) \\ & \hspace{25em} \text{(Lemma 39 and } \mathbb{E}_k[G_{k, s, a, h}] = x_k(s, a, h) \bar{Q}_k(s, a, h) \text{)} \end{aligned}$$

$$\leq \sum_{k=1}^K \sum_{s,h \leq H} q^*(s,h) \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{L'(\bar{x}_k(s,a,h) - \underline{x}_k(s,a,h) + \theta)}{\bar{x}_k(s,a,h) + \theta} + \tilde{\mathcal{O}}(T_\star).$$

For the second term, first note that $G_{k,s,a,h} \leq L' m_k(s,a,h)$ and

$$\begin{aligned} \text{Var}_k \left[\left\langle \pi_k(\cdot|s,h), \tilde{Q}_k(s, \cdot, h) \right\rangle \right] &\leq \mathbb{E}_k \left[\left\langle \pi_k(\cdot|s,h), \tilde{Q}_k(s, \cdot, h) \right\rangle^2 \right] \\ &\leq \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{\mathbb{E}_k[G_{k,s,a,h}^2]}{(\bar{x}_k(s,a,h) + \theta)^2} \leq \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{L'^2}{\bar{x}_k(s,a,h) + \theta}. \end{aligned}$$

(Cauchy-Schwarz inequality and $\mathbb{E}_k[G_{k,s,a,h}] \leq L' x_k(s,a,h)$)

Therefore, by [Lemma 50](#), with probability at least $1 - \delta$,

$$\begin{aligned} &\sum_{k=1}^K \sum_{s,h \leq H} q^*(s,h) \sum_{a \in \mathcal{A}} \pi_k(a|s,h) (\mathbb{E}_k[\tilde{Q}_k(s,a,h)] - \tilde{Q}_k(s,a,h)) \\ &= \tilde{\mathcal{O}} \left(\sum_{s,h \leq H} q^*(s,h) \left(\sqrt{\sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{L'^2}{\bar{x}_k(s,a,h) + \theta}} + \frac{L'}{\theta} \right) \right) \\ &\leq \sum_{s,h \leq H} q^*(s,h) \sum_{k=1}^K \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{\theta L'}{\bar{x}_k(s,a,h) + \theta} + \tilde{\mathcal{O}} \left(\frac{T_\star L'}{\theta} \right). \quad (\text{AM-GM inequality}) \end{aligned}$$

Summing these two terms, we have:

$$\text{BIAS}_1 \leq \sum_{k=1}^K \sum_{s,h \leq H} q^*(s,h) \sum_{a \in \mathcal{A}} \pi_k(a|s,h) \frac{L'(\bar{x}_k(s,a,h) - \underline{x}_k(s,a,h)) + 2\theta L'}{\bar{x}_k(s,a,h) + \theta} + \tilde{\mathcal{O}} \left(\frac{T_\star L'}{\theta} \right).$$

This completes the proof. ■

Lemma 43 *With probability at least $1 - \delta$, $\text{BIAS}_2 = \tilde{\mathcal{O}}(T_\star L'/\theta)$.*

Proof By [Lemma 44](#) with $Z_k(s,a,h) = G_{k,s,a,h}/L'$ and $\mathbb{E}_k[G_{k,s,a,h}] = x_k(s,a,h)\bar{Q}_k(s,a,h)$, we have with probability at least $1 - \delta$:

$$\sum_{k=1}^K \tilde{Q}_k(s,a,h) - \bar{Q}_k(s,a,h) = \tilde{\mathcal{O}}(L'/\theta),$$

for any $(s,a) \in \mathcal{S} \times \mathcal{A}, h \leq H$. Therefore,

$$\text{BIAS}_2 = \sum_{k=1}^K \sum_{s,h \leq H} q^*(s,h) \sum_{a \in \mathcal{A}} \pi^*(a|s,h) (\tilde{Q}_k(s,a,h) - \bar{Q}_k(s,a,h)) = \tilde{\mathcal{O}}(T_\star L'/\theta). \quad \blacksquare$$

Lemma 44 For any random variable $Z_k(s, a, h)$ depending on interaction before episode k such that $Z_k(s, a, h) \in [0, 1]$, $\mathbb{E}_k[Z_k(s, a, h)] = z_k(s, a, h) \leq x_k(s, a, h)$, we have with probability at least $1 - \delta$:

$$\sum_{k=1}^K \left(\frac{Z_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta} - \frac{z_k(s, a, h)}{\bar{x}_k(s, a, h)} \right) \leq \frac{\ln \frac{1}{\delta}}{2\theta}.$$

Proof The statement is clearly true when $x_k(s, a, h) = 0$. When $x_k(s, a, h) > 0$, we also have $\bar{x}_k(s, a, h) > 0$. By $\frac{z}{1+z/2} \leq \ln(1+z)$ for $z \geq 0$, we have:

$$\begin{aligned} \frac{2\theta Z_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta} &\leq \frac{2\theta Z_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta Z_k(s, a, h)} = \frac{2\theta Z_k(s, a, h)/\bar{x}_k(s, a, h)}{1 + \theta Z_k(s, a, h)/\bar{x}_k(s, a, h)} \\ &\leq \ln(1 + 2\theta Z_k(s, a, h)/\bar{x}_k(s, a, h)). \end{aligned}$$

This gives

$$\begin{aligned} \mathbb{E}_k \left[\exp \left(\frac{2\theta Z_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta} \right) \right] &\leq \mathbb{E}_k \left[1 + \frac{2\theta Z_k(s, a, h)}{\bar{x}_k(s, a, h)} \right] = 1 + \frac{2\theta z_k(s, a, h)}{\bar{x}_k(s, a, h)} \\ &\leq \exp(2\theta z_k(s, a, h)/\bar{x}_k(s, a, h)). \quad (1 + z \leq e^z) \end{aligned}$$

Therefore, by Markov inequality,

$$\begin{aligned} P \left(\sum_{k=1}^K \frac{2\theta Z_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta} - \frac{2\theta z_k(s, a, h)}{\bar{x}_k(s, a, h)} > \ln \frac{1}{\delta} \right) \\ \leq \delta \cdot \mathbb{E} \left[\exp \left(\sum_{k=1}^K \frac{2\theta Z_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta} - \frac{2\theta z_k(s, a, h)}{\bar{x}_k(s, a, h)} \right) \right] \leq \delta. \end{aligned}$$

Thus, with probability at least $1 - \delta$, $\sum_{k=1}^K \frac{Z_k(s, a, h)}{\bar{x}_k(s, a, h) + \theta} - \frac{z_k(s, a, h)}{\bar{x}_k(s, a, h)} \leq \frac{\ln \frac{1}{\delta}}{2\theta}$. ■

D.3. Dilated Bonus in SDA

Below we present lemmas related to dilated bonus in $\mathring{\mathcal{M}}$. We first show that a form of dilated value function is well-defined.

Lemma 45 For some policy π in $\mathring{\mathcal{M}}$, transition $P \in \Lambda_{\mathcal{M}}$, and bonus function $b : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, \rho]$ for some $\rho > 0$, define $B(s, a, h) = b(s, a, h) + (1 + \frac{1}{H'}) P_{s, a, h} B$, $B(s, h) = \sum_a \pi(a|s, h) B(s, a, h)$ and $B(g) = B(s, a, H+1) = 0$. Then, $\max_{s, a} B(s, a, h) \leq \frac{15\rho(H-h+1)}{1-\gamma}$.

Proof Define $\gamma' = (1 + \frac{1}{H'})\gamma$ and recall that $H' = \frac{8(H+1)\ln(2K)}{1-\gamma}$. Now note that $\frac{1}{1-\gamma'} \leq \frac{1+\frac{1}{H'}}{1-\gamma}$ by simple algebra. Finally, define $\bar{b}(s, a, h) = (1 + \frac{1}{H'}) \langle P_{s, a, h}(\cdot, h+1), B(\cdot, h+1) \rangle$ for $h \leq H$, and $P'_{s, a, h}(s') = (1 + \frac{1}{H'}) P_{s, a, h}(s', h)$.

We prove that B is well defined and the statement holds by induction on $h = H + 1, \dots, 1$. The base case is true by definition $B(s, a, H + 1) = 0$. For $h \leq H$ we have:

$$\begin{aligned} B(s, a, h) &= b(s, a, h) + \left(1 + \frac{1}{H'}\right) (\langle P_{s,a,h}(\cdot, h), B(\cdot, h) \rangle + \langle P_{s,a,h}(\cdot, h+1), B(\cdot, h+1) \rangle) \\ &= b(s, a, h) + \bar{b}(s, a, h) + P'_{s,a,h} B(\cdot, h). \end{aligned}$$

Therefore, $B(\cdot, \cdot, h)$ can be treated as the action-value function in an SSP with cost $(b + \bar{b})(\cdot, \cdot, h)$ and transition function P' (thus well defined). By $\sum_{s'} P'_{s,a,h}(s', h) \leq \gamma'$, we have the expected hitting time of any policy starting from any state in an SSP with transition P' is upper bounded by $\frac{1}{1-\gamma'} \leq \frac{1+\frac{1}{H}}{1-\gamma}$. Let $R(h) = \max_{s,a} B(s, a, h)$ and note that $R(H + 1) = 0$. Since $b(s, a, h) \leq \rho$ and $\bar{b}(s, a, h) \leq (1 + \frac{1}{H'}) (1 - \gamma) R(h + 1)$ by $\sum_{s'} P_{s,a,h}(s', h + 1) \leq 1 - \gamma$, we have:

$$\begin{aligned} R(h) &\leq \frac{\rho + (1 + \frac{1}{H'}) (1 - \gamma) R(h + 1)}{1 - \gamma'} \leq \frac{\rho}{1 - \gamma'} + \left(1 + \frac{1}{H'}\right) \left(1 + \frac{1}{H}\right) R(h + 1) \\ &\leq \frac{\rho(1 + \frac{1}{H})}{1 - \gamma} + \left(1 + \frac{1}{H'}\right) \left(1 + \frac{1}{H}\right) R(h + 1), \end{aligned}$$

where the two last inequalities follow because $\frac{1}{1-\gamma'} \leq \frac{1+\frac{1}{H}}{1-\gamma}$. The proof is now finished by solving the recursion and obtaining:

$$R(h) \leq \frac{\rho}{1 - \gamma} \sum_{i=0}^{H-h} \left(1 + \frac{1}{H'}\right)^i \left(1 + \frac{1}{H}\right)^{i+1},$$

which implies that $R(h) \leq \frac{15\rho(H-h+1)}{1-\gamma}$ since $(1 + \frac{1}{H})^{H+1} (1 + \frac{1}{H'})^H \leq 2e^2 \leq 15$. \blacksquare

Lemma 46 *Let π be a policy in $\mathring{\mathcal{M}}$ and b be a non-negative cost function in $\mathring{\mathcal{M}}$ such that $b(s, a, H + 1) = 0$ and $b(s, a, h) \leq \rho$. Moreover, let $\hat{P} \in \Lambda_{\mathcal{M}}$ be an optimistic transition so that*

$$B(s, a, h) = b(s, a, h) + \left(1 + \frac{1}{H'}\right) \hat{P}_{s,a,h} B \geq b(s, a, h) + \left(1 + \frac{1}{H'}\right) P_{s,a,h} B,$$

where $B(s, h) = \sum_{a \in \mathcal{A}} \pi(a|s, h) B(s, a, h)$ and $B(g) = B(s, H + 1) = 0$. Then,

$$\begin{aligned} \sum_{s,h} q^*(s, h) \sum_{a \in \mathcal{A}} (\pi(a|s, h) - \hat{\pi}^*(a|s, h)) B(s, a, h) + \frac{1}{H'} \sum_{s,h} q^*(s, h) B(s, h) \\ + \sum_{s,a,h} q^*(s, a, h) b(s, a, h) \leq 3V^{\pi, \hat{P}, b}(s_{\text{init}}, 1) + \tilde{\mathcal{O}}\left(\frac{H\rho}{K(1-\gamma)}\right). \end{aligned}$$

Proof By the optimism property of \hat{P} , we have:

$$\sum_{s,h} q^*(s, h) \sum_{a \in \mathcal{A}} (\pi(a|s, h) - \hat{\pi}^*(a|s, h)) B(s, a, h)$$

$$\begin{aligned}
 & + \frac{1}{H'} \sum_{s,h} q^*(s,h) \sum_{a \in \mathcal{A}} \pi(a|s,h) B(s,a,h) + \sum_{s,a,h} q^*(s,a,h) b(s,a,h) \\
 \leq & \left(1 + \frac{1}{H'}\right) \sum_{s,h} q^*(s,h) \sum_{a \in \mathcal{A}} \pi(a|s,h) B(s,a,h) + \sum_{s,a,h} q^*(s,a,h) b(s,a,h) \\
 & - \sum_{s,a,h} q^*(s,a,h) \left(b(s,a,h) + \left(1 + \frac{1}{H'}\right) \sum_{s',h'} P_{s,a,h}(s',h') B(s',h') \right) \\
 = & \left(1 + \frac{1}{H'}\right) \sum_{s',h'} \left(q^*(s',h') - \sum_{s,a,h} q^*(s,a,h) P_{s,a,h}(s',h') \right) B(s',h') \\
 = & \left(1 + \frac{1}{H'}\right) B(s_{\text{init}}, 1). \tag{18}
 \end{aligned}$$

The last relation is by $q^*(s,h) - \sum_{s',a',h'} q^*(s',a',h') P_{s',a',h'}(s,h) = \mathbb{I}\{(s,h) = (s_{\text{init}}, 1)\}$ (see [Rosenberg and Mansour, 2021](#), Appendix B.1)).

Let J be the number of steps until the goal state g is reached in \mathcal{M} , and $n = \frac{8H}{1-\gamma} \ln(2K)$. Now note that for any policy, the expected hitting time in an SSP with transition \hat{P} is upper bounded by $\frac{H}{1-\gamma} + 1$ by $\hat{P} \in \Lambda_{\mathcal{M}}$. Therefore, by [Lemma 31](#), $P(J \geq n) \leq \frac{1}{K}$, and

$$\begin{aligned}
 B(s,h) &= \mathbb{E} \left[\sum_{t=1}^J \left(1 + \frac{1}{H'}\right)^{t-1} b(s_t, a_t, h_t) \middle| \pi, \hat{P}, (s_1, h_1) = (s, h) \right] \\
 &= \mathbb{E} \left[\sum_{t=1}^n \left(1 + \frac{1}{H'}\right)^{t-1} b(s_t, a_t, h_t) + \left(1 + \frac{1}{H'}\right)^n B(s_{t+1}, h_{t+1}) \middle| \pi, \hat{P}, (s_1, h_1) = (s, h) \right] \\
 &\leq \left(1 + \frac{1}{H'}\right)^{n-1} V^{\pi, \hat{P}, b}(s, h) + \tilde{O} \left(\frac{H\rho}{K(1-\gamma)} \right). \tag{Lemma 45}
 \end{aligned}$$

Plugging this back into [Eq. \(18\)](#) and by $(1 + 1/H')^n \leq e < 3$, we get the desired result. \blacksquare

D.4. Computation of B_k

We study an operator on value function, from which B_k can be computed as a fixed point. For any policy π , cost function c , transition confidence set $\mathcal{P} \subseteq \Lambda_{\mathcal{M}}$, and interest factor $\rho \geq 0$, we define the dilated Bellman operator \mathcal{T}_ρ that maps any value function $V : \mathring{S}_+ \rightarrow \mathbb{R}_+$ to another value function $\mathcal{T}_\rho V : \mathring{S}_+ \rightarrow \mathbb{R}_+$, such that:

$$\begin{aligned}
 (\mathcal{T}_\rho V)(s, h) &= \sum_a \pi(a|s, h) \left(c(s, a, h) + (1 + \rho) \max_{P \in \mathcal{P}} P_{s,a,h} V \right), \\
 (\mathcal{T}_\rho V)(g) &= 0, \quad (\mathcal{T}_\rho V)(s, H+1) = \max_a c(s, a, H+1). \tag{19}
 \end{aligned}$$

In this work, we have $\mathcal{P} \in \{\mathcal{P}_k\}_{k=1}^K$, and $\mathcal{P}_k = \bigcap_{s,a,h} \mathcal{P}_{k,s,a,h}$, where $\mathcal{P}_{k,s,a,h}$ is a convex set that specifies constraints on $((s, h), a)$. In other words, \mathcal{P}_k is a product of constraints on each $((s, h), a)$

(note that $\Lambda_{\mathcal{M}}$ can also be decomposed into shared constraints on $P_{s,a,H+1}$ and independent constraints on each $s, a, h \leq H$). Thus, there exists $P' \in \mathcal{P}$ that satisfies $P' = \operatorname{argmax}_{P \in \mathcal{P}} P_{s,a,h} V$ in Eq. (19) for all $((s, h), a)$ simultaneously. Moreover, finding such P' can be done by linear programming for each $((s, h), a)$ independently. Now we show that iteratively applying \mathcal{T}_ρ to some initial value function converges to a fixed point sufficiently fast.

Lemma 47 *Define value function $V^0 : \mathring{S}_+ \rightarrow \mathbb{R}_+$ such that $V^0(s, h) = V^0(g) = 0$ for any $(s, h) \in \mathcal{S} \times [H]$ and $V^0(s, H+1) = \max_a c(s, a, H+1)$. Then for any $\rho \geq 0$ such that $\gamma' = (1 + \rho)\gamma < 1$, the limit $V_\rho = \lim_{n \rightarrow \infty} \mathcal{T}_\rho^n V^0$ exists. Moreover, when $n \geq Hl$ with $l = \lceil \frac{\ln \frac{1}{\epsilon}}{1-\gamma'} \rceil$ for some $\epsilon > 0$, we have $\|\mathcal{T}_\rho^n V^0 - V_\rho\|_\infty \leq H \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^{H-1} \kappa \epsilon$, where $\kappa = \sum_{j=0}^{H-1} \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^j \frac{\|c\|_\infty}{1-\gamma'}$.*

Proof Define a sequence of value functions $\{V^i\}_{i=0}^\infty$ such that $V^{i+1} = \mathcal{T}_\rho V^i$. We first show that $\|V^i(\cdot, h)\|_\infty \leq \sum_{j=0}^{H-h} \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^j \frac{\|c\|_\infty}{1-\gamma'}$ for $i \geq 0$ and $h \leq H$. We prove this by induction on i . Note that this is clearly true when $i = 0$. For $i > 0$, by $\mathcal{P} \subseteq \Lambda_{\mathcal{M}}$ and Eq. (19), we have:

$$\begin{aligned} V^i(s, h) &= (\mathcal{T}_\rho V^{i-1})(s, h) \leq \|c\|_\infty + \gamma' \|V^{i-1}(\cdot, h)\|_\infty + (1+\rho)(1-\gamma) \|V^{i-1}(\cdot, h+1)\|_\infty \\ &\leq \|c\|_\infty + \gamma' \sum_{j=0}^{H-h} \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^j \frac{\|c\|_\infty}{1-\gamma'} + \sum_{j=1}^{H-h} \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^j \|c\|_\infty \\ &\leq \sum_{j=0}^{H-h} \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^j \frac{\|c\|_\infty}{1-\gamma'}. \end{aligned}$$

Therefore, $\|V^i\|_\infty \leq \kappa$. We now show that $\{V^i\}_i$ converges to a fixed point. Specifically, we show that for some $\epsilon > 0$ and any $i, j \in \mathbb{N}$, when $n \geq (H-h+1)l$, we have $\|(\mathcal{T}_\rho^n V^i)(\cdot, h) - (\mathcal{T}_\rho^n V^j)(\cdot, h)\|_\infty \leq (H-h+1) \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^{H-h} \kappa \epsilon$ (note that $\frac{(1+\rho)(1-\gamma)}{1-\gamma'} > 1$). Therefore, when $n \geq Hl$, we have $\|\mathcal{T}_\rho^n V^i - \mathcal{T}_\rho^n V^j\|_\infty \leq H \left(\frac{(1+\rho)(1-\gamma)}{1-\gamma'} \right)^{H-1} \kappa \epsilon$. Setting $\epsilon \rightarrow 0$, the statement above implies that for any $\mathring{s} \in \mathring{S}$, $\{V^i(\mathring{s})\}_{i=1}^\infty$ is a Cauchy sequence and thus converges. Moreover, letting $j \rightarrow \infty$ implies that $\{V^i\}_i$ converges to V_ρ with the rate shown above. We prove the statement above by induction on $h = H, \dots, 1$. First note that for any $s \in \mathcal{S}, h \in [H]$:

$$\begin{aligned} |(\mathcal{T}_\rho V^i)(s, h) - (\mathcal{T}_\rho V^j)(s, h)| &= (1+\rho) \left| \sum_a \pi(a|s, h) \left(\max_{P \in \mathcal{P}} P_{s,a,h} V^i - \max_{P \in \mathcal{P}} P_{s,a,h} V^j \right) \right| \\ &\leq (1+\rho) \sum_a \pi(a|s, h) \max_{P \in \mathcal{P}} |P_{s,a,h} (V^i - V^j)| \\ &\leq \gamma' \|V^i(\cdot, h) - V^j(\cdot, h)\|_\infty + (1+\rho)(1-\gamma) \|V^i(\cdot, h+1) - V^j(\cdot, h+1)\|_\infty, \end{aligned} \quad (20)$$

where the last inequality is by $\sum_{s'} P_{s,a,h}(s', h) \leq \gamma$, $\sum_{s'} P_{s,a,h}(s', h+1) \leq 1-\gamma$, and $P_{s,a,h}(s', h') = 0$ for $h' \notin \{h, h+1\}$, for any $P \in \Lambda_{\mathcal{M}}$. Now for the base case $h = H$, Eq. (20) implies $\|(\mathcal{T}_\rho V^i)(\cdot, H) - (\mathcal{T}_\rho V^j)(\cdot, H)\|_\infty \leq \gamma' \|V^i(\cdot, H) - V^j(\cdot, H)\|_\infty$. Thus for $n \geq l$, $\|(\mathcal{T}_\rho^n V^i)(\cdot, H) - (\mathcal{T}_\rho^n V^j)(\cdot, H)\|_\infty \leq \gamma'^n \cdot \kappa \leq \kappa \epsilon$. For the induction step $h < H$, if $n \geq (H-h+1)l$, then Eq. (20) implies:

$$|(\mathcal{T}_\rho^n V^i)(s, h) - (\mathcal{T}_\rho^n V^j)(s, h)|$$

$$\begin{aligned}
 &\leq \gamma'^l \left\| (\mathcal{T}_\rho^{n-l} V^i)(s, h) - (\mathcal{T}_\rho^{n-l} V^j)(s, h) \right\|_\infty + (1 - \gamma') \left(\frac{(1 + \rho)(1 - \gamma)}{1 - \gamma'} \right)^{H-h} \sum_{i=0}^{l-1} \gamma'^i (H - h) \kappa \epsilon \\
 &\hspace{15em} \text{(by the induction assumption)} \\
 &\leq (H - h + 1) \left(\frac{(1 + \rho)(1 - \gamma)}{1 - \gamma'} \right)^{H-h} \kappa \epsilon.
 \end{aligned}$$

This completes the proof of the statement above. \blacksquare

Now note that B_k is a fixed point of \mathcal{T}_ρ with $\pi = \pi_k$, $\mathcal{P} = \mathcal{P}_k$, $c = b_k$, and $\rho = 1/H'$. Thus, B_k can be approximated efficiently.

D.5. Computation of \bar{x}_k and \underline{x}_k

Note that $\bar{x}_k(s, a, h)$ can be computed by solving the following linear program (it is straightforward to verify that the constraints on π_q and P_q are linear):

$$\begin{aligned}
 &\max_{q \in \mathbb{R}_{\geq 0}^{S \times A \times [H] \times \hat{S}_+}} \sum_{s' \in \hat{S}_+} q(s, a, h, s') \\
 \text{s.t.} \quad &\sum_{a' \in \mathcal{A}, s' \in \mathcal{S}_+} q(s', a', h', s') \\
 &\quad - \sum_{(s'', h'') \in \hat{S}, a'' \in \mathcal{A}} q(s'', a'', h'', (s', h')) = \mathbb{I}\{(s', h') = (s_{\text{init}}, 1)\}, \forall (s', h') \\
 &\pi_q = \pi_k, \quad P_q \in \bigcap_{(s', a', h') \in (\mathcal{S} \times \mathcal{A} \times [H]) \setminus \{(s, a, h)\}} \mathcal{P}_{k, s', a', h'}, \quad P_{q, s, a, h}(g) = 1
 \end{aligned}$$

That is, we try to compute the occupancy measure that maximizes the number of visits to (s, a, h) in an augmented MDP, where the transition lies in \mathcal{P}_k except that taking action a at state (s, h) directly transits to the goal state (so that the number of visits to (s, a, h) is at most 1 and the occupancy measure at (s, a, h) is the probability of visiting (s, a, h)). The computation of $\underline{x}_k(s, a, h)$ is similar. Thus, both \bar{x}_k and \underline{x}_k can be computed efficiently (in a weakly polynomial time).

Appendix E. Learning without Some Parameters

In this section, we discuss the achievable regret guarantee without knowing some of the parameters assumed to be known. For simplicity, we only describe the high level ideas. We first describe the general ideas of dealing with each parameter being unknown, which are applicable under all types of feedback.

- **Unknown D and unknown fast policy:** we can simply follow the ideas in (Chen and Luo, 2021) to estimate D and fast policy. For unknown fast policy, we maintain an instance of Bernstein-SSP (Cohen et al., 2020) \mathcal{B}_f . When we need to switch to the fast policy, we simply involve \mathcal{B}_f as if this is a new episode for this algorithm, follow its decision until reaching g , and always feed cost 1 for all state-action pairs. Following the arguments in (Chen and Luo, 2021, Lemma 1), the scheme above only incurs constant extra regret. For unknown D , we maintain an estimate of it and update the algorithm's parameters whenever the estimate

is updated. Specifically, we separate the state space into known states and unknown states. A state is known if the number of visits to it is more than some threshold, and it is unknown otherwise. Whenever the learner visits an unknown state, it involves a Bernstein-SSP instance to approximate the behavior of fast policy until reaching g . When an unknown state s becomes known, we update the diameter estimate by incorporating an estimate of $T^{\pi^f}(s)$, and then updates the algorithm's parameters with respect to the new estimate. In terms of regret, this approach does not affect the transition estimation error, but brings an extra \sqrt{S} factor in the regret from policy optimization due to at most S updates to the algorithm's parameters.

- **Unknown B_\star :** We can estimate B_\star following the procedure in (Cohen et al., 2021, Appendix C). The main idea is pretty similar to the unknown D case: we again maintain an estimate of B_\star and separate states into known states and unknown states based on how many times a state has been visited. The learner updates algorithm's parameters whenever the estimate of B_\star is updated. Similarly, this approach brings an extra \sqrt{S} factor in the regret from policy optimization.
- **Unknown T_\star :** We can replace T_\star in parameters by B_\star/c_{\min} in stochastic costs setting and D/c_{\min} in other settings since $T_\star \leq B_\star/c_{\min}$ (or $T_\star \leq D/c_{\min}$). How to estimate D or B_\star is discussed above.
- **Unknown T_{\max} :** Similar to (Chen and Luo, 2021), we simply replace T_{\max} in parameters by K^p for some $p \in (0, \frac{1}{2})$.

Next, we describe under each setting, what regret guarantee we can achieve with each parameter being unknown by applying the corresponding method above.

Stochastic Costs In this setting, we need the knowledge of D , B_\star and T_{\max} .

- **Unknown D :** Since the regret from policy optimization is a lower order term, the dominating term of the final regret remains to be $\tilde{O}(B_\star S \sqrt{AK})$.
- **Unknown B_\star :** Since the regret from policy optimization is a lower order term, the dominating term of the final regret remains to be $\tilde{O}(B_\star S \sqrt{AK})$.
- **Unknown T_{\max} :** We replace T_{\max} in parameters by $K^{1/12}$. If $K^{1/12} \leq T_{\max}$, then clearly the regret is of order $\tilde{O}(LK) = \tilde{O}(T_{\max}^{13})$. Otherwise, by Theorem 5 we have $R_K = \tilde{O}(B_\star S \sqrt{AK} + S^4 A^{2.5} K^{1/3})$.

Stochastic Adversary In this setting, we need the knowledge of D , T_\star , and T_{\max} . We consider the following cases:

- **Unknown D :** Since the regret from policy optimization is a lower order term, the dominating term of the final regret remains to be $\tilde{O}(\sqrt{DT_\star K} + DS\sqrt{AK})$ in the full information setting, and $\tilde{O}(\sqrt{DT_\star SAK} + DS\sqrt{AK})$ in the bandit feedback setting.
- **Unknown T_\star :** Ignoring the lower order terms, we have $R_K = \tilde{O}(D\sqrt{K/c_{\min}} + DS\sqrt{AK})$ in the full information setting by Theorem 6, and $R_K = \tilde{O}(D\sqrt{SAK/c_{\min}} + DS\sqrt{AK})$ in the bandit feedback setting by Theorem 7.

- **Unknown T_{\max} :** We replace T_{\max} in parameters by $K^{1/13}$. If $K^{1/13} \leq T_{\max}$, then R_K is of order $\tilde{O}(LK) = \tilde{O}(T_{\max}^{14})$. Otherwise, we have $R_K = \tilde{O}(\sqrt{DT_{\star}K} + DS\sqrt{AK} + (S^2A^3)^{1/4}K^{25/52} + S^4A^{2.5}K^{4/13})$ in the full information setting by [Theorem 6](#), and $R_K = \tilde{O}(\sqrt{SADT_{\star}K} + DS\sqrt{AK} + SA^{5/4}K^{25/52} + S^4A^{2.5}K^{4/13})$ in the bandit feedback setting by [Theorem 7](#).

Adversarial Costs, Full Information In this setting, we need the knowledge of D , T_{\star} , and T_{\max} . We consider the following cases:

- **Unknown D :** With an extra \sqrt{S} factor in the policy optimization term, we have $R_K = \tilde{O}(T_{\star}\sqrt{SDK} + \sqrt{S^2ADT_{\star}K})$ ignoring the lower order terms.
- **Unknown T_{\star} :** Ignoring the lower order terms, we have $R_K = \tilde{O}(\frac{D^{1.5}}{c_{\min}}\sqrt{K} + D\sqrt{S^2AK/c_{\min}})$.
- **Unknown T_{\max} :** We replace T_{\max} in parameters by $K^{1/11}$. If $K^{1/11} \leq T_{\max}$, then clearly the regret is of order $\tilde{O}(LK) = \tilde{O}(T_{\max}^{12})$. Otherwise, by [Theorem 9](#) we have $R_K = \tilde{O}(T_{\star}\sqrt{DK} + \sqrt{S^2ADT_{\star}K} + S^4A^2K^{5/11})$.

Adversarial Costs, Bandit Feedback In this setting, we need the knowledge of D and T_{\max} . We consider the following cases:

- **Unknown D :** Tracing the proof of [Theorem 10](#), the regret from policy optimization is of order $\tilde{O}(\sqrt{SAT_{\max}^4K})$. With an extra \sqrt{S} factor in the policy optimization term, we still have $R_K = \tilde{O}(\sqrt{S^2AT_{\max}^5K})$ ignoring the lower order terms.
- **Unknown T_{\max} :** We replace T_{\max} in parameters by K^p for any $p \in (0, \frac{1}{5})$. If $K^p \leq T_{\max}$, then clearly the regret is of order $\tilde{O}(LK) = \tilde{O}(T_{\max}^{1+1/p})$. Otherwise, by [Theorem 10](#) we have $R_K = \tilde{O}(\sqrt{S^2AK^{1+5p}} + S^{5.5}A^{3.5}K^{5p})$.

Appendix F. Auxiliary Lemma

Lemma 48 *If $x \leq (a\sqrt{x} + b)\ln^p(cx)$ for some $a, b, c > 0$ and absolute constant $p \geq 0$, then $x = \tilde{O}(a^2 + b)$. Specifically, $x \leq a\sqrt{x} + b$ implies $x \leq (a + \sqrt{b})^2 \leq 2a^2 + 2b$.*

Lemma 49 ([Luo et al., 2021](#), Lemma A.4) *Let $\eta > 0$, $\pi_k \in \Delta(A)$, and $\ell_k \in \mathbb{R}^A$ satisfy the following for all $k \in [K]$ and $a \in A$:*

$$\pi_1(a) = \frac{1}{A}, \quad \pi_{k+1}(a) \propto \pi_k(a) \exp(-\eta \ell_k(a)), \quad |\eta \ell_k(a)| \leq 1.$$

Then for any $\pi^ \in \Delta(A)$, $\sum_{k=1}^K \langle \pi_k - \pi^*, \ell_k \rangle \leq \frac{\ln A}{\eta} + \eta \sum_{k=1}^K \sum_{a \in A} \pi_k(a) \ell_k^2(a)$.*

Lemma 50 ([Chen et al., 2021b](#), Lemma 38) *Let $\{X_i\}_{i=1}^{\infty}$ be a martingale difference sequence adapted to the filtration $\{\mathcal{F}_i\}_{i=0}^{\infty}$ and $|X_i| \leq B$ for some $B > 0$. Then with probability at least $1 - \delta$, for all $n \geq 1$ simultaneously,*

$$\left| \sum_{i=1}^n X_i \right| \leq 3 \sqrt{\sum_{i=1}^n \mathbb{E}[X_i^2 | \mathcal{F}_{i-1}] \ln \frac{4B^2 n^3}{\delta}} + 2B \ln \frac{4B^2 n^3}{\delta}.$$

Lemma 51 (*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 and $X_n \in [0, B]$ almost surely. Then with probability at least $1 - \delta$, for any $n \geq 1$:

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

Lemma 52 (*Cohen et al., 2020, Lemma D.4*) and (*Cohen et al., 2021, Lemma C.2*) Let $\{X_i\}_{i=1}^\infty$ be a sequence of random variables w.r.t to the filtration $\{\mathcal{F}_i\}_{i=0}^\infty$ and $X_i \in [0, B]$ almost surely. Then with probability at least $1 - \delta$, for all $n \geq 1$ simultaneously:

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}[X_i | \mathcal{F}_{i-1}] &\leq 2 \sum_{i=1}^n X_i + 4B \ln \frac{4n}{\delta}, \\ \sum_{i=1}^n X_i &\leq 2 \sum_{i=1}^n \mathbb{E}[X_i | \mathcal{F}_{i-1}] + 8B \ln \frac{4n}{\delta}. \end{aligned}$$