
Tiered Reinforcement Learning: Pessimism in the Face of Uncertainty and Constant Regret

Jiawei Huang^{1*} Li Zhao² Tao Qin² Wei Chen² Nan Jiang¹ Tie-Yan Liu²

¹ Department of Computer Science, University of Illinois at Urbana-Champaign
{jiaweih, nanjiang}@illinois.edu

² Microsoft Research Asia
{lizo, taoqin, weic, tyliu}@microsoft.com

Abstract

We propose a new learning framework that captures the tiered structure of many real-world user-interaction applications, where the users can be divided into two groups based on their different tolerance on exploration risks and should be treated separately. In this setting, we simultaneously maintain two policies π^O and π^E : π^O (“O” for “online”) interacts with more risk-tolerant users from the first tier and minimizes regret by balancing exploration and exploitation as usual, while π^E (“E” for “exploit”) exclusively focuses on exploitation for risk-averse users from the second tier utilizing the data collected so far. An important question is whether such a separation yields advantages over the standard online setting (i.e., $\pi^E = \pi^O$) for the risk-averse users. We individually consider the gap-independent vs. gap-dependent settings. For the former, we prove that the separation is indeed not beneficial from a minimax perspective. For the latter, we show that if choosing Pessimistic Value Iteration as the exploitation algorithm to produce π^E , we can achieve a constant regret for risk-averse users independent of the number of episodes K , which is in sharp contrast to the $\Omega(\log K)$ regret for any online RL algorithms in the same setting, while the regret of π^O (almost) maintains its online regret optimality and does not need to compromise for the success of π^E .

1 Introduction

Reinforcement learning (RL) has been applied to many real-world user-interaction applications to provide users with better services, such as in recommendation systems [Afsar et al., 2021] and medical treatment [Yu et al., 2021, Lipsky and Sharp, 2001]. In those scenarios, the users take the role of the environments and the interaction strategies (e.g. recommendation or medical treatment) correspond to the agents in RL. In the theoretical study of such problems, most of the existing literature adopts the online interaction protocol, where in each episode $k \in [K]$, the learning agent executes a policy π_k to interact with users (i.e. environments), receives new data to update the policy, and moves on to the next episode. While this formulation *treats each user equivalently* when optimizing the regret, many scenarios have a special “**Tiered Structure**”²: *users can be divided into multiple groups depending on their different preference and tolerance about the risk that results from the necessary exploration to improve the policy*, and such grouping is available to the learner in advance so it would be better to treat them separately. As a concrete example, in medical treatment, after a new treatment plan comes out, some courageous patients or paid volunteers (denoted as G^O ; “O” for “Online”) may prefer it given the potential risks, while some conservative patients (denoted as G^E ; “E” for “Exploit”) may tend to receive mature and well-tested plans, even if the new one is promising to be more effective.

*Work done during the internship at Microsoft Research Asia.

²We consider the cases with two tiers in this paper.

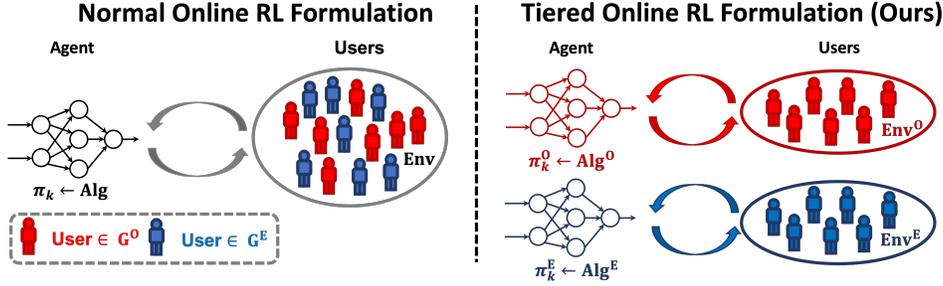


Figure 1: Comparison between the standard setting and our tiered RL setting (#Tiers = 2), where we use red and blue to color users from different groups. The main difference is that, in the standard setting (LHS), the learner does not distinguish users from different groups and treats them equivalently with a single policy π_k produced by algorithm Alg, while in our setting (RHS), we leverage the tier information and interact with different groups with different policies π^E and π^O .

As another example, companies offering recommendation services may recruit paid testers or use bonus to attract customers (G^O) to interact with the system to shoulder the majority of the exploration risk during policy improvement, which may result in better service (low regret) for the remaining customers (G^E). Moreover, many online platforms have free service open for everyone (G^O), while some users are willing to pay for enhanced service (G^E). If we follow the traditional online setting and treat the users in these two groups equivalently, then in expectation each group will suffer the same regret and risk. In contrast, if we leverage the group information by using policies with different risk levels to interact with different groups, it is potentially possible to transfer some exploration risks from users in G^E to G^O , while the additional risks suffered by G^O will be compensated in other forms (such as payment, the users' inherent motivation, or the free service itself).

To make our objective more clear, we abstract the problem setting into Frw. 1 and compare it with the standard online setting in Fig. 1, where we use Alg^O and Alg^E to denote the two algorithms producing policies π^E and π^O to interact with users in G^O and G^E , respectively. To enable theoretical analyses, we do adopt a few simplification assumptions while still modelling the core challenges in the aforementioned scenarios: firstly, at each iteration of Frw. 1, the algorithms will interact with and collect one trajectory from each group, which assumes that users from two groups will come to seek for service in pair with the same frequency. In practice, usually, the users come in random order and the frequencies from different groups are not the same; see Appx. B for why our abstraction is still a valid surrogate and how our results can be generalized. Secondly, for convenience, *we only use the samples generated from G^O* , because Alg^E is expected to best exploit the available information and not encouraged to perform intelligent exploration. Nonetheless, our results hold with minor modifications if one also chooses to use trajectories from G^E . Thirdly, we assume that the dynamics and rewards during the interactions with users in different tiers can behave differently, and we leave the relaxation of such an assumption to future work.

Framework 1: The Tiered RL Framework

- 1 Input: Env^O and Env^E // Note that $\text{Env}^O = \text{Env}^E$
 - 2 Initialize $D_1 \leftarrow \{\}$.
 - 3 **for** $k = 1, 2, \dots, K$ **do**
 - 4 $\pi_k^O \leftarrow \text{Alg}^O(D_k); \pi_k^E \leftarrow \text{Alg}^E(D_k)$.
 - 5 π_k^O interacts with customers/users/patients in G^O (i.e. Env^O), and collect data τ_k^O .
 - 6 π_k^E interacts with customers/users/patients in G^E (i.e. Env^E), and collect data τ_k^E .
 - 7 $D_{k+1} = D_k \cup \{\tau_k^O\}$. // We do not consider to use τ_k^E in this paper.
 - 8 **end**
-

Similar to the online setting, we use the expected pseudo-regret to measure the performance of the algorithm, which is formalized in Def. 2.1. The key problem we would like to investigate is provable benefits of leveraging the tiered structure by Frw. 1 comparing with the standard online setting:

Is it possible for $\text{Regret}(\text{Alg}^E)$ to be strictly lower than any online learning algorithms in certain scenarios, while keeping $\text{Regret}(\text{Alg}^O)$ near-optimal?

Note that we still expect regret of Alg^O to enjoy near-optimal regret guarantees, which is a reasonable requirement as the experience of users in G^O also matters in many of our motivating applications. We regard the above problem formulation as **our first contribution**, which is mainly conceptual.

As **our second contribution**, Sec. 3 shows that Alg^E has the same minimax gap-independent lower bound as online learning algorithms. This result reveals the difficulty to leverage tiered structure in standard tabular MDPs, and motivates us to investigate the benefits under the gap-dependent setting, which is frequently considered in the Multi-Armed Bandit (MAB) [Lattimore and Szepesvári, 2020, Rouyer and Seldin, 2020] and RL literature [Xu et al., 2021, Simchowitz and Jamieson, 2019].

As **our third contribution** and our main technical results, Sec. 4 establishes provable benefits of Frw. 1 by proposing a new algorithmic framework and showing $\text{Regret}(\text{Alg}^E)$ is constant and independent of the number of episodes K , which is in sharp contrast with the $\Omega(\log K)$ regret lower bound for any algorithms in the standard online setting that do not leverage the tiered structure. Specifically, we use Pessimistic Value Iteration (PVI) as Alg^E for exploitation to interact with G^E , while Alg^O can be arbitrary online algorithms with near-optimal regret. Concretely, we first study stochastic MABs as a warm-up, where we choose Alg^E to be LCB (Lower Confidence Bonus), a degenerated version of PVI in bandits, and choose UCB (Upper Confidence Bonus) as Alg^O for a concrete case study. We prove that Alg^E can achieve constant pseudo-regret $\Theta\left(\frac{1}{\Delta_1} - \frac{\Delta_i}{\Delta_1^2}\right)$ with A being the number of actions and Δ_i 's being the gaps with $\Delta_1 \geq \dots \geq \Delta_{A-1} \geq \Delta_A = 0$, while Alg^O is near-optimal due to the regret guarantee of UCB. After that, Sec. 4.2 extends the success of PVI to tabular MDPs, and establishes results that apply to a wide range of online algorithms Alg^O with near-optimal regret. Although the benefits of pessimism have been widely recognized in offline RL [Jin et al., 2021], to our knowledge, we are the first to study PVI in a gap-dependent online setting. We also contribute several novel techniques for overcoming the difficulties in achieving constant regret, and defer their summary to Sec. 4.2. Moreover, in Appx.H, we report experiment results to demonstrate the advantage of leveraging tiered structure as predicted by theory.

Closest Related Work Due to space limit, we only discuss the closest related work here and defer the rest to Appx. A. To our knowledge, there is no previous works on leveraging tiered structure in MDPs. In the bandit setting, there is a line of related works studying decoupling exploration and exploitation [Avner et al., 2012, Rouyer and Seldin, 2020], where [Rouyer and Seldin, 2020] studied “best of both worlds” methods and reported a similar constant regret. First, in stochastic bandits, there are many cases when our result is tighter than theirs, (see a detailed comparison in Sec. 4.1), and more importantly, our methods can naturally extend to RL (i.e., MDPs), whereas a similar extension of their techniques can run into serious difficulties: they relied on *importance sampling* to provide unbiased off-policy estimation for policy value, which incurs the infamous “curse of horizon”, a.k.a., a sample complexity *exponential* in the planning horizon H in long-horizon RL (see examples in Sec. 2 in [Liu et al., 2018]). Our approach overcomes this difficulty by developing a pessimism-based learning framework, which is fundamentally different from their approach and requires several novel techniques in the analyses. Second, they did not provide any guarantee for the regret of exploration algorithm, whereas in our results the regret of Alg^O can be near-optimal, which we believe is also important as the experience of users in G^O also matters in many of our motivating applications. Third, their bandit results require a unique best arm, whereas we allow the optimal arms/policies to be non-unique, which can cause non-trivial difficulties in the analyses as we will discuss in Sec. 4.2.3.

2 Preliminary and Problem formulation

Stochastic Multi-Armed Bandits (MABs) The MAB model consists of a set of arms $\mathcal{A} = \{1, 2, \dots, A\}$. When sampling an arm $i \in \mathcal{A}$, the agent observes a random variable $r_i \in [0, 1]$. We use $\mu_i = \mathbb{E}[r_i]$ to denote the mean value for arm i for each $i \in \mathcal{A}$. We allow the optimal arms to be non-unique. For simplicity of notation, we assume the arms are ordered such that $\mu_1 \leq \mu_2 \dots \leq \mu_A$.

Finite-Horizon Tabular Markov Decision Processes (MDPs) For the reinforcement learning (RL) setting, we consider the episodic tabular MDPs denoted by $M(\mathcal{S}, \mathcal{A}, H, P, r)$, where \mathcal{S} is the finite

state space, \mathcal{A} is the finite action space, H is the horizon length, and $\mathbb{P} = \{\mathbb{P}_h\}_{h=1}^H$ and $r = \{r_h\}_{h=1}^H$ are the time-dependent transition and reward functions, respectively. We assume all steps share the state and action space (i.e. $\mathcal{S}_1 = \mathcal{S}_2 \dots = \mathcal{S}_H = \mathcal{S}$, $\mathcal{A}_1 = \mathcal{A}_2 \dots = \mathcal{A}_H = \mathcal{A}$) while the transition and reward functions can be different. At the beginning of each episode, the environment will start from a fixed initial state s_1 (w.l.o.g.). Then, for each time step $h \in [H]$, the agent selects an action $a_h \in \mathcal{A}$ based on the current state s_h , receives the reward $r_h(s_h, a_h)$, and observes the system transition to the next state s_{h+1} , until s_{H+1} is reached. W.l.o.g., we assume the reward function r is deterministic and our results can be easily extended to handle stochastic rewards.

A time-dependent policy is specified as $\pi = \{\pi_1, \pi_2, \dots, \pi_H\}$ with $\pi_h : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ for all $h \in [H]$. Here $\Delta(\mathcal{A})$ denotes the probability simplex over the action space. With a slight abuse of notation, when π_h is a deterministic policy, we use $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$ to refer to a deterministic mapping. $V_h^\pi(s)$ and $Q_h^\pi(s, a)$ denote the value function and Q-function at step $h \in [H]$, which are defined as: $V_h^\pi(s) = \mathbb{E}[\sum_{h^0=h}^H r_{h^0}(s_{h^0}, a_{h^0}) | s_h = s, \pi]$, $Q_h^\pi(s, a) = \mathbb{E}[\sum_{h^0=h}^H r_{h^0}(s_{h^0}, a_{h^0}) | s_h = s, a_h = a, \pi]$.

We use $V_h^*(\cdot) := \max_\pi V_h^\pi(\cdot)$ and $Q_h^*(\cdot, \cdot) := \max_\pi Q_h^\pi(\cdot, \cdot)$ to refer to the optimal state/action-value functions, and $\Pi^*(s_h) := \{a_h | Q_h^*(s_h, a_h) = V_h^*(s_h)\}$ to denote the collection of all optimal actions at state s_h . With an abuse of notation, we define $\Pi^* := \{\pi : V_1^\pi(s_1) = V_1^*(s_1)\}$, i.e., the set of policies that maximize the total expected return. In this paper, when we say that the MDP has “unique optimal (deterministic) policy”, it is up to the occupancy measure, that is, all policies in Π^* share the same state-action occupancy $d^\pi(s_h, a_h) := \Pr(S_h = s_h, A_h = a_h | S_1 = s_1, \pi)$ for all $h \in [H]$, $s_h \in \mathcal{S}_h$, $a_h \in \mathcal{A}_h$. In the following, we use $|\Pi^*| = 1$ to refer to the case of unique optimal (deterministic) policy, where the cardinality of Π^* is counted up to the equivalence of occupancies. Besides, for any function $V : \mathcal{S} \rightarrow \mathbb{R}$, we denote $\mathbb{P}_h V(s_h, a_h) := \mathbb{E}_{s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, a_h)}[V(s_{h+1})]$.

Gap-Dependent Setting We follow the standard formulation of gap-dependent setting in previous bandits [Lattimore and Szepesvári, 2020] and RL literature [Simchowitz and Jamieson, 2019, Xu et al., 2021, Dann et al., 2021]. In bandits, the gap w.r.t. arm i is defined as $\Delta_i := \max_{j \in [A]} \mu_j - \mu_i, \forall i \in [A]$, and we assume that there exists a strictly positive value Δ_{\min} such that, either $\Delta_i = 0$ or $\Delta_i \geq \Delta_{\min}$. For tabular RL setting, we define $\Delta_h(s_h, a_h) := V^*(s_h) - Q^*(s_h, a_h), \forall h \in [H], s_h \in \mathcal{S}_h, a_h \in \mathcal{A}_h$. We use the same notation Δ_{\min} to refer to the minimal gap in tabular setting and assume that either $\Delta_h(s_h, a_h) = 0$ or $\Delta_h(s_h, a_h) \geq \Delta_{\min}$.

Performance Measure We use Pseudo-Regret defined below to measure the performance of Alg^O and Alg^E . In the following, we will also use “exploitation regret” to refer to $\text{Regret}_K(\text{Alg}^E)$.

Definition 2.1 (Pseudo-Regret). *We define the regret of Alg^O and Alg^E to be:*

$$\text{Regret}_K(\text{Alg}^O) := \mathbb{E} \sum_{k=1}^K V_1^*(s_1) - V_1^{\pi_k^O}(s_1) \quad ; \quad \text{Regret}_K(\text{Alg}^E) := \mathbb{E} \sum_{k=1}^K V_1^*(s_1) - V_1^{\pi_k^E}(s_1) \quad ,$$

where π_k^O and π_k^E are generated according to the procedure in Framework 1 and the expectation is taken over the randomness in data generation and algorithms.

3 Lower Bound of Regret(Alg^E) without Gap Assumption

In this section, we show that, in normal tabular RL setting, for arbitrary algorithm pair $(\text{Alg}^O, \text{Alg}^E)$, even if we do not constrain Alg^O to be near-optimal, the regret of Alg^E has the same minimax lower bound as algorithms in online setting. We defer the formal statement and proof to Appendix C.1.

Theorem 3.1. [Lower Bound for Alg^E without Gap Assumption] *There exist positive constants $c, \varepsilon_0, \delta_0$, such that, for arbitrary $S \geq 4, A \geq 2, H \geq 2, K \geq \frac{c}{\varepsilon_0^2} H^3 S A$, and arbitrary algorithm pair $(\text{Alg}^O, \text{Alg}^E)$, there must exist a hard tabular MDP $M_{\text{hard}}, \mathbb{E}_{(\text{Alg}^O, \text{Alg}^E), M_{\text{hard}}} \sum_{k=1}^K V^* - V^{\pi_k^E} \geq \delta_0 \sqrt{c H^3 S A K}$, where the expectation is taken over the randomness of algorithms and MDP.*

The theorem above is stating that, comparing with the regret lower bound for online algorithms $\Theta(\sqrt{H^3 S A K})$ in Theorem 9 of Domingues et al. [2021], the exploitation algorithm cannot reduce the dependence on any of parameters H, S, A, K in hard MDPs, even if we allow Alg^O to sacrifice its performance to gather the best possible data for Alg^E . Also, the lower bound would still hold

even if we allow both Alg^{O} and Alg^{E} to additionally use the data τ^{E} generated by π^{E} . This negative result implies that without any further assumptions, the separation is not beneficial from a minimax optimality perspective, and we can simply choose both Alg^{E} and Alg^{O} to be the same near-optimal online algorithm as without worrying about separating them.

However, in the next section, we will show that, in tabular MDPs with strictly positive gaps, in contrast with the $\Omega(\log K)$ lower bound for online algorithms, we can have Alg^{E} such that its regret is constant and independent on the number of time horizon K , which reveals the fundamental differences between the pure online setting and the Tiered RL setting considered in this paper.

4 Pessimism in the Face of Uncertainty and Constant Regret

In this section, we consider the gap-dependent setting and contribute to identifying the possibility to achieve constant regret by using pessimistic algorithms for Alg^{E} . Intuitively, the main reason why PVI can lead to a constant regret is that the quality of the policy returned by PVI is positively correlated to the accumulation of optimal trajectories in the dataset D , which is directly connected with $\text{Regret}(\text{Alg}^{\text{O}})$. As a result, on the one hand, the regret minimization objective of Alg^{E} coincidentally aligns with the optimality constraint of Alg^{O} . On the other hand, thanks to the positive gap assumption, π^{E} will gradually converge to the optimal policy with high probability when Alg^{E} is PVI, so there will be no regret after that. In Sec. 4.1, we start with stochastic MAB as a warm-up, and in Sec. 4.2 we extend our success to tabular RL setting. We defer the proofs in this section to Appendix D.

4.1 Warm-Up: Gap-Dependent Regret Bound for Stochastic Multi-Armed Bandits

Algorithm 2: UCB-Exploration-LCB-Exploitation

```

1 Initialize:  $\alpha > 1$ ;  $N_i(1) \leftarrow 0$ ,  $\hat{\mu}_i(1) \leftarrow 0$ ,  $\forall i \in \mathcal{A}$ ;  $f(k) := 1 + 16A^2(k+1)^2$ 
2 for  $k = 1, 2, \dots, K$  do
3    $\pi_k^{\text{O}} \leftarrow \arg \max_i \hat{\mu}_i(k) + \mathcal{O}\left(\frac{2\alpha \log f(k)}{N_i(k)}\right)$ ,  $\pi_k^{\text{E}} \leftarrow \arg \max_i \hat{\mu}_i(k) - \mathcal{O}\left(\frac{2\alpha \log f(k)}{N_i(k)}\right)$ .
4   Interact with  $G^{\text{E}}$  and  $G^{\text{O}}$  by  $\pi_k^{\text{E}}$  and  $\pi_k^{\text{O}}$ , and observe reward  $r(\pi_k^{\text{E}})$  and  $r(\pi_k^{\text{O}})$ , respectively.
5   for  $i = 1, 2, \dots, A$  do
6      $N_i(k+1) \leftarrow N_i(k) + \mathbb{1}[\pi_k^{\text{O}} = i]$ ;  $\hat{\mu}_i(k+1) \leftarrow \hat{\mu}_i(k) \frac{N_i(k)}{N_i(k+1)} + r(\pi_k^{\text{O}}) \frac{\mathbb{1}[\pi_k^{\text{O}} = i]}{N_i(k+1)}$ .
7   end
8 end

```

Our main algorithm for bandit setting is shown in Alg 1, where we consider the UCB algorithm [Lattimore and Szepesvári, 2020] as Alg^{O} and choose the LCB as Alg^{E} , which flips the sign of the bonus term in UCB. We use $N_i(k)$ to denote the number of times that arm i was pulled previous to step k , and use $\hat{\mu}_i$ to record the empirical average of arm i . Besides, we assume $1/N_i(\cdot) = +\infty$ if $N_i(\cdot) = 0$, which implies that at the first $|\mathcal{A}|$ steps the algorithm will pull each arm one by one. Moreover, as we will show later, the choice of $\alpha > 1$ is crucial to avoiding dependence on K in $\text{Regret}(\text{Alg}^{\text{E}})$ with our techniques. For Alg. 2, we have the following guarantee:

Theorem 4.1. [Exploitation Regret] *In Algorithm 2, by choosing arbitrary $\alpha > 1$, there exists an absolute constant c , such that, for arbitrary $K \geq 1$, the pseudo-regret of Alg^{E} is upper bounded by:*

$$\text{Regret}_K(\text{Alg}^{\text{E}}) \leq \Theta\left(\frac{A}{\alpha-1} + \alpha \prod_{\Delta_i > 0} (A-i) \frac{1}{\Delta_i} - \frac{\Delta_i}{\Delta_i^2 - 1}\right) \quad \text{where } \Delta_0 := \infty \text{ so } \frac{\Delta_1}{\Delta_0^2} = 0.$$

Our result implies that by choosing PVI as Alg^{E} , we can achieve constant regret while keeping Alg^{O} near-optimal. Besides the advantages discussed in the related work paragraph in Sec. 1, our guarantee is also more favorable in certain cases compared to the $O\left(\frac{A}{\Delta_{\min}} \prod_{\Delta_i > 0} \frac{1}{\Delta_i}\right)$ result in Rouyer and Seldin [2020]: while it is not easy to verify whether our guarantee dominates theirs, in many cases ours can be strictly better (or at least no worse) than theirs. For example, consider the following two representative cases: $\Delta_1 = \Delta_2 = \dots \Delta_{A-1} = \Delta_{\min}$ (uniform gap) and $\Delta_1 = \Delta_2 = \dots = \Delta_{A-2} \gg \Delta_{A-1} = \Delta_{\min}$ (small last gap); our result achieves $\Theta\left(\frac{A}{\Delta_{\min}}\right)$ and $\Theta\left(\frac{1}{\Delta_{\min}}\right)$, respectively, in contrast to their $\Theta\left(\frac{A}{\Delta_{\min}}\right)$ and $\Theta\left(\frac{\sqrt{A}}{\Delta_{\min}}\right)$.

Proof Sketch: The proof consists of two novel technique lemmas with a carefully chosen failure rate $\delta_k \sim O(1/k^{\Theta(\alpha)})$ so that the accumulative failure probability $\sum_{k=1}^{\infty} \delta_k < +\infty$. The first one is Lem. 4.2, where we show that w.p. $1 - \delta_k$, LCB will not prefer i with $\Delta_i > 0$ as long as another better arm has been visited enough times in the dataset. The second step is to identify a key property of UCB algorithm as stated in Lem. 4.3, where we provide a high probability upper bound that $N_i(k) \leq k/\lambda$ if $k \geq \Theta(\lambda/\Delta_i^2)$ for arbitrary $\lambda \in [1, 4A]$, and it serves to indicate that the condition required by the success of LCB is achievable as long as k is large enough³.

Lemma 4.2. *[Blessing of Pessimism] With the choice that $f(k) = 1 + 16A^2(k+1)^2$, for arbitrary i with $\Delta_i > 0$, for the LCB algorithm in Alg 2, and arbitrary j satisfying $\Delta_j < \Delta_i$, we have:*

$$\Pr \{i = \pi_k^E\} \cap \{\Delta_j < \Delta_i\} \cap \{N_j(k) \geq \frac{8\alpha \log f(k)}{(\Delta_j - \Delta_i)^2}\} \leq \frac{2}{k^2}.$$

Lemma 4.3. *[Property of UCB] With the choice that $f(k) = 1 + 16A^2(k+1)^2$, there exists a constant c , for arbitrary i with $\Delta_i > 0$ and arbitrary $\lambda \in [1, 4A]$, in UCB algorithm, we have:*

$$\Pr(N_i(k) \geq \frac{k}{\lambda}) \leq \frac{2}{k^2 - 1}, \quad \forall k \geq \lambda + c \cdot \frac{\alpha \lambda}{\Delta_i^2} \log(1 + \frac{\alpha A}{\Delta_{\min}}).$$

Directly combining the above two results, we can obtain an upper bound for $\text{Regret}(\text{Alg}^E)$ of order $\Theta(A/\Delta_{\min}^{-2})$, which is already independent of K . To achieve better dependence on Δ_{\min} in the regret, we conduct a finer analysis. For each arm i with $\Delta_i > 0$, we separate all the arms including i into two groups based on whether its gap exceeds $\Delta_i/2$: $G_i^{\text{lower}} = \{j : \Delta_j > \Delta_i/2\}$ and $G_i^{\text{upper}} = \{j : \Delta_j \leq \Delta_i/2\}$. As a result of Lem. 4.2, we know that π_k^E will not prefer arm i as long as there exists $j \in G_i^{\text{upper}}$ such that $N_j(k) = \Omega(4\Delta_i^{-2}) = \Omega(\Delta_i^{-2})$. Based on Lem. 4.3, we know it is true with high probability, as long as $k \geq \Theta(A \cdot \Delta_i^{-2})$, since at that time $N_l(k) \leq k/A$ holds for arbitrary $l \in G_i^{\text{lower}}$, which directly implies that $\max_{j:j \in G_i^{\text{upper}}} N_j(k) \geq \Omega(\Delta_i^{-2})$. Then, combining Lem. 4.2, with high probability, the regret resulting from taken arm i cannot be higher than $\Theta(A \cdot \Delta_i^{-2}) \cdot \Delta_i = \Theta(A \cdot \Delta_i^{-1})$, which results in a $\Theta(\sum_{\Delta_i > 0} A/\Delta_i)$ regret bound. As for the techniques leading to the further improvement in our final result, please refer to Lem. D.1 and the proof of Thm. 4.1 in Appx. D.

4.2 Constant Regret of Alg^E in Tabular MDPs

In this section, we establish constant regret of Alg^E based on realistic conditions for Alg^O and Alg^E . We highlight the key steps of our analysis and our technical contributions here.

First of all, in Sec.4.2.1, we propose the concrete PVI algorithm, and inspired by the clipping trick used for optimistic online algorithms [Simchowitz and Jamieson, 2019], we develop a high-probability gap-dependent upper bound for the sub-optimality of π^E , which is related to the accumulation of the optimal trajectories in dataset D_k . Secondly, in Sec. 4.2.2, we first introduce a general condition (Cond. 4.6) for the choice of Alg^O , based on which we quantify the accumulation of optimal trajectories in D_k with the regret of Alg^O , and connect the exploration by Alg^O and the optimality of Alg^E . We also supplement some details about how to relax such a condition and inherit the guarantees by the doubling-trick in Appx. G, which may be of independent interest. In Sec. 4.2.3, i.e. the last part of analysis, we bring the above two steps together and complete the proof. However, there is an additional challenge when the tabular MDP has multiple deterministic optimal policies, which is possible when there are non-unique optimal actions at some states. We overcome this difficulty by Thm. 4.8 about policy coverage. To our knowledge, the only paper that runs into a similar challenge is [Papini et al., 2021], and they bypass the difficulty by assuming the uniqueness of optimal policy. Finally, Section 4.2.4 provide some interpretation and implications of our results.

4.2.1 Pessimistic Value Iteration as Alg^E and its Property

The full details of our algorithm for tiered RL setting is provided in Alg. 3, where we use PVI as Alg^E . Here we do not specify a concrete **Bonus** function, but provide general results for a range of

³Comparing with results in Thm. 8.1 of [Lattimore and Szepesvári, 2020], although our upper bounds of $N_i(k)$ is linear w.r.t. k rather than \log scale, we want to highlight that ours hold with high probability $O(1 - k^{-\Theta(\cdot)})$ while [Lattimore and Szepesvári, 2020] only upper bounded the expectation.

Algorithm 3: Tiered-RL Algorithm with Pessimistic Value Iteration as Alg^E

```

1 Input: Episode number  $K$ ; Confidence level  $\{\delta_k\}_{k=1}^K$ ; Bonus function  $\mathbf{Bonus}(\cdot, \cdot)$ 
2 for  $k = 1, 2, \dots, K$  do
3    $\{b_{k,1}(\cdot, \cdot), b_{k,2}(\cdot, \cdot), \dots, b_{k,H}(\cdot, \cdot)\} \leftarrow \mathbf{Bonus}(D_k, \delta_k)$ . //Compute bonus function for PVI.
4   for  $h = H, H-1, \dots, 1$  do
5     for  $s_h \in \mathcal{S}_h, a_h \in \mathcal{A}_h$  do
6        $N_{k,h}(s_h, a_h) \leftarrow$  the number of times  $s_h, a_h$  occurs in the dataset  $D_k$ .
7        $N_{k,h}(s_h, a_h, s_{h+1}) \leftarrow$  the number of times  $(s_h, a_h, s_{h+1})$  occurs in the dataset  $D_k$ .
8        $\mathfrak{p}_{k,h}(\cdot | s_h, a_h) \leftarrow \begin{cases} 0, & \text{if } N_{k,h}(s_h, a_h) = 0; \\ \frac{N_{k,h}(s_h, a_h, \cdot)}{N_{k,h}(s_h, a_h)}, & \text{otherwise.} \end{cases}$ 
9     end
10     $\mathfrak{Q}_{k,h}(\cdot, \cdot) \leftarrow \max\{R(\cdot, \cdot) + \mathfrak{p}_{k,h} \mathfrak{V}_{k,h+1}(\cdot, \cdot) - b_{k,h}(\cdot, \cdot), 0\}$ .
11     $\mathfrak{V}_{k,h}(\cdot) = \max_{a_h \in \mathcal{A}} \mathfrak{Q}_{k,h}(\cdot, a_h)$ ,  $\pi_{k,h}^{\text{PVI}}(\cdot) \leftarrow \arg \max_a \mathfrak{Q}_{k,h}(\cdot, a)$ .
12  end
13   $\pi_k^{\text{E}} \leftarrow \{\pi_{k,1}^{\text{PVI}}, \pi_{k,2}^{\text{PVI}}, \dots, \pi_{k,H}^{\text{PVI}}\}$ 
14  // Step 2: Use AlgO satisfying Cond. 4.6 to compute  $\pi_k^{\text{O}}$  for  $G^{\text{O}}$ 
15   $\pi_k^{\text{O}} \leftarrow \text{Alg}^{\text{O}}(D_k)$ .
16  // Step 3: Sample trajectories and collect new data
17  Interact with  $G^{\text{E}}$  and  $G^{\text{O}}$  by  $\pi_k^{\text{E}}$  and  $\pi_k^{\text{O}}$ , and observe  $\tau_k^{\text{E}}$  and  $\tau_k^{\text{O}}$ , respectively.
18   $D_{k+1} \leftarrow D_k \cup \{\tau_k^{\text{O}}\}$ .
19 end

```

qualified bonus functions satisfying Cond. 4.4 below. Cond. 4.4 can be satisfied by many bonus term considered in online literatures, and we briefly discuss some examples in Appx. F.1.

Condition 4.4 (Condition on Bonus Term for Alg^E). *We define the following event at iteration $k \in [K]$ during the running of Alg. 3: $\mathcal{E}_{\text{Bonus},k} := \bigcap_{h \in [H]} \bigcap_{s_h \in \mathcal{S}_h, a_h \in \mathcal{A}_h} \left\{ \left| \mathfrak{p}_{k,h} \mathfrak{V}_{k,h+1}(s_h, a_h) - \mathfrak{p}_h \mathfrak{V}_{k,h+1}(s_h, a_h) \right| < b_{k,h}(s_h, a_h) \right\} \cap \left\{ b_{k,h}(s_h, a_h) \leq B_1 \frac{\log(B_2/\delta_k)}{N_{k,h}(s_h, a_h)} \right\}$ where B_1 and B_2 are parameters depending on S, A, H and Δ but independent of δ_k, k .⁴ We assume that, **Bonus** function satisfies that, in Alg. 3, given arbitrary sequence $\{\delta_k\}_{k=1}^K$ with $\delta_1, \delta_2, \dots, \delta_K \in (0, 1/2)$, at arbitrary iteration $k \in [K]$, we have $\Pr(\mathcal{E}_{\text{Bonus},k}) \geq 1 - \delta_k$.*

Next, we provide an upper bound for the sub-optimality gap of π_k^{PVI} with the clipping operator $\text{Clip}[x|\varepsilon] := x \cdot \mathbb{1}[x \geq \varepsilon]$. Previous upper bounds of PVI [e.g., Theorem 4.4 of Jin et al., 2021] do not leverage the strictly positive gap and can be much looser when $N_{k,h}$ is large, and directly applying those results to our analysis would result in a regret scaling with \sqrt{K} .

Theorem 4.5. *By running Algorithm 3 with confidence level δ_k , a function **Bonus** satisfying Condition 4.4, and a dataset $D = \{\tau_1, \dots, \tau_k\}$ consisting of k complete trajectories generated by executing a sequence of policies π_1, \dots, π_k , on the event $\mathcal{E}_{\text{Bonus}}$ defined in Condition 4.4:*

$$V_1^*(s_1) - V_1^{\pi_k^{\text{PVI}}}(s_1) \leq 2\varepsilon_{\pi} \prod_{h=1}^H \text{Clip} \left(\min_{s_h \in \mathcal{S}_h, a_h \in \mathcal{A}_h} \frac{\log(B_2/\delta_k)}{N_{k,h}(s_h, a_h)} \right) \varepsilon_{\text{Clip}} \quad (1)$$

where π^* can be an arbitrary optimal policy, $\varepsilon_{\text{Clip}} := \frac{\Delta_{\min}}{2H+2}$ if $|\Pi^*| = 1$ and $\varepsilon_{\text{Clip}} := \frac{d_{\min} \Delta_{\min}}{2SAH}$ if $|\Pi^*| > 1$, where $d_{\min} := \min_{\pi \in \Pi, s_h \in \mathcal{S}_h, a_h \in \mathcal{A}_h} d^{\pi}(s_h, a_h)$ subject to $d^{\pi}(s_h, a_h) > 0$.

4.2.2 Choice and Analysis of Alg^O

Next, we introduce our general condition for Alg^O that the Alg^O can achieve $O(\log k)$ -regret with high probability. It is worth noting that many existing near-optimal online RL algorithms for gap-dependent settings may not directly satisfy the condition [Simchowitz and Jamieson, 2019, Xu et al.,

⁴Note that we do not require the knowledge of π^* 's to compute $b_{k,h}$.

2021, Dann et al., 2021] since they use a fixed confidence interval δ . In Appx. G, we will introduce a more realistic abstraction of those algorithms in Cond. G.1, and discuss in more details about how to close this gap with an algorithm framework inspired by the doubling trick.

Condition 4.6 (Condition on Alg^O). Alg^O is an algorithm which returns deterministic policies at each iteration, and for arbitrary $k \geq 2$, we have: $\Pr \sum_{\mathfrak{k}=1}^k V_1^*(s_1) - V_1^{\pi_{\mathfrak{k}}} (s_1) > C_1 + \alpha C_2 \log k \leq \frac{1}{k}$, where C_1, C_2 are parameters only depending on S, A, H and gap Δ and independent of k .

Implication of Condition 4.6 for Alg^O Intuitively, low regret implies high accumulation of optimal trajectories in the dataset collected by Alg^O . We formalize this intuition in Thm. 4.7 by establishing the relationship between the regret of Alg^O , d^{π} and $\sum_{\mathfrak{k}=1}^k d^{\pi_{\mathfrak{k}}}(s_h, a_h)$ (the expectation of $N_{k,h}$).

Theorem 4.7. For an arbitrary sequence of deterministic policies $\pi_1, \pi_2, \dots, \pi_k$, there must exist a sequence of deterministic optimal policies $\pi_1^*, \pi_2^*, \dots, \pi_k^*$, such that $\forall h \in [H], s_h \in \mathcal{S}_h, a_h \in \mathcal{A}_h$:

$$\sum_{\mathfrak{k}=1}^k d^{\pi_{\mathfrak{k}}}(s_h, a_h) \geq \sum_{\mathfrak{k}=1}^k d^{\pi_{\mathfrak{k}}^*}(s_h, a_h) - \frac{1}{\Delta_{\min}} \sum_{\mathfrak{k}=1}^k V_1^*(s_1) - V_1^{\pi_{\mathfrak{k}}}(s_1) .$$

4.2.3 Main Results and Analysis

The main analysis is based on our discussion about the properties of Alg^E and Alg^O in previous sub-sections. In the following, we first discuss the proof sketch for the case when $|\Pi^*| = 1$. The main idea is to show that the unique π^* will be “well-covered” by dataset, where we say a policy π^* is “well-covered” if for each $(s_h, a_h) \in \mathcal{S}_h \times \mathcal{A}_h$ with $d^{\pi^*}(s_h, a_h) > 0$, $N_{k,h}(s_h, a_h)$ can strictly increase so that the RHS of Eq.(1) in Thm. 4.5 will gradually decay to zero (e.g. $N_{k,h}(s_h, a_h) \geq \Theta(k)$). To show this, the key observation is that, with high probability, $N_{k,h}(s_h, a_h)$ will not deviate too much from its expectation $\sum_{\mathfrak{k}=1}^k d^{\pi_{\mathfrak{k}}}(s_h, a_h)$ (Lem. F.8), and can be lower bounded by $\sum_{\mathfrak{k}=1}^k d^{\pi_{\mathfrak{k}}}(s_h, a_h) - O(\log k) = k d^{\pi^*}(s_h, a_h) - O(\log k)$ as a result of Thm. 4.7. As a result, the clipping operator in Eq.(1) will take effects as long as k is large enough, and π_k^{PVI} will converge to the optimal policy with no regret. All that remains is to show the regret under failure events is also at the constant level because we choose a gradually decreasing failure rate $O(\frac{1}{k})$, and $\lim_{K \rightarrow \infty} \sum_{k=1}^K O(\frac{1}{k}) < \infty$ as long as $\alpha > 1$.

However, when $|\Pi^*| > 1$, the analysis becomes more challenging. The main difficulty is that, when the optimal policy is not unique, it is not obvious about the existence of “well-covered” π^* , since it is not guarantee that how much similarity is shared by the sequence of policies π_1^*, \dots, π_k^* , especially when $|\Pi^*|$ is exponentially large (e.g. $|\Pi^*| = \Omega((SA)^H)$). We overcome this difficulty by proving the existence of “well-covered” policy in the theorem stated below:

Theorem 4.8. [The existence of well-covered optimal policy] Given an arbitrary tabular MDP, and an arbitrary sequence of deterministic optimal policies $\pi_1^*, \pi_2^*, \dots, \pi_k^*$ (π_i^* may not equal to π_j^* for arbitrary $1 \leq i < j \leq k$ when there are multiple deterministic optimal policies), there exists a (possibly stochastic) policy π_{cover}^* such that $\forall h \in [H], \forall (s_h, a_h) \in \mathcal{S}_h \times \mathcal{A}_h$ with $d^{\pi_{\text{cover}}}(s_h, a_h) > 0$:

$$\sum_{\mathfrak{k}=1}^k d^{\pi_{\mathfrak{k}}}(s_h, a_h) \geq \frac{k}{2} \cdot d^{\pi_{\text{cover}}}(s_h, a_h), \text{ with } d^{\pi_{\text{cover}}}(\cdot, \cdot) := \max \frac{d_{h,\min}^*(\cdot, \cdot)}{(|\mathcal{Z}_{h,\text{div}}^*| + 1)H}, d^{\pi_{\text{cover}}}(\cdot, \cdot) .$$

where $\mathcal{Z}_{h,\text{div}}^* := \{(s_h, a_h) \in \mathcal{S}_h \times \mathcal{A}_h | \exists \pi^*, \pi^* \in \Pi^*, \text{ s.t. } d^{\pi^*}(s_h) > 0, d^{\pi^*}(s_h) = 0\}$, and $d_{h,\min}^*(s_h, a_h) := \min_{\pi \in \Pi} d^{\pi}(s_h, a_h)$ subject to $d^{\pi}(s_h, a_h) > 0$.

Here we provide some explanation to the above result. According to the definition, $\mathcal{Z}_{h,\text{div}}^*$ is the set including the state-action pairs which can be covered by some deterministic policies but is not reachable by some other deterministic policies, and therefore $|\mathcal{Z}_{h,\text{div}}^*| \leq SA$ (or even $|\mathcal{Z}_{h,\text{div}}^*| \ll SA$). Besides, $d_{h,\min}^*(s_h, a_h)$ denotes the minimal occupancy over all possible deterministic optimal policies which can hit s_h, a_h , and therefore, is no less than d_{\min} defined in Thm. 4.5. As a result, we know there exists a “well-covered” π_{cover}^* , since the accumulative density of its arbitrary reachable states can be lower bounded by $O(k)$. Then, following a similar discussion as the case $|\Pi^*| = 1$, we can finish the proof. We summarize our main result below.

Theorem 4.9. *By running an Algorithm satisfying Condition 4.6 as Alg^O , running Alg 3 as Alg^E with a bonus term function **Bonus** satisfying Condition 4.4 and $\delta_k = 1/k^\alpha$, for some constant $\alpha > 1$, for arbitrary $K \geq 1$, the exploitation regret of Alg^E can be upper bounded by:*

(i) *When $|\Pi^*| = 1$ (unique optimal deterministic policy):*

$$\text{Regret}_K(\text{Alg}^E) \leq O \prod_{h=1}^H \times_{\substack{s_h, a_h: \\ d(s_h, a_h) > 0}} \frac{C_1 + C_2}{\Delta_{\min}} \log \frac{\text{SAH}(C_1 + C_2)}{d^\pi(s_h, a_h) \Delta_{\min}} + \frac{B_1 H}{\Delta_{\min}} \log \frac{B_2 H}{d^\pi(s_h, a_h) \Delta_{\min}}$$

(ii) *When $|\Pi^*| > 1$ (non-unique optimal deterministic policies):*

$$\text{Regret}_K(\text{Alg}^E) \leq O \prod_{h=1}^H \times_{\substack{s_h, a_h: \\ d_{\text{cover}}(s_h, a_h) > 0}} \frac{C_1 + C_2}{\Delta_{\min}} \log \frac{\text{SAH}(C_1 + C_2)}{\mathcal{E}^{\pi_{\text{cover}}}(s_h, a_h) \Delta_{\min}} + \frac{B_1 \text{SAH}}{d_{\min} \Delta_{\min}} \log \frac{B_2 \text{SAH}}{d_{\min} \Delta_{\min}}$$

where π_{cover}^* and $\mathcal{E}^{\pi_{\text{cover}}}(s_h, a_h)$ are introduced in Theorem 4.8.

4.2.4 Interpretation of Results in Tabular RL

Recall our objective in Sec. 1 is to establish the benefits of leveraging the tiered structure by showing $\text{Regret}_K(\text{Alg}^E)$ is constant. This contrasts the lower bound of online algorithms that continuously increases with the episode number K , which corresponds to the regret suffered by users in G^E without leveraging the tiered structure, while $\text{Regret}_K(\text{Alg}^O)$ keeps (near-)optimal as before. In Appx. H, we also provide some simulation results as a verification of our theoretical discovery.

One limitation of our results is that our bounds have additional dependence on d^π (or even $1/d^\pi$) compared to most of the regret bounds in the online setting, although similar dependence on $\log d^\pi$ also appeared in a few recent works [e.g., λ_h^+ in Thms. 8 and 9 of Papini et al., 2021]. Besides, according to the lower and upper bound of online RL in gap-dependent settings [Simchowitz and Jamieson, 2019], $C_1 + C_2$ in Cond. 4.6 have dependence on $O(\Delta_{\min}^{-1})$, which implies that in the regret bound in Thm. 4.9, the dependence on Δ_{\min} would be $O(\Delta_{\min}^{-2})$. For the former, in Appx. C.2, we prove a lower bound, showing that $\log \frac{1}{d}$ is unavoidable when Alg^O is allowed to behave adversarially without violating Cond. 4.6; for the latter, we note that in the analysis of MAB setting (Sec.4.1), specifying the detailed behavior of Alg^O can help tighten the bound. Therefore, we conjecture that our results can be improved by putting more constraints on the behavior of Alg^O , which we leave to future work.

5 Conclusion

In this paper, we identify the tiered structure in many real-world applications and study the potential advantages of leveraging it by interacting with users from different groups with different strategies. Under the gap-dependent setting, we provide theoretical evidence of benefits by deriving constant regret for the exploitation policy while maintaining the optimality of the online learning policy.

As for the future work, we propose several potentially interesting directions. **(i)** As we mentioned in Section 4.2.4, it is worth investigating the possibility of improving the regret bound of Alg^E by considering a more concrete choice of Alg^O , or maybe other choices for Alg^E . **(ii)** It would be interesting to relax our constraint on the optimality of Alg^O by introducing the notion of budget C as the tolerance on the sub-optimality of Alg^O . As a result, our setting and the decoupling exploration and exploitation setting can be regarded as special cases of a more general framework when $C = 0$ and $C = \infty$. **(iii)** We assume that the users from different groups share the same transition and reward function, and it would also be interesting to extend our results to more general settings, where the group ID serves as context and will affect the dynamics [Abbasi-Yadkori and Neu, 2014, Modi et al., 2018]. **(iv)** We only consider the setting with two tiers, and it may be worth studying the possibility and potential benefits under the setting with multiple tiers.

Acknowledgements

JH’s research activities on this work were conducted during his internship at MSRA. NJ’s last involvement was in December 2021. NJ also acknowledges funding support from ARL Cooperative Agreement W911NF-17-2-0196, NSF IIS-2112471, NSF CAREER award, and Adobe Data Science Research Award. The authors thank Yuanying Cai for valuable discussion.

References

- Yasin Abbasi-Yadkori and Gergely Neu. Online learning in mdps with side information. *arXiv preprint arXiv:1406.6812*, 2014.
- M Mehdi Afsar, Trafford Crump, and Behrouz Far. Reinforcement learning based recommender systems: A survey. *arXiv preprint arXiv:2101.06286*, 2021.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- Orly Avner, Shie Mannor, and Ohad Shamir. Decoupling exploration and exploitation in multi-armed bandits. *arXiv preprint arXiv:1205.2874*, 2012.
- 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.
- Jacob Buckman, Carles Gelada, and Marc G Bellemare. The importance of pessimism in fixed-dataset policy optimization. *arXiv preprint arXiv:2009.06799*, 2020.
- Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning, 2017.
- Christoph Dann, Teodor Vanislavov Marinov, Mehryar Mohri, and Julian Zimmert. Beyond value-function gaps: Improved instance-dependent regret bounds for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Omar Darwiche Domingues, Pierre Ménard, Emilie Kaufmann, and Michal Valko. Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*, pages 578–598. PMLR, 2021.
- Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *arXiv preprint arXiv:2106.06860*, 2021.
- Jiafan He, Dongruo Zhou, and Quanquan Gu. Logarithmic regret for reinforcement learning with linear function approximation. In *International Conference on Machine Learning*, pages 4171–4180. PMLR, 2021.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(4), 2010.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? *arXiv preprint arXiv:1807.03765*, 2018.
- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *arXiv preprint arXiv:2006.04779*, 2020.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.

- Martin S Lipsky and Lisa K Sharp. From idea to market: the drug approval process. *The Journal of the American Board of Family Practice*, 14(5):362–367, 2001.
- Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. *Advances in Neural Information Processing Systems*, 31, 2018.
- Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Provably good batch reinforcement learning without great exploration. *arXiv preprint arXiv:2007.08202*, 2020.
- Aditya Modi, Nan Jiang, Satinder Singh, and Ambuj Tewari. Markov decision processes with continuous side information. In *Algorithmic Learning Theory*, pages 597–618. PMLR, 2018.
- Matteo Papini, Andrea Tirinzoni, Aldo Pacchiano, Marcello Restelli, Alessandro Lazaric, and Matteo Pirodda. Reinforcement learning in linear mdps: Constant regret and representation selection. *Advances in Neural Information Processing Systems*, 34, 2021.
- Chloé Rouyer and Yevgeny Seldin. Tsallis-inf for decoupled exploration and exploitation in multi-armed bandits. In *Conference on Learning Theory*, pages 3227–3249. PMLR, 2020.
- Max Simchowitz and Kevin G Jamieson. Non-asymptotic gap-dependent regret bounds for tabular mdps. *Advances in Neural Information Processing Systems*, 32:1153–1162, 2019.
- Aleksandrs Slivkins. Introduction to multi-armed bandits. *arXiv preprint arXiv:1904.07272*, 2019.
- Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under partial coverage. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=tyrJsbKAe6>.
- Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent pessimism for offline reinforcement learning. *arXiv preprint arXiv:2106.06926*, 2021.
- Haike Xu, Tengyu Ma, and Simon S Du. Fine-grained gap-dependent bounds for tabular mdps via adaptive multi-step bootstrap. *arXiv preprint arXiv:2102.04692*, 2021.
- Kunhe Yang, Lin Yang, and Simon Du. Q-learning with logarithmic regret. In *International Conference on Artificial Intelligence and Statistics*, pages 1576–1584. PMLR, 2021.
- Ming Yin and Yu-Xiang Wang. Towards instance-optimal offline reinforcement learning with pessimism. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. *Advances in Neural Information Processing Systems*, 33:14129–14142, 2020.
- Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*, pages 7304–7312. PMLR, 2019.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]

 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Detailed Related Work

Online RL The online RL/MAB is the most basic framework studying the trade-off between exploration and exploitation [Auer et al., 2002, Slivkins, 2019, Lattimore and Szepesvári, 2020], where the agent targets at exploring the MDP to identify good actions as fast as possible to minimize the accumulative regrets. In tabular MDPs [Jaksch et al., 2010, Dann et al., 2017, Jin et al., 2018], the regret lower bounds for non-stationary MDPs [Domingues et al., 2021] have been achieved by Azar et al. [2017], Zanette and Brunskill [2019]. Recently, there has been interests in studying the gap-dependent regret [Simchowitz and Jamieson, 2019, Xu et al., 2021, Dann et al., 2021, Yang et al., 2021, He et al., 2021], where the agent can achieve dependence on the number of episodes under additional dependence on (the inverse of) the minimal gap. Simchowitz and Jamieson [2019] reports that, similar to stochastic MABs [Lattimore and Szepesvári, 2020], the regret in the gap-dependent setting must scale as $\log K$, which implies that the $\log K$ upper bound is asymptotically tight. However, all of these works treat the customers equivalently and ignore the opportunities of leveraging the tiered structure. Recently, [Papini et al., 2021] achieved similar constant regret in online setting with linear function approximation. Comparing with ours, they investigated the benefits of good features, while we focus on the benefits of considering a different learning protocol. Besides, although their results were established on a more general linear setting, their assumptions on the feature and the uniqueness of optimal policy are quite restrictive in tabular setting.

Offline RL Offline RL considers how to learn a good policy with a fixed dataset [Levine et al., 2020]. Without the requirement of exploration, offline RL prefers algorithms with strong guarantees for exploitation and safety, and Pessimism in the Face of Uncertainty (PFU) becomes a major principle for achieving this both theoretically and empirically [Yin and Wang, 2021, Uehara and Sun, 2022, Liu et al., 2020, Xie et al., 2021, Buckman et al., 2020, Kumar et al., 2020, Fujimoto and Gu, 2021, Yu et al., 2020]. Similar to the offline setting, we choose Alg^E to be a pessimistic algorithm. However, we still consider to interact the environment with Alg^E, although we ignore the data collected by Alg^E for now and leave the investigation of its value to future work. As another difference, offline RL assumes the dataset is fixed and only the final performance matters, whereas we evaluate the accumulative regret of Alg^F.

B More Discussion about Framework. 1 and Motivating Examples

In this section, we try to justify that our Frw. 1 is an appropriate abstraction for our motivating examples and user-interaction real-world applications.

In the standard online learning protocol, at iteration k , the algorithm Alg compute a policy π_k based on previous exploration data, the environment samples a user from G^O and G^E according to the probability where $P(u_k \in G^O)$ and $P(u_k \in G^E)$, respectively (note that $P(u_k \in G^O) + P(u_k \in G^E) = 1$). After that, π_k will interact with u_k and obtain a new trajectory for the future learning, while u_k suffers loss $V_1 - V^k$, and the expected accumulative loss suffered by users from two groups till step K is

$$\text{Regret}_K(\text{Alg}) := E\left[\sum_{k=1}^K V_1(s_1) - V^k(s_1)\right]$$

Now, we consider a realistic assumption about the probability p_k from different groups:

Assumption 1 (Assumption on the Ratio between Users from Different Groups) We assume that:

$$\frac{P(u_k \in G^E)}{P(u_k \in G^O)} = C; \quad \forall k \in [1, K]$$

for some constant C .

Based on the assumption above, if we do not leverage the tiered structure, and treat the users from different groups equivalently, then the loss suffered by each group will be proportional to the size of

that group. Therefore, even if we assume Alg is near-optimal, the expected loss for each group will scale with $\log K$, i.e.:

$$\begin{aligned} \text{Loss}_K(G^O) &:= E\left[\sum_{k=1}^K \ell[u_k \in G^O](V_1(s_1) - V^k(s_1))\right] \\ &= \frac{1}{1+C} E\left[\sum_{k=1}^K V_1(s_1) - V^k(s_1)\right] = O\left(\frac{\log K}{1+C}\right) \end{aligned} \quad (2)$$

$$\begin{aligned} \text{Loss}_K(G^E) &:= E\left[\sum_{k=1}^K \ell[u_k \in G^E](V_1(s_1) - V^k(s_1))\right] \\ &= \frac{C}{1+C} E\left[\sum_{k=1}^K V_1(s_1) - V^k(s_1)\right] = O\left(\frac{C \log K}{1+C}\right) \end{aligned} \quad (3)$$

Besides, we can also leverage the tiered structure, and consider an alternative protocol below:

Algorithm 4: Online Interaction Protocol after Leveraging Tiered Structure

```

1 Initialize:  $D_1 = f, g$ ;  $k = 1$ ;  $\mathcal{O}_1 = \text{Alg}^O(D_1)$ ;  $\mathcal{E}_1 = \text{Alg}^E(D_1)$ 
2 for  $k = 1; 2; \dots$  do
3   User  $u_k$  comes.
4   if  $u_k \in G^O$  then
5     Use  $\mathcal{O}_k$  to interact with  $u_k$ , and collect data  $\mathcal{O}_k$ .
6      $D_{k+1} = D_k [f_{\mathcal{O}_k}, g_{\mathcal{O}_k}]$ ;  $\mathcal{O}_{k+1} = \text{Alg}^O(D_k)$ ;  $\mathcal{E}_{k+1} = \text{Alg}^E(D_k)$ 
7   end
8   else
9      $u_k \in G^E$  interacts with  $u_k$ , and collect data  $\mathcal{E}_k$ . // We do not use  $\mathcal{E}_k$  for now.
10     $D_{k+1} = D_k [f_{\mathcal{E}_k}, g_{\mathcal{E}_k}]$ ;  $\mathcal{O}_{k+1} = \mathcal{O}_k$ ;  $\mathcal{E}_{k+1} = \mathcal{E}_k$ 
11  end
12 end

```

In another word, in this new protocol, we use two policies at different exploitation level to interact with users from different groups, and only update policies if user comes from G^O . Note that in expectation $\ell[u_k \in G^O]$ will happen for $\frac{K}{1+C}$ times, and therefore we have:

$$\begin{aligned} \text{Loss}_K^O(G^O) &:= E\left[\sum_{k=1}^K \ell[u_k \in G^O](V_1(s_1) - V^k(s_1))\right] \\ &= \text{Regret}_{k=\frac{K}{1+C}}(\text{Alg}^O) \\ \text{Loss}_K^O(G^E) &:= E\left[\sum_{k=1}^K \ell[u_k \in G^E](V_1(s_1) - V^k(s_1))\right] \\ &= C \text{Regret}_{k=\frac{K}{1+C}}(\text{Alg}^E) \end{aligned}$$

where $\text{Regret}_{k=\frac{K}{1+C}}(\text{Alg}^O)$ and $\text{Regret}_{k=\frac{K}{1+C}}(\text{Alg}^E)$ are originally defined in Def. 2.1, and they are exactly the metric we used to measure the performance of Alg^O under our Frw. 1.

Based on our results in Sec. 4.1 and 4.2, we know that under our framework, it is possible to achieve that:

$$\text{Loss}_K^O(G^O) = \text{Regret}_{k=\frac{K}{1+C}}(\text{Alg}^O) = O\left(\log \frac{K}{1+C}\right) \quad (4)$$

$$\text{Loss}_K^O(G^E) = C \text{Regret}_{k=\frac{K}{1+C}}(\text{Alg}^E) = C \cdot \text{constant} \quad (5)$$

where constant means independence of K but may include dependence on other parameters such as $S; A; H; \dots$. Comparing with Eq(2), (3), (4), and (5), we can see that users from G^O will suffer less regret than before because we "transfer" some the regret \mathcal{E} to G^O , and the additional regret suffered by G^E can be compensated in other forms as we discussed in Sec. 1.

Remark Besides, our methods and results can be applied to those scenarios suggested by the decoupling setting [Avner et al., 2012, Rouyer and Seldin, 2020], where Alg^E does not necessarily interact with the environment, and we omit the discussion here.

C Lower Bounds

C.1 Regret Lower Bounds for Tabular MDP without Strictly Positive Gap Assumption

We first recall a Theorem from [Dann et al., 2017]:

Theorem C.1 (Theorem C.1 in [Dann et al., 2017]) There exist positive constants $c_0 > 0, \epsilon_0 > 0$, such that for every $\epsilon \in (0, \epsilon_0)$, $S \geq 4, A \geq 2$ and for every algorithm Alg and $n \geq \frac{c_0 A S H^3}{\epsilon^2}$ there is a fixed-horizon episodic MDP $\mathcal{M}_{\text{hard}}$ with time-dependent transition probabilities and S states and A actions so that returning an ϵ -optimal policy after n episodes is at most ϵ_0 .

Theorem 3.1. [Lower Bound for Alg^E without Gap Assumption] There exist positive constants $c, \epsilon_0, \epsilon_1$, such that, for arbitrary $S \geq 4; A \geq 2; H \geq 2; K \geq \frac{c}{\epsilon_0} H^3 S A$, and arbitrary algorithm pair $(\text{Alg}^O; \text{Alg}^E)$, there must exist a hard tabular MDP $\mathcal{M}_{\text{hard}}, E_{(\text{Alg}^O; \text{Alg}^E); \mathcal{M}_{\text{hard}}} \left[\sum_{k=1}^K V_{\text{opt}} - V_k^E \right] \geq \frac{\epsilon_1}{c} H^3 S A K$; where the expectation is taken over the randomness of algorithms and MDP.

Proof. Suppose we have a pair algorithm $(\text{Alg}^O; \text{Alg}^E)$, we can construct a PAC algorithm with Alg^O and Alg^E , in the following way:

- Input: K .
- For $k = 1; 2; \dots; K$, run Alg^O to collect data and run Alg^E to generate a sequence of policies $\pi_1^E; \dots; \pi_K^E$.
- Uniformly randomly select an index from $1; 2; \dots; K$, and denote it as k_{PAC} .
- Output $\pi_{k_{\text{PAC}}}^E$.

In the following, we denote such an algorithm as Alg_{PAC} . Then, for an arbitrary MDP \mathcal{M} , we must have:

$$E_{\text{Alg}_{\text{PAC}}; \mathcal{M}} [V_{\text{opt}} - V_{k_{\text{PAC}}}^E] = \frac{1}{K} E_{(\text{Alg}^O; \text{Alg}^E); \mathcal{M}} \left[\sum_{k=1}^K V_{\text{opt}} - V_k^E \right]$$

As a result of Markov inequality and that $V_{\text{opt}} \geq 0$ for arbitrary \mathcal{M} , for arbitrary $\epsilon > 0$, we have:

$$\Pr(V_{\text{opt}} - V_{k_{\text{PAC}}}^E \geq \epsilon) \leq \frac{E_{\text{Alg}_{\text{PAC}}; \mathcal{M}} [V_{\text{opt}} - V_{k_{\text{PAC}}}^E]}{\epsilon} \quad (6)$$

Since the above holds for arbitrary \mathcal{M} , by choosing $\mathcal{M} := \frac{\epsilon}{c} \frac{c}{\epsilon} H^3 S A = K$, since $K \geq \frac{c}{\epsilon} H^3 S A$, we have $\epsilon < \epsilon_0$ and

$$\Pr(V_{\text{opt}} - V_{k_{\text{PAC}}}^E \geq \epsilon) \leq \frac{E_{\text{Alg}_{\text{PAC}}; \mathcal{M}} [V_{\text{opt}} - V_{k_{\text{PAC}}}^E]}{\frac{\epsilon}{c} H^3 S A = K} \quad (7)$$

Because $K \geq \frac{c}{\epsilon} H^3 S A = \frac{c}{\epsilon} H^3 S A$, Thm. C.1 implies that, for Alg_{PAC} , for arbitrary $S \geq 4; A \geq 2; H \geq 1$, there must exist a hard MDP $\mathcal{M}_{\text{hard}}$, such that:

$$\frac{E_{\text{Alg}_{\text{PAC}}; \mathcal{M}_{\text{hard}}} [V_{\text{opt}} - V_{k_{\text{PAC}}}^E]}{\frac{\epsilon}{c} H^3 S A = K} \geq \epsilon_0$$

and it is equivalent to:

$$E_{(\text{Alg}^O; \text{Alg}^E); \mathcal{M}_{\text{hard}}} \left[\sum_{k=1}^K V_{\text{opt}} - V_k^E \right] \geq K \epsilon_0 = \frac{\epsilon_0}{c} H^3 S A K \quad (8)$$

which finishes the proof. \square

Remark C.2 (Regret lower bound when F is used in Framework 1) Our techniques can be extended to the case when F are used by Alg^O and Alg^E , and establish the same $\tilde{O}(\frac{1}{\epsilon} \sqrt{H^3 S A K})$ lower bound. Because the only difference would be in this new setting, at iteration k , Alg^O and Alg^E will use $2k$ trajectories to compute \bar{Q}_k^O and \bar{Q}_k^E , which only double the sample size comparing with trajectories in Framework 1. Therefore, with the same techniques (and choosing $\frac{1}{\epsilon} \sqrt{H^3 S A K} = (2K)$), one can obtain a lower bound which differs from $\tilde{O}(\frac{1}{\epsilon} \sqrt{H^3 S A K})$ by constant.

C.2 Lower Bound for the Dependence on $\log d_{\min}$ when $j = 1$

In this section, because we will conduct discussion on multiple different MDPs, and in order to distinguish them, we will introduce a subscript M to highlight which MDP we are discussing. Therefore, we revise some key notations and re-introduce them here.

Notation Given arbitrary tabular MDP $M = (S; A; P; r; H; g)$. We use \mathcal{M} to denote the set of deterministic optimal policies. For each deterministic optimal policy $\mu \in \mathcal{M}$, we define d_{\min}^M to be the minimal non-zero occupancy of the reachable state-action pair.

$$d_{\min}^M := \min_{h; s_h; a_h} d^M(s_h; a_h); \quad s.t.: d^M(s_h; a_h) > 0 \quad (9)$$

Then, we define:

$$d_{M; \min} := \min_{M \in \mathcal{M}} d_{\min}^M; \quad M; d_{\min} := \arg \min_{M \in \mathcal{M}} d_{\min}^M; \quad (10)$$

Different from regret analysis in online setting [Simchowitz and Jamieson, 2019, Xu et al., 2021, Dann et al., 2021], our regret bound in Thm. 4.9 has additional dependence on d_{\min} , even if the optimal policy is unique. In Thm. C.3, we show that the dependence on d_{\min} is unavoidable if we do not have other assumptions about the behavior of Alg^O besides Cond. 4.6. In another word, even if constrained by satisfies Cond. 4.6, Alg^O can be arbitrarily adversarial so that d_{\min} exists in the lower bound. We defer the proof to Appx. C.2

Theorem C.3. For arbitrary $S; A; H \geq 3$, arbitrary $d_{\min} > 0$ and $d_{M; \min} > 0$, if there exists an MDP $M = (S; A; P; r; H; g)$ such that $|S| = S$, $|A| = A$, $d_{M; \min} = d_{\min}$ and the minimal gap is lower bounded by d_{\min} , then there exists a hard MDP $M^+ = (S^+; A^+; P^+; r^+; H; g)$ with $|S^+| = S + 1$, minimal gap lower bounded by $d_{\min} = 4$ and $d_{M^+; \min} = d_{M; \min} = 4$, and an adversarial choice of Alg^O satisfying Cond. 4.6, such that when K is large enough, the expected Pseudo-Regret of Alg^E is lower bounded by:

$$E_{\text{Alg}^O; M; \text{Alg}^E} \left[\sum_{k=1}^K V - V^E \right] = \tilde{O} \left((C_1 + C_2) \log \frac{C_1 + C_2}{d_{M^+; \min} d_{\min}} \right)$$

Proof. The proof is divided into three steps.

Step 1: Construction of the Hard MDP Instance Now, we construct a hard MDP instance $M^+ := (S^+; A^+; P^+; r^+; H; g)$ based on M by expanding the state space with an absorbing state s_{absorb} for layer $h = 2$ (we use $s_{h; \text{absorb}}$ to distinguish the absorbing state at different time step), and define the transition and reward function by:

$$\begin{aligned} 8 a_1 \in A_1; s_2 \in S_2; \quad P^+(s_{2; \text{absorb}} | s_1; a_1) &= \frac{d_{M; \min}}{4}; \\ P^+(s_2 | s_1; a_1) &= \left(1 - \frac{d_{M; \min}}{4}\right) P(s_2 | s_1; a_1) \\ r^+(s_1; a_1) &= \left(1 - \frac{d_{M; \min}}{4}\right) r(s_1; a_1) \\ 8 h = 2; s_h \in S_h; a_h \in A_h; \quad P^+(j | s_h; a_h) &= P(j | s_h; a_h); \quad r^+(s_h; a_h) = r(s_h; a_h) \\ 8 h = 2; a_h \in A_h; \quad P^+(s_{h+1; \text{absorb}} | s_h; a_h) &= 1 \\ 8 H = h = 2; a_h \in A_h; \quad r(s_{h; \text{absorb}}; a_h) &= \frac{1}{\min} \mathbb{I}[a_h = a_h] \end{aligned}$$

Briefly speaking, at the initial state, by taking arbitrary action, with probability $\frac{1}{2}$, it will transit to absorbing state at layer 2, and the agent can not escape from the absorbing state till the end of the episodes. Besides, at the absorbing states, for each layer H , there always exists an optimal action a_H with reward r_{\min} and taking any the other actions will lead to 0 reward. Moreover, M^+ agrees with M for all the transition and rewards when $h \geq 2$.

Easy to see that:

$$V_{M^+}(s_1) - Q_{M^+}(s_1; a_1) = (1 - \frac{d_{M^+; \min}}{4})(V_M(s_1) - Q_M(s_1; a_1))$$

Therefore, if $V_M(s_1) - Q_M(s_1; a_1) > 0$, we still have:

$$V_{M^+}(s_1) - Q_{M^+}(s_1; a_1) \geq \frac{3}{4}(V_M(s_1) - Q_M(s_1; a_1)) \geq \frac{3}{4} r_{\min}$$

Combining with the transition and reward functions in absorbing states, we can conclude that the gap of M^+ is still $O(r_{\min})$.

Step 2: Construction of Adversarial Alg^O. Let's use \mathcal{M}^+ to denote the set of deterministic optimal policies at MDRM⁺. It's easy to see that, for arbitrary $\pi_{M^+} \in \mathcal{M}^+$, there must exists an optimal policy $\pi_M \in \mathcal{M}$ agrees with π_{M^+} at all non-absorbing states (and vice versa), i.e.

$$\pi_{M^+}(s_h) = \pi_M(s_h); \quad \forall h \geq 2, s_h \in S_h$$

Then, for arbitrary $\pi_{M^+} \in \mathcal{M}^+$, we have:

$$d_{M^+}(s_{h; \text{absorb}}) = \frac{d_{M^+; \min}}{4} \left(1 - \frac{d_{M^+; \min}}{4}\right) d_{M^+; \min} < d_{M^+}(s_{h^0}); \quad \forall h^0 \geq 2, s_{h^0} \in S_{h^0}$$

which implies that

$$d_{M^+; \min} = \frac{d_{M^+; \min}}{4}$$

and $s_{h; \text{absorb}}$ are the hardest state to reach for all deterministic optimal policies. In the following, we randomly choose an optimal deterministic policy π_{M^+} from \mathcal{M}^+ , and randomly select one action a_H from A_H with $a_H \in A_H$ and discuss them in the following discussion.

Based on the definition above, we define a deterministic policy π_{M^+} , which agree with π_{M^+} for all states except $s_{H; \text{absorb}}$:

$$\pi_{M^+}(s_h) = \begin{cases} \pi_{M^+}(s_h); & \text{if } s_h \in S_{H; \text{absorb}}; \\ a_H; & \text{if } s_h = s_{H; \text{absorb}}; \end{cases}$$

Now, we are ready to design the adversarial choice Alg^P satisfying the condition 4.6. We consider the following algorithm:

$$\text{Alg}^O(k) = \begin{cases} \pi_{M^+}; & \text{if } k \leq k_{\text{sup}}; \\ \pi_{M^+}; & \text{if } k > k_{\text{sup}}; \end{cases}$$

where k_{sup} is defined to be:

$$k_{\text{sup}} := \sup_{k \in \mathbb{N}^+} : f(k) \leq \frac{1}{d_{M^+; \min}} (C_1 + C_2 \log k) g \leq O\left(\frac{C_1 + C_2}{d_{M^+; \min}} \log \frac{C_1 + C_2}{d_{M^+; \min}}\right)$$

We can easily verify that Cond. 4.6 will not be violated, since

$$\begin{aligned} 8k &\leq 1; \quad \sum_{k=1}^k V - V^{\circ} \\ &= d_{M^+; \min} (V(s_{2; \text{absorb}}) - V^{\circ}(s_{2; \text{absorb}})) \min_{k; k_{\text{sup}}} \\ &= d_{M^+; \min} (V(s_{H; \text{absorb}}) - V^{\circ}(s_{H; \text{absorb}})) \frac{1}{d_{M^+; \min}} (C_1 + C_2 \log \min_{k; k_{\text{sup}}}) \\ &= d_{M^+; \min} \min_{k; k_{\text{sup}}} \frac{1}{d_{M^+; \min}} (C_1 + C_2 \log k) \\ &\quad C_1 + C_2 \log k \end{aligned}$$

Step 3: Lower Bound of Alg^E under the Choice of Adversarial Alg^O Now, we can derive an lower bound for Alg^E . Since in the k_{sup} steps, Alg^E can only observe what happens if action a_H is taken at $s_{H;\text{absorb}}$, and therefore, it has no idea about which action among a_H is the optimal action a_H . We use M^+ to denote a set of MDPs by permuting the position of a_H in M^+ . Since $|A_H| = A$, we have $|M^+| = A^{-1}$ and $M^+ \subseteq M^+$.

Then, we uniformly sample an MDP from M^+ and run the adversarial Alg^O above to generate the data for Alg^E to learn. We use M_i^+ with $i = 1; 2; \dots; A^{-1}$ to refer to the MDPs in M^+ and use index i to refer to the position of the optimal actions a_H in each MDP. For the simplicity of the notation, we use A as the index to refer to the position of a_H .

Because Alg^E do not have prior knowledge about which MDP in M^+ is sampled, we have:

$$\begin{aligned}
 & E_{M^+; \text{Alg}^O; \text{Alg}^E} \left[\sum_{k=1}^{k_{\text{sup}}} V_{k-1} - V_k^E \right] \\
 & E_{M^+; \text{Alg}^O; \text{Alg}^E} \left[\sum_{k=1}^{k_{\text{sup}}} V_{k-1} - V_k^E \right] \\
 &= \frac{1}{A^{-1}} \sum_{i \in [A^{-1}]} E_{M_i^+; \text{Alg}^O; \text{Alg}^E} \left[\sum_{k=1}^{k_{\text{sup}}} V_{k-1} - V_k^E \right] \\
 &= \frac{d_{M; \min} \min}{A^{-1}} \sum_{k=1}^{k_{\text{sup}}} \sum_{i \in [A^{-1}]} \sum_{j \in [A]} \Pr_{\text{Alg}^O; M_i} (E_k(s_{H;\text{absorb}}) = j) \\
 & \quad \text{(Drop the probability that } E_k \text{ is sub-optimal at non-absorbing states)} \\
 &= \frac{d_{M; \min} \min}{A^{-1}} \sum_{k=1}^{k_{\text{sup}}} \sum_{i \in [A^{-1}]} \sum_{j \in [A]} \Pr_{\text{Alg}^O; M_i} (E_k(s_{H;\text{absorb}}) = j) \\
 & \quad + \sum_{i \in [A^{-1}]} \sum_{j \in [A]} \Pr_{\text{Alg}^O; M_i} (E_k(s_{H;\text{absorb}}) = j) \\
 &= \frac{d_{M; \min} \min}{A^{-1}} \sum_{k=1}^{k_{\text{sup}}} \sum_{i \in [A^{-1}]} \Pr_{\text{Alg}^O; M_i} (E_k(s_{H;\text{absorb}}) = i) \\
 & \quad + \sum_{i \in [A^{-1}]} \sum_{j \in [A]} \Pr_{\text{Alg}^O; M_i} (E_k(s_{H;\text{absorb}}) = j) \\
 & \quad \text{(Alg}^E \text{ can not distinguish between } M_i^+ \text{)} \\
 &= \frac{d_{M; \min} \min}{A^{-1}} \sum_{k=1}^{k_{\text{sup}}} \sum_{i \in [A^{-1}]} \sum_{j \in [A]} \Pr_{\text{Alg}^O; M_i} (E_k(s_{H;\text{absorb}}) = j) \\
 &= \frac{A^{-2}}{A^{-1}} d_{M; \min} \min \sum_{k=1}^{k_{\text{sup}}} \\
 &= O(C_1 + C_2) \log \frac{C_1 + C_2}{d_{M^+; \min} \min} = O(C_1 + C_2) \log \frac{C_1 + C_2}{d_{M; \min} \min}
 \end{aligned}$$

□

$$2=f(k) \quad 2=k^2$$

where the last but two step is because of the Azuma-Hoeffding's inequality. \square

Lemma 4.3. [Property of UCB] With the choice that $\delta(k) = 1 + 16 A^2(k + 1)^2$, there exists a constant c_i for arbitrary i with $c_i > 0$ and arbitrary $\epsilon \in [1/4, 1]$, in UCB algorithm, we have: $\Pr(N_i(k) \leq \frac{k}{\epsilon}) \leq \frac{2}{k^2 - 1} + 8k^{-\epsilon} + c_i \frac{1}{k} \log(1 + \frac{A}{\epsilon})$:

Proof. We choose c_i defined in Eq.(11) to be the constant in this Lemma.

The key idea of the proof is that, because $N_i(k) \leq k$ for all k , if $N_i(k) \leq k = \epsilon$, there must exists an iteration R between $k = \epsilon - 1$ and k , such that $N_i(R) = k = \epsilon - 1$ and $N_i(R + 1) = k = \epsilon$ (i.e. R is the time step that UCB takes arm for the $k = \epsilon$ -th time). Therefore, for arbitrary fixed $\epsilon \in [1/4, 1]$, when $k = \epsilon$, we have:

$$\begin{aligned} & \Pr(N_i(k) \leq k = \epsilon) = \Pr(N_i(k) \leq k = \epsilon) \\ &= \sum_{R=k=\epsilon-1}^{k-1} \Pr(f N_i(R) = k = \epsilon - 1; N_i(R + 1) = k = \epsilon) \leq \sum_{R=k=\epsilon-1}^{k-1} \left(\Pr(f b_i(R) \leq \frac{2 \log f(R)}{N_i(R)} b_i(R) + \frac{2 \log f(R)}{N_i(R)} g) \right) \end{aligned}$$

(Union bound.)

$$\begin{aligned} &= \sum_{R=k=\epsilon-1}^{k-1} \Pr(f N_i(R) = k = \epsilon - 1; N_i(R + 1) = k = \epsilon) \\ & \leq \sum_{R=k=\epsilon-1}^{k-1} \left(\Pr(f b_i(R) \leq \frac{2 \log f(R)}{N_i(R)} b_i(R) + \frac{2 \log f(R)}{N_i(R)} g) \right) \end{aligned}$$

(Subtract $\frac{2 \log f(R)}{N_i(R)} g$ at both sides)

$$\begin{aligned} &= \sum_{R=k=\epsilon-1}^{k-1} \left(\Pr(f N_i(R) = k = \epsilon - 1; N_i(R + 1) = k = \epsilon) \leq \Pr(f b_i(R) \leq \frac{2 \log f(R)}{N_i(R)} b_i(R) + \frac{2 \log f(R)}{N_i(R)} g) \right) \\ &+ \Pr(f N_i(R) = k = \epsilon - 1; N_i(R) = k = \epsilon + 1) \leq \Pr(f 0 \leq \frac{2 \log f(R)}{N_i(R)} g) \end{aligned}$$

(12)

$$\begin{aligned} &= \sum_{R=k=\epsilon-1}^{k-1} \left(\Pr(f b_i(R) \leq \frac{2 \log f(R)}{N_i(R)} b_i(R) + \frac{2 \log f(R)}{N_i(R)} g) \right) \\ &+ \Pr(f N_i(R) = k = \epsilon - 1; N_i(R + 1) = k = \epsilon) \leq \Pr(f 0 \leq \frac{2 \log f(R)}{N_i(R)} g) \end{aligned}$$

(Under our choice of $\delta(k)$, and $N_i(R) = k = \epsilon - 1$, $\frac{\log f(R)}{N_i(R)} \leq \frac{\log f(k)}{k-1} \leq \frac{2}{8(k-1)}$)

$$\begin{aligned} &= \sum_{R=k=\epsilon-1}^{k-1} \left(\Pr(f b_i(R) \leq \frac{2 \log f(R)}{N_i(R)} b_i(R) + \frac{2 \log f(R)}{N_i(R)} g) \right) \\ &+ \Pr(f 0 \leq \frac{2 \log f(R)}{N_i(R)} g) \end{aligned}$$

(Azuma-Hoeffding Inequality)

$$\frac{2k}{f(k-1)} = \frac{2k}{(16A^2k^2 - 2 + 1)} \leq \frac{2}{k^2 - 1} \quad (\epsilon \in [1/4, 1])$$

where the step (12) is because:

$$f b_i(R) \leq \frac{2 \log f(R)}{N_i(R)} b_i(R) + \frac{2 \log f(R)}{N_i(R)} g$$

$$2f 0 < b_i(\mathbb{R}) \quad i \quad i + \frac{s}{N_i(\mathbb{R})} \log \left[f b_i(\mathbb{R}) \quad i + \frac{s}{N_i(\mathbb{R})} \right] \leq 0g$$

□

Lemma D.1. Given an arm i , we separate all the arms into two parts depending on whether its gap is larger than $\frac{1}{2} \frac{A}{\min}$ and define $G_i^{\text{lower}} := \{j \mid \Delta_j > \frac{1}{2} \frac{A}{\min}\}$ and $G_i^{\text{upper}} := \{j \mid \Delta_j \leq \frac{1}{2} \frac{A}{\min}\}$. With the choice that $N_j(k) = 1 + 16 A^2 (k+1)^2$, there is a constant c , such that for arbitrary j with $\Delta_j > 0$, for the LCB algorithm in Alg 2, we have:

$$\Pr(i = \hat{k}) \leq \frac{c}{k^2} + 2A = k^2 \quad 1; \quad 8c \quad k_i := 8c \quad \sum_{j \in G_i^{\text{lower}}} \frac{1}{2} + \frac{4jG_i^{\text{upper}}}{2} \log\left(1 + \frac{A}{\min}\right) \quad (13)$$

where c is the constant considered in Lem. 4.3 (defined in Eq(11)).

Proof. We want to remark that the constants in the definition of (i.e. $8c$ in “ $8c$ ” and 4 in “ $4jG_i^{\text{upper}}$ ”) can be replaced by others, but we choose them carefully in order to make sure some steps in the proof of this Lemma and Thm. 4.1 can go through.

The main idea of the proof is to use Lem. 4.3 to show that, for those arms j with $\Delta_j > 0$, when k is large, $N_j(k)$ will be small for those $j \in G_i^{\text{lower}}$. As a result, there must exist an arm $j \in G_i^{\text{upper}}$, such that $N_j(k)$ is large than the threshold considered in Lem. 4.2 and therefore, with high probability, arm j will not be preferred.

First, we try to apply Lem. 4.3 to upper bound the quantity $N_j(k)$ for those arms $j \in G_i^{\text{lower}}$. For each $j \in G_i^{\text{lower}}$, we define the following quantity, which measures the magnitude of Δ_j with k :

$$k_{j,k} := \frac{k}{\frac{8c}{j} \log\left(1 + \frac{A}{\min}\right)}$$

We only consider $k_{j,k} > 1$, where we always have $k_{j,k} > 1$ based on the definition of $k_{j,k}$.

Next, we separately consider two cases depending on whether $k_{j,k} > 2A$ or not.

Case 1: $k_{j,k} > 2A$: In this case, j is relatively large (or say more sub-optimal) comparing with iteration k . For arbitrary $k > k_i$, we have

$$k = k_{j,k} \frac{8c}{j} \log\left(1 + \frac{A}{\min}\right) \geq 2A \frac{8c}{j} \log\left(1 + \frac{A}{\min}\right) \geq 2A + \frac{2cA}{j} \log\left(1 + \frac{A}{\min}\right):$$

which implies that k satisfying the condition of applying Lemma 4.3 with $\Delta_j = 2A$, and we can conclude that:

$$\Pr(N_k(j) \leq \frac{k}{2A}) \leq \frac{2}{k^2 - 1}:$$

Case 2: $k_{j,k} \leq 2A$: Note that

$$\frac{4c}{j} \log\left(1 + \frac{A}{\min}\right) = \frac{k}{2k_{j,k}}:$$

Since $k_{j,k}$ locates in the interval $[1, 4A]$ and:

$$k = k_{j,k} \frac{8c}{j} \log\left(1 + \frac{A}{\min}\right) \leq 2k_{j,k} + 2k_{j,k} \frac{c}{j} \log\left(1 + \frac{A}{\min}\right)$$

which satisfies the condition of applying Lem. 4.3 with $\Delta_j = 2k_{j,k}$. Therefore, we have:

$$\Pr(N_k(j) \leq \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right)) = \Pr(N_k(j) \leq \frac{k}{2k_{j,k}}) \leq \frac{2}{k^2 - 1} \quad (14)$$

Combining the above two cases, we can conclude that, for arbitrary $j \in G_i^{\text{lower}}$,

$$\Pr(N_j(k) \leq \frac{k}{2A} + \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right)) \leq \min\left\{\Pr(N_j(k) \leq \frac{k}{2A}); \Pr(N_j(k) \leq \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right))\right\} \leq \frac{2}{k^2 - 1}$$

which reflects that with high probability, $\sum_{j \in 2G_i^{\text{lower}}} N_j(k)$ is small:

$$\Pr\left(\sum_{j \in 2G_i^{\text{lower}}} N_j(k) \leq \frac{k}{2} + \sum_{j \in 2G_i^{\text{lower}}} \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right)\right)$$

$$\Pr\left(N_j(k) \leq \frac{k}{2jG_i^{\text{lower}j}} + \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right)\right)$$

$$\Pr\left(N_j(k) \leq \frac{k}{2A} + \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right)\right)$$

$$\frac{2jG_i^{\text{lower}j}}{k^2 - 1} \leq \frac{2A}{k^2 - 1}$$

Since $\sum_{j \in 2G_i^{\text{upper}}} N_j(k) + \sum_{j \in 2G_i^{\text{lower}}} N_j(k) = k$, and note that,

$$k - \left(\frac{k}{2} + \sum_{j \in 2G_i^{\text{lower}}} \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right)\right) = \frac{k}{2} - \sum_{j \in 2G_i^{\text{lower}}} \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right) = \frac{16c \sum_{j \in 2G_i^{\text{upper}}} \log\left(1 + \frac{A}{\min}\right)}{j}$$

we have:

$$\Pr\left(\sum_{j \in 2G_i^{\text{upper}}} N_j(k) \geq \frac{16c \sum_{j \in 2G_i^{\text{upper}}} \log\left(1 + \frac{A}{\min}\right)}{j}\right)$$

$$= \Pr\left(\sum_{j \in 2G_i^{\text{lower}}} N_j(k) \geq \frac{k}{2} + \sum_{j \in 2G_i^{\text{lower}}} \frac{4c}{j} \log\left(1 + \frac{A}{\min}\right)\right)$$

$$\frac{2A}{k^2 - 1}$$

Therefore, w.p.1 $\frac{2A}{k^2 - 1}$, there exists $\sum_{j \in 2G_i^{\text{upper}}} N_j(k)$, such that

$$\sum_{j \in 2G_i^{\text{upper}}} N_j(k) \geq \frac{1}{\sum_{j \in 2G_i^{\text{upper}}} j} \sum_{j \in 2G_i^{\text{upper}}} N_j(k) \geq \frac{16c}{j} \log\left(1 + \frac{A}{\min}\right):$$

Recall our choice of ϵ (Eq.(11)), the above implies that:

$$N_j(k) \geq \frac{32 \log f(k)}{j} \tag{15}$$

Therefore,

$$\Pr(f_i = \frac{E}{k}g \setminus f(k) \leq k_i g) = \Pr(f_i = \frac{E}{k}g \setminus f(k) \leq k_i g \setminus \sum_{j \in 2G_i^{\text{upper}}} N_j(k) \geq \frac{32 \log f(k)}{j}g)$$

$$+ \Pr(f_i = \frac{E}{k}g \setminus f(k) \leq k_i g \setminus \sum_{j \in 2G_i^{\text{upper}}} N_j(k) < \frac{32 \log f(k)}{j}g)$$

$$\Pr(f_i = \frac{E}{k}g \setminus f(k) \leq k_i g \setminus \sum_{j \in 2G_i^{\text{upper}}} N_j(k) \geq \frac{32 \log f(k)}{j}g) + \frac{2A}{k^2 - 1}$$

(Eq.(15))

$$\Pr(f_i = \frac{E}{k}g \setminus f(k) \leq k_i g \setminus \sum_{j \in 2G_i^{\text{upper}}} N_j(k) < \frac{32 \log f(k)}{j}g) + \frac{2A}{k^2 - 1}$$

(Lem. 4.2)

□

Lemma D.2 (Integral Lemma) For arbitrary $k_0 \geq 1$ and $\epsilon > 1$, we have:

$$\sum_{k=k_0+1}^{\infty} \frac{1}{k} \leq \int_{k_0}^{\infty} \frac{1}{x} dx \leq \frac{1}{(k_0 - 1)k_0}$$

Theorem 4.1. [Exploitation Regret] In Algorithm 2, by choosing arbitrary $\epsilon > 0$, there exists an absolute constant c , such that, for arbitrary $K \geq 1$, the pseudo-regret of Alg^E is upper bounded by: $\text{Regret}_K(\text{Alg}^E) \leq \frac{A}{1-\epsilon} + \sum_{i>0} (A - \epsilon^i) \frac{1}{i} \frac{\epsilon^i}{1-\epsilon^i}$ where $\epsilon := 1 - \frac{1}{2K}$ so $\frac{1}{2K} = 0$.

Proof. Recall the definition of k_i in Eq.(13) in Lem. D.1 above. For $i \geq 2$, if $j_i \neq j_{i-1}$, we have:

$$\begin{aligned} k_i - k_{i-1} &= c \sum_{2G_i^{\text{lower}}} \frac{8}{2} \sum_{2G_{i-1}^{\text{lower}}} \frac{8}{2} + \frac{32j_i^{\text{upper}}}{2} \frac{32j_{i-1}^{\text{upper}}}{2} \log\left(1 + \frac{A}{\min}\right) \\ &= c \left((j_i^{\text{lower}} - j_{i-1}^{\text{lower}}) \frac{32}{2} + \frac{32j_i^{\text{upper}}}{2} \frac{32j_{i-1}^{\text{upper}}}{2} \log\left(1 + \frac{A}{\min}\right) \right) \\ &\quad \left(8 \sum_{2G_i^{\text{lower}}} \sum_{2G_{i-1}^{\text{lower}}} \right), \text{ we have } \frac{1}{2} = \frac{1}{2} - \frac{1}{2} = \frac{1}{2} \\ &= c \left((j_i^{\text{upper}} - j_{i-1}^{\text{upper}}) \frac{32}{2} + \frac{32j_i^{\text{upper}}}{2} \frac{32j_{i-1}^{\text{upper}}}{2} \log\left(1 + \frac{A}{\min}\right) \right) \\ &\quad \left(j_i^{\text{upper}} + j_{i-1}^{\text{lower}} = j_i^{\text{upper}} + j_{i-1}^{\text{lower}} = A \right) \\ &= 32c \left(j_i^{\text{upper}} - j_{i-1}^{\text{upper}} \right) \frac{1}{2} \frac{1}{1-\epsilon^i} \log\left(1 + \frac{A}{\min}\right); \end{aligned}$$

and if $j_i = j_{i-1}$, we also have:

$$k_i - k_{i-1} = 0 - 32c \left(j_i^{\text{upper}} - j_{i-1}^{\text{upper}} \right) \frac{1}{2} \frac{1}{1-\epsilon^i} \log\left(1 + \frac{A}{\min}\right);$$

Moreover, for $i = 1$, with the extended definition that $\epsilon_0 = 1$ (so that $\frac{1}{2} = \frac{1}{2} = 0$) and $j_0^{\text{upper}} = A$, we also have:

$$\begin{aligned} k_1 &:= 8c \sum_{2G_1^{\text{lower}}} \frac{1}{2} + \frac{4j_1^{\text{upper}}}{2} \log\left(1 + \frac{A}{\min}\right) \\ &= 8c \sum_{2G_1^{\text{lower}}} \frac{4}{2} + \frac{4j_1^{\text{upper}}}{2} \log\left(1 + \frac{A}{\min}\right) \\ &= \frac{32cA}{2} \log\left(1 + \frac{A}{\min}\right) \\ &= 32c \left(j_0^{\text{upper}} - j_1^{\text{upper}} \right) \left(\frac{1}{2} - \frac{1}{2} \right) \log\left(1 + \frac{A}{\min}\right); \end{aligned}$$

Therefore, we have (we denote $\epsilon := 1 - \frac{1}{2K}$ and $k_0 := 0$):

$$\begin{aligned} \lim_{K \rightarrow \infty} \text{Regret}_K(\text{Alg}^E) &= \sum_{k=1}^{\infty} \sum_{j: j > 0} \Pr(j = \frac{E}{k}) \\ &= \sum_{i=1}^{\infty} \sum_{k=k_{i-1}+1}^{\infty} \sum_{j > 0} \Pr(j = \frac{E}{k}) \\ &= \sum_{i=1}^{\infty} \sum_{k=k_{i-1}+1}^{\infty} \sum_{j < i-1} \Pr(j = \frac{E}{k}) + \sum_{j < i-1} \Pr(j = \frac{E}{k}) \\ &= \sum_{i=1}^{\infty} \sum_{k=k_{i-1}+1}^{\infty} \sum_{j < i-1} \left(\frac{2}{k^2} + \frac{2A}{k^2-1} \right) \Pr(j = \frac{E}{k}) \\ &\quad \text{(Lemma D.1)} \\ &= \sum_{i=1}^{\infty} \sum_{k=k_{i-1}+1}^{\infty} \sum_{j < i-1} \left(\frac{2}{k^2} + \frac{2A}{k^2-1} \right) \Pr(j = \frac{E}{k}) \\ &\quad \text{(Lemma D.1)} \end{aligned}$$

$$\sum_{i=1}^X \sum_{k=k_i+1}^X \frac{2}{k^2} + \frac{2A}{k^2-1} + \sum_{i=1}^X \sum_{k=k_{i-1}+1}^X \Pr(j = \frac{E}{k}) \cdot j$$

$$\sum_{i=1}^X \sum_{k=k_i+1}^X \frac{2(A+1)}{k^2-1} + \sum_{i=1}^X \sum_{k=k_{i-1}+1}^X j \cdot \Pr(j = \frac{E}{k})$$

(Second term is maximized when $\Pr(j = \frac{E}{k}) = 1$.)

$\Theta(\frac{A}{1}) + \sum_{i=1}^X \sum_{k=k_i+1}^X (k_i - k_{i-1})$
 (First term: Lemma D.2 and some simplification; Second term: Definition of k_i)

$$= \Theta(\frac{A}{1}) + \sum_{i: i > 0} \sum_{k=k_i+1}^X (k_i - k_{i-1})$$

$$\Theta(\frac{A}{1}) + \sum_{i > 0} \sum_{k=k_i+1}^X 32c \cdot j \cdot G_{i-1}^{\text{upper}} \cdot \frac{1}{i} \cdot \frac{i}{2^{i-1}} \log(1 + \frac{A}{\min})$$

$$\Theta(\frac{A}{1}) + \sum_{i > 0} \sum_{k=k_i+1}^X j \cdot G_{i-1}^{\text{upper}} \cdot \frac{1}{i} \cdot \frac{i}{2^{i-1}}$$

According to the definition, we always have $G_{i-1}^{\text{upper}} \leq A \cdot i + 1$, therefore,

$$\lim_{K \rightarrow \infty} \text{Regret}_K(\text{Alg}^E) = \Theta(\frac{A}{1}) + \sum_{i > 0} \sum_{k=k_i+1}^X (A - i) \cdot \frac{1}{i} \cdot \frac{i}{2^{i-1}}$$

□

E Behavior Analysis of Optimistic Algorithm

Definition E.1 (Definition of Events)

$$E_{k,h} := f_{k,h}(s_h) \notin g_h(s_h); \quad \mathbb{E}_{k,h} := E_{k,h} \setminus \bigcap_{h^0=1}^h E_{k,h^0-1}; \quad E_k := \bigcap_{h=1}^H E_{k,h};$$

$$E_{k,h} := f_{k,h}(s_h) \notin g_h(s_h); \quad \mathbb{E}_{k,h} := E_{k,h} \setminus \bigcap_{h^0=1}^h E_{k,h^0-1}; \quad E_k := \bigcap_{h=1}^H E_{k,h}$$

In another word, $E_{k,h}$ means k disagrees with π at state s_h which occurs at step h , $\mathbb{E}_{k,h}$ means the first disagreement between k and π and occurs at step h , and E_k denotes the event that there exists one state s_h at some time step $h \in [H]$ such that k agrees with π at s_h .

Besides, $E_{k,h}$ denotes the events that $k_h(s_h)$ will not be taken by any optimal policy. Note that here we use $g_h(s_h)$ to denote the set of all possible optimal actions at state s_h . Given a deterministic optimal policy π , in general $\mathbb{E}_{k,h} \subseteq E_{k,h}$ when there are multiple optimal actions at one state.

Lemma E.2. For arbitrary reward function R , given a fixed deterministic policy, we have:

$$V_1(s_1) - V_1^k(s_1) = E_k \left[\sum_{h=1}^H I[\mathbb{E}_{k,h}] (V_h(s_h) - V_h^k(s_h)) \right]$$

Proof.

$$V_1(s_1) - V_1^k(s_1) = I[\mathbb{E}_{k,1}] Q_1(s_1; k) - Q_1^k(s_1; k) + I[\mathbb{E}_{k,1}] V_1(s_1) - V_1^k(s_1)$$

$$= E_k [I[\mathbb{E}_{k,1}] V_2(s_2) - V_2^k(s_2)] + I[\mathbb{E}_{k,1}] V_1(s_1) - V_1^k(s_1)$$

($\mathbb{E}_{k,1} = E_{k,1}$ by definition)

$$\begin{aligned}
&= E_k [I[E_{k;2} \setminus E_{k;1}] Q_2(s_2; k) - Q_2^k(s_2; k)] \\
&\quad + E_k [I[E_{k;2} \setminus E_{k;1}] V_2(s_2) - V_2^k(s_2)] + I[E_{k;1}] V_1(s_1) - V_1^k(s_1) \\
&= E_k [I[E_{k;2} \setminus E_{k;1}] Q_2(s_2; k) - Q_2^k(s_2; k)] \\
&\quad + E_k [I[E_{k;2}] V_2(s_2) - V_2^k(s_2)] + I[E_{k;1}] V_1(s_1) - V_1^k(s_1) \\
&= \dots \\
&= E_k \left[\sum_{h=1}^H I[E_{k;h}] (V_h(s_h) - V_h^k(s_h)) \right]
\end{aligned}$$

□

Lemma E.3 (Relationship between Density Difference and Policy Disagreement Probability)

$$d^k(s_h; a_h) - d(s_h; a_h) \leq \min_j \Pr(E_k; j; k); d(s_h; a_h) \leq \sum_{s_h \in S_h, a_h \in A_h, h \in [H]} \Pr(E_k; j; k)$$

where we use $\Pr(E_k; j; k)$ as a short note of $\Pr(E_{s_1; a_1; s_2; a_2; \dots; s_H; a_H} = E_k; j; k)$.

Proof. By applying Lemma E.2 with $r_{s_h; a_h} := I[S_h = s_h; A_h = a_h]$ as reward function, we have:

$$\begin{aligned}
d(s_h; a_h) - d^k(s_h; a_h) &= V_1(s_1; s_h; a_h) - V_1^k(s_1; s_h; a_h) \\
&= E_k \left[\sum_{h^0=1}^H I[E_{k;h^0}] (V_{h^0}(s_{h^0}; s_h; a_h) - V_{h^0}^k(s_{h^0}; s_h; a_h)) \right] \\
&\quad (V_{h^0} = V_{h^0}^k = 0 \text{ for all } h^0 \neq h + 1) \\
&= E_k \left[\sum_{h^0=1}^H I[E_{k;h^0}] V_{h^0}(s_{h^0}; s_h; a_h) \right] - (V_{h^0}^k(s_{h^0}^0; s_h; a_h) - 0) \\
&= E_k \left[\sum_{h^0=1}^H I[E_{k;h^0}] \right] - (V_{h^0}^k(s_{h^0}^0; s_h; a_h) - 1) \\
&= E_{s_1; a_1; s_2; a_2; \dots; s_H; a_H} [E_k; j; k] = \Pr(E_k; j; k)
\end{aligned}$$

which implies that,

$$d^k(s_h; a_h) - d(s_h; a_h) \leq \Pr(E_k; j; k)$$

Combining with $d^k \geq 0$, we finish the proof. □

Definition E.4 (Conversion to Optimal Deterministic Policy) Given arbitrary deterministic policy $\pi = f_1; \dots; f_H$, we use $\pi^* = f_1^*; \dots; f_H^*$ to denote an optimal deterministic policy, such that:

$$f_h^*(s_h) = \begin{cases} f_h(s_h); & \text{if } f_h(s_h) \in \pi_h(s_h); \\ \text{Select}_h(\pi_h(s_h)); & \text{otherwise} \end{cases}$$

where Select is a function which returns the first optimal action from $\pi_h(s_h)$.

In another word, π^* agrees with π if $f_h(s_h)$ is one of the optimal action at state s_h . Otherwise, π^* takes one of a fixed optimal action from $\pi_h(s_h)$. In order to make sure π^* is a deterministic mapping, we assume function Select only choose the first optimal action in $\pi_h(s_h)$ (ordered by index of action).

Theorem 4.7. For an arbitrary sequence of deterministic policies $\pi_2; \dots; \pi_k$, there must exist a sequence of deterministic optimal policies $\pi_2^*; \dots; \pi_k^*$, such that $\forall h \in [H]; s_h \in S_h; a_h \in A_h$:

$$d^k(s_h; a_h) - d^{\pi^*}(s_h; a_h) \leq \frac{1}{\min_{\pi \in \Pi} \sum_{h=1}^H V_1(s_1) - V_1^{\pi}(s_1)}$$

Proof. For each k , we construct an optimal deterministic policy $\pi_k := \pi_k$, where π_k is defined in Def. E.4. By applying Lemma E.2 with the reward function in MDP, and π_k , we have:

$$\begin{aligned} V_1^k(s_1) - V_1^k(s_1) &= E_k \left[\sum_{h=1}^H \mathbb{I}[\mathbb{E}_{k;h; \pi_k}] (V_h^k(s_h) - V_h^k(s_h)) \right] \\ &= E_k \left[\sum_{h=1}^H \mathbb{I}[\mathbb{E}_{k;h; \pi_k}] (V_h^k(s_h) - Q_h^k(s_h; \pi_k(s_h))) \right] \\ &= E_k \left[\sum_{h=1}^H \mathbb{I}[\mathbb{E}_{k;h; \pi_k}] \min \right] = \min \Pr(\mathbb{E}_{k; \pi_k} \geq k) \end{aligned}$$

Therefore, we have:

$$\Pr(\mathbb{E}_{k; \pi_k} \geq k) \leq \frac{1}{\min} (V_1^k(s_1) - V_1^k(s_1))$$

By applying Lemma E.3, we have:

$$d^k(s_h; a_h) - d^k(s_h; a_h) \leq \frac{1}{\min} |V_1(s_1) - V_1^k(s_1)|; \quad \forall s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h; h \in [H]$$

After the same discussion for all $k \in [K]$, and the above inequality of each together, we have:

$$\sum_{k=1}^K d^k(s_h; a_h) - \sum_{k=1}^K d^k(s_h; a_h) \leq \frac{1}{\min} \sum_{k=1}^K |V_1(s_1) - V_1^k(s_1)|$$

□

Corollary E.5 (Unique Optimal Policy) When $j = 1$, Thm. 4.7 implies that:

$$\sum_{k=1}^K d^k(s_h; a_h) - K d^k(s_h; a_h) \leq \frac{1}{\min} \sum_{k=1}^K |V_1(s_1) - V_1^k(s_1)|$$

Theorem 4.8. [The existence of well-covered optimal policy] Given an arbitrary tabular MDP, and an arbitrary sequence of deterministic optimal policies π_1, \dots, π_K (π_i may not equal to π_j for arbitrary $1 \leq i < j \leq K$ when there are multiple deterministic optimal policies), there exists a (possibly stochastic) policy π_{cover} such that $\forall s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h$ with $d^{\text{cover}}(s_h; a_h) > 0$:

$$\sum_{k=1}^K d^k(s_h; a_h) \leq \frac{K}{2} d^{\text{cover}}(s_h; a_h); \text{ with } d^{\text{cover}}(s_h; a_h) := \max_{\pi} \frac{d_{h; \min}(\pi; s_h; a_h)}{(j Z_{h; \text{div}} + 1) H}; d^{\text{cover}}(s_h; a_h) > 0$$

where $Z_{h; \text{div}} := \sum_{a_h \in \mathcal{A}_h} f(s_h; a_h) \in \mathcal{S}_h; a_h \in \mathcal{A}_h; \forall s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h$; $d^e(s_h; a_h) := \min_{\pi} d^e(s_h; a_h)$ subject to $d^e(s_h; a_h) > 0$.

Proof. For arbitrary $j \in [H]$, we define:

$$N_{I_K}(s_h; a_h) := \sum_{k=1}^K \mathbb{I}[d^k(s_h; a_h) > 0]$$

In another word $N_{I_K}(s_h; a_h)$ denotes the number of optimal policies in the sequence, which can hit states s_h and take actions a_h .

Next, we define Z_h and Z_h^{insuff} as

$$\begin{aligned} Z_h &:= \sum_{a_h \in \mathcal{A}_h} f(s_h; a_h) \in \mathcal{S}_h; a_h \in \mathcal{A}_h; \forall s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h \\ Z_h^{\text{insuff}} &:= \sum_{a_h \in \mathcal{A}_h} f(s_h; a_h) \in \mathcal{S}_h; a_h \in \mathcal{A}_h; \forall s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h; \text{ s.t. } d^k(s_h; a_h) < \frac{K}{2(j Z_{h; \text{div}} + 1) H} \\ I_h^{\text{insuff}} &:= \sum_{k \in [K]} \mathbb{I}[d^k(s_h; a_h) > 0] \end{aligned}$$

In a word, Z_h is the collection of states actions reachable by at least one optimal policy, Z_h^{insuff} is a collection of "insufficiently hit" states actions at step h which are only covered by a small portion of optimal policies in the sequence, and $I_{1:H}^{insuff}$ is a collection of the index of the optimal policies in the sequence, which cover at least one state action pair.

Note that we must have $Z_h^{insuff} \subseteq Z_h$, because if one state action pair, (s_h, a_h) is reachable by arbitrary deterministic policy, then $N_k(s_h; a_h) = K$. Then, we have:

$$|I_{1:H}^{insuff}| < |Z_h^{insuff}| \frac{K}{2(|Z_h^{insuff}| + 1)H} \leq |Z_h^{insuff}| \frac{K}{2(|Z_h^{insuff}| + 1)H} \leq \frac{K}{2H}$$

We define $I_{1:H}^{suff} := I_K \cap \bigcap_{h=1}^H I_h^{insuff}$. Intuitively, $I_{1:H}^{suff}$ is the set including the indices of optimal policies in the sequence only hitting those states which are covered by most of the other optimal policies. In fact, $I_{1:H}^{suff}$ is non-empty since:

$$|I_{1:H}^{suff}| \geq K - \frac{K}{2H} \cdot H = \frac{K}{2}$$

We use $\pi_{1:H}^{suff}$ to denote the average mixture policy over $\{i \in I_{1:H}^{suff}\}$, a direct result is that:

$$\sum_{k=1}^K d^k(s_h; a_h) \pi_{1:H}^{suff}(s_h; a_h) = |I_{1:H}^{suff}| \sum_{i \in I_{1:H}^{suff}} \frac{1}{|I_{1:H}^{suff}|} d^i(s_h; a_h) \geq \frac{K}{2} \sum_{i \in I_{1:H}^{suff}} \frac{1}{|I_{1:H}^{suff}|} d^i(s_h; a_h)$$

On the other hand, for (s_h, a_h) such that $d^i(s_h; a_h) > 0$, we must have $(s_h, a_h) \in Z_h^{insuff}$, and therefore:

$$\sum_{k=1}^K d^k(s_h; a_h) \leq \frac{K}{2(|Z_h^{insuff}| + 1)H} d_{h, \min}(s_h; a_h)$$

Combining the above two inequalities, we finish the proof. \square

F Analysis of Pessimistic Value Iteration

In this section, we provide analysis for Alg. 3. Our analyses base on an extension of the Clipping Trick in [Simchowitz and Jamieson, 2019] into our setting.

F.1 Underestimation and Some Concrete Choices of Bonus Term

Lemma F.1 (Underestimation) Given a Bonus satisfying Cond. 4.4, for arbitrary data \mathcal{D}_k consisting of k trajectories by a sequence of policies $\{\pi_k; \dots; \pi_1\}$, by running Alg 3 with \mathcal{D}_k and the bonus term $b(\cdot; \cdot)$ returned by $\text{Bonus}(\mathcal{D}_k; k)$, on the event $\mathcal{E}_{\text{Bonus}}$ defined in Cond. 4.4:

$$|Z_h| \geq 2H; \quad \forall (s_h, a_h) \in Z_h, \quad Q_h(s_h; a_h) \geq Q_{h+1}^\phi(s_h; a_h) - Q_h(s_h; a_h) \quad (16)$$

where we use $Q_h = \{Q_{h,1}; \dots; Q_{h,H}\}$ to denote the greedy policy with Q_h .

Proof. We only prove the first inequality holds, since the second one holds directly because of the definition of optimal policy.

First of all, $V_{H+1} = 0 = V_{H+1}^\phi$ holds directly, which implies that Eq.(16) holds at step $h = H$ as a result of the deterministic reward function.

Now, we conduct the induction. Suppose Eq.(16) already holds for $h+1$, which implies that:

$$V_{h+1}(s_{h+1}) = Q_{h+1}^\phi(s_{h+1}; Q_{h+1}(s_{h+1})) - Q_{h+1}(s_{h+1}; Q_{h+1}(s_{h+1})) = V_{h+1}^\phi(s_{h+1}) \quad (17)$$

then, at step h , we have:

$$\begin{aligned} Q_h(s_h; a_h) - Q_h^\phi(s_h; a_h) &= P_h V_{h+1}(s_h; a_h) - b_h(s_h; a_h) - P_h V_{h+1}^\phi(s_h; a_h) \\ &= \underbrace{(P_h - P_h^\phi) V_{h+1}(s_h; a_h)}_{\text{part 1}} - b_h(s_h; a_h) + \underbrace{P_h (V_{h+1} - V_{h+1}^\phi)(s_h; a_h)}_{\text{part 2}} \end{aligned}$$

As we can see, part 1 is non-positive with probability 1 as a result of Cond. 4.4, while part 2 is also less than or equal to zero because of the induction condition in Eq.(17). \square

Choice 1: Naive Bound According to Hoeffding inequality, with probability $1 - \delta$, we have the following holds for each $s_h; a_h; h$

$$|V_h(s_h) - \mathbb{E}[V_h(s_h)]| \leq \sqrt{\frac{\log(2/\delta)}{N(s_h; a_h)}}$$

which implies that condition 4.4 holds with probability $1 - \delta$ as long as:

$$N(s_h; a_h) \geq \frac{2 \log(2/\delta)}{\epsilon^2}$$

Choice 2: Adaptive Bonus Term based on the ‘‘Bernstein Trick’’ One can also consider an analogue of the bonus term functions in Alg. 3 of [Simchowitz and Jamieson, 2019], which is originally designed for optimistic algorithms. We omit the discussions here.

F.2 Definition of ‘‘Surplus’’ in Pessimistic Algorithms and the Clipping Trick

We consider the pessimistic algorithm, and denote the estimation of value function as \hat{V}_h . We assume they are pessimistic estimation, i.e.:

$$\hat{V}_h(s_h) = Q_h(s_h; a_h) - \beta_h V_h^k(s_h) \quad \hat{Q}_h(s_h; a_h) = Q_h(s_h; a_h) - \beta_h V_h^k(s_h):$$

Definition F.2 (Definition of Surplus in Pessimistic Algorithm setting) We define the surplus in Pessimistic Algorithm setting:

$$E_{k;h}(s_h; a_h) = r(s_h; a_h) + P_h \hat{V}_{k;h+1}(s_h; a_h) - \hat{Q}_{k;h}(s_h; a_h):$$

Because of the underestimation, different from the surplus in overestimation cases [Simchowitz and Jamieson, 2019], here we flip the role between β and $r + PV$ to make sure the quantity is non-negative (with high probability).

Based on our definition, we have the following lemma:

Lemma F.3. Under the same condition as Lemma F.1, for arbitrary s_h , the policy π_k^{PVI} returned by Alg.3 satisfying:

$$V_h^{\pi_k^{PVI}}(s_h) - \hat{V}_{k;h}(s_h) = \mathbb{E}_k^{\pi_k^{PVI}} \left[\sum_{h^0=h}^H E_{k;h^0}(s_{h^0}; a_{h^0}) | s_h \right]$$

Moreover, for arbitrary optimal deterministic or non-deterministic policy, we have:

$$V_h(s_h) - \hat{V}_{k;h}(s_h) = V_h(s_h) - \hat{Q}_{k;h}(s_h; a_h) = \mathbb{E}_k \left[\sum_{h^0=h}^H E_{k;h^0}(s_{h^0}; a_{h^0}) \right]$$

Proof.

$$\begin{aligned} & V_h^{\pi_k^{PVI}}(s_h) - \hat{V}_{k;h}(s_h) \\ &= \mathbb{E}_{a_h}^{\pi_k^{PVI}} [r(s_h; a_h) + P_h V_{h+1}^{\pi_k^{PVI}}(s_h; a_h) - \hat{Q}_{k;h}(s_h; a_h) - P_h \hat{V}_{k;h+1}(s_h; a_h)] \\ &= \mathbb{E}_k^{\pi_k^{PVI}} [r(s_h; a_h) + P_h V_{h+1}^{\pi_k^{PVI}}(s_h; a_h) - \hat{Q}_{k;h}(s_h; a_h) + P_h (V_{h+1}^{\pi_k^{PVI}} - \hat{V}_{k;h+1})(s_h; a_h)] \\ &= \mathbb{E}_k^{\pi_k^{PVI}} \left[\sum_{h^0=h}^H E_{k;h^0}(s_{h^0}; a_{h^0}) | s_h \right] \end{aligned}$$

Besides, given arbitrary optimal policy π_k , we have:

$$V_h(s_h) - \hat{V}_{k;h}(s_h) = V_h(s_h) - \hat{Q}_{k;h}(s_h; a_h) = (V_h^{\pi_k} - \hat{Q}_{k;h})(s_h; a_h) \quad (\pi_k^{PVI} \text{ is greedy policy w.r.t. } \hat{Q}_{k;h})$$

$$= E_{a_h} [r(s_h; a_h) + P_h \mathbb{V}_{k;h+1}(s_h; a_h) \mathbb{Q}_{k;h}(s_h; a_h) + P_h (V_{h+1} \mathbb{V}_{k;h+1})(s_h; a_h)]$$

$$E_{h^0=h} [E_{k;h^0}(s_{h^0}; a_{h^0})]$$

□

Lemma F.4. Under the same condition as Lemma F.1, we have:

$$E_{k;h} \min_f H_{h+1; 2B_1} \frac{s \overline{\log(B_{2^k})}}{N_{k;h}(s_h; a_h)} g:$$

Proof.

$$E_{k;h} := r(s_h; a_h) + P_h \mathbb{V}_{k;h+1}(s_h; a_h) \mathbb{Q}_{k;h}(s_h; a_h)$$

$$= P_h \mathbb{V}_{k;h+1} \mathbb{P}_{k;h} \mathbb{V}_{k;h+1} + \mathbb{b}_{k;h}(s_h; a_h) \mathbb{2b}_{k;h}(s_h; a_h) \mathbb{2B}_1 \frac{s \overline{\log(B_{2^k})}}{N_{k;h}(s_h; a_h)}:$$

On the other hand, because the reward function is always local and \mathbb{Q} is always larger than zero, we have $E_{k;h}(s_h; a_h) \geq H_{h+1} - H$. □

In the following, we define

$$\tilde{E}_{k;h}(s_h; a_h) := \text{clip}[E_{k;h}(s_h; a_h)]_{\text{Clip}}:$$

where $\text{clip} := \frac{\min}{2H+2}$, and $\text{Clip}[x] := x - |x|$. Then, we recursively define

$$\mathbb{Q}_{k;h}(s_h; a_h) = E_{h^0=h} [r(s_h; a_h) - \tilde{E}_{k;h}(s_h; a_h) + P_h \mathbb{V}_{k;h+1}(s_h; a_h) | s_h; a_h]; \quad \mathbb{V}_{k;h}(s_h) := \mathbb{Q}_{k;h}(s_h; a_h)$$

Note that although different optimal policies and \mathbb{V}^e have the same optimal value, \mathbb{V} may no longer equal to \mathbb{V}^e because they may have different state occupancy. \mathbb{V} depends on E . Therefore, in the following, when we consider the for optimal policies, we will always specify which optimal policy we are referring to.

Lemma F.5 (Relationship between \mathbb{V} , V_k^{PVI} and $\mathbb{V}_{k;h}$). Under the same condition as Lemma F.1, for arbitrary optimal policy π , we have:

$$\mathbb{V}_{k;h}(s_h) \leq \mathbb{V}_{k;h}(s_h) + (H - h + 1) \text{Clip} \leq V_k^{\text{PVI}}(s_h) + (H - h + 1) \text{Clip}$$

Proof. Note that:

$$\tilde{E}_{k;h}(s_h; a_h) \leq E_{k;h}(s_h; a_h) \text{Clip}$$

Therefore,

$$V_h(s_h) - \mathbb{V}_h(s_h) = E_{h^0=h} [E_{k;h^0}(s_{h^0}; a_{h^0}) | s_h]$$

$$E_{h^0=h} [E_{k;h^0}(s_{h^0}; a_{h^0}) \text{Clip} | s_h]$$

$$V_h(s_h) \leq \min_f \mathbb{Q}_{k;h}(s_h; \cdot); \mathbb{V}_{k;h}(s_h) g + (H - h + 1) \text{Clip} \quad (\text{Lemma F.3})$$

$$V_h(s_h) \leq \min_f Q_h^{\text{PVI}}(s_h; \cdot); V_h^{\text{PVI}}(s_h) g + (H - h + 1) \text{Clip} \quad (\text{Underestimation (Lemma F.1)})$$

Therefore,

$$\mathbb{V}_{k;h}(s_h) \leq \mathbb{V}_{k;h}(s_h) + (H - h + 1) \text{Clip} \leq V_h^{\text{PVI}}(s_h) + (H - h + 1) \text{Clip}$$

□

F.3 Additional Lemma for the Analysis of the Regret of Alg^F when Optimal Deterministic Policies are non-unique

We first introduce a useful Lemma related to the clipping operator from [Simchowitz and Jamieson, 2019]

Lemma F.6 (Lemma B.3 in [Simchowitz and Jamieson, 2019]) Let $M \geq 2$, $a_1, \dots, a_m \geq 0$ and $\epsilon \in [0, 1]$. $\text{Clip}[\frac{1}{M} \sum_{i=1}^m a_i] \geq \frac{1}{M} \sum_{i=1}^m \text{Clip}[a_i]_{\frac{1}{2M}}$.

Next, based on the definition of d_{\min} in Eq.(10), we have the following Lemma:

Lemma F.7. Given arbitrary deterministic policy, if $\epsilon \geq 2^{-2}$, we have:

$$V_1(s_1) - V_1^*(s_1) \leq d_{\min} \cdot \min_{h \in [H]} \epsilon_h$$

Proof. We use π^* to denote the converted deterministic optimal policy, where is defined in Def. E.4. As a direct application of Lemma E.2, we have:

$$\begin{aligned} V_1(s_1) - V_1^*(s_1) &= V_1(s_1) - V_1^*(s_1) = E \left[\sum_{h=1}^{H-1} \mathbb{I}[\mathcal{E}_h] (V_h(s_h) - V_h^*(s_h)) \right] \\ &\leq \sum_{h=1}^{H-1} E \left[\mathbb{I}[\mathcal{E}_h] \right] \\ &\leq \sum_{h=1}^{H-1} \Pr(\mathcal{E}_{h_{\text{init}}}) \\ &\leq d_{\min} \end{aligned}$$

where we use \mathcal{E}_h to denote the event that at step h first disagrees with π^* , or equivalently, π^* take non-optimal action; in the second inequality, we define $\epsilon_h := \min_{s \in \mathcal{S}_h} |V_h(s) - V_h^*(s)|$; $s^* := \arg \min_{s \in \mathcal{S}_h} |V_h(s) - V_h^*(s)| > 0$. Besides, the last inequality is because:

$$\begin{aligned} \Pr(\mathcal{E}_{h_{\text{init}}}) &= \sum_{s \in \mathcal{S}_h} \mathbb{I}[|V_h(s) - V_h^*(s)| \geq \epsilon_h] d(s) \\ &= \sum_{s \in \mathcal{S}_h} \mathbb{I}[|V_h(s) - V_h^*(s)| \geq \epsilon_h] d(s) \\ &\leq \sum_{s \in \mathcal{S}_h} \mathbb{I}[|V_h(s) - V_h^*(s)| \geq \epsilon_h] d_{\min} \\ &\leq d_{\min} \end{aligned}$$

where the last step is because, according to the definition of ϵ_h , there is at least one $s \in \mathcal{S}_h$ such that $|V_h(s) - V_h^*(s)| = \epsilon_h$. \square

F.4 Upper Bound for the Regret of Alg^F

Theorem 4.5. By running Algorithm 3 with confidence level δ , a function B satisfying Condition 4.4, and a dataset $\mathcal{D} = \{f_1, \dots, f_k\}$ consisting of k complete trajectories generated by executing a sequence of policies π_1, \dots, π_k , on the event $\mathcal{E}_{\text{Bonus}}$ defined in Condition 4.4:

$$V_1(s_1) - V_1^{\text{PVI}}(s_1) \leq 2E \sum_{h=1}^{H-1} \text{Clip} \left(\min_{j \in [k]} \sum_{s \in \mathcal{S}_h} \frac{\log(B_{2^k})}{N_{k,h}(s; a_h)} \right) \cdot \epsilon_{\text{Clip}} \quad (1)$$

where π^* can be an arbitrary optimal policy, $\epsilon_{\text{Clip}} := \frac{\min_{j \in [k]} \epsilon_j}{2H+2}$ if $j = 1$ and $\epsilon_{\text{Clip}} := \frac{d_{\min} \cdot \min_{s \in \mathcal{S}_h} |V_h(s) - V_h^*(s)|}{2SAH}$ if $j > 1$, where $d_{\min} := \min_{h \in [H]} \min_{s \in \mathcal{S}_h} \sum_{a \in \mathcal{A}_h} d(s; a_h)$ subject to $d(s; a_h) > 0$.

Proof. We separately discuss the cases when there are unique or multiple deterministic optimal policies.

Case 1: Unique Deterministic Optimal Policy For arbitrary s_h , suppose $V_k^{PVI}(s_h) \geq V_h(s_h)$, we have:

$$\begin{aligned} V_h(s_h) - V_{k,h}(s_h) &= V_h(s_h) - \min_{a_h} \{ V_h(s_h) + \gamma \mathbb{E}_{k,h}(V_h(s_h) + Q_h(s_h; a_h)) - \min_{a_h} \{ V_h(s_h) + \gamma \mathbb{E}_{k,h}(V_h(s_h) + Q_h(s_h; a_h)) \} \} \\ &= \frac{1}{2} V_h(s_h) - V_{k,h}(s_h) + \frac{1}{2} V_h(s_h) - V_h^{PVI}(s_h) - \frac{\min}{2} \\ &= \frac{1}{2} V_h(s_h) - V_{k,h}(s_h) + \frac{h(s_h; V_k^{PVI}(s_h))}{2} - \frac{\min}{2} \\ &= \frac{1}{2} V_h(s_h) - V_{k,h}(s_h) \end{aligned}$$

Recall the definition of Events in Def.E.1, and note that when the optimal policy is unique, the events $E_{k,h}^{\leq}; E_{k,h}^{\geq}$ collapse to $E_{k,h}^{\leq}; E_{k,h}^{\geq}$, respectively. For arbitrary optimal policy, we have:

$$\begin{aligned} V_1(s_1) - V_1^{PVI}(s_1) &= V_1(s_1) - V_1^{PVI}(s_1) \\ &= I[E_{k,1}^{\leq}] V_1(s_1) - V_1^{PVI}(s_1) + I[E_{k,1}^{\geq}] V_1(s_1) - V_1^{PVI}(s_1) \\ &= I[E_{k,1}^{\leq}] \frac{1}{2} V_1(s_1) - V_{k,1}(s_1) + I[E_{k,1}^{\geq}] P(V_2 - V_2^{PVI})(s_1; a_1) \\ &\vdots \\ &= \frac{1}{2} E \left[\sum_{h=1}^H I[E_{k,h}^{\leq}] (V_h(s_h) - V_{k,h}(s_h)) \right] \end{aligned}$$

Besides, on the other hand,

$$\begin{aligned} V_1(s_1) - V_1^{PVI}(s_1) &= I[E_{k,1}^{\leq}] (V_1(s_1) - V_1^{PVI}(s_1)) + I[E_{k,1}^{\geq}] (V_1(s_1) - V_1^{PVI}(s_1)) \\ &= I[E_{k,1}^{\leq}] (V_1(s_1) - V_1^{PVI}(s_1)) + I[E_{k,1}^{\geq}] P_1(V_2 - V_2^{PVI})(s_1; a_1) \\ &= \vdots \\ &= E_k^{PVI} \left[\sum_{h=1}^H I[E_{k,h}^{\leq}] (V_h(s_h) - V_h^{PVI}(s_h)) \right] \\ &= E_k^{PVI} \left[\sum_{h=1}^H I[E_{k,h}^{\geq}] (V_h(s_h) - V_{k,h}(s_h)) \right] \end{aligned}$$

Combining the above two results and Lemma F.4, we finish the discussion for Case 1.

Case 2: Non-unique Optimal Deterministic Policies From Lemma F.3, we know that,

$$V_1(s_1) - V_1^{PVI}(s_1) = V_1(s_1) - V_{k,1}(s_1) - E \left[\sum_{h=1}^H E_{k,h}(s_h; a_h) \right]$$

where π can be arbitrary optimal policy. Combining with Lemma F.7, we know that:

$$\begin{aligned} V_1(s_1) - V_1^{PVI}(s_1) &\leq \text{Clip} \left[E \left[\sum_{h=1}^H E_{k,h}(s_h; a_h) \right] \right] d_{\min} - \min \\ &= \frac{1}{2} \sum_{h=1}^H \sum_{s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h} \text{Clip} [d(s_h; a_h) E_{k,h}(s_h; a_h)] \frac{d_{\min} - \min}{2SAH} \\ &= \frac{1}{2} E \left[\sum_{h=1}^H \text{Clip} [E_{k,h}(s_h; a_h)] \frac{d_{\min} - \min}{2SAH} \right] \end{aligned} \quad (\text{Lemma F.6})$$

where the last inequality is because $\text{Clip}[x] \leq \text{Clip}[x] + \epsilon$ as long as $\epsilon < 1$. Combining with Lemma F.4, we finish the proof. \square

Next, we introduce a useful Lemma from [Dann et al., 2017]:

Lemma F.8 (Lemma 7.4 in [Dann et al., 2017]) Let F_i for $i = 1, \dots, n$ be a filtration and X_1, \dots, X_n be a sequence of Bernoulli random variables with $\Pr(X_i = 1 | F_{i-1}) = P_i$ with P_i being F_{i-1} -measurable and X_i being F_i -measurable. It holds that

$$\Pr\left(\sum_{i=1}^n X_i < \sum_{i=1}^n P_i - W\right) \leq e^{-W}$$

Definition of Good Events We first introduce some notations about good events which holds with high probability. We override the definition in Cond. 4.4 by assigning $\beta_k = 1/k$ at iteration k , i.e.:

$$E_{\text{Bonus}_k} := \bigcap_{h \in [H]; s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h} \bigcap_{j=1}^n \mathbb{P}_{k,h} V_{k,h+1}(s_h; a_h) - \mathbb{P}_h V_{k,h+1}(s_h; a_h) < b_{k,h}(s_h; a_h) g$$

$$\forall b_{k,h}(s_h; a_h) \leq B_1 \frac{\log(B_2/k)}{N_{k,h}(s_h; a_h)};$$

with $b_k = (b_{k,1}, \dots, b_{k,H}) \in \text{Bonus}(D_k; 1/k)$:

Besides, we use $E_{\text{Con};k}$ to denote the concentration event that

$$E_{\text{Con};k} := \bigcap_{h \in [H]; s_h \in \mathcal{S}_h; a_h \in \mathcal{A}_h} \mathbb{P}_{k,h}(s_h; a_h) \leq \frac{1}{2} \sum_{k^0=1}^k d_{k^0}(s_h; a_h) \log(\text{SAH}_k) g$$

Finally, we use $E_{\text{Alg}^0;k}$ to denote the good events that the regret of Alg only at the level k :

$$E_{\text{Alg}^0;k} := \sum_{R=1}^k V_1(s_1) - V_1^R(s_1) < C_1 + C_2 \log k$$

Based on Cond. 4.4, Cond. 4.6 and Lemma F.8, we have:

$$\Pr(E_{\text{Alg}^0;k}) \geq 1 - \frac{1}{k}; \quad \Pr(E_{\text{Bonus}_k}) \geq 1 - \frac{1}{k}$$

$$\Pr(E_{\text{Con};k}) \geq 1 - \text{SAH} \exp(-\log(\text{SAH}_k)) = 1 - \frac{\text{SAH}}{(\text{SAH}_k)} \geq 1 - \frac{1}{k}$$

Lemma F.9. [One Step Sub-optimality Gap Conditioning on Good Events] At iteration k on the good events $E_{\text{Bonus}_k}; E_{\text{Con};k}$ and $E_{\text{Alg}^0;k}$, the sub-optimality gap of F_k can be upper bounded by:

(i) when $j = 1$ (i.e. the optimal deterministic policy is unique):

$$V_1(s_1) - V_k^E(s_1) \leq 2E_{h=1}^h \text{Clip} \left[\sum_{h=1}^h \mathbb{P}_{k,h} V_{k,h+1}(s_h; a_h) - \mathbb{P}_h V_{k,h+1}(s_h; a_h) \right] H + I[k \leq \text{SAH}_k] B_1 \frac{\log(B_2/k)}{k d(s_h; a_h)} \quad \text{ii}$$

(ii) when $j > 1$ (i.e. there are multiple optimal deterministic policy):

$$V_1(s_1) - V_k^E(s_1) \leq 2E_{\text{cover}}^h \text{Clip} \left[\sum_{h=1}^h \mathbb{P}_{k,h} V_{k,h+1}(s_h; a_h) - \mathbb{P}_h V_{k,h+1}(s_h; a_h) \right] H + I[k \leq \text{SAH}_k^{\text{cover}}] B_1 \frac{\log(B_2/k)}{k \mathfrak{C}^{\text{cover}}(s_h; a_h)} \quad \text{ii}$$

where $\text{Clip} := \frac{\min}{2H+2}$ and $\text{Clip}^0 := \frac{d_{\min} - \min}{2\text{SAH}}$; cover and $\mathfrak{C}^{\text{cover}}(s_h; a_h)$ are defined in Thm. 4.8; besides,

$$c_{s_h; a_h} := c \frac{(C_1 + C_2)}{d(s_h; a_h)_{\min}} \log \frac{\text{SAH} (C_1 + C_2)}{d(s_h; a_h)_{\min}}; \quad c_{s_h; a_h}^{\text{cover}} := c^0 \frac{(C_1 + C_2)}{\mathfrak{C}^{\text{cover}}(s_h; a_h)_{\min}} \log \frac{\text{SAH} (C_1 + C_2)}{\mathfrak{C}^{\text{cover}}(s_h; a_h)_{\min}}$$

for some constant c, c^0 .

Proof. We first discuss the case when $j = 1$.

Case 1: unique optimal deterministic policy As a result of Thm. 4.5, on the event $\mathcal{E}_{\text{Con};k}$, we show that the sub-optimality gap of V_1^k can be upper bounded by:

$$V_1(s_1) - V_1^k(s_1) \leq 2 \mathbb{E} \left[\sum_{h=1}^H \mathbb{E}_{k,h}(s_h; a_h) \right] = \sum_{h=1}^H \mathbb{E}_{s_h; a_h} [d(s_h; a_h) \mathbb{E}_{k,h}(s_h; a_h)]$$

Because of Lemma F.4, the above further implies that:

$$V_1(s_1) - V_1^k(s_1) \leq 2 \mathbb{E} \left[\min_{h=1, \dots, H} \left(\frac{1}{2} \sum_{k^0=1}^k d^{k^0}(s_h; a_h) \log(\text{SAH}k) + \frac{1}{\min} (C_1 + C_2 \log k) \right) \right]$$

Because of Thm. 4.7, on the event $\mathcal{E}_{\text{Con};k}$ and $\mathcal{E}_{\text{Alg}^0;k}$, we further have:

$$\begin{aligned} N_{k,h}(s_h; a_h) &\leq \frac{1}{2} \sum_{k^0=1}^k d^{k^0}(s_h; a_h) \log(\text{SAH}k) \\ &\leq \frac{1}{2} \sum_{k^0=1}^k d^{k^0}(s_h; a_h) \log(\text{SAH}k) + \frac{1}{\min} (C_1 + C_2 \log k) \\ &\leq \frac{k}{2} d(s_h; a_h) \log(\text{SAH}k) + \frac{1}{\min} (C_1 + C_2 \log k) \end{aligned}$$

Now, we define that,

$$c_{s_h; a_h} := \inf_{t: 8t^0} \frac{1}{t} d(s_h; a_h) \log(\text{SAH}t) + \frac{1}{\min} (C_1 + C_2 \log t)g$$

there must exists a constant independent with $C_1; C_2; ; d(s_h; a_h)$ and \min , such that:

$$8h \geq 2[H]; s_h \geq 2 S_h; a_h \geq 2 A_h; \quad c_{s_h; a_h} := c \frac{(C_1 + C_2)}{d(s_h; a_h) \min} \log \frac{\text{SAH} (C_1 + C_2)}{d(s_h; a_h) \min}.$$

Easy to check that, for arbitrary $s_h; a_h$, on the good events, we can verify that

$\frac{k}{4} d(s_h; a_h) - \frac{s_h; a_h}{4} d(s_h; a_h) - \frac{c}{4} > 0$, and as a result, we have:

$$\begin{aligned} V_1(s_1) - V_1^k(s_1) &\leq \sum_{h=1}^H \mathbb{E} \left[\text{Clip} \left(\min_{h=1, \dots, H} \left(\frac{1}{2} \sum_{k^0=1}^k d^{k^0}(s_h; a_h) \log(\text{SAH}k) + \frac{1}{\min} (C_1 + C_2 \log k) \right) \right) \right] \\ &\leq \sum_{h=1}^H \mathbb{E} \left[\text{Clip} \left(I[k < s_h; a_h]g + I[k \geq s_h; a_h] B_1 \frac{\log(B_2k)}{kd(s_h; a_h)} \right) \right] \end{aligned}$$

Case 2: multiple optimal deterministic policies The discussion are similar. As a result of Thm. 4.8, on the event $\mathcal{E}_{\text{Con};k}$ and $\mathcal{E}_{\text{Alg}^0;k}$, we further have:

$$\begin{aligned} N_{k,h}(s_h; a_h) &\leq \frac{k}{4} \max_f \frac{d_{h; \min}(s_h; a_h)}{(jZ_h; \text{div}j + 1)H}; d^{\text{cover}}(s_h; a_h)g + \log(\text{SAH}k) + \frac{1}{\min} \sum_{k=1}^k V_1(s_1) - V_1^k(s_1) \\ &\leq \frac{k}{4} d^{\text{cover}}(s_h; a_h) \log(\text{SAH}k) + \frac{1}{\min} (C_1 + C_2 \log k) \end{aligned}$$

Similarly, we define that,

$$c_{s_h; a_h}^{\text{cover}} := \inf_{t: 8t^0} \frac{t}{8} d^{\text{cover}}(s_h; a_h) \log(\text{SAH}t) + \frac{1}{\min} (C_1 + C_2 \log t)g$$

there must exists a constant independent with $C_1; C_2; ; d^{\text{cover}}(s_h; a_h)$ and \min , such that:

$$8h \geq 2[H]; s_h \geq 2 S_h; a_h \geq 2 A_h; \quad c_{s_h; a_h}^{\text{cover}} := c^0 \frac{(C_1 + C_2)}{d^{\text{cover}}(s_h; a_h) \min} \log \frac{\text{SAH} (C_1 + C_2)}{d^{\text{cover}}(s_h; a_h) \min}.$$

Then, we have:

$$\begin{aligned}
 & \lim_{K \rightarrow \infty} \sum_{k=1}^K \sum_{h=1}^h \mathbb{E} \left[\text{Clip} \left[\frac{1}{k} \sum_{k=1}^k \sum_{h=1}^h \log(B_2 k) \right] \right] \\
 &= 2 \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right] + 2 \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right] \\
 &= 2 \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right] + 2 \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right] \\
 &= 2 \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right] + 2 \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right]
 \end{aligned}$$

For the first part, we have:

$$\begin{aligned}
 \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right] &= \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \\
 &= \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \\
 &= \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k)
 \end{aligned}$$

For the second part, we have:

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \right] \\
 &= \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \\
 &= \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \\
 &= \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k)
 \end{aligned}$$

(Lemma F.10)

$$\sum_{h=1}^h \sum_{k=1}^k \log(B_2 k)$$

where in the last step, we drop the term $\frac{1}{k} \sum_{k=1}^k \log(B_2 k)$, and c_2 is a constant.

Combining the above results, we have:

$$\begin{aligned}
 & \mathbb{E}_{\text{Alg}^0; M; \text{Alg}^E} \left[\sum_{k=1}^K V_1(s_1) - V_1^E(s_1) \right] \\
 &= \sum_{k=1}^K \frac{3H}{k} + c \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k) \\
 &+ c_2 \sum_{h=1}^h \sum_{k=1}^k \log(B_2 k)
 \end{aligned}$$

$$\frac{3H}{1} + C_{Alg^E} \sum_{h=1}^H \sum_{\substack{s_h; a_h: \\ d(s_h; a_h) > 0}} X \frac{(C_1 + C_2)}{\min} \log \frac{SAH (C_1 + C_2)}{d(s_h; a_h) \min} + \frac{B_1 H}{\min} \log \frac{B_2 H}{d(s_h; a_h) \min}$$

where C_{Alg^E} is some constant.

Case 2: $j > 1$ (Non-Unique Optimal Policy) Similar to the discussion above, we define:

$$e_{s_h; a_h}^{cover, 0, Clip} := \inf_{t; 8t^0} f_{B_1} \frac{8 \log(B_2 t)}{t e^{cover(s_h; a_h)}} < e_{Clip}^0$$

Recall that $e_{Clip}^0 := d_{\min} \min = (2SAH)$, it's easy to verify that, there exists a constant, such that,

$$e_{s_h; a_h}^{cover, 0, Clip} := C_{Clip} \frac{(SAH)^2}{e^{cover(s_h; a_h)} (d_{\min} \min)^2} \log \frac{B_2 SAH}{e^{cover(s_h; a_h)} d_{\min} \min}$$

Following a similar discussion, we have:

$$\begin{aligned} \lim_{K \uparrow} \sum_{k=1}^K 2E_{cover} \sum_{h=1}^H I[k < e_{s_h; a_h}^{cover}] g_H + I[k > e_{s_h; a_h}^{cover}] B_1 \frac{8 \log(B_2 k)}{k e^{cover(s_h; a_h)}} e_{Clip}^0 \\ 2E_{cover} \sum_{h=1}^H \sum_{k=1}^K [H] + 2 \sum_{h=1}^H \sum_{\substack{s_h; a_h: \\ d(s_h; a_h) > 0}} X \frac{B_1^p}{8 d^{cover(s_h; a_h)}} \int_{x = e_{s_h; a_h}^{cover}}^r \frac{\log(B_2 x)}{x} dx \\ C_0 \frac{H (C_1 + C_2)}{\min} \sum_{h=1}^H \sum_{\substack{s_h; a_h: \\ d(s_h; a_h) > 0}} X \log \frac{SAH (C_1 + C_2)}{e^{cover(s_h; a_h)} \min} \\ + C_0 \frac{B_1 SAH}{d_{\min} \min} \sum_{h=1}^H \sum_{\substack{s_h; a_h: \\ d(s_h; a_h) > 0}} X \log \frac{B_2 SAH}{d_{\min} \min} \quad (\text{Note that } e^{cover} d^{cover}) \end{aligned}$$

Therefore, we have:

$$\begin{aligned} E_{Alg^0; M; Alg^E} \sum_{k=1}^K V_1(s_1) - V_1^E(s_1) \\ \frac{3H}{1} + C_{Alg^E} \sum_{h=1}^H \sum_{\substack{s_h; a_h: \\ d(s_h; a_h) > 0}} X \frac{(C_1 + C_2)}{\min} \log \frac{SAH (C_1 + C_2)}{e^{cover(s_h; a_h)} \min} + \frac{B_1 SAH}{d_{\min} \min} \log \frac{B_2 SAH}{d_{\min} \min} \end{aligned}$$

□

Lemma F.10 (Computation of Integral) Suppose $1, b, a, e = p$ then we have:

$$\int_a^b \frac{\log px}{x} dx = 2 \left(\frac{1}{b \log pb} - \frac{1}{a \log pa} \right)$$

Proof.

$$\int_a^b \frac{\log px}{x} dx = \int_a^b \frac{1}{x \log px} + \int_a^b \frac{\log px}{x} dx = 2 \int_a^b \left(\frac{1}{x \log px} \right)^0 = 2 \left(\frac{1}{b \log pb} - \frac{1}{a \log pa} \right)$$

□

G Doubling Trick for Alg^O Satisfying Cond. G.1

As we briefly mentioned in Sec.4.2.2, Cond.4.6 may not holds for some algorithms with near-optimal regret guarantees. For example, in [Simchowitx and Jamieson, 2019, Xu et al., 2021, Dann et al., 2021], although these algorithms are anytime, they require a confidence interval δ as input at the beginning of the algorithm and fix it during the running, which we abstract into the Cond.G.1 below:

Condition G.1 (Alternative Condition of Alg^O). *Alg^O is an algorithm which returns a deterministic policies $\pi_{\mathbb{k}}^O$ at each iteration \mathbb{k} , and for arbitrary fixed $k \geq 2$, with probability $1 - \delta$, we have the following holds:*

$$\prod_{\mathbb{k}=1}^k V_1^*(s_1) - V_1^{\pi_{\mathbb{k}}^O}(s_1) \leq C_1 + C_2 \log \frac{k}{\delta}$$

where C_1, C_2 are some parameters depending on S, A, H and $\Delta_h(s_h, a_h)$ and independent with k .

As a result, no matter how small δ is chosen at the beginning, when $k \geq \lceil 1/\delta \rceil$, the Cond. 4.6 can not be directly guaranteed. To overcome this issue, we present a new framework in Alg 5 inspired by doubling trick.

Algorithm 5: Tiered RL Algorithm with Doubling Trick

```

1 Input:  $\alpha > 1$ .
2  $K_0 = 1, \quad k = 1, \quad \pi_{1,1}^E \leftarrow \text{Alg}^E(\{\})$ .
3 for  $n = 1, 2, \dots$  do
4    $K_n \leftarrow 2K_{n-1}, \quad \delta_{n-1} = 1/K_n^\alpha, \quad D_{n,1} \leftarrow \{\}$ 
5   for  $k = 1, \dots, K_n$  do
6     // Here we do not update  $\pi^E$ 
7      $\pi_{n,k+1}^O \leftarrow \text{Alg}^O(D_{n,k}, \delta_n)$ .
8      $\pi_{n,k+1}^E = \begin{cases} \pi_{n-1, K_{n-1}/2 + \lceil k/2 \rceil}^E, & \text{If } k \leq K_n/2, \\ \text{Alg}^E(D_{n,k}, 1/k^\alpha), & \text{Otherwise.} \end{cases}$ 
9      $\tau_{k+1} \sim \pi_{n,k+1}^O$ 
10     $D_{n,k+1} = D_{n,k} \cup \tau_{n,k+1}$ 
11  end
12 end

```

The basic idea is to iteratively run Alg^O satisfying Cond.G.1 from scratch while gradually doubling the number of iterations (i.e. K_n) and shrinking the confidence level δ_n rather than running with a fixed δ forever. Besides, another crucial part is the computation of π_k^E . Instead of continuously updating π^E with the data collected before, we only update the exploitation policy when $k \geq K_n/2$ for each outer loop n . As we will discuss in Lemma G.2, Alg^E _{n,k} will behave as if the dataset is generated by another online algorithm satisfying Cond. 4.6, and therefore, the analysis based on Cond. 4.6 can be adapted here, which we summarize to Thm. G.3 below.

Lemma G.2. *By running an algorithm satisfying Cond. G.1 in Alg. 5 as Alg^O, for arbitrary $n \geq 1$ and $K_n/2 + 1 \leq k < K_n/2$, we have:*

$$\Pr\left(\prod_{k=1}^k V^* - V^{\pi_{n:k}^O} > C'_1 + \alpha C'_2 \log k\right) \leq 1/k^\alpha$$

with $C'_1 = C_1 + (\alpha + 1)C_2 \log 2$ and $C'_2 = \frac{\alpha+1}{\alpha}C_2$.

Proof. Based on Cond. G.1, we know that:

$$\Pr\left(\prod_{k^O=1}^k V_1^*(s_1) - V_1^{\pi_{k^O}^O}(s_1) > C_1 + C_2 \log \frac{k}{\delta_n}\right) \leq \delta_n$$

Since $\delta_n = 1/K_n^\alpha$ and $k \geq K_n/2$, we have:

$$\begin{aligned}
& \Pr(\prod_{k^0=1}^{\times^k} V_1^*(s_1) - V_1^{\pi_{k^0}}(s_1) > C_1 + (1 + \alpha)C_2 \log 2k) \\
&= \Pr(\prod_{k^0=1}^{\times^k} V_1^*(s_1) - V_1^{\pi_{k^0}}(s_1) > C_1 + C_2 \log(2k)^{1+\alpha}) \\
&\leq \Pr(\prod_{k^0=1}^{\times^k} V_1^*(s_1) - V_1^{\pi_{k^0}}(s_1) > C_1 + C_2 \log \frac{k}{\delta_n}) \quad ((2k)^{1+\alpha} \geq 2kK_n^\alpha > k/\delta_n) \\
&\leq \delta_n \leq 1/k^\alpha
\end{aligned}$$

□

Now, we are ready to upper bound the regret of Alg^E :

Theorem G.3. *By choosing an arbitrary algorithm satisfying Cond. G.1 as Alg^O , choosing Alg. 3 as Alg^E and choosing a bonus function satisfying Cond. 4.4 as **Bonus**, the Pseudo regret of $\pi_{n,k}^E$ in Alg. 5 can be upper bounded by:*

(i) $|\Pi^*| = 1$ (unique optimal deterministic policy):

$$\begin{aligned}
& \mathbb{E}[\prod_{n=1}^{\times^N} \prod_{k=1}^{\times^k} V_1^*(s_1) - V^{\pi_{n,k}^E}] \leq 2H + \frac{9\alpha H}{\alpha - 1} \\
&+ 3c_{\text{Alg}^E} \cdot \prod_{h=1}^{\times^H} \prod_{\substack{s_h, a_h: \\ d(s_h, a_h) > 0}}^{\times} \frac{\alpha(C'_1 + C'_2)}{\Delta_{\min}} \log \frac{\alpha \text{SAH}(C'_1 + C'_2)}{d^\pi(s_h, a_h) \Delta_{\min}} + \frac{\alpha B_1 H}{\Delta_{\min}} \log \frac{\alpha B_2 H}{d^\pi(s_h, a_h) \Delta_{\min}}
\end{aligned}$$

(ii) $|\Pi^*| > 1$ (non-unique optimal deterministic policies):

$$\begin{aligned}
& \mathbb{E}[\prod_{n=1}^{\times^N} \prod_{k=1}^{\times^k} V_1^*(s_1) - V^{\pi_{n,k}^E}] \leq 2H + \frac{9\alpha H}{\alpha - 1} \\
&+ 3c'_{\text{Alg}^E} \cdot \prod_{h=1}^{\times^H} \prod_{\substack{s_h, a_h: \\ d_{\text{cover}}(s_h, a_h) > 0}}^{\times} \frac{\alpha(C'_1 + C'_2)}{\Delta_{\min}} \log \frac{\alpha \text{SAH}(C'_1 + C'_2)}{d^{\pi_{\text{cover}}}(s_h, a_h) \Delta_{\min}} + \frac{\alpha B_1 \text{SAH}}{d_{\min} \Delta_{\min}} \log \frac{\alpha B_2 \text{SAH}}{d_{\min} \Delta_{\min}}
\end{aligned}$$

where $C'_1 = C_1 + (\alpha + 1)C_2 \log 2$ and $C'_2 = \frac{\alpha+1}{\alpha}C_2$.

Remark G.4 ($O(\log^2 K)$ -Regret of Alg^O). *Although the regret of Alg^E stays constant under this framework, it is easy to verify that the pseudo-regret of Alg^O will be $O(\log^2 K)$ as a result of the doubling trick, which is worse than $O(\log K)$ up to a factor of $\log K$. Therefore, more rigorously speaking, the regret of Alg^O will be almost near-optimal.*

Proof. The key observation is that one can decompose the total expected regret into two parts:

$$\begin{aligned}
\mathbb{E}[\prod_{n=1}^{\times^N} \prod_{k=1}^{\times^k} V^* - V^{\pi_{n,k}^E}] &= \mathbb{E}[\prod_{n=1}^{\times^N} \prod_{k=1}^{\times^{K_n/2}} V^* - V^{\pi_{n,k}^E}] + \mathbb{E}[\prod_{n=1}^{\times^N} \prod_{k=K_n/2+1}^{\times^k} V^* - V^{\pi_{n,k}^E}] \\
&= 2\mathbb{E}[\prod_{n=0}^{\times^{N-1}} \prod_{k=K_n/2+1}^{\times^k} V^* - V^{\pi_{n,k}^E}] + \mathbb{E}[\prod_{n=1}^{\times^N} \prod_{k=K_n/2+1}^{\times^k} V^* - V^{\pi_{n,k}^E}] \\
&\leq 2K_0 H + 3\mathbb{E}[\prod_{n=1}^{\times^N} \prod_{k=K_n/2+1}^{\times^k} V^* - V^{\pi_{n,k}^E}] \tag{18}
\end{aligned}$$

Therefore, all we need to do is to upper bound the second part of Eq.(18). As a result of Lemma G.2, we can apply Lemma F.9 to upper bound the regret of the policy sequence $\{\{\pi_{n,k}^E\}_{n=1}^N\}_{k=K_n/2+1}^{K_n}$,

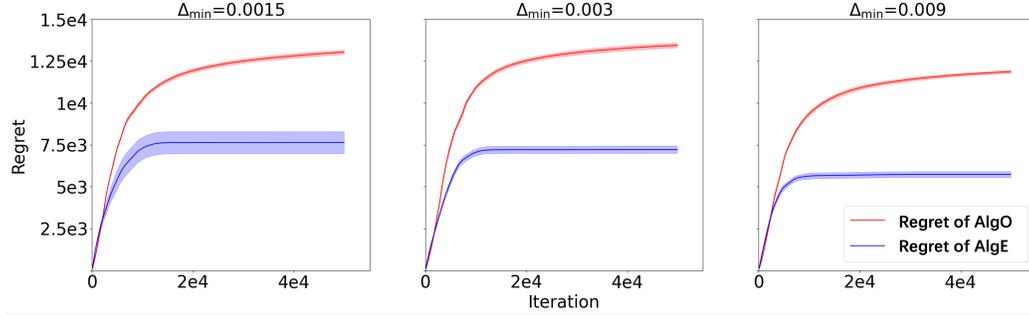


Figure 2: Simulation results with $S = A = H = 5$ and different Δ_{\min} , averaged over 10 different random seeds. Error bars show double the standard errors, which correspond to 95% confidence intervals. Our choice of Alg^E can achieve constant regret as predicted by theory. We can also see the tendency that larger Δ_{\min} will result in smaller accumulative regret.

since the $\mathcal{P}\text{Cond. 4.6}$ is satisfied when generating those policies. Therefore, the Pseudo regret $\mathbb{E}[\sum_{n=1}^N \sum_{k=K_n/2+1}^{K_n} V^* - V^{\pi_{n:k}^E}]$ can be upper bounded by extending the results in Thm. 4.9 here, and we finish the proof. \square

H Experiments

H.1 Experiment Setup

Environment We test our algorithms in tabular MDPs with randomly generated transition and rewards functions. To generate the MDP, for each layer h and each state action pair (s_h, a_h) , we first sample a random vector $\mathcal{P}(\cdot|s_h, a_h)$, where each element is uniformly sampled from $\{1, 2, 3, \dots, 10\}$, and then normalize it to a valid probability vector. Besides, the reward function is set to $\xi/10$ where ξ is randomly generated from $\{1, 2, \dots, 10\}$ to make sure it locates in $[0, 1]$.

Algorithm We implement the StrongEuler algorithm in [Simchowitz and Jamieson, 2019] as Alg^O and construct the same adaptive bonus term (Alg. 3 in [Simchowitz and Jamieson, 2019]) for Alg^E to match Cond. 4.4. Although for the convenience of analysis, in our Framework 1, we do not consider to use the data generated by Alg^E , in experiments, we use both τ^O and τ^E , which slightly improves the performance. Besides, in practice, we observe that the bonus term is quite loose, and it will take a long time before the estimated Q/V value fallen in the interval $[0, H]$, which is the value range of true value functions. Therefore, we introduce a multiplier α and adjust the bonus term from $b_{k,h}$ to $\alpha \cdot b_{k,h}$, and set $\alpha = 0.25$ for both Alg^O and Alg^E .

H.2 Results

We test the algorithms in tabular MDPs with $S = A = H = 5^5$. Although the minimal gap Δ_{\min} is hard to control since we generate the MDP in a random way, we filter out three random seeds in MDP construction, which correspond to minimal gaps (approximately) equal to 0.0015, 0.003 and 0.009, respectively. We report the accumulative regret in Fig. 2.

As predicted by our theory, Alg^E can indeed achieve constant regret in contrast with the continuously increasing regret of Alg^O , which demonstrates the advantage of leveraging tiered structure.

⁵The code can be find in <https://github.com/jiaweihsuang/Tiered-RL-Experiments>.