

REINFORCEMENT LEARNING WITH SEGMENT FEEDBACK

Anonymous authors

Paper under double-blind review

ABSTRACT

Classic reinforcement learning (RL) assumes that an agent can observe a reward for each state-action pair. However, in practical applications, it is often difficult and costly to collect a reward for each state-action pair. While there have been several works considering RL with trajectory feedback, it is unclear if trajectory feedback is inefficient for learning when trajectories are long. In this work, we consider a model named RL with segment feedback, which offers a general paradigm filling the gap between per-state-action feedback and trajectory feedback seamlessly. In this model, we consider an episodic Markov decision process (MDP), where each episode is equally divided into m segments, and the agent observes reward feedback only at the end of each segment. Under this model, we study two popular feedback settings: binary feedback and sum feedback, where the agent observes a binary outcome and a reward sum according to the underlying reward function, respectively. To investigate the impact of the number of segments m on learning performance, we design efficient algorithms and establish regret upper and lower bounds for both feedback settings. Our theoretical and empirical results show that: under binary feedback, increasing the number of segments m decreases the regret at an exponential rate; in contrast, surprisingly under sum feedback, increasing m does not reduce the regret significantly.

1 INTRODUCTION

Reinforcement learning (RL) is a class of sequential decision-making algorithms, where an agent interacts with an unknown environment through time with the goal of maximizing the obtained reward. RL is applied to extensive fields such as robotics, autonomous driving and game playing.

In classic RL, when the agent takes an action in a state, the environment will provide a reward for this state-action pair. However, in real-world applications, it is often difficult and costly to collect a reward for each state-action pair. For example, in robotics, when we instruct a robot to scramble eggs, it is hard to specify a reward for each individual action. In autonomous driving, it is difficult and onerous to evaluate each action, considering multiple criteria including safety, comfort and speed.

Motivated by this fact, there have been several works that consider RL with trajectory feedback (Efroni et al., 2021; Chatterji et al., 2021). In these works, the agent observes a reward signal only at the end of each episode, instead of at each step, with the signal indicating the quality of the trajectory generated during the episode. While these works mitigate the issue of impractical per-step reward feedback in classic RL, the relationship between the frequency of feedback and the performance of RL algorithms is unknown. In particular, if for example we get feedback twice in each trajectory, does that significantly improve performance over once per trajectory feedback?

To answer this question, we study a general model called RL with segment feedback, which bridges the gap between per-state-action feedback in classic RL (Sutton & Barto, 2018) and trajectory feedback in recent works (Efroni et al., 2021; Chatterji et al., 2021). In this model, we consider an episodic Markov decision process (MDP), where an episode is equally divided into m segments. In each episode, at each step, the agent first observes the current state, and takes an action, and then transitions to a next state according to the transition distribution. The agent *observes a reward signal at the end of each segment*. Under this model, we consider two reward feedback settings: binary feedback and sum feedback. In the binary feedback setting, the agent observes a binary outcome (e.g., thumbs up/down) generated by a sigmoid function of the reward on this segment. In the sum

054 feedback setting, the agent observes the sum of the rewards over this segment. In our model, the
 055 agent needs to learn the underlying reward function (i.e., the expected reward as a function of states
 056 and actions) from binary or sum segment feedback, and maximize the expected reward achieved.
 057 (Reviewer NUmy) While Tang et al. (2024) also studied this segment model before (they called it
 058 bagged reward), their work is mostly an empirical work, and did not provide theoretical guarantees
 059 for algorithms and rigorously reveal the influence of segments on learning.

060 This model is applicable to many scenarios involving human queries. For instance, in autonomous
 061 driving, a driving trajectory is often divided into several segments, and human annotators are asked to
 062 provide feedback for each segment, e.g., thumbs up/down. Compared to state-action pairs or whole
 063 trajectories, segments are easier and more efficient to evaluate, since human annotators can focus on
 064 and rate behaviors in each segment, e.g., passing interactions, reversing the car and park.

065 In this segment model, there is an interesting balance between the number of segments (queries to
 066 humans) and the collected observations, i.e., we desire more observations, but we also want to reduce
 067 the number of queries. Therefore, in this problem, it is critical to investigate the trade-off between
 068 the benefits brought by segments and the increase of queries, which essentially comes down to a
 069 question: *How does the number of segments m impact learning performance?*

070 To answer this question, we design efficient algorithms for binary and sum feedback settings in both
 071 known and unknown transition cases. Regret upper and lower bounds are provided to rigorously show
 072 the influence of the number of segments on learning performance. We also present experimental
 073 results to validate our theoretical findings.

074 (Reviewer NUmy) (Reviewer R7pS) Note that studying RL with equal segments is an important
 075 start point and serves as a foundation for further investigation on more general models and analysis for
 076 RL with unequal segments. Even for the equal segment case, this problem is already very challenging:
 077 (i) This problem cannot be solved by applying prior trajectory feedback works, e.g., (Efroni et al.,
 078 2021), since they use the martingale property of subsequent trajectories in analysis, while subsequent
 079 segments are not a martingale due to the dependence of segments within a trajectory. (ii) In prior
 080 trajectory feedback works (Efroni et al., 2021; Chatterji et al., 2021), there exists a gap between upper
 081 and lower bounds for sum feedback, and there is no lower bound for binary feedback. This fact poses
 082 a great challenge for figuring out the essential influence of the number of segments m on learning
 083 performance, since one cannot get too many hints from prior works.

084 Our work resolves the above challenges and makes the following contributions:

- 086 1. We study a general model called RL with segment feedback, which bridges the gap between
 087 per-state-action feedback in classic RL and trajectory feedback seamlessly. Under this
 088 model, we consider two feedback settings: binary feedback and sum feedback.
- 089 2. For binary feedback, we design computationally and sample efficient algorithms SegBiTS
 090 and SegBiTS-Tran for known and unknown transitions, respectively. We provide regret up-
 091 per and lower bounds which depend on $\exp(\frac{Hr_{\max}}{2m})$, where H is the length of each episode,
 092 and r_{\max} is the universal upper bound of rewards. Our results exhibit that under binary
 093 feedback, increasing the number of segments m significantly helps accelerate learning.
- 094 3. For sum feedback, we devise algorithms E-LinUCB and LinUCB-Tran for known and un-
 095 known transitions, respectively, which achieve near-optimal regrets in terms of H and m .
 096 We also establish lower bounds to validate the optimality, and show that optimal regrets do
 097 not depend on m . Our results reveal that surprisingly, under sum feedback, increasing the
 098 number of segments m does not help expedite learning much.
- 099 4. We develop novel techniques which can be of independent interest, including the KL
 100 divergence analysis to derive an exponential lower bound under binary feedback, and the
 101 use of E-optimal experimental design in algorithm E-LinUCB to refine the eigenvalue of the
 102 covariance matrix and reduce the regret.

105 2 RELATED WORK

106 In this section, we briefly review prior works related to ours.
 107

The algorithm design and theoretical analysis for classic RL were well studied in the literature (Sutton & Barto, 2018; Jaksch et al., 2010; Azar et al., 2017; Jin et al., 2018; Zanette & Brunskill, 2019). (Reviewer NUmY) Tang et al. (2024) proposed the RL with segment feedback problem (they called it RL with bagged reward) and designed a transformer-based algorithm. However, their work is mostly an empirical work and only considered the sum feedback type. They did not provide theoretical results for their algorithm and rigorously quantified the influence of segments on learning. There are two recent theoretical works (Efroni et al., 2021; Chatterji et al., 2021) studying RL with trajectory feedback, which are most related to our work. Efroni et al. (2021) investigated RL with sum trajectory feedback, and designed upper confidence bound (UCB)-type and Thompson sampling (TS)-type algorithms with regret guarantees. Chatterji et al. (2021) considered RL with binary trajectory feedback, and developed algorithms based on UCB value iteration. For binary trajectory feedback, Chatterji et al. (2021) provided regret upper bounds that exponentially depend on the scale of rewards.

Different from (Efroni et al., 2021; Chatterji et al., 2021), we study RL with segment feedback, which allows feedback from multiple segments within a trajectory, with per-state-action feedback and trajectory feedback as the two extremes. (Reviewer NUmY) Under sum feedback, we improve the regret bound in (Efroni et al., 2021) by \sqrt{H} using experimental design, and demonstrate the optimality of our result by establishing a lower bound. Under binary feedback, we propose TS-style algorithms which are computationally efficient, and also build a lower bound to validate the inevitability of the exponential dependence in the regret bound, which is new to the literature.

Our work is also related to linear bandits (Abbasi-Yadkori et al., 2011) and logistic bandits (Filippi et al., 2010; Faury et al., 2020; Russac et al., 2021), and uses analytical techniques from that literature.

3 FORMULATION

In this section, we present the formulation of RL with binary and sum segment feedback.

We consider an episodic MDP denoted by $\mathcal{M}(\mathcal{S}, \mathcal{A}, H, r, p, \rho)$. Here \mathcal{S} is the state space, and \mathcal{A} is the action space. H is the length of each episode. (Reviewer R7pS) $r : \mathcal{S} \times \mathcal{A} \rightarrow [-r_{\max}, r_{\max}]$ is an unknown reward function, where $r_{\max} > 0$ is a universal constant and used to prevent the input of binary feedback (the sigmoid function) from being too large. Define the reward parameter $\theta^* := [r(s, a)]_{(s, a) \in \mathcal{S} \times \mathcal{A}} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$. $p : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ is the transition distribution. For any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, $p(s'|s, a)$ is the probability of transitioning to s' if action a is taken in state s . $\rho \in \Delta_{\mathcal{S}}$ is the initial state distribution.

A policy $\pi : \mathcal{S} \times [H] \rightarrow \mathcal{A}$ is defined as a mapping from the state space and step indices to the action space, so that $\pi_h(s)$ specifies what action to take in state s at step h . For any policy π , $h \in [H]$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$, let $V_h^\pi(s)$ be the state value function, and $Q_h^\pi(s, a)$ be the state-action value function, which denote the cumulative expected reward obtained under policy π till the end of an episode, starting from s and (s, a) at step h , respectively. Formally,

$$V_h^\pi(s) := \mathbb{E} \left[\sum_{t=h}^H r(s_t, a_t) | s_h = s, \pi \right], \quad Q_h^\pi(s, a) := \mathbb{E} \left[\sum_{t=h}^H r(s_t, a_t) | s_h = s, a_h = a, \pi \right].$$

The optimal policy is defined as $\pi^* = \operatorname{argmax}_{\pi} V_h^\pi(s)$ for all $s \in \mathcal{S}$ and $h \in [H]$. For any $s \in \mathcal{S}$ and $h \in [H]$, denote $V_h^*(s) := V_h^{\pi^*}(s)$.

The process of RL with segment feedback is as follows. In each episode k , the agent chooses a policy π^k at the beginning of this episode, and starts from $s_1^k \sim \rho$. At each step $h \in [H]$, the agent first observes the current state s_h^k , and takes an action $a_h^k = \pi_h^k(s_h^k)$ according to her policy, and then transitions to a next state $s_{h+1}^k \sim p(\cdot | s_h^k, a_h^k)$.

Each episode is equally divided into m segments, and each segment is of length $\frac{H}{m}$. For convenience, assume that H is divisible by m . For any $k > 0$ and $i \in [m]$, let $\tau^k = (s_1^k, a_1^k, \dots, s_h^k, a_h^k)$ denote the trajectory in episode k , and $\tau_i^k = (s_{\frac{H}{m} \cdot (i-1) + 1}^k, a_{\frac{H}{m} \cdot (i-1) + 1}^k, \dots, s_{\frac{H}{m} \cdot i}^k, a_{\frac{H}{m} \cdot i}^k)$ denote the i -th segment of the trajectory in episode k .

For any trajectory or trajectory segment τ , $\phi^\tau \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ denotes the vector where each entry $\phi^\tau(s, a)$ is the number of times (s, a) is visited in τ . For any policy π , $\phi^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ denotes the vector where

Algorithm 1: SegBiTS

Input: $\delta, \delta' := \frac{\delta}{3}, \alpha := \exp(\frac{Hr_{\max}}{m}) + \exp(-\frac{Hr_{\max}}{m}) + 2, \lambda$.

1 **for** $k = 1, \dots, K$ **do**

2 $\hat{\theta}_{k-1} \leftarrow \operatorname{argmin}_{\theta} -(\sum_{k'=1}^{k-1} \sum_{i=1}^m (y_i^{k'} \cdot \log(\mu((\phi^{\tau_i^{k'}})^\top \theta)) + (1 - y_i^{k'}) \cdot \log(1 - \mu((\phi^{\tau_i^{k'}})^\top \theta))) - \frac{1}{2} \lambda \|\theta\|_2^2$;

3 $\Sigma_{k-1} \leftarrow \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \alpha \lambda I$;

4 **Sample** $\xi_k \sim \mathcal{N}(0, \alpha \cdot \nu(k-1)^2 \cdot \Sigma_{k-1}^{-1})$, where $\nu(k-1)$ is defined in Eq. (1);

5 $\hat{\theta}_k \leftarrow \hat{\theta}_{k-1} + \xi_k$;

6 $\pi^k \leftarrow \operatorname{argmax}_{\pi} (\phi^\pi)^\top \hat{\theta}_k$, where ϕ^π is defined in Eq. (3);

7 **Play episode** k with policy π^k . **Observe trajectory** τ^k and binary segment feedback $\{y_i^k\}_{i=1}^m$;

each entry $\phi^\pi(s, a)$ is the expected number of times (s, a) is visited in an episode under policy π , i.e., $\phi^\pi(s, a) := \mathbb{E}[\sum_{h=1}^H \mathbb{1}\{s_h = s, a_h = a\} | \pi]$.

In our model, the agent *observes reward feedback only at the end of each segment*, instead of each step as in classic RL. We consider two reward feedback settings as follows.

Binary Segment Feedback. Denote the sigmoid function by $\mu(x) := \frac{1}{1 + \exp(-x)}$ for any $x \in \mathbb{R}$. In the binary segment feedback setting, in each episode k , at the end of each segment $i \in [m]$, the agent observes a binary outcome

$$y_i^k = \begin{cases} 1, & \text{w.p. } \mu\left(-\sum_{t=\frac{H}{m} \cdot (i-1) + 1}^{\frac{H}{m} \cdot i} r(s_t^k, a_t^k)\right) = \mu((\phi^{\tau_i^k})^\top \theta^*), \\ 0, & \text{w.p. } 1 - \mu\left(-\sum_{t=\frac{H}{m} \cdot (i-1) + 1}^{\frac{H}{m} \cdot i} r(s_t^k, a_t^k)\right) = 1 - \mu((\phi^{\tau_i^k})^\top \theta^*). \end{cases}$$

Sum Segment Feedback. In the sum segment feedback setting, in each episode k , at each step h , the environment generates an underlying random reward $R_h^k = r(s_h^k, a_h^k) + \varepsilon_h^k$, where ε_h^k is a zero-mean and 1-Sub-Gaussian noise, and independent of the transition distribution. At the end of each segment $i \in [m]$, the agent observes the sum of random rewards

$$R_i^k = \sum_{t=\frac{H}{m} \cdot (i-1) + 1}^{\frac{H}{m} \cdot i} R(s_t^k, a_t^k) = (\phi^{\tau_i^k})^\top \theta^* + \sum_{t=\frac{H}{m} \cdot (i-1) + 1}^{\frac{H}{m} \cdot i} \varepsilon_t^k.$$

In the sum feedback setting, when $m = H$, our model degenerates to classic RL (Azar et al., 2017; Sutton & Barto, 2018). When $m = 1$, the above two settings reduce to the problems of RL with binary (Chatterji et al., 2021) and sum trajectory feedback (Efroni et al., 2021), respectively.

In our model, the agent needs to infer the reward function from sparse and implicit reward feedback. Let K denote the number of episodes played. The goal of the agent is to minimize the cumulative regret, which is defined as $\mathcal{R}(K) := \sum_{k=1}^K (V_1^*(s_1) - V_1^{\pi^k}(s_1))$.

4 REINFORCEMENT LEARNING WITH BINARY SEGMENT FEEDBACK

In this section, we investigate RL with binary segment feedback. To isolate the effect of segment feedback from transition model learning, we first design a computationally-efficient and sample-efficient algorithm SegBiTS for the known transition case with a regret guarantee, and establish a nearly matching regret lower bound. Then, we further develop an algorithm SegBiTS-Tran with carefully-designed transition bonuses for the unknown transition case, and provide a regret analysis.

4.1 ALGORITHM SegBiTS AND REGRET UPPER BOUND FOR KNOWN TRANSITION

Building upon the Thompson sampling algorithm (Thompson, 1933), SegBiTS adopts the maximum likelihood estimator (MLE) to learn rewards from binary feedback, and performs posterior sampling to compute the optimal policy. (Reviewer NUmY) Different from prior trajectory feedback

algorithms (Chatterji et al., 2021) which are either computationally inefficient or have an $O(K^{\frac{3}{5}})$ regret bound, SegBiTS is both computationally efficient and has an $O(\sqrt{K})$ regret bound.

Algorithm 1 presents the procedure of SegBiTS. Specifically, in each episode k , SegBiTS first employs MLE with past binary reward observations to obtain the estimated reward parameter $\hat{\theta}_{k-1}$ (Line 2). Then, SegBiTS calculates the feature covariance matrix of past segments Σ_{k-1} (Line 3). After that, SegBiTS samples a noise ξ_k from Gaussian distribution $\mathcal{N}(0, \alpha \cdot \nu(k-1)^2 \cdot \Sigma_{k-1}^{-1})$ (Line 4). Here α represents the universal upper bound of the inverse of the sigmoid function’s derivative. For any $k > 0$, we define

$$\nu(k) := \frac{m\sqrt{\lambda}}{H} \left(1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \sqrt{1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m}} \omega(k) + \frac{H^2}{m^2\lambda} \cdot \omega(k)^2 \right)^{\frac{3}{2}}, \quad (1)$$

and

$$\omega(k) := \sqrt{\lambda} \left(r_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|} + \frac{1}{2} \right) + \frac{|\mathcal{S}||\mathcal{A}|}{\sqrt{\lambda}} \log \left(\frac{4}{\delta'} \left(1 + \frac{H^2k}{4|\mathcal{S}||\mathcal{A}|\lambda m} \right) \right). \quad (2)$$

(Reviewer 1qnB) $\nu(k)$ stands for the confidence radius of the MLE estimate $\hat{\theta}_k$. With high probability, we have $|\phi^\top \theta^* - \phi^\top \hat{\theta}_k| \leq \sqrt{\alpha} \cdot \nu(k) \|\phi\|_{\Sigma_k^{-1}}$, where ϕ is the visitation indicator of any trajectory.

Adding noise ξ_k to $\hat{\theta}_{k-1}$, SegBiTS obtains a posterior reward estimate $\tilde{\theta}_k$ (Line 5). Then, it computes the optimal policy π^k under reward $\tilde{\theta}_k$ (Line 6). Note that the step in Line 6 is computationally efficient, which can be easily solved by any MDP planning algorithm, e.g., value iteration, with reward $\tilde{\theta}_k$. After obtaining π^k , SegBiTS plays episode k , and observes trajectory τ^k and binary feedback $\{y_i^k\}_{i=1}^m$ on each segment (Line 7).

Now we provide a regret upper bound for algorithm SegBiTS.

Theorem 1. (Reviewer R7pS) *With probability at least $1 - \delta$, for any $K > 0$, the regret of algorithm SegBiTS is bounded by*

$$\begin{aligned} \mathcal{R}(K) = & \tilde{O} \left(\exp \left(\frac{Hr_{\max}}{2m} \right) \cdot \left(\sqrt{Km|\mathcal{S}||\mathcal{A}|} \max \left\{ \frac{H^2}{m\alpha\lambda}, 1 \right\} + H\sqrt{\frac{K}{\alpha\lambda}} \right) \cdot \frac{m\sqrt{\lambda}|\mathcal{S}||\mathcal{A}|}{H} \right. \\ & \left(1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \sqrt{1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m}} \left(\sqrt{\lambda} \left(r_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|} + \frac{1}{2} \right) + \frac{|\mathcal{S}||\mathcal{A}|}{\sqrt{\lambda}} \right) \right. \\ & \left. \left. + \frac{H^2}{m^2\lambda} \left(\sqrt{\lambda} \left(r_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|} + \frac{1}{2} \right) + \frac{|\mathcal{S}||\mathcal{A}|}{\sqrt{\lambda}} \right)^2 \right)^{\frac{3}{2}} \right). \end{aligned}$$

(Reviewer 1qnB) Theorem 1 exhibits that the regret of algorithm SegBiTS depends on $\exp(\frac{Hr_{\max}}{2m})$, which is usually the dominating factor. This implies that as the number of segments m increases, the regret decays rapidly. Thus, under binary feedback, increasing the number of segments significantly helps accelerate learning. (Reviewer 1qnB) The intuition behind this is that when the reward scale $x = \frac{Hr_{\max}}{m}$ is large, the binary feedback is generated from the range where the sigmoid function $\mu(x) = \frac{1}{1+\exp(-x)}$ is flat, i.e., the derivative of the sigmoid function $\mu'(x)$ is small. Then, the generated binary feedback is likely always 0 or always 1, and it is hard to distinguish between a good action and a bad action, leading to a higher regret; On the contrary, when the reward scale $x = \frac{Hr_{\max}}{m}$ is small, the binary feedback is generated from the range where the sigmoid function $\mu(x)$ is steep, i.e., $\mu'(x)$ is large. Then, the generated binary feedback is more dispersed to be 0 or 1, and it is easier to distinguish between a good action and a bad action, leading to a lower regret. In other words, the regret bound depends on the inverse of the sigmoid function’s derivative $\mu'(x) = \frac{1}{\exp(x)+\exp(-x)+2}$.

(Reviewer R7pS) In Theorem 1, the dependence on $|\mathcal{S}|, |\mathcal{A}|, H$ are $|\mathcal{S}|^3, |\mathcal{A}|^3, \exp(\frac{Hr_{\max}}{2m})H^2$, respectively. Since the exponential dependence on $\exp(\frac{Hr_{\max}}{m})$ usually dominates the bound, here we mainly aim to reveal such exponential influence on the regret, and are not pursuing the absolute tightness of every polynomial factor. Below we establish a lower bound to demonstrate the inevitability of this exponential dependence.

4.2 REGRET LOWER BOUND FOR KNOWN TRANSITION

The lower bound for RL with binary segment feedback and known transition is as follows.

Theorem 2. *Consider the problem of RL with binary segment feedback and known transition. There exists a distribution of instances where for any $c_0 \in (0, \frac{1}{2})$, when $K \geq \exp(\frac{Hr_{\max}}{m}) \frac{4|\mathcal{S}||\mathcal{A}|m}{H^2 r_{\max}^2 c_0^2}$, the regret of any algorithm must be*

$$\Omega \left(\exp \left(\left(\frac{1}{2} - c_0 \right) \frac{Hr_{\max}}{m} \right) \sqrt{|\mathcal{S}||\mathcal{A}|mK} \right).$$

Theorem 2 shows that the exponential dependence on $\frac{Hr_{\max}}{m}$ in the regret is indispensable, and the $\exp(\frac{Hr_{\max}}{2m})$ factor in Theorem 1 nearly matches the exponential dependence in the lower bound up to an arbitrarily small factor c_0 in $\exp(\cdot)$. Together with Theorem 1, Theorem 2 reveals that when the number of segments m increases, the regret decreases at an exponential rate. (Reviewer NUmy) In addition, this regret lower bound also holds for the unknown transition case, by constructing the same problem instance as in this lower bound proof.

This lower bound and its analysis is novel to the literature. In analysis, we calculate the KL divergence of Bernoulli distributions with the sigmoid function as parameters. Then, we employ Pinsker’s inequality and the fact that $\mu'(x) = \mu(x)(1 - \mu(x))$ to build a connection between the calculated KL divergence and $\mu'(\frac{Hr_{\max}}{m})$. Since $\mu'(\frac{Hr_{\max}}{m})$ contains an exponential factor, we can finally derive an exponential dependence in the lower bound.

4.3 ALGORITHM SegBiTS-Tran AND REGRET UPPER BOUND FOR UNKNOWN TRANSITION

In the following, we extend our results to the unknown transition case. We develop an efficient algorithm SegBiTS-Tran for binary segment feedback and unknown transition. SegBiTS-Tran includes a transition bonus p_{k-1}^{pv} in posterior reward estimate $\tilde{\theta}_k^b$, and replaces visitation indicator ϕ^π by its estimate $\hat{\phi}_{k-1}^\pi$. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $\hat{\phi}_{k-1}^\pi(s, a)$ is the expected number of times (s, a) is visited in an episode under policy π on empirical MDP \hat{p}_{k-1} , where \hat{p}_{k-1} is the empirical estimate of transition distribution p . Then, SegBiTS-Tran finds the optimal policy via $\arg\max_\pi (\hat{\phi}_{k-1}^\pi)^\top \tilde{\theta}_k^b$, which can be efficiently solved by any MDP planning algorithm with transition \hat{p}_{k-1} and reward $\tilde{\theta}_k^b$. We defer the pseudo-code and details of SegBiTS-Tran to Appendix C.3.

The regret performance of algorithm SegBiTS-Tran is stated below.

Theorem 3. *With probability at least $1 - \delta$, for any $K > 0$, the regret of algorithm SegBiTS-Tran is bounded by*

$$\begin{aligned} & \tilde{O} \left(\exp \left(\frac{Hr_{\max}}{2m} \right) \nu(K) \sqrt{|\mathcal{S}||\mathcal{A}|} \left(\sqrt{Km|\mathcal{S}||\mathcal{A}| \max \left\{ \frac{H^2}{m\alpha\lambda}, 1 \right\}} + H \sqrt{\frac{K}{\alpha\lambda}} \right) \right. \\ & \left. + \left(\nu(K) \sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{\lambda}} + Hr_{\max} \right) |\mathcal{S}|^2 |\mathcal{A}|^{\frac{3}{2}} H^{\frac{3}{2}} \sqrt{K} \right). \end{aligned}$$

Similar to algorithm SegBiTS (Theorem 1), the regret of algorithm SegBiTS-Tran also has a dependence on $\exp(\frac{Hr_{\max}}{2m})$. When the number of segments m increases, the regret of SegBiTS-Tran significantly decreases. Compared to SegBiTS, the regret of SegBiTS-Tran has an additional term polynomial in $|\mathcal{S}|$, $|\mathcal{A}|$, H and \sqrt{K} , which is incurred by learning the unknown transition distribution.

5 REINFORCEMENT LEARNING WITH SUM SEGMENT FEEDBACK

In this section, we turn to RL with sum segment feedback. (Reviewer NUmy) Different from prior trajectory feedback algorithm (Efroni et al., 2021) which directly uses the least squares estimator and has a suboptimal regret bound, we develop an algorithm E-LinUCB for the known transition case, which employs experimental design to perform an initial exploration and achieves a near-optimal regret with respect to H and m . To validate the optimality, we further establish a regret lower bound.

Algorithm 2: E-LinUCB

Input: $\delta, \delta' := \frac{\delta}{3}, \lambda := \frac{H}{r_{\max}^2 m}$, rounding procedure ROUND, rounding approximation parameter

$$\gamma := \frac{1}{10}, \beta(k) := \sqrt{\frac{H|\mathcal{S}||\mathcal{A}|}{m} \log(1 + \frac{kH^2}{\lambda|\mathcal{S}||\mathcal{A}|m}) + 2 \log(\frac{1}{\delta'}) + r_{\max} \sqrt{\lambda|\mathcal{S}||\mathcal{A}|}, \forall k > 0.$$

1 Let $w^* \in \Delta_{\Pi}$ and z^* be the optimal solution and optimal value of the optimization:

$$\min_{w \in \Delta_{\Pi}} \left\| \left(\sum_{\pi \in \Pi} w(\pi) \left(\sum_{i=1}^m \mathbb{E}_{\tau_i \sim \pi} [\phi^{\tau_i} (\phi^{\tau_i})^{\top}] \right) \right)^{-1} \right\| \quad (3)$$

2 $K_0 \leftarrow \lceil \max\{26(1 + \gamma)^2 (z^*)^2 H^4 \log(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'}), \frac{|\mathcal{S}||\mathcal{A}|}{\gamma^2}\} \rceil$;

3 $(\pi^1, \dots, \pi^{K_0}) \leftarrow \text{ROUND}(\{\sum_{i=1}^m \mathbb{E}_{\tau_i \sim \pi} [\phi^{\tau_i} (\phi^{\tau_i})^{\top}]\}_{\pi \in \Pi}, w^*, \gamma, K_0)$;

4 Play K_0 episodes with policies π^1, \dots, π^{K_0} . Observe trajectories $\tau^1, \dots, \tau^{K_0}$ and rewards

$$\{R_i^1\}_{i=1}^m, \dots, \{R_i^{K_0}\}_{i=1}^m;$$

5 **for** $k = K_0 + 1, \dots, K$ **do**

$$6 \quad \hat{\theta}_{k-1} \leftarrow (\lambda I + \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^{\top})^{-1} \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau_i^{k'}} R_i^{k'};$$

$$7 \quad \Sigma_{k-1} \leftarrow \lambda I + \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^{\top};$$

$$8 \quad \pi^k \leftarrow \operatorname{argmax}_{\pi \in \Pi} ((\phi^{\pi})^{\top} \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\phi^{\pi}\|_{(\Sigma_{k-1})^{-1}});$$

9 Play episode k with policy π^k . Observe trajectory τ^k and sum segment feedback $\{R_i^k\}_{i=1}^m$;

Moreover, we design an algorithm LinUCB-Tran equipped with a variance-aware transition bonus to handle the unknown transition scenario.

5.1 ALGORITHM E-LinUCB AND REGRET UPPER BOUND FOR KNOWN TRANSITION

If we regard visitation indicators $\phi^{\tau_i^k}$ as feature vectors and θ^* as the reward parameter, the problem of RL with sum segment feedback and known transition is similar to linear bandits.

Building upon the classic linear bandit algorithm LinUCB (Abbasi-Yadkori et al., 2011), algorithm E-LinUCB adopts the E-optimal design (Pukelsheim, 2006) to conduct an initial exploration to ensure sufficient coverage of the covariance matrix and further sharpen the norm under the inverse of the covariance matrix, resulting in a near-optimal regret.

As shown in Algorithm 2, the procedure of E-LinUCB is as follows. E-LinUCB first performs the E-optimal design to compute an optimal policy distribution w^* , which maximizes the minimum eigenvalue of the feature covariance matrix (Line 1). In Line 1, $\sum_{\pi \in \Pi} w(\pi) (\sum_{i=1}^m \mathbb{E}_{\tau_i \sim \pi} [\phi^{\tau_i} (\phi^{\tau_i})^{\top}])$ is the feature covariance matrix of segments under policy distribution w , and we assume that this matrix is invertible. Then, E-LinUCB calculates the number of samples K_0 for initial exploration according to the optimal value z^* of the E-optimal design (Line 2).

Then, in Line 3, E-LinUCB calls a rounding procedure ROUND (Allen-Zhu et al., 2021) to transform sampling distribution w^* into discrete sampling sequence $(\pi^1, \dots, \pi^{K_0})$ which satisfies (see Appendix B for more details of ROUND)

$$\left\| \left(\sum_{k=1}^{K_0} \left(\sum_{i=1}^m \mathbb{E}_{\tau_i \sim \pi^k} [\phi^{\tau_i} (\phi^{\tau_i})^{\top}] \right) \right)^{-1} \right\| \leq (1 + \gamma) \left\| \left(K_0 \sum_{\pi \in \Pi} w^*(\pi) \left(\sum_{i=1}^m \mathbb{E}_{\tau_i \sim \pi} [\phi^{\tau_i} (\phi^{\tau_i})^{\top}] \right) \right)^{-1} \right\|.$$

After that, E-LinUCB plays K_0 episodes with $(\pi^1, \dots, \pi^{K_0})$ to perform initial exploration (Line 4). Owing to the E-optimal design, the covariance matrix of initial exploration Σ_{K_0} has an optimized minimum eigenvalue, and then $\|\phi^{\pi}\|_{(\Sigma_{k-1})^{-1}}$ has a sharpened upper bound for any $k > K_0$. This is the key to the optimality of E-LinUCB with respect to H and m .

After initial exploration, in each episode $k > K_0$, E-LinUCB first calculates the least squares reward estimate $\hat{\theta}_{k-1}$ with past reward observations and covariance matrix Σ_{k-1} (Lines 6-7). Then, it computes the optimal policy with reward estimate $\hat{\theta}_{k-1}$ and bonus $\|\phi^{\pi}\|_{(\Sigma_{k-1})^{-1}}$ (Line 8). With the computed optimal policy π^k , E-LinUCB plays episode k , and collects trajectory τ^k and reward observation on each segment $\{R_i^k\}_{i=1}^m$ (Line 9).

Below we present a regret upper bound for algorithm E-LinUCB.

Theorem 4. *With probability at least $1 - \delta$, for any $K > 0$, the regret of algorithm E-LinUCB is bounded by*

$$O\left(|\mathcal{S}||\mathcal{A}|\sqrt{HK} \log\left(\left(1 + \frac{KHr_{\max}}{|\mathcal{S}||\mathcal{A}|m}\right)\frac{1}{\delta}\right) + (z^*)^2 H^5 \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{\delta}\right) + |\mathcal{S}||\mathcal{A}|H\right).$$

(Reviewer IqnB) Surprisingly, under sum feedback, when the number of segments m increases, the regret bound does not decrease significantly as people might expect, e.g., at a rate of $\frac{1}{\sqrt{m}}$ or $\frac{1}{m}$. While this looks surprising at the first glance, we discover an *intuition* through analysis: The objective of RL measures the expected reward sum of an episode, namely, we only need to accurately estimate the expected reward sum of an episode. When the number of segments m increases, although we obtain more observations, the segment features $\phi^{\tau_i^{k'}}$ contributed to covariance matrix Σ_k shrink, which makes the reward estimation uncertainty $\|\phi^\pi\|_{(\Sigma_k)^{-1}}$ inflate. When we focus on the estimation performance of the expected reward sum of an episode, these two effects cancel out with each other, and the regret result is not influenced by m distinctly.

(Reviewer NUmy) When $m = 1$, our problem reduces to RL with sum trajectory feedback (Efroni et al., 2021), and our result improves theirs by a factor of \sqrt{H} . This improvement comes from that we conduct the E-optimal design and perform an initial exploration to guarantee that $\|\phi^\pi\|_{(\Sigma_{k-1})^{-1}} \leq 1$, instead of $\|\phi^\pi\|_{(\Sigma_{k-1})^{-1}} \leq \frac{H}{\sqrt{\lambda}}$ as used in (Efroni et al., 2021).

Next, we investigate the lower bound to see if the number of segments m does not influence the regret bound distinctly in essence.

5.2 REGRET LOWER BOUND FOR KNOWN TRANSITION

We establish a lower bound for RL with sum segment feedback and known transition as follows.

Theorem 5. *Consider the problem of RL with sum segment feedback and known transition. There exists a distribution of instances where the regret of any algorithm must be*

$$\Omega\left(\sqrt{|\mathcal{S}||\mathcal{A}|HK}\right).$$

Theorem 5 demonstrates that our regret upper bound for algorithm E-LinUCB (Theorem 4) is optimal with respect to H and m when ignoring logarithmic factors. In addition, this lower bound corroborates that the number of segments m does not impact the regret result essentially.

5.3 ALGORITHM LinUCB-Tran AND REGRET UPPER BOUND FOR UNKNOWN TRANSITION

Now we study RL with sum segment feedback in the unknown transition scenario. For unknown transition, we design an algorithm LinUCB-Tran, which constructs a variance-aware uncertainty bound for the estimated visitation indicator $\hat{\phi}_k^\pi$, and takes into account this uncertainty bound in exploration bonuses. In analysis, we handle the estimation error of visitation indicators $\|\hat{\phi}_k^\pi - \phi^\pi\|_1$ by this variance-aware uncertainty bound, which enables us to achieve a near-optimal regret in terms of H . The pseudo-code and details of LinUCB-Tran are deferred to Appendix D.3.

In the following, we state the regret performance of algorithm LinUCB-Tran.

Theorem 6. *With probability at least $1 - \delta$, for any $K > 0$, the regret of algorithm LinUCB-Tran is bounded by*

$$\tilde{O}\left((1 + r_{\max})|\mathcal{S}|^{\frac{5}{2}}|\mathcal{A}|^2 H\sqrt{K}\right).$$

Theorem 6 shows that similar to algorithm E-LinUCB, the regret of LinUCB-Tran does not depend on the number of segments m when ignoring logarithmic factors. The heavier dependence on $|\mathcal{S}|$, $|\mathcal{A}|$ and H is due to the estimation of the unknown transition distribution. We also provide a lower bound for the unknown transition case, which demonstrates that the optimal regret indeed does not depend on m and our upper bound is near-optimal in terms of H (see Appendix D.5).

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

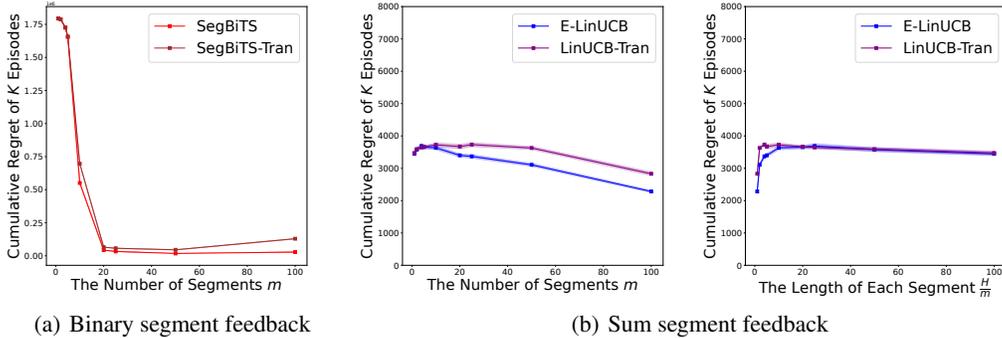


Figure 1: Experimental results for RL with binary or sum segment feedback.

6 EXPERIMENTS

Below we present experiments for RL with segment feedback to validate our theoretical results.

For the binary feedback setting, we evaluate our algorithms SegBiTS and SegBiTS-Tran in known and unknown transition environments, respectively, and we set $|\mathcal{S}| = 9$, $|\mathcal{A}| = 5$ and $K = 30000$. For the sum feedback setting, similarly, we run our algorithms E-LinUCB and LinUCB-Tran in known and unknown transition environments, respectively. Since E-LinUCB and LinUCB-Tran are computationally inefficient (mainly designed to exhibit the optimal dependence on m), we use a small MDP with $|\mathcal{S}| = 3$ and $|\mathcal{A}| = 5$, and set $K = 1000$. The details of the instances used in our experiments are described in Appendix A. In both settings, we set $r_{\max} = 0.5$, $\delta = 0.005$, $H = 100$ and $m \in \{1, 2, 4, 5, 10, 20, 25, 50, 100\}$. We perform 20 independent runs for each algorithm, and plot the average cumulative regret up to episode K across runs with a 90% confidence interval.

Figure 1(a) draws the regrets of algorithms SegBiTS and SegBiTS-Tran under binary feedback. One sees that as the number of segments m increases, the regret decreases rapidly. Specifically, when m decreases from 20 to 1, i.e., $\frac{H}{2m}$ increases from $\exp(2.5)$ to $\exp(50)$, the regret grows explosively. This matches our theoretical results (Theorems 1 and 3) which show a dependence on $\exp(\frac{Hr_{\max}}{2m})$.

Figure 1(b) plots the regrets of algorithms E-LinUCB and LinUCB-Tran under sum feedback. To see the impact of segments on regrets clearly, here we show the regrets with respect to the number of segments m and the length of each segment $\frac{H}{m}$ in the left and right subfigures, respectively. In the left subfigure, when m increases, the regrets almost keep the same for small m and slightly decrease for large m . To see the dependence on m more clearly, we turn to the right subfigure: When the length of each segment $\frac{H}{m}$ increases, the regrets slightly increase in a logarithmic trend. This also matches our theoretical bounds, which do not depend on m except for the $\log(\frac{H}{m})$ factor (Theorems 4 and 6).

7 CONCLUSION

In this work, we formulate a model named RL with segment feedback, which offers a general paradigm for feedback, bridging the gap between per-state-action feedback in classic RL and trajectory feedback. In the binary feedback setting, we design efficient algorithms SegBiTS and SegBiTS-Tran, and provide regret upper and lower bounds which show a dependence on $\exp(\frac{Hr_{\max}}{2m})$. These results reveal that under binary feedback, increasing the number of segment m greatly helps expedite learning. In the sum feedback setting, we develop near-optimal algorithms E-LinUCB and LinUCB-Tran in terms of H and m , where the regret results do not depend on m when ignoring logarithmic factors. These results exhibit that under sum feedback, increasing m does not help accelerate learning much.

There are several interesting directions worth further investigation. One direction is to consider segments of unequal lengths and study how to divide segments to optimize learning. Another direction is to generalize the results to the function approximation setting.

REFERENCES

- 486
487
488 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic
489 bandits. In *Advances in Neural Information Processing Systems*, volume 24, 2011.
- 490
491 Zeyuan Allen-Zhu, Yuanzhi Li, Aarti Singh, and Yining Wang. Near-optimal discrete optimization for
492 experimental design: A regret minimization approach. *Mathematical Programming*, 186:439–478,
493 2021.
- 494
495 Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed
496 bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- 497
498 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for rein-
499 forcement learning. In *International Conference on Machine Learning*, pp. 263–272. PMLR,
2017.
- 500
501 P Borjesson and C-E Sundberg. Simple approximations of the error function $q(x)$ for communications
502 applications. *IEEE Transactions on Communications*, 27(3):639–643, 1979.
- 503
504 Niladri Chatterji, Aldo Pacchiano, Peter Bartlett, and Michael Jordan. On the theory of reinforcement
505 learning with once-per-episode feedback. In *Advances in Neural Information Processing Systems*,
volume 34, pp. 3401–3412, 2021.
- 506
507 Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds
508 for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*,
509 volume 30, 2017.
- 510
511 Yonathan Efroni, Nadav Merlis, and Shie Mannor. Reinforcement learning with trajectory feedback.
512 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 7288–7295, 2021.
- 513
514 Louis Faury, Marc Abeille, Clément Calauzènes, and Olivier Fercoq. Improved optimistic algorithms
515 for logistic bandits. In *International Conference on Machine Learning*, pp. 3052–3060. PMLR,
2020.
- 516
517 Sarah Filippi, Olivier Cappé, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The
518 generalized linear case. In *Advances in Neural Information Processing Systems*, volume 23, 2010.
- 519
520 Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the
521 sample complexity of reinforcement learning with a generative model. *Machine Learning*, 91:
325–349, 2013.
- 522
523 Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement
524 learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- 525
526 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient?
527 In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- 528
529 Tor Lattimore and Marcus Hutter. Pac bounds for discounted mdps. In *International Conference on
530 Algorithmic Learning Theory*, pp. 320–334. Springer, 2012.
- 531
532 Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection.
533 *Annals of Statistics*, pp. 1302–1338, 2000.
- 534
535 Pierre Ménard, Omar Darwiche Domingues, Anders Jonsson, Emilie Kaufmann, Edouard Leurent,
and Michal Valko. Fast active learning for pure exploration in reinforcement learning. In *Internation-
536 al Conference on Machine Learning*, pp. 7599–7608. PMLR, 2021.
- 537
538 Rémi Munos and Andrew Moore. Influence and variance of a markov chain: Application to adaptive
539 discretization in optimal control. In *Proceedings of the IEEE Conference on Decision and Control*,
volume 2, pp. 1464–1469. IEEE, 1999.
- Friedrich Pukelsheim. *Optimal design of experiments*. SIAM, 2006.

540 Yoan Russac, Louis Faury, Olivier Cappé, and Aurélien Garivier. Self-concordant analysis of
541 generalized linear bandits with forgetting. In *International Conference on Artificial Intelligence
542 and Statistics*, pp. 658–666. PMLR, 2021.

543 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

544 Yuting Tang, Xin-Qiang Cai, Yao-Xiang Ding, Qiyu Wu, Guoqing Liu, and Masashi Sugiyama.
545 Reinforcement learning from bagged reward. In *ICML 2024 Workshop: Aligning Reinforcement
546 Learning Experimentalists and Theorists*, 2024.

547 William R Thompson. On the likelihood that one unknown probability exceeds another in view of
548 the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.

549 Joel A Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in
550 Machine Learning*, 8(1-2):1–230, 2015.

551 Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement
552 learning without domain knowledge using value function bounds. In *International Conference on
553 Machine Learning*, pp. 7304–7312. PMLR, 2019.

554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

APPENDIX

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

CONTENTS

1	Introduction	1
2	Related Work	2
3	Formulation	3
4	Reinforcement Learning with Binary Segment Feedback	4
4.1	Algorithm SegBiTS and Regret Upper Bound for Known Transition	4
4.2	Regret Lower Bound for Known Transition	6
4.3	Algorithm SegBiTS-Tran and Regret Upper Bound for Unknown Transition	6
5	Reinforcement Learning with Sum Segment Feedback	6
5.1	Algorithm E-LinUCB and Regret Upper Bound for Known Transition	7
5.2	Regret Lower Bound for Known Transition	8
5.3	Algorithm LinUCB-Tran and Regret Upper Bound for Unknown Transition	8
6	Experiments	9
7	Conclusion	9
A	Details of the Experimental Setup	13
B	Rounding Procedure ROUND	14
C	Proofs for RL with Binary Segment Feedback	14
C.1	Proof for the Regret Upper Bound with Known Transition	14
C.2	Proof for the Regret Lower Bound with Known Transition	24
C.3	Pseudo-code and Detailed Description of Algorithm SegBiTS-Tran	27
C.4	Proof for the Regret Upper Bound with Unknown Transition	28
D	Proofs for RL with Sum Segment Feedback	34
D.1	Proof for the Regret Upper Bound with Known Transition	34
D.2	Proof for the Regret Lower Bound with Known Transition	38
D.3	Pseudo-code and Detailed Description of Algorithm LinUCB-Tran	40
D.4	Proof for the Regret Upper Bound with Unknown Transition	41
D.5	A Lower Bound for Unknown Transition and its Proof	50
E	Technical Tools	52

A DETAILS OF THE EXPERIMENTAL SETUP

In this section, we detail the instances used in our experiments.

For the binary segment feedback setting, we consider an MDP as in Figure 1(a): There are 9 states and 5 actions. For any $a \in \mathcal{A}$, we have $r(s_0, a) = 0$, $r(s_i, a) = r_{\max}$ for any $i \in \{1, 3, 5, 7\}$ (called good states), and $r(s_i, a) = -r_{\max}$ for any $i \in \{2, 4, 6, 8\}$ (called bad states). There is an optimal action a^* and four suboptimal actions a^{sub} for all states. The agent starts from an initial state s_0 . For any $0 \leq i \leq 6$, in state s_i , under the optimal action a^* , the agent transitions to the good state and bad state at the next horizon with probabilities 0.9 and 0.1, respectively; Under the suboptimal action a^{sub} , the agent transitions to the good state and bad state at the next horizon with probabilities 0.1 and 0.9, respectively. In s_7 or s_8 , under the optimal action a^* , the agent transitions to s_1 and s_2 with probabilities 0.9 and 0.1, respectively; Under the suboptimal action a^{sub} , the agent transitions to s_1 and s_2 with probabilities 0.1 and 0.9, respectively.

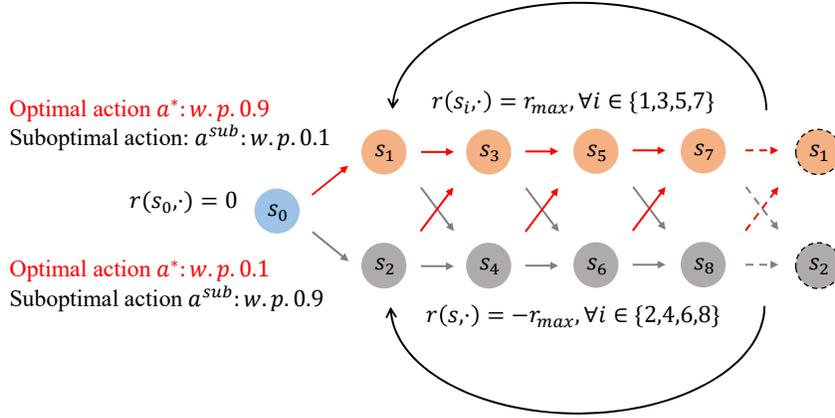


Figure 2: Instance used in the experiment for RL with binary segment feedback.

For the sum segment feedback setting, since algorithms E-LinUCB and LinUCB-Tran are computationally inefficient (which are mainly designed for revealing the optimal dependency on H and m), we consider a smaller MDP as in Figure 1(b): There are 3 states and 5 actions. For any $a \in \mathcal{A}$, we have $r(s_0, a) = 0$, $r(s_1, a) = r_{\max}$ (called a good state), and $r(s_2, a) = -r_{\max}$ (called a bad state). There is an optimal action a^* and four suboptimal actions a^{sub} for all states. The agent starts from an initial state s_0 . In any state $s \in \mathcal{S}$, under the optimal action a^* , the agent transitions to s_1 and s_2 with probabilities 0.9 and 0.1, respectively; Under the suboptimal action a^{sub} , the agent transitions to s_1 and s_2 with probabilities 0.1 and 0.9, respectively.

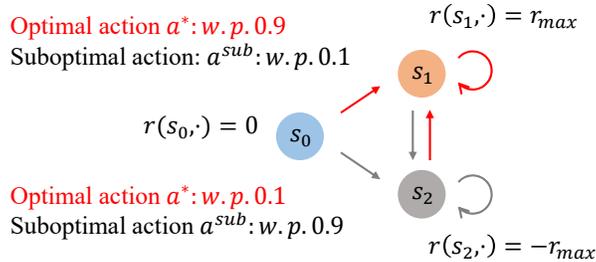


Figure 3: Instance used in the experiment for RL with sum segment feedback.

B ROUNDING PROCEDURE ROUND

Algorithm E-LinUCB calls a rounding procedure ROUND (Allen-Zhu et al., 2021) in the experimental design literature. Taking $X_1, \dots, X_n \in \mathbb{S}_+^d$, distribution $w \in \Delta_{\{X_1, \dots, X_n\}}$, rounding approximation error $\gamma > 0$ and the number of samples $T \geq \frac{d}{\gamma^2}$ as inputs, ROUND rounds sampling distribution w into a discrete sampling sequence $(Y_1, \dots, Y_T) \in \{X_1, \dots, X_n\}^T$ that satisfies

$$\left\| \left(\sum_{t=1}^T Y_t \right)^{-1} \right\| \leq (1 + \gamma) \left\| \left(T \sum_{i \in [n]} w(X_i) X_i \right)^{-1} \right\|.$$

In implementation, we can regard xx^\top in (Allen-Zhu et al., 2021) as $\sum_{i=1}^m \mathbb{E}_{\tau_i \sim \pi} [\phi^{\tau_i} (\phi^{\tau_i})^\top]$, and regard sampling weight on x as the sampling weight on π in our work.

C PROOFS FOR RL WITH BINARY SEGMENT FEEDBACK

In this section, we present the proofs for RL with binary segment feedback.

C.1 PROOF FOR THE REGRET UPPER BOUND WITH KNOWN TRANSITION

First, we prove the regret upper bound (Theorem 1) of algorithm SegBiTS for known transition.

For any $k > 0$ and $\theta \in \Theta$, define

$$Z_k := \sum_{k'=1}^k \sum_{i=1}^m \varepsilon_{k',i} \cdot \phi^{\tau_i^{k'}},$$

$$g_k(\theta) := \sum_{k'=1}^k \sum_{i=1}^m \mu((\phi^{\tau_i^{k'}})^\top \theta) \cdot \phi^{\tau_i^{k'}} + \lambda \theta, \quad (4)$$

$$\Lambda_k(\theta) := \sum_{k'=1}^k \sum_{i=1}^m \mu'((\phi^{\tau_i^{k'}})^\top \theta) \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I. \quad (5)$$

Lemma 1. For any $k > 0$ and $\theta \in \Theta$, we have

$$\det(\Lambda_k(\theta)) \leq \left(\frac{H^2 \mu'_{\max} k}{|\mathcal{S}||\mathcal{A}|m} + \lambda \right)^{|\mathcal{S}||\mathcal{A}|}.$$

Proof. For any $k > 0$, we have

$$\begin{aligned} \det(\Lambda_k(\theta)) &\leq \left(\frac{\text{tr}(\Lambda_k(\theta))}{|\mathcal{S}||\mathcal{A}|} \right)^{|\mathcal{S}||\mathcal{A}|} \\ &\leq \left(\frac{1}{|\mathcal{S}||\mathcal{A}|} \cdot \left(km \cdot \mu'_{\max} \cdot \left(\frac{H}{m} \right)^2 + \lambda |\mathcal{S}||\mathcal{A}| \right) \right)^{|\mathcal{S}||\mathcal{A}|} \\ &= \left(\frac{H^2 \mu'_{\max} k}{|\mathcal{S}||\mathcal{A}|m} + \lambda \right)^{|\mathcal{S}||\mathcal{A}|}. \end{aligned}$$

□

For any $k > 0$, let F_k denote the filtration that includes all events up to the end of episode k , and \tilde{F}_k denote the filtration that includes all events before playing π^k in episode k . Then, π^k is \tilde{F}_k -measurable.

For any $k > 0$ and $i \in [m]$, let $\varepsilon_{k,i} := y_i^k - (\phi^{\tau_i^k})^\top \theta^*$ denote the noise of binary feedback, and $v_{k,i}^2 := \mathbb{E}[\varepsilon_{k,i}^2 | \tilde{F}_k] = (\phi^{\tau_i^k})^\top \theta^* \cdot (1 - (\phi^{\tau_i^k})^\top \theta^*) = \mu'((\phi^{\tau_i^k})^\top \theta^*)$ denote the variance of $\varepsilon_{k,i}$ conditioning on \tilde{F}_k .

Then, we have

$$\begin{aligned}\Lambda_k(\theta^*) &:= \sum_{k'=1}^k \sum_{i=1}^m \mu'((\phi^{\tau_i^{k'}})^\top \theta^*) \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I \\ &= \sum_{k'=1}^k \sum_{i=1}^m v_{k',i}^2 \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I.\end{aligned}$$

Lemma 2 (Concentration of Noises under Binary Feedback). *With probability at least $1 - \delta'$, for any $k > 0$,*

$$\left\| \sum_{k'=1}^k \sum_{i=1}^m \varepsilon_{k',i} \cdot \phi^{\tau_i^{k'}} \right\|_{\Lambda_k^{-1}(\theta^*)} \leq \frac{\sqrt{\lambda}}{2} + \frac{|\mathcal{S}||\mathcal{A}|}{\sqrt{\lambda}} \log \left(\frac{4}{\delta'} \cdot \left(1 + \frac{H^2 \mu'_{\max} k}{|\mathcal{S}||\mathcal{A}| m \lambda} \right) \right).$$

Proof. According to Theorem 1 in (Faury et al., 2020), we have that with probability at least $1 - \delta'$, for any $k > 0$,

$$\begin{aligned}\left\| \sum_{k'=1}^k \sum_{i=1}^m \varepsilon_{k',i} \cdot \phi^{\tau_i^{k'}} \right\|_{\Lambda_k^{-1}(\theta^*)} &\leq \frac{\sqrt{\lambda}}{2} + \frac{2}{\sqrt{\lambda}} \log \left(\frac{\det(\Lambda_k(\theta^*))^{\frac{1}{2}} \cdot \lambda^{-\frac{|\mathcal{S}||\mathcal{A}|}{2}}}{\delta'} \right) + \frac{2}{\sqrt{\lambda}} |\mathcal{S}||\mathcal{A}| \log(2) \\ &\stackrel{(a)}{\leq} \frac{\sqrt{\lambda}}{2} + \frac{2}{\sqrt{\lambda}} \log \left(\frac{1}{\delta'} \left(1 + \frac{H^2 \mu'_{\max} k}{|\mathcal{S}||\mathcal{A}| m \lambda} \right)^{\frac{|\mathcal{S}||\mathcal{A}|}{2}} \right) + \frac{2}{\sqrt{\lambda}} |\mathcal{S}||\mathcal{A}| \log(2) \\ &\leq \frac{\sqrt{\lambda}}{2} + \frac{|\mathcal{S}||\mathcal{A}|}{\sqrt{\lambda}} \log \left(\frac{1}{\delta'} \left(1 + \frac{H^2 \mu'_{\max} k}{|\mathcal{S}||\mathcal{A}| m \lambda} \right) \right) + \frac{2}{\sqrt{\lambda}} |\mathcal{S}||\mathcal{A}| \log(2) \\ &\leq \frac{\sqrt{\lambda}}{2} + \frac{|\mathcal{S}||\mathcal{A}|}{\sqrt{\lambda}} \log \left(\frac{4}{\delta'} \left(1 + \frac{H^2 \mu'_{\max} k}{|\mathcal{S}||\mathcal{A}| m \lambda} \right) \right),\end{aligned}$$

where (a) uses Lemma 1. □

Define event

$$\mathcal{E} := \left\{ \left\| g_k(\hat{\theta}_k) - g_k(\theta^*) \right\|_{\Lambda_k^{-1}(\theta^*)} \leq \omega(k), \forall k > 0 \right\}.$$

Lemma 3. *It holds that*

$$\Pr[\mathcal{E}] \geq 1 - \delta'.$$

Proof. This proof is similar to that for Lemma 8 in (Faury et al., 2020).

Define

$$\mathcal{L}_k(\theta) := - \left(\sum_{k'=1}^k \sum_{i=1}^m \left(y_i^{k'} \cdot \log \left(\mu((\phi^{\tau_i^{k'}})^\top \theta) \right) + (1 - y_i^{k'}) \cdot \log \left(1 - \mu((\phi^{\tau_i^{k'}})^\top \theta) \right) \right) - \frac{1}{2} \lambda \|\theta\|_2^2 \right).$$

Recall that $\hat{\theta}_k = \operatorname{argmin}_\theta \mathcal{L}_k(\theta)$. Using the facts that $\nabla \mathcal{L}_k(\hat{\theta}_k) = 0$ and $\mu'(x) = \mu(x)(1 - \mu(x))$, we have

$$\underbrace{\sum_{k'=1}^k \sum_{i=1}^m \mu((\phi^{\tau_i^{k'}})^\top \hat{\theta}_k) \cdot \phi^{\tau_i^{k'}}}_{g_k(\hat{\theta}_k)} + \lambda \hat{\theta}_k = \sum_{k'=1}^k \sum_{i=1}^m y_i^{k'} \cdot \phi^{\tau_i^{k'}}.$$

Hence, we have

$$g_k(\hat{\theta}_k) - g_k(\theta^*) = \sum_{k'=1}^k \sum_{i=1}^m y_i^{k'} \cdot \phi^{\tau_i^{k'}} - \left(\sum_{k'=1}^k \sum_{i=1}^m \mu((\phi^{\tau_i^{k'}})^\top \theta^*) \cdot \phi^{\tau_i^{k'}} + \lambda \theta^* \right)$$

$$= \sum_{k'=1}^k \sum_{i=1}^m \varepsilon_{k',i} \cdot \phi^{\tau_i^{k'}} - \lambda \theta^*. \quad (6)$$

Then, using Lemma 2, we have that with probability at least $1 - \delta'$, for any $k > 0$,

$$\begin{aligned} \left\| g_k(\hat{\theta}_k) - g_k(\theta^*) \right\|_{\Lambda_k^{-1}(\theta^*)} &\leq \left\| \sum_{k'=1}^k \sum_{i=1}^m \varepsilon_{k',i} \cdot \phi^{\tau_i^{k'}} \right\|_{\Lambda_k^{-1}(\theta^*)} + r_{\max} \sqrt{\lambda |\mathcal{S}| |\mathcal{A}|} \\ &\leq \frac{\sqrt{\lambda}}{2} + \frac{|\mathcal{S}| |\mathcal{A}|}{\sqrt{\lambda}} \log \left(\frac{4}{\delta'} \cdot \left(1 + \frac{H^2 \mu'_{\max} k}{|\mathcal{S}| |\mathcal{A}| m \lambda} \right) \right) + r_{\max} \sqrt{\lambda |\mathcal{S}| |\mathcal{A}|} \\ &= \omega(k). \end{aligned}$$

□

For any $\phi \in \mathbb{R}^{|\mathcal{S}| |\mathcal{A}|}$ and $\theta_1, \theta_2 \in \Theta$, define

$$b(\phi, \theta_1, \theta_2) := \int_{z=0}^1 \mu'((1-z) \cdot \phi^\top \theta_1 + z \cdot \phi^\top \theta_2) dz.$$

For any $k > 0$ and $\theta_1, \theta_2 \in \Theta$, define

$$\Gamma_k(\theta_1, \theta_2) := \sum_{k'=1}^k \sum_{i=1}^m b(\phi, \theta_1, \theta_2) \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I.$$

In the definitions of $b(\phi, \theta_1, \theta_2)$ and $\Gamma_k(\theta_1, \theta_2)$, θ_1 and θ_2 have the same roles and can be interchanged.

Recall that

$$\alpha := \exp\left(\frac{H r_{\max}}{m}\right) + \exp\left(-\frac{H r_{\max}}{m}\right) + 2.$$

Then, we have

$$\sup_{\tau^{\text{seg}}, \theta} \frac{1}{\mu'((\phi^{\tau^{\text{seg}}})^\top \theta)} \leq \alpha,$$

where τ^{seg} denotes the visitation indicator of any possible trajectory segment.

Lemma 4. For any $k \geq 1$ and $\theta \in \Theta$, we have

$$\Sigma_k \preceq \alpha \Lambda_k(\theta).$$

Proof. We have

$$\frac{1}{\alpha} = \inf_{\tau^{\text{seg}}, \theta} \mu'((\phi^{\tau^{\text{seg}}})^\top \theta).$$

Then, it holds that

$$\begin{aligned} \Sigma_k &= \sum_{k'=1}^k \sum_{i=1}^m \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \alpha \lambda I \\ &= \alpha \left(\sum_{k'=1}^k \sum_{i=1}^m \frac{1}{\alpha} \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I \right) \\ &\leq \alpha \left(\sum_{k'=1}^k \sum_{i=1}^m \mu'((\phi^{\tau_i^{k'}})^\top \theta) \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I \right) \\ &= \alpha \Lambda_k(\theta). \end{aligned}$$

□

Lemma 5. For any $\phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\theta_1, \theta_2 \in \Theta$, we have

$$\mu(\phi^\top \theta_1) - \mu(\phi^\top \theta_2) = b(\phi, \theta_2, \theta_1) \cdot \phi^\top (\theta_1 - \theta_2).$$

In addition, for any $k > 0$ and $\theta_1, \theta_2 \in \Theta$, we have

$$\|\theta_1 - \theta_2\|_{\Gamma_k(\theta_2, \theta_1)} = \|g_k(\theta_1) - g_k(\theta_2)\|_{\Gamma_k^{-1}(\theta_2, \theta_1)}.$$

Proof. The first statement follows from the mean-value theorem.

Then, using the first statement, we have that for any $k > 0$,

$$\begin{aligned} g_k(\theta_1) - g_k(\theta_2) &= \sum_{k'=1}^k \sum_{i=1}^m \left(\mu((\phi^{\tau_{i}^{k'}})^\top \theta_1) - \mu((\phi^{\tau_{i}^{k'}})^\top \theta_2) \right) \cdot \phi^{\tau_{i}^{k'}} + \lambda(\theta_1 - \theta_2) \\ &= \sum_{k'=1}^k \sum_{i=1}^m b(\phi^{\tau_{i}^{k'}}, \theta_2, \theta_1) \cdot \phi^{\tau_{i}^{k'}} (\phi^{\tau_{i}^{k'}})^\top (\theta_1 - \theta_2) + \lambda(\theta_1 - \theta_2) \\ &= \Gamma_k(\theta_2, \theta_1) \cdot (\theta_1 - \theta_2), \end{aligned}$$

and thus

$$\begin{aligned} \|\theta_1 - \theta_2\|_{\Gamma_k(\theta_2, \theta_1)} &= \sqrt{(\theta_1 - \theta_2)^\top \cdot \Gamma_k(\theta_2, \theta_1) \cdot (\theta_1 - \theta_2)} \\ &= \sqrt{(\theta_1 - \theta_2)^\top \cdot \Gamma_k(\theta_2, \theta_1) \cdot \Gamma_k^{-1}(\theta_2, \theta_1) \cdot \Gamma_k(\theta_2, \theta_1) \cdot (\theta_1 - \theta_2)} \\ &= \|g_k(\theta_1) - g_k(\theta_2)\|_{\Gamma_k^{-1}(\theta_2, \theta_1)}, \end{aligned}$$

which gives the second statement. \square

Recall that for any $k > 0$, $Z_k := \sum_{k'=1}^k \sum_{i=1}^m \varepsilon_{k',i} \cdot \phi^{\tau_{i}^{k'}}$.

Lemma 6. For any $k > 0$, we have

$$\begin{aligned} \Gamma_k(\theta^*, \hat{\theta}_k) &\succeq \left(1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \right)^{-1} \Lambda_k(\theta^*), \\ \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} &\leq \sqrt{1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m}} \|Z_k\|_{\Lambda_k^{-1}(\theta^*)} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Lambda_k^{-1}(\theta^*)}^2. \end{aligned}$$

Furthermore, assuming that event \mathcal{E} holds, we have

$$\|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \leq \sqrt{1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m}} \cdot \omega(k) + \frac{H}{m\sqrt{\lambda}} \cdot \omega(k)^2.$$

Proof. This proof follows the analysis of Proposition 6 and Corollary 5 in (Russac et al., 2021).

From Eq. (6), we have that for any $k > 0$,

$$g_k(\hat{\theta}_k) - g_k(\theta^*) = Z_k - \lambda\theta^*.$$

Using Lemma 34, we have that for any $\phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ such that $\|\phi\|_2 \leq L_\phi$,

$$\begin{aligned} b(\phi, \theta^*, \hat{\theta}_k) &\geq \left(1 + \left| \phi^\top (\theta^* - \hat{\theta}_k) \right| \right)^{-1} \mu'(\phi^\top \theta^*) \\ &= \left(1 + \left| \phi^\top \Gamma_k^{-1}(\theta^*, \hat{\theta}_k) \cdot (g_k(\theta^*) - g_k(\hat{\theta}_k)) \right| \right)^{-1} \mu'(\phi^\top \theta^*) \\ &\geq \left(1 + \|\phi\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \|g_k(\theta^*) - g_k(\hat{\theta}_k)\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \right)^{-1} \mu'(\phi^\top \theta^*) \\ &\geq \left(1 + \frac{L_\phi}{\sqrt{\lambda}} \|g_k(\theta^*) - g_k(\hat{\theta}_k)\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \right)^{-1} \mu'(\phi^\top \theta^*) \end{aligned}$$

$$\begin{aligned}
&= \left(1 + \frac{L_\phi}{\sqrt{\lambda}} \|Z_k - \lambda\theta^*\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}\right)^{-1} \mu'(\phi^\top \theta^*) \\
&\geq \left(1 + L_\phi r_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|} + \frac{L_\phi}{\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}\right)^{-1} \mu'(\phi^\top \theta^*).
\end{aligned}$$

Using the above equation with $\phi = \phi^{\tau_i^{k'}}$ and $L_\phi = \frac{H}{m}$, we have

$$\begin{aligned}
\Gamma_k(\theta^*, \hat{\theta}_k) &:= \sum_{k'=1}^k \sum_{i=1}^m b(\phi^{\tau_i^{k'}}, \theta^*, \hat{\theta}_k) \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I \\
&\succeq \sum_{k'=1}^k \sum_{i=1}^m \left(1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}\right)^{-1} \mu'(\phi^\top \theta^*) \cdot \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \lambda I \\
&= \left(1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}\right)^{-1} \Lambda_k(\theta^*).
\end{aligned}$$

This implies

$$\|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}^2 \leq \left(1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}\right) \|Z_k\|_{\Lambda_k^{-1}(\theta^*)}^2,$$

which is equivalent to

$$\|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}^2 - \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Lambda_k^{-1}(\theta^*)}^2 \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} - \left(1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m}\right) \|Z_k\|_{\Lambda_k^{-1}(\theta^*)}^2 \leq 0.$$

By analysis of quadratic functions, we have

$$\begin{aligned}
\|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} &\leq \sqrt{1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m}} \|Z_k\|_{\Lambda_k^{-1}(\theta^*)} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Lambda_k^{-1}(\theta^*)}^2 \\
&\leq \sqrt{1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m}} \cdot \omega(k) + \frac{H}{m\sqrt{\lambda}} \cdot \omega(k)^2.
\end{aligned}$$

□

Lemma 7 (Concentration of $\phi^\top \hat{\theta}_k$ under Binary Feedback). *Assume that event \mathcal{E} holds. Then, for any $k > 0$ and $\phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$,*

$$|\phi^\top \theta^* - \phi^\top \hat{\theta}_k| \leq \sqrt{\alpha} \cdot \nu(k) \|\phi\|_{\Sigma_k^{-1}}.$$

Proof. We have

$$\begin{aligned}
&|\phi^\top \theta^* - \phi^\top \hat{\theta}_k| \\
&= \|\phi\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \|\theta^* - \hat{\theta}_k\|_{\Gamma_k(\theta^*, \hat{\theta}_k)} \\
&\stackrel{(a)}{\leq} \sqrt{1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}} \|\phi\|_{\Lambda_k^{-1}(\theta^*)} \|g_k(\theta^*) - g_k(\hat{\theta}_k)\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \\
&= \sqrt{1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}} \|\phi\|_{\Lambda_k^{-1}(\theta^*)} \|Z_k - \lambda\theta^*\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \\
&\leq \sqrt{1 + \frac{Hr_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)}} \|\phi\|_{\Lambda_k^{-1}(\theta^*)} \left(\|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} + r_{\max} \sqrt{\lambda|\mathcal{S}||\mathcal{A}|}\right)
\end{aligned}$$

$$\begin{aligned}
&= \frac{m\sqrt{\lambda}}{H} \sqrt{1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \|\phi\|_{\Lambda_k^{-1}(\theta^*)}} \\
&\quad \left(\frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m} \right) \\
&\leq \frac{m\sqrt{\lambda}}{H} \left(1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \|Z_k\|_{\Gamma_k^{-1}(\theta^*, \hat{\theta}_k)} \right)^{\frac{3}{2}} \|\phi\|_{\Lambda_k^{-1}(\theta^*)} \\
&\stackrel{(b)}{\leq} \frac{m\sqrt{\alpha\lambda}}{H} \left(1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m} + \frac{H}{m\sqrt{\lambda}} \left(\sqrt{1 + \frac{Hr_{\max}\sqrt{|\mathcal{S}||\mathcal{A}|}}{m}} \omega(k) + \frac{H}{m\sqrt{\lambda}} \omega(k)^2 \right) \right)^{\frac{3}{2}} \|\phi\|_{\Sigma_k^{-1}} \\
&= \sqrt{\alpha} \cdot \nu(k) \|\phi\|_{\Sigma_k^{-1}}.
\end{aligned}$$

where inequality (a) is due to Lemmas 5 and 6, and inequality (b) follows from Lemmas 4 and 6. \square

Lemma 8 (Gaussian Anti-Concentration). *Assume that event \mathcal{E} holds. Then, for any $k > 0$ and F_{k-1} -measurable random variable $X \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, we have*

$$\Pr \left[X^\top \tilde{\theta}_k > X^\top \theta^* \mid F_{k-1} \right] \geq \frac{1}{2\sqrt{2\pi e}}.$$

Proof. This proof is originated from the analysis of Lemma 11 in (Efroni et al., 2021).

Using Lemma 7, we have that for any $k > 0$,

$$\|X^\top \theta^* - X^\top \hat{\theta}_{k-1}\| \leq \sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}.$$

It holds that

$$\begin{aligned}
&\Pr \left[X^\top \tilde{\theta}_k > X^\top \theta^* \mid F_{k-1} \right] \\
&= \Pr \left[\frac{X^\top \tilde{\theta}_k - X^\top \hat{\theta}_{k-1}}{\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}} > \frac{X^\top \theta^* - X^\top \hat{\theta}_{k-1}}{\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}} \mid F_{k-1} \right].
\end{aligned}$$

Here given F_{k-1} , $X^\top \tilde{\theta}_k - X^\top \hat{\theta}_{k-1} = X^\top \xi_k$ is a Gaussian random variable with mean 0 and standard deviation $\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}$.

Since when event \mathcal{E} holds,

$$\frac{X^\top \theta^* - X^\top \hat{\theta}_{k-1}}{\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}} \leq \frac{\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}}{\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}} = 1,$$

we have

$$\begin{aligned}
\Pr \left[X^\top \tilde{\theta}_k > X^\top \theta^* \mid F_{k-1} \right] &\geq \Pr \left[\frac{X^\top \tilde{\theta}_k - X^\top \hat{\theta}_{k-1}}{\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}} > 1 \mid F_{k-1} \right] \\
&= \Pr \left[\frac{X^\top \xi_k}{\sqrt{\alpha} \cdot \nu(k-1) \|X\|_{\Sigma_{k-1}^{-1}}} > 1 \mid F_{k-1} \right] \\
&\stackrel{(a)}{\geq} \frac{1}{2\sqrt{2\pi e}},
\end{aligned}$$

where inequality (a) comes from that if $Z \sim \mathcal{F}_{\text{UTran}}^{\text{B}}(0, 1)$, $\Pr[Z > z] \geq \frac{1}{\sqrt{2\pi}} \cdot \frac{z}{1+z^2} e^{-\frac{z^2}{2}}$ (Borjesson & Sundberg, 1979). \square

Lemma 9. *Let $\xi_k, \xi'_k \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ be i.i.d. random variables given F_{k-1} . Let \tilde{p} be a F_{k-1} -measurable transition model, and $x_{k-1} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ be a F_{k-1} -measurable random variable. For any policy π ,*

1026 denote the visitation indicator under policy π on MDP \tilde{p} by $\tilde{\phi}^\pi$. Let $\tilde{\pi}^k := \operatorname{argmax}_\pi (\tilde{\phi}^\pi)^\top (x_{k-1} +$
 1027 $\xi_k)$. Then, we have

$$1028 \mathbb{E} \left[\left((\tilde{\phi}^{\tilde{\pi}^k})^\top (x_{k-1} + \xi_k) - \mathbb{E} \left[(\tilde{\phi}^{\tilde{\pi}^k})^\top (x_{k-1} + \xi_k) \mid F_{k-1} \right] \right)^+ \mid F_{k-1} \right]$$

$$1029 \leq \mathbb{E} \left[|(\tilde{\phi}^{\tilde{\pi}^k})^\top \xi_k| + |(\tilde{\phi}^{\tilde{\pi}^k})^\top \xi'_k| \mid F_{k-1} \right].$$

1030
 1031
 1032
 1033 *Proof.* This proof is originated from Lemma 12 in (Efroni et al., 2021).

1034 First, using the definition of $\tilde{\pi}^k$ and the fact that ξ_k and ξ'_k follow the same distribution, we have

$$1035 \mathbb{E} \left[(\tilde{\phi}^{\tilde{\pi}^k})^\top (x_{k-1} + \xi_k) \mid F_{k-1} \right] = \mathbb{E} \left[\max_\pi (\tilde{\phi}^\pi)^\top (x_{k-1} + \xi'_k) \mid F_{k-1} \right]. \quad (7)$$

1036 Then, since given F_{k-1} , ξ_k and ξ'_k are independent, we have

$$1037 \mathbb{E} \left[\max_\pi (\tilde{\phi}^\pi)^\top (x_{k-1} + \xi'_k) \mid F_{k-1} \right] = \mathbb{E} \left[\max_\pi (\tilde{\phi}^\pi)^\top (x_{k-1} + \xi'_k) \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right]$$

$$1038 \geq \mathbb{E} \left[(\phi^{\tilde{\pi}_k})^\top (x_{k-1} + \xi'_k) \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right]. \quad (8)$$

1039 Hence, combining Eqs. (7) and (8), we have

$$1040 \mathbb{E} \left[\left((\phi^{\tilde{\pi}_k})^\top (x_{k-1} + \xi_k) - \mathbb{E} \left[(\phi^{\tilde{\pi}_k})^\top (x_{k-1} + \xi_k) \mid F_{k-1} \right] \right)^+ \mid F_{k-1} \right]$$

$$1041 \leq \mathbb{E} \left[\left((\phi^{\tilde{\pi}_k})^\top (x_{k-1} + \xi_k) - \mathbb{E} \left[(\phi^{\tilde{\pi}_k})^\top (x_{k-1} + \xi'_k) \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right] \right)^+ \mid F_{k-1} \right]$$

$$1042 = \mathbb{E} \left[\left(\mathbb{E} \left[(\phi^{\tilde{\pi}_k})^\top (x_{k-1} + \xi_k) \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right] - \mathbb{E} \left[(\phi^{\tilde{\pi}_k})^\top (x_{k-1} + \xi'_k) \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right] \right)^+ \mid F_{k-1} \right]$$

$$1043 = \mathbb{E} \left[\left(\mathbb{E} \left[(\phi^{\tilde{\pi}_k})^\top \xi_k - (\phi^{\tilde{\pi}_k})^\top \xi'_k \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right] \right)^+ \mid F_{k-1} \right]$$

$$1044 \leq \mathbb{E} \left[\left| \mathbb{E} \left[(\phi^{\tilde{\pi}_k})^\top \xi_k - (\phi^{\tilde{\pi}_k})^\top \xi'_k \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right] \right| \mid F_{k-1} \right]$$

$$1045 \leq \mathbb{E} \left[\mathbb{E} \left[|(\phi^{\tilde{\pi}_k})^\top \xi_k - (\phi^{\tilde{\pi}_k})^\top \xi'_k| \mid F_{k-1}, \xi_k, \tilde{\pi}_k \right] \mid F_{k-1} \right]$$

$$1046 = \mathbb{E} \left[|(\phi^{\tilde{\pi}_k})^\top \xi_k - (\phi^{\tilde{\pi}_k})^\top \xi'_k| \mid F_{k-1} \right]$$

$$1047 \leq \mathbb{E} \left[|(\phi^{\tilde{\pi}_k})^\top \xi_k| \mid F_{k-1} \right] + \mathbb{E} \left[|(\phi^{\tilde{\pi}_k})^\top \xi'_k| \mid F_{k-1} \right].$$

1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

For any $k > 0$ and $\delta_k \in (0, 1)$, define event

$$\mathcal{M}_k(\delta_k) := \left\{ \forall \phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} : |\phi^\top \xi_k| \leq \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 2\sqrt{\log \left(\frac{1}{\delta_k} \right)} \right) \|\phi\|_{\Sigma_{k-1}^{-1}} \right\}.$$

Lemma 10. For any $k > 0$ and $\delta_k \in (0, 1)$, we have

$$\Pr[\mathcal{M}_k(\delta_k) \mid F_{k-1}] \geq 1 - \delta_k.$$

In addition, for a random variable $X \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ such that $\|X\|_{\Sigma_{k-1}^{-1}} \leq L_X$, we have

$$\mathbb{E} \left[|X^\top \xi_k| \mid F_{k-1} \right] \leq \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 2\sqrt{\log \left(\frac{1}{\delta_k} \right)} \right) \mathbb{E} \left[\|X\|_{\Sigma_{k-1}^{-1}} \mid F_{k-1} \right]$$

$$+ \sqrt{\alpha} \cdot \nu(k-1) \cdot L_X \sqrt{|\mathcal{S}||\mathcal{A}| \delta_k}.$$

Proof. This proof is similar to the analysis of Lemma 13 in (Efroni et al., 2021).

First, we prove the first statement.

For any $\phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, we have

$$\begin{aligned}
|\phi^\top \xi_k| &= |\phi^\top \Sigma_{k-1}^{-\frac{1}{2}} \Sigma_{k-1}^{\frac{1}{2}} \xi_k| \\
&\leq \left\| \Sigma_{k-1}^{-\frac{1}{2}} \phi \right\|_2 \left\| \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2 \\
&= \sqrt{\alpha} \cdot \nu(k-1) \|\phi\|_{\Sigma_{k-1}^{-1}} \left\| \frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2. \tag{9}
\end{aligned}$$

Since given F_{k-1} , $\frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ is a vector with each entry being a standard Gaussian random variable, we have that $\left\| \frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2$ is chi-distributed with parameter $|\mathcal{S}||\mathcal{A}|$.

Then, using Lemma 1 in (Laurent & Massart, 2000), we have that with probability at least $1 - \delta_k$,

$$\begin{aligned}
\left\| \frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2 &\leq \sqrt{|\mathcal{S}||\mathcal{A}| + 2\sqrt{|\mathcal{S}||\mathcal{A}| \log\left(\frac{1}{\delta_k}\right)} + 2\log\left(\frac{1}{\delta_k}\right)} \\
&= \sqrt{\left(\sqrt{|\mathcal{S}||\mathcal{A}|} + \sqrt{\log\left(\frac{1}{\delta_k}\right)}\right)^2 + \log\left(\frac{1}{\delta_k}\right)} \\
&\leq \sqrt{|\mathcal{S}||\mathcal{A}|} + 2\sqrt{\log\left(\frac{1}{\delta_k}\right)}.
\end{aligned}$$

Next, we prove the second statement.

For a random variable $X \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, we have

$$\begin{aligned}
\mathbb{E}[|X^\top \xi_k| | F_{k-1}] &= \Pr[\mathcal{M}_k(\delta_k)] \cdot \mathbb{E}[|X^\top \xi_k| | F_{k-1}, \mathcal{M}_k(\delta_k)] \\
&\quad + \Pr[\bar{\mathcal{M}}_k(\delta_k)] \cdot \mathbb{E}[|X^\top \xi_k| | F_{k-1}, \bar{\mathcal{M}}_k(\delta_k)] \\
&\leq \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 2\sqrt{\log\left(\frac{1}{\delta_k}\right)} \right) \mathbb{E}[\|X\|_{\Sigma_{k-1}^{-1}} | F_{k-1}] \\
&\quad + \sqrt{\Pr[\bar{\mathcal{M}}_k(\delta_k)] \cdot \mathbb{E}[|X^\top \xi_k|^2 | F_{k-1}, \bar{\mathcal{M}}_k(\delta_k)]} \\
&\stackrel{(a)}{\leq} \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 2\sqrt{\log\left(\frac{1}{\delta_k}\right)} \right) \mathbb{E}[\|X\|_{\Sigma_{k-1}^{-1}} | F_{k-1}] \\
&\quad + \sqrt{\alpha} \cdot \nu(k-1) \sqrt{\delta_k \mathbb{E} \left[\|X\|_{\Sigma_{k-1}^{-1}}^2 \cdot \left\| \frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2^2 \mid F_{k-1} \right]} \\
&\leq \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 2\sqrt{\log\left(\frac{1}{\delta_k}\right)} \right) \mathbb{E}[\|X\|_{\Sigma_{k-1}^{-1}} | F_{k-1}] \\
&\quad + \sqrt{\alpha} \cdot \nu(k-1) \sqrt{\delta_k L_X^2 \mathbb{E} \left[\left\| \frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2^2 \mid F_{k-1} \right]} \\
&\stackrel{(b)}{\leq} \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 2\sqrt{\log\left(\frac{1}{\delta_k}\right)} \right) \mathbb{E}[\|X\|_{\Sigma_{k-1}^{-1}} | F_{k-1}] \\
&\quad + \sqrt{\alpha} \cdot \nu(k-1) \cdot L_X \sqrt{|\mathcal{S}||\mathcal{A}| \delta_k}.
\end{aligned}$$

Here inequality (a) follows from the Cauchy-Schwarz inequality. Inequality (b) is due to the fact that given F_{k-1} , $\left\| \frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2$ is chi-distributed with parameter $|\mathcal{S}||\mathcal{A}|$, and then $\mathbb{E}[\left\| \frac{1}{\sqrt{\alpha} \cdot \nu(k-1)} \Sigma_{k-1}^{\frac{1}{2}} \xi_k \right\|_2^2 | F_{k-1}] = |\mathcal{S}||\mathcal{A}|$. \square

1134 Define event

$$1135 \mathcal{F}_{\text{KTran}}^{\text{B}} := \left\{ \left| \sum_{k'=1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| \leq 4H \sqrt{\frac{k}{\alpha\lambda}} \log\left(\frac{4k}{\delta'}\right), \right. \\ 1137 \left. \left| \sum_{k'=1}^k \left(\mathbb{E} \left[(\phi^{\pi^{k'}})^\top \theta^* |F_{k'-1}] - (\phi^{\pi^{k'}})^\top \theta^* \right) \right| \leq 4Hr_{\max} \sqrt{k \log\left(\frac{4k}{\delta'}\right)}, \forall k > 0 \right\}. \quad (10)$$

1141 **Lemma 11.** *It holds that*

$$1142 \Pr[\mathcal{F}_{\text{KTran}}^{\text{B}}] \geq 1 - 2\delta'.$$

1145 *Proof.* We prove the first inequality as follows.

1147 For any $k' \geq 1$, we have that $\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \leq \frac{H}{\sqrt{\alpha\lambda}}$, and then $|\mathbb{E}_{\tau \sim \pi^{k'}}[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}}| \leq \frac{2H}{\sqrt{\alpha\lambda}}$.

1149 Using the Azuma-Hoeffding inequality, we have that for any fixed $k > 0$, with probability at least $1 - \frac{\delta'}{2k^2}$,

$$1152 \left| \sum_{k'=1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| \leq \sqrt{2 \cdot \frac{4H^2}{\alpha\lambda} \cdot k \log\left(\frac{4k^2}{\delta'}\right)}.$$

1156 Since $\sum_{k=1}^{\infty} \frac{\delta'}{2k^2} \leq \delta'$, by a union bound over k , we have that with probability at least δ' , for any $k \geq 1$,

$$1158 \left| \sum_{k'=1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| \leq \sqrt{2 \cdot \frac{4H^2}{\alpha\lambda} \cdot k \log\left(\frac{4k^2}{\delta'}\right)} \\ 1160 \leq 4H \sqrt{\frac{k}{\alpha\lambda}} \log\left(\frac{4k}{\delta'}\right).$$

1164 The second inequality can be obtained by a similar argument and the fact that $|(\phi^{\pi^k})^\top \theta^*| \leq Hr_{\max}$ for any $k > 0$. \square

1167 **Lemma 12.** *For any $K \geq 1$, we have*

$$1168 \sum_{k=1}^K \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}} \leq \sqrt{2Km|\mathcal{S}||\mathcal{A}| \cdot \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \cdot \log\left(1 + \frac{KH^2}{\alpha\lambda|\mathcal{S}||\mathcal{A}|m}\right)}.$$

1172 *Proof.* We have

$$1174 \sum_{k=1}^K \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}} \leq \sqrt{Km \sum_{k=1}^K \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2} \\ 1175 \stackrel{(a)}{\leq} \sqrt{2Km \cdot \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \cdot \sum_{k=1}^K \log\left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2\right)} \\ 1176 = \sqrt{2Km \cdot \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \cdot \log\left(\frac{\det(\Sigma_K)}{\det(\alpha\lambda I)}\right)} \\ 1177 \leq \sqrt{2Km|\mathcal{S}||\mathcal{A}| \cdot \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \cdot \log\left(1 + \frac{KH^2}{\alpha\lambda|\mathcal{S}||\mathcal{A}|m}\right)}, \quad (11)$$

1181 where inequality (a) is due to that for any $x \in [0, c]$ with constant $c \geq 0$, it holds that $x \leq 2 \max\{c, 1\} \cdot \log(1 + x)$. \square

1188 *Proof of Theorem 1.* Letting $\delta' = \frac{\delta}{3}$, we have $\Pr[\mathcal{E} \cap \mathcal{F}_{K\text{Tran}}^{\text{B}}] \leq 1 - \delta$. Then, to prove this theorem,
 1189 it suffices to prove the regret bound when event $\mathcal{E} \cap \mathcal{F}_{K\text{Tran}}^{\text{B}}$ holds.
 1190

1191 Assume that event $\mathcal{E} \cap \mathcal{F}_{K\text{Tran}}^{\text{B}}$ holds. Then, we have

$$\begin{aligned}
 1192 \mathcal{R}(K) &= \sum_{k=1}^K \left((\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* \right) \\
 1193 &= \sum_{k=1}^K \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] + \mathbb{E} \left[(\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] - (\phi^{\pi^k})^\top \theta^* \right) \\
 1194 &\leq \sum_{k=1}^K \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \right) + 4Hr_{\max} \sqrt{K \log \left(\frac{4K}{\delta'} \right)}. \quad (12)
 \end{aligned}$$

1201 For the first term, we have

$$\begin{aligned}
 1202 \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \\
 1203 = \sum_{k=1}^K \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] + \mathbb{E} \left[(\phi^{\pi^k})^\top \tilde{\theta}_k - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \right). \quad (13)
 \end{aligned}$$

1204 In the following, we prove

$$1205 \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \leq 2\sqrt{2\pi e} \cdot \mathbb{E} \left[\left((\phi^{\pi^k})^\top \tilde{\theta}_k - \mathbb{E} \left[(\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \right)^+ | F_{k-1} \right]. \quad (14)$$

1206 If $\mathbb{E}[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1}] < 0$, then Eq. (14) trivially holds.

1207 Otherwise, letting $z := \mathbb{E}[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1}]$, we have

$$\begin{aligned}
 1208 &\mathbb{E} \left[\left((\phi^{\pi^k})^\top \tilde{\theta}_k - \mathbb{E} \left[(\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \right)^+ | F_{k-1} \right] \\
 1209 &\geq z \Pr \left[(\phi^{\pi^k})^\top \tilde{\theta}_k - \mathbb{E} \left[(\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \geq z | F_{k-1} \right] \\
 1210 &\geq \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \right) \cdot \Pr \left[(\phi^{\pi^k})^\top \tilde{\theta}_k \geq (\phi^{\pi^*})^\top \theta^* | F_{k-1} \right] \\
 1211 &\stackrel{(a)}{\geq} \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \right) \cdot \Pr \left[(\phi^{\pi^*})^\top \tilde{\theta}_k \geq (\phi^{\pi^*})^\top \theta^* | F_{k-1} \right] \\
 1212 &\stackrel{(b)}{\geq} \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \right) \cdot \frac{1}{2\sqrt{2\pi e}},
 \end{aligned}$$

1213 where inequality (a) uses the definition of π^k , and inequality (b) follows from Lemma 8. Thus, we
 1214 complete the proof of Eq. (14).
 1215

1216 Let $\xi'_k \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ be a random variable that is i.i.d. with ξ given F_{k-1} . Then, using Lemma 9 with
 1217 $p' = p$, $x_{k-1} = \hat{\theta}_{k-1}$ and $\tilde{\pi}^k = \pi^k$, we have

$$\begin{aligned}
 1218 \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] &\leq 2\sqrt{2\pi e} \cdot \mathbb{E} \left[\left((\phi^{\pi^k})^\top \tilde{\theta}_k - \mathbb{E} \left[(\phi^{\pi^k})^\top \tilde{\theta}_k | F_{k-1} \right] \right)^+ | F_{k-1} \right] \\
 1219 &\leq 2\sqrt{2\pi e} \cdot \mathbb{E} \left[|\phi(\pi^k)^\top \xi_k| + |\phi(\pi^k)^\top \xi'_k| | F_{k-1} \right].
 \end{aligned}$$

1220 Plugging the above inequality into Eq. (13) and using Lemma 10 with $\delta_k = \frac{1}{k^4}$ and $L_X = \frac{H}{\sqrt{\alpha\lambda}}$, we
 1221 have

$$1222 \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right]$$

$$\begin{aligned}
&= \sum_{k=1}^K \left(2\sqrt{2\pi e} \mathbb{E} \left[|(\phi^{\pi^k})^\top \xi_k| + |(\phi^{\pi^k})^\top \xi'_k| \mid F_{k-1} \right] + \mathbb{E} \left[(\phi^{\pi^k})^\top (\hat{\theta}_{k-1} + \xi_k) - (\phi^{\pi^k})^\top \theta^* \mid F_{k-1} \right] \right) \\
&= \sum_{k=1}^K \left((2\sqrt{2\pi e} + 1) \cdot \mathbb{E} \left[|(\phi^{\pi^k})^\top \xi_k| \mid F_{k-1} \right] + 2\sqrt{2\pi e} \cdot \mathbb{E} \left[|(\phi^{\pi^k})^\top \xi'_k| \mid F_{k-1} \right] \right. \\
&\quad \left. + \mathbb{E} \left[(\phi^{\pi^k})^\top \hat{\theta}_{k-1} - (\phi^{\pi^k})^\top \theta^* \mid F_{k-1} \right] \right) \\
&\stackrel{(a)}{\leq} \sum_{k=1}^K \left((4\sqrt{2\pi e} + 2) \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 4\sqrt{\log(k)} \right) \cdot \mathbb{E} \left[\left\| \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} \mid F_{k-1} \right] \right. \\
&\quad \left. + (4\sqrt{2\pi e} + 1) \sqrt{\alpha} \cdot \nu(k-1) \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{k^2} \cdot \frac{H}{\sqrt{\alpha\lambda}} \right), \tag{15}
\end{aligned}$$

where inequality (a) uses Lemmas 7 and 10.

Here according to the definition of event $\mathcal{F}_{K\text{Tran}}^B$ and Lemma 12, we have

$$\begin{aligned}
\sum_{k=1}^K \mathbb{E} \left[\left\| \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} \mid F_{k-1} \right] &= \sum_{k=1}^K \left(\mathbb{E} \left[\left\| \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} \mid F_{k-1} \right] - \left\| \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} \right) + \sum_{k=1}^K \left\| \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} \\
&\leq 4H \sqrt{\frac{K}{\alpha\lambda} \log\left(\frac{4K}{\delta'}\right)} \\
&\quad + \sqrt{2Km|\mathcal{S}||\mathcal{A}| \cdot \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \cdot \log\left(1 + \frac{KH^2}{\alpha\lambda|\mathcal{S}||\mathcal{A}|m}\right)}. \tag{16}
\end{aligned}$$

Therefore, plugging the above two equations into Eq. (12), we have

$$\begin{aligned}
\mathcal{R}(K) &\leq (4\sqrt{2\pi e} + 2) \sqrt{\alpha} \cdot \nu(K) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 4\sqrt{\log(K)} \right) \cdot \\
&\quad \left(4H \sqrt{\frac{K}{\alpha\lambda} \log\left(\frac{4K}{\delta'}\right)} + \sqrt{2Km|\mathcal{S}||\mathcal{A}| \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \log\left(1 + \frac{KH^2}{\alpha\lambda|\mathcal{S}||\mathcal{A}|m}\right)} \right) \\
&\quad + 2(4\sqrt{2\pi e} + 1) H \cdot \nu(K) \sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{\lambda}} + 4Hr_{\max} \sqrt{K \log\left(\frac{4K}{\delta'}\right)} \\
&\stackrel{(a)}{=} \tilde{O} \left(\exp\left(\frac{Hr_{\max}}{2m}\right) \cdot \nu(K) \sqrt{|\mathcal{S}||\mathcal{A}|} \left(\sqrt{Km|\mathcal{S}||\mathcal{A}| \cdot \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\}} + H \sqrt{\frac{K}{\alpha\lambda}} \right) \right),
\end{aligned}$$

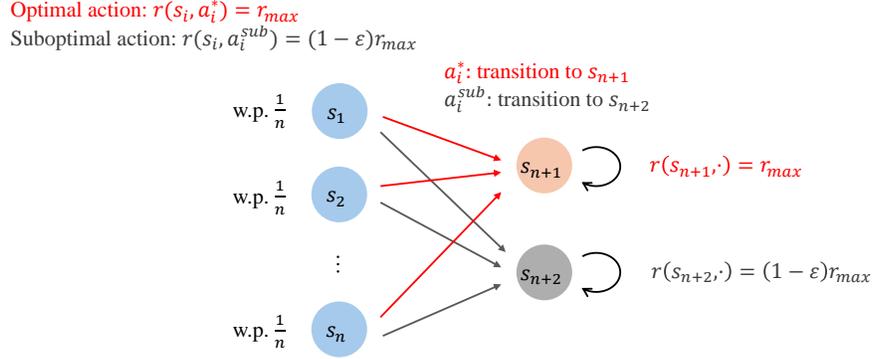
where in equality (a), the last two terms are absorbed into $\tilde{O}(\cdot)$. \square

C.2 PROOF FOR THE REGRET LOWER BOUND WITH KNOWN TRANSITION

In the following, we prove the regret lower bound (Theorem 2) for RL with binary segment feedback and known transition.

Proof of Theorem 2. We construct a random instance \mathcal{I} as follows. As shown in Figure 4, there are n bandit states s_1, \dots, s_n (i.e., there is an optimal action and multiple suboptimal actions), a good absorbing state s_{n+1} and a bad absorbing state s_{n+2} . The agent starts from s_1, \dots, s_n with equal probability $\frac{1}{n}$. For any $i \in [n]$, in state s_i , one action a_J is uniformly chosen from \mathcal{A} as the optimal action. In state s_i , under the optimal action a_J , the agent transitions to s_{n+1} deterministically, and $r(s_i, a_J) = r_{\max}$; Under any suboptimal action $a \in \mathcal{A} \setminus \{s_J\}$, the agent transitions to s_{n+2} deterministically, and $r(s_i, a) = (1 - \varepsilon)r_{\max}$, where $\varepsilon \in (0, \frac{1}{2})$ is a parameter specified later. For all actions $a \in \mathcal{A}$, $r(s_{n+1}, a) = r_{\max}$ and $r(s_{n+2}, a) = (1 - \varepsilon)r_{\max}$.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308



1309

Figure 4: Instance for the lower bound under binary segment feedback and known transition.

1311

In this proof, we will also use an alternative uniform instance $\mathcal{I}_{\text{unif}}$. The only difference between $\mathcal{I}_{\text{unif}}$ and \mathcal{I} is that for any $i \in [n]$, in state s_i , under all actions $a \in \mathcal{A}$, the agent transitions to s_{n+2} deterministically, and $r(s_i, a) = (1 - \varepsilon)r_{\max}$.

1312

1313

1314

1315

1316

1317

1318

Fix an algorithm \mathbb{A} . Let $\mathbb{E}_{\text{unif}}[\cdot]$ denote the expectation with respect to $\mathcal{I}_{\text{unif}}$. Let $\mathbb{E}_*[\cdot]$ denote the expectation with respect to \mathcal{I} . For any $i \in [n]$ and $j \in [|\mathcal{A}|]$, let $\mathbb{E}_{i,j}[\cdot]$ denote the expectation with respect to the case where a_j is the optimal action in state s_i , and $N_{i,j}$ denote the number of episodes where algorithm \mathbb{A} chooses a_j in state s_i , i.e., $N_{i,j} = \sum_{k=1}^K \mathbb{1}\{\pi_1^k(s_i) = a_j\}$.

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

The KL divergence of binary observations if taking a_j in s_i in each episode between $\mathcal{I}_{\text{unif}}$ and \mathcal{I} is

$$\begin{aligned} & \sum_{i=1}^m \text{KL} \left(\text{Ber} \left(\mu \left((1 - \varepsilon)r_{\max} \cdot \frac{H}{m} \right) \right) \parallel \text{Ber} \left(\mu \left(r_{\max} \cdot \frac{H}{m} \right) \right) \right) \\ & \stackrel{(a)}{\leq} m \cdot \frac{\left(\mu \left((1 - \varepsilon)r_{\max} \cdot \frac{H}{m} \right) - \mu \left(r_{\max} \cdot \frac{H}{m} \right) \right)^2}{\mu' \left(r_{\max} \cdot \frac{H}{m} \right)} \\ & \stackrel{(b)}{\leq} m \cdot \frac{\mu' \left((1 - \varepsilon) \frac{Hr_{\max}}{m} \right)^2 \left(\varepsilon \cdot \frac{Hr_{\max}}{m} \right)^2}{\mu' \left(\frac{Hr_{\max}}{m} \right)}, \end{aligned}$$

1329

1330

1331

1332

where inequality (a) uses the fact that $\text{KL}(\text{Ber}(p) \parallel \text{Ber}(q)) \leq \frac{(p-q)^2}{q(1-q)}$, and inequality (b) is due to that $\mu'(x)$ is monotonically decreasing when $x > 0$.

1333

1334

1335

1336

1337

1338

1339

1340

1341

In addition, the agent has probability only $\frac{1}{n}$ to arrive at (observe) state s_i .

Thus, using Lemma A.1 in (Auer et al., 2002), we have that for any $i \in [n]$, in state s_i ,

$$\begin{aligned} \mathbb{E}_{i,j}[N_{i,j}] & \leq \mathbb{E}_{\text{unif}}[N_{i,j}] + \frac{K}{2} \sqrt{\frac{1}{n} \cdot \mathbb{E}_{\text{unif}}[N_{i,j}] \cdot m \cdot \frac{\mu' \left((1 - \varepsilon) \frac{Hr_{\max}}{m} \right)^2 \left(\varepsilon \cdot \frac{Hr_{\max}}{m} \right)^2}{\mu' \left(\frac{Hr_{\max}}{m} \right)}} \\ & = \mathbb{E}_{\text{unif}}[N_{i,j}] + \frac{K}{2} \cdot \varepsilon \cdot \frac{Hr_{\max}}{m} \sqrt{\frac{m}{n} \cdot \mathbb{E}_{\text{unif}}[N_{i,j}] \cdot \frac{\mu' \left((1 - \varepsilon) \frac{Hr_{\max}}{m} \right)^2}{\mu' \left(\frac{Hr_{\max}}{m} \right)}}. \end{aligned}$$

1342

1343

1344

1345

1346

1347

1348

1349

Summing over $j \in [|\mathcal{A}|]$, using the Cauchy-Schwarz inequality and the fact that $\sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{\text{unif}}[N_{i,j}] = K$, we have

$$\begin{aligned} \sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{i,j}[N_{i,j}] & \leq K + \frac{KHr_{\max}\varepsilon}{2} \sqrt{\frac{|\mathcal{A}|K}{mn} \cdot \frac{\mu' \left((1 - \varepsilon) \frac{Hr_{\max}}{m} \right)^2}{\mu' \left(\frac{Hr_{\max}}{m} \right)}} \\ & \leq K + \frac{KHr_{\max}\varepsilon}{2} \sqrt{\frac{|\mathcal{A}|K}{mn} \cdot \frac{\mu' \left((1 - c_0) \frac{Hr_{\max}}{m} \right)^2}{\mu' \left(\frac{Hr_{\max}}{m} \right)}}, \end{aligned}$$

where $c_0 \in (0, \frac{1}{2})$ is a constant which satisfies $c_0 \geq \varepsilon$. We will specify how to make $c_0 \geq \varepsilon$ to satisfy this condition later.

Then, we have

$$\begin{aligned} \mathcal{R}(K) &= \sum_{k=1}^K \mathbb{E}_* [V^* - V^{\pi^k}] \\ &= r_{\max} H K - \frac{1}{n} \sum_{i=1}^n \left((1 - \varepsilon) r_{\max} H K + \varepsilon r_{\max} H \cdot \frac{1}{|\mathcal{A}|} \sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{i,j} [N_{i,j}] \right) \\ &\geq \varepsilon r_{\max} H \left(K - \frac{K}{|\mathcal{A}|} - \frac{K H r_{\max} \varepsilon}{2} \sqrt{\frac{K}{|\mathcal{A}| m n} \cdot \frac{\mu' \left((1 - c_0) \frac{H r_{\max}}{m} \right)^2}{\mu' \left(\frac{H r_{\max}}{m} \right)}} \right). \end{aligned}$$

Let

$$\varepsilon = \frac{1}{2 H r_{\max}} \sqrt{\frac{|\mathcal{A}| m n}{K} \cdot \frac{\mu' \left(\frac{H r_{\max}}{m} \right)}{\mu' \left((1 - c_0) \frac{H r_{\max}}{m} \right)^2}}.$$

Then, the constant c_0 should satisfy

$$\varepsilon = \frac{1}{2 H r_{\max}} \sqrt{\frac{|\mathcal{A}| m n}{K} \cdot \frac{\mu' \left(\frac{H r_{\max}}{m} \right)}{\mu' \left((1 - c_0) \frac{H r_{\max}}{m} \right)^2}} \leq c_0.$$

Since

$$\begin{aligned} \frac{\mu' \left(\frac{H r_{\max}}{m} \right)}{\mu' \left((1 - c_0) \frac{H r_{\max}}{m} \right)^2} &= \frac{(\exp \left((1 - c_0) \frac{H r_{\max}}{m} \right) + \exp \left(-(1 - c_0) \frac{H r_{\max}}{m} \right) + 2)^2}{\exp \left(\frac{H r_{\max}}{m} \right) + \exp \left(-\frac{H r_{\max}}{m} \right) + 2} \\ &\leq \frac{(4 \exp \left((1 - c_0) \frac{H r_{\max}}{m} \right))^2}{\exp \left(\frac{H r_{\max}}{m} \right)} \\ &= 16 \exp \left(\left(1 - 2c_0 \right) \frac{H r_{\max}}{m} \right), \end{aligned}$$

it suffices to let c_0 satisfy

$$\frac{1}{2 H r_{\max}} \sqrt{\frac{|\mathcal{A}| m n}{K} \cdot 16 \exp \left((1 - 2c_0) \frac{H r_{\max}}{m} \right)} \leq c_0,$$

which is equivalent to $K \geq \frac{4|\mathcal{A}| m n}{H^2 r_{\max}^2 c_0^2} \exp \left((1 - 2c_0) \frac{H r_{\max}}{m} \right)$.

It suffices to let

$$K \geq \frac{4|\mathcal{A}| m n}{H^2 r_{\max}^2 c_0^2} \exp \left(\frac{H r_{\max}}{m} \right),$$

and then c_0 can be any constant in $(0, \frac{1}{2})$.

Let $|\mathcal{S}| \geq 3$, $|\mathcal{A}| \geq 2$, $c_0 \in (0, \frac{1}{2})$ and $K \geq \frac{4|\mathcal{A}| m n}{H^2 r_{\max}^2 c_0^2} \exp \left(\frac{H r_{\max}}{m} \right)$. Since

$$\begin{aligned} \frac{\mu' \left(\frac{H r_{\max}}{m} \right)}{\mu' \left((1 - c_0) \frac{H r_{\max}}{m} \right)^2} &= \frac{(\exp \left((1 - c_0) \frac{H r_{\max}}{m} \right) + \exp \left(-(1 - c_0) \frac{H r_{\max}}{m} \right) + 2)^2}{\exp \left(\frac{H r_{\max}}{m} \right) + \exp \left(-\frac{H r_{\max}}{m} \right) + 2} \\ &\geq \frac{(\exp \left((1 - c_0) \frac{H r_{\max}}{m} \right))^2}{4 \exp \left(\frac{H r_{\max}}{m} \right)} \end{aligned}$$

Algorithm 3: SegBiTS-Tran**Input:** $\delta, \delta' := \frac{\delta}{8}, \lambda$.1 **for** $k = 1, \dots, K$ **do**2 $\hat{\theta}_{k-1} \leftarrow \operatorname{argmin}_{\theta} -(\sum_{k'=1}^{k-1} \sum_{i=1}^m (y_i^{k'} \cdot \log(\mu((\phi^{\tau_i^{k'}})^\top \theta)) + (1 - y_i^{k'}) \cdot \log(1 - \mu((\phi^{\tau_i^{k'}})^\top \theta))) - \frac{1}{2} \lambda \|\theta\|_2^2$;3 $\Sigma_{k-1} \leftarrow \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau_i^{k'}} (\phi^{\tau_i^{k'}})^\top + \alpha \lambda I$;4 Draw a noise $\xi_k \sim \mathcal{N}(0, \alpha \cdot \nu(k-1)^2 \cdot \Sigma_{k-1}^{-1})$, where $\nu(k-1)$ is defined in Eq. (1);5 $b_{k-1}^{pv}(s, a) \leftarrow \min\{2Hr_{\max} \sqrt{\frac{\log(\frac{KH|S||\mathcal{A}|}{\delta'})}{n_{k-1}(s,a)}}, Hr_{\max}\}$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$;6 $\tilde{\theta}_k^b \leftarrow \hat{\theta}_{k-1} + \xi_k + b_{k-1}^{pv}$;7 $\pi^k \leftarrow \operatorname{argmax}_{\pi} (\hat{\phi}_{k-1}^\pi)^\top \tilde{\theta}_k^b$, where $\hat{\phi}_{k-1}^\pi$ is defined in Eq. (17);8 Play episode k with policy π^k . Observe τ^k and binary segment feedback $\{y_i^k\}_{i=1}^m$;

$$= \frac{1}{4} \exp\left(\left(1 - 2c_0\right) \frac{Hr_{\max}}{m}\right),$$

we have

$$\begin{aligned} \mathcal{R}(K) &\geq \frac{1}{2Hr_{\max}} \sqrt{\frac{|\mathcal{A}|mn}{K} \cdot \frac{\mu' \left(\frac{Hr_{\max}}{m}\right)}{\mu' \left(\left(1 - c_0\right) \frac{Hr_{\max}}{m}\right)^2} \cdot r_{\max} H \left(K - \frac{K}{|\mathcal{A}|} - \frac{K}{4}\right)} \\ &= \Omega \left(\sqrt{\exp\left(\left(1 - 2c_0\right) \frac{Hr_{\max}}{m}\right)} |\mathcal{S}| |\mathcal{A}| m K \right) \\ &= \Omega \left(\exp\left(\left(\frac{1}{2} - c_0\right) \frac{Hr_{\max}}{m}\right) \sqrt{|\mathcal{S}| |\mathcal{A}| m K} \right). \end{aligned}$$

□

C.3 PSEUDO-CODE AND DETAILED DESCRIPTION OF ALGORITHM SegBiTS-Tran

Algorithm 3 illustrates the procedure of SegBiTS-Tran. In episode k , similar to SegBiTS, SegBiTS-Tran first uses MLE with past binary segment observations to obtain a reward estimate $\hat{\theta}_{k-1}$, and calculates the covariance matrix of past observations Σ_{k-1} (Lines 2-3). After that, SegBiTS-Tran samples a Gaussian noise ξ_k using Σ_{k-1} (Line 3).

For any $k > 0$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$, let $\hat{p}_k(\cdot|s, a)$ denote the empirical estimate of $p(\cdot|s, a)$, and $n_k(s, a)$ denote the number of times (s, a) was visited at the end of episode k . Then, SegBiTS-Tran constructs a transition bonus $b_{k-1}^{pv}(s, a)$, which represents the uncertainty on transition estimation. Incorporating the MLE estimate $\hat{\theta}_{k-1}$, noise ξ_k and transition bonus $b_{k-1}^{pv}(s, a)$, SegBiTS-Tran constitutes a posterior estimate of the reward parameter $\tilde{\theta}_k$ (Line 6).

For any policy π , $k > 0$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$, we define

$$\hat{\phi}_k^\pi(s, a) := \mathbb{E}_{\hat{p}_k} \left[\sum_{h=1}^H \mathbb{1}\{s_h = s, a_h = a\} | \pi \right], \quad (17)$$

which denotes the expected number of times (s, a) is visited in an episode under policy π on the empirical MDP \hat{p}_k . In addition, let $\hat{\phi}_k^\pi := [\hat{\phi}_k^\pi(s, a)]_{(s,a) \in \mathcal{S} \times \mathcal{A}} \in \mathbb{R}^{|\mathcal{S}| |\mathcal{A}|}$.

Then, SegBiTS-Tran finds the optimal policy via $\operatorname{argmax}_{\pi} (\hat{\phi}_{k-1}^\pi)^\top \tilde{\theta}_k^b$, which can be efficiently solved by any MDP planning algorithm with transition \hat{p}_{k-1} and reward $\tilde{\theta}_k^b$ (Line 7). With the computed optimal policy π^k , SegBiTS-Tran plays episode k , and observes a trajectory and binary feedback on each segment (Line 8).

1458 C.4 PROOF FOR THE REGRET UPPER BOUND WITH UNKNOWN TRANSITION
 1459

1460 In the following, we prove the regret upper bound (Theorem 3) of algorithm SegBiTS-Tran for
 1461 unknown transition.

1462 Define event

$$1463 \mathcal{G}_{\text{Hoeff}} := \left\{ \left| \hat{p}_{k-1}(\cdot|s, a)^\top V_{h+1}^* - p(\cdot|s, a)^\top V_{h+1}^* \right| \leq \left(2Hr_{\max} \sqrt{\frac{\log\left(\frac{KH|\mathcal{S}||\mathcal{A}|}{\delta'}\right)}{n_{k-1}(s, a)}} \wedge Hr_{\max} \right), \right. \\ 1464 \left. \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall k > 0 \right\}.$$

1470 **Lemma 13.** *It holds that*

$$1471 \Pr[\mathcal{G}_{\text{Hoeff}}] \geq 1 - 2\delta'.$$

1474 *Proof.* This lemma follows from the Hoeffding inequality and a union bound over $n_{k-1}(s, a) \in$
 1475 $[KH]$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$. \square

1477 **Lemma 14** (Optimism of Thompson Sampling with Unknown Transition). *Assume that event \mathcal{E} and*
 1478 $\mathcal{G}_{\text{Hoeff}}$ *holds. Then, for any $k > 0$, we have*

$$1479 \Pr\left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b > (\phi^{\pi^*})^\top \theta^* \mid F_{k-1}\right] \geq \frac{1}{2\sqrt{2\pi e}}.$$

1483 *Proof.* This proof follows the analysis of Lemma 17 in (Efroni et al., 2021).

1484 Using the value difference lemma (see Lemma 35), we have

$$1485 \hat{\phi}_{k-1}(\pi^*)^\top \tilde{\theta}_k^b - (\phi^{\pi^*})^\top \theta^* \\ 1486 = \mathbb{E}_{\hat{p}_{k-1}, \pi^*} \left[\sum_{h=1}^H \left(\tilde{\theta}_k^b(s_h, a_h) - \theta^*(s_h, a_h) + (\hat{p}_{k-1}(\cdot|s_h, a_h) - p(\cdot|s_h, a_h))^\top V_{h+1}^* \right) \right] \\ 1487 \\ 1488 = \mathbb{E}_{\hat{p}_{k-1}, \pi^*} \left[\sum_{h=1}^H \left(\tilde{\theta}_k(s_h, a_h) - \theta^*(s_h, a_h) + b_{k-1}^{pv}(s_h, a_h) + (\hat{p}_{k-1}(\cdot|s_h, a_h) - p(\cdot|s_h, a_h))^\top V_{h+1}^* \right) \right] \\ 1489 \\ 1490 \stackrel{(a)}{\geq} \mathbb{E}_{\hat{p}_{k-1}, \pi^*} \left[\sum_{h=1}^H \left(\tilde{\theta}_k(s_h, a_h) - \theta^*(s_h, a_h) + b_{k-1}^{pv}(s_h, a_h) - b_{k-1}^{pv}(s_h, a_h) \right) \right] \\ 1491 \\ 1492 = \mathbb{E}_{\hat{p}_{k-1}, \pi^*} \left[\sum_{h=1}^H \left(\tilde{\theta}_k(s_h, a_h) - \theta^*(s_h, a_h) \right) \right] \\ 1493 \\ 1494 = \hat{\phi}_{k-1}(\pi^*)^\top \tilde{\theta}_k - \hat{\phi}_{k-1}(\pi^*)^\top \theta^*, \\ 1495 \\ 1496 \\ 1497 \\ 1498 \\ 1499 \\ 1500$$

1501 where inequality (a) uses the definition of event $\mathcal{G}_{\text{Hoeff}}$.

1502 Thus, by the definition of π^k , we have

$$1503 \Pr\left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b > (\phi^{\pi^*})^\top \theta^* \mid F_{k-1}\right] \stackrel{(a)}{\geq} \Pr\left[\hat{\phi}_{k-1}(\pi^*)^\top \tilde{\theta}_k^b > (\phi^{\pi^*})^\top \theta^* \mid F_{k-1}\right] \\ 1504 = \Pr\left[\hat{\phi}_{k-1}(\pi^*)^\top \tilde{\theta}_k^b - (\phi^{\pi^*})^\top \theta^* > 0 \mid F_{k-1}\right] \\ 1505 \\ 1506 \geq \Pr\left[\hat{\phi}_{k-1}(\pi^*)^\top \tilde{\theta}_k - \hat{\phi}_{k-1}(\pi^*)^\top \theta^* > 0 \mid F_{k-1}\right] \\ 1507 \\ 1508 \stackrel{(b)}{\geq} \frac{1}{2\sqrt{2\pi e}}, \\ 1509 \\ 1510 \\ 1511$$

where inequality (a) is due to the definition of π^k , and inequality (b) follows from Lemma 8. \square

1512 Define event

$$1513 \mathcal{G}_{\text{KL}} := \left\{ \text{KL}(\hat{p}_{k-1}(\cdot|s, a), p(\cdot|s, a)) \leq \frac{L}{n_{k-1}(s, a)}, \forall k > 0, \forall (s, a) \in \mathcal{S} \times \mathcal{A} \right\}. \quad (18)$$

1516 **Lemma 15** (Concentration of Transition). *It holds that*

$$1517 \Pr[\mathcal{G}_{\text{KL}}] \geq 1 - \delta'.$$

1519 *Proof.* This lemma can be obtained by Theorem 3 and Lemma 3 in (Ménard et al., 2021). \square

1521 Recall that for any $k > 0$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$, $n_k(s, a)$ denotes the cumulative number of times
1522 that (s, a) is visited at the end of episode k . For any $k > 0$, $h \in [H]$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$, let
1523 $w_{k,h}(s, a)$ denote the probability that (s, a) is visited at step h in episode k , and let $w_k(s, a) :=$
1524 $\sum_{h=1}^H w_{k,h}(s, a)$.

1525 Define event

$$1527 \mathcal{H} := \left\{ n_k(s, a) \geq \frac{1}{2} \sum_{k'=1}^k w_{k'}(s, a) - H \log \left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'} \right), \forall k > 0, \forall (s, a) \in \mathcal{S} \times \mathcal{A} \right\}. \quad (19)$$

1530 **Lemma 16** (Concentration of the Number of Visitations). *It holds that*

$$1531 \Pr[\mathcal{H}] \geq 1 - \delta'.$$

1533 *Proof.* This lemma can be obtained from Lemma F.4 in (Dann et al., 2017) and summing over
1534 $h \in [H]$. \square

1536 Define event

$$1537 \mathcal{F}_{\text{UTran}}^{\text{B}} := \left\{ \left| \sum_{k'=1}^k \left(\mathbb{E} \left[(\phi^{\pi^{k'}})^{\top} b_{k'-1}^{pv} | F_{k'-1} \right] - (\phi^{\pi^{k'}})^{\top} b_{k'-1}^{pv} \right) \right| \leq 4H^2 r_{\max} \sqrt{k \log \left(\frac{4k}{\delta'} \right)}, \right. \\ 1541 \left. \left| \sum_{k'=1}^k \left(\mathbb{E} \left[\left\| \hat{\phi}_{k'-1}(\pi^{k'}) - \phi(\pi^{k'}) \right\|_1 | F_{k'-1} \right] - \left\| \hat{\phi}_{k'-1}(\pi^{k'}) - \phi(\pi^{k'}) \right\|_1 \right) \right| \right. \\ 1544 \left. \leq 8H \sqrt{k \log \left(\frac{4k}{\delta'} \right)}, \forall k > 0 \right\}.$$

1546 **Lemma 17.** *It holds that*

$$1547 \Pr[\mathcal{F}_{\text{UTran}}^{\text{B}}] \geq 1 - 2\delta'.$$

1549 *Proof.* This lemma can be obtained by a similar analysis as Lemma 11, and the facts that
1550 $|(\phi^{\pi^k})^{\top} b_{k-1}^{pv}| \leq H^2 r_{\max}$ and $\|\hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k}\|_1 \leq 2H$ for any $k \geq 1$. \square

1552 **Lemma 18.** *Assume that event $\mathcal{F}_{\text{UTran}}^{\text{B}} \cap \mathcal{G}_{\text{KL}} \cap \mathcal{H}$ holds. Then, we have*

$$1554 \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right] \leq 24e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H^{\frac{3}{2}} \sqrt{KL \log(2KH)} \\ 1556 + 192e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log \left(\frac{2KH|\mathcal{S}||\mathcal{A}|}{\delta'} \right).$$

1559 *Proof.* First, from Lemmas 29 and 30, we have

$$1561 \sum_{k=1}^K \left\| \hat{\phi}_{k-1}(\pi) - \phi(\pi) \right\|_1 \\ 1562 \leq e^{12} |\mathcal{S}||\mathcal{A}| \sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} w_h^{\pi^k}(s, a) \left(8H \sqrt{\frac{L}{n_{k-1}(s, a)}} + \frac{46H^2 L}{n_{k-1}(s, a)} \right)$$

$$\begin{aligned}
& + e^{12} |\mathcal{S}| |\mathcal{A}| H \sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \notin D_k} w_h^{\pi^k}(s, a) \\
& \leq 8e^{12} |\mathcal{S}| |\mathcal{A}| H \sqrt{L} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} w_h^{\pi^k}(s, a)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} \frac{w_h^{\pi^k}(s, a)}{n_{k-1}(s, a)}} \\
& \quad + 46e^{12} |\mathcal{S}| |\mathcal{A}| H^2 L \sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} \frac{w_h^{\pi^k}(s, a)}{n_{k-1}(s, a)} + 8e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 \log \left(\frac{|\mathcal{S}| |\mathcal{A}| H}{\delta'} \right) \\
& \leq 16e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H^{\frac{3}{2}} \sqrt{KL \log(2KH)} + 184e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log(2KH) \\
& \quad + 8e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 \log \left(\frac{|\mathcal{S}| |\mathcal{A}| H}{\delta'} \right) \\
& \leq 16e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H^{\frac{3}{2}} \sqrt{KL \log(2KH)} + 192e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log \left(\frac{2KH |\mathcal{S}| |\mathcal{A}|}{\delta'} \right).
\end{aligned}$$

Next, we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right] \\
& \leq \sum_{k=1}^K \left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 + \sum_{k=1}^K \left(\mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right] - \left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 \right) \\
& \leq 16e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H^{\frac{3}{2}} \sqrt{KL \log(2KH)} + 192e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log \left(\frac{2KH |\mathcal{S}| |\mathcal{A}|}{\delta'} \right) \\
& \quad + 8H \sqrt{K \log \left(\frac{4K}{\delta'} \right)} \\
& \leq 24e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H^{\frac{3}{2}} \sqrt{KL \log(2KH)} + 192e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log \left(\frac{2KH |\mathcal{S}| |\mathcal{A}|}{\delta'} \right).
\end{aligned}$$

□

Lemma 19. Assume that event $\mathcal{F}_{\text{UTran}}^{\text{B}}$ holds. Then, we have

$$\sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^k})^\top b_{k-1}^{pv} | F_{k-1} \right] \leq 20 |\mathcal{S}| |\mathcal{A}| H^2 r_{\max} \sqrt{K} \log \left(\frac{4KH |\mathcal{S}| |\mathcal{A}|}{\delta'} \right).$$

Proof. It holds that

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^k})^\top b_{k-1}^{pv} | F_{k-1} \right] \\
& = \sum_{k=1}^K (\phi^{\pi^k})^\top b_{k-1}^{pv} + \sum_{k=1}^K \left(\mathbb{E} \left[(\phi^{\pi^k})^\top b_{k-1}^{pv} | F_{k-1} \right] - (\phi^{\pi^k})^\top b_{k-1}^{pv} \right) \\
& \leq \sum_{k=1}^K \sum_{h=1}^H \sum_{s,a} w_h^{\pi^k}(s, a) \left(2Hr_{\max} \sqrt{\frac{\log \left(\frac{KH |\mathcal{S}| |\mathcal{A}|}{\delta'} \right)}{n_{k-1}(s, a)}} \wedge Hr_{\max} \right) + 4H^2 r_{\max} \sqrt{K \log \left(\frac{4K}{\delta'} \right)} \\
& \leq 2Hr_{\max} \sqrt{\log \left(\frac{KH |\mathcal{S}| |\mathcal{A}|}{\delta'} \right)} \sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} \frac{w_h^{\pi^k}(s, a)}{\sqrt{n_{k-1}(s, a)}}
\end{aligned}$$

$$\begin{aligned}
& + Hr_{\max} \sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \notin D_k} w_h^{\pi^k}(s,a) + 4H^2 r_{\max} \sqrt{K \log \left(\frac{4K}{\delta'} \right)} \\
& \leq 2Hr_{\max} \sqrt{\log \left(\frac{KH|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} \cdot \sqrt{KH} \cdot \sqrt{\sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} \frac{w_h^{\pi^k}(s,a)}{n_{k-1}(s,a)}} \\
& \quad + 8|\mathcal{S}||\mathcal{A}|H^2 r_{\max} \log \left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'} \right) + 4H^2 r_{\max} \sqrt{K \log \left(\frac{4K}{\delta'} \right)} \\
& \leq 2Hr_{\max} \sqrt{\log \left(\frac{KH|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} \cdot \sqrt{KH} \cdot \sqrt{4|\mathcal{S}||\mathcal{A}| \log(2KH)} \\
& \quad + 8|\mathcal{S}||\mathcal{A}|H^2 r_{\max} \log \left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'} \right) + 4H^2 r_{\max} \sqrt{K \log \left(\frac{4K}{\delta'} \right)} \\
& \leq 16|\mathcal{S}||\mathcal{A}|H^2 r_{\max} \sqrt{K} \log \left(\frac{4KH|\mathcal{S}||\mathcal{A}|}{\delta'} \right).
\end{aligned}$$

□

Proof of Theorem 3. Letting $\delta' = \frac{\delta}{8}$, we have $\Pr[\mathcal{E} \cap \mathcal{F}_{\text{KTran}}^{\text{B}} \cap \mathcal{G}_{\text{Hoeff}} \cap \mathcal{G}_{\text{KL}} \cap \mathcal{H} \cap \mathcal{F}_{\text{UTran}}^{\text{B}}] \leq 1 - \delta$. Then, to prove this theorem, it suffices to prove the regret bound when event $\mathcal{E} \cap \mathcal{F}_{\text{KTran}}^{\text{B}} \cap \mathcal{G}_{\text{Hoeff}} \cap \mathcal{G}_{\text{KL}} \cap \mathcal{H} \cap \mathcal{F}_{\text{UTran}}^{\text{B}}$ holds.

Assume that event $\mathcal{E} \cap \mathcal{F}_{\text{KTran}}^{\text{B}} \cap \mathcal{G}_{\text{Hoeff}} \cap \mathcal{G}_{\text{KL}} \cap \mathcal{H} \cap \mathcal{F}_{\text{UTran}}^{\text{B}}$ holds. Then, we have

$$\begin{aligned}
\mathcal{R}(K) &= \sum_{k=1}^K \left((\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* \right) \\
&= \sum_{k=1}^K \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] + \mathbb{E} \left[(\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] - (\phi^{\pi^k})^\top \theta^* \right) \\
&= \sum_{k=1}^K \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \right) + 4Hr_{\max} \sqrt{K \log \left(\frac{4K}{\delta'} \right)}. \tag{20}
\end{aligned}$$

For the first term, we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \\
&= \sum_{k=1}^K \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - \hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] + \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \right). \tag{21}
\end{aligned}$$

In the following, we prove

$$\begin{aligned}
& \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - \hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \\
& \leq 2\sqrt{2\pi e} \cdot \mathbb{E} \left[\left(\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b - \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \right)^+ | F_{k-1} \right]. \tag{22}
\end{aligned}$$

If $\mathbb{E}[(\phi^{\pi^*})^\top \theta^* - \hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1}] < 0$, then Eq. (22) trivially holds.

Otherwise, letting $z := \mathbb{E}[(\phi^{\pi^*})^\top \theta^* - \hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1}]$, we have

$$\mathbb{E} \left[\left(\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b - \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \right)^+ | F_{k-1} \right]$$

$$\begin{aligned}
1674 & \geq z \Pr \left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b - \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \geq z | F_{k-1} \right] \\
1675 & \geq \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - \hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \right) \cdot \Pr \left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b \geq (\phi^{\pi^*})^\top \theta^* | F_{k-1} \right] \\
1676 & \stackrel{(a)}{\geq} \left(\mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - \hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \right) \cdot \frac{1}{2\sqrt{2\pi e}}, \\
1677 & \\
1678 & \\
1679 &
\end{aligned}$$

1680 where inequality (a) uses Lemma 14. Thus, we complete the proof of Eq. (22).

1681 Let $\xi'_k \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ be an i.i.d. random variable with ξ given F_{k-1} . Then, using Lemma 9 with
1682 $p' = \hat{p}_{k-1}$, $x_{k-1} = \hat{\theta}_{k-1} + b_{k-1}^{pv}$ and $\tilde{\pi}^k = \pi^k$, we have

$$\begin{aligned}
1683 & \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - \hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \\
1684 & \leq 2\sqrt{2\pi e} \cdot \mathbb{E} \left[\left(\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b - \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \tilde{\theta}_k^b | F_{k-1} \right] \right)^+ | F_{k-1} \right] \\
1685 & \leq 2\sqrt{2\pi e} \cdot \mathbb{E} \left[|\hat{\phi}_{k-1}(\pi^k)^\top \xi_k| + |\hat{\phi}_{k-1}(\pi^k)^\top \xi'_k| | F_{k-1} \right]. \\
1686 & \\
1687 & \\
1688 & \\
1689 & \\
1690 &
\end{aligned}$$

1691 Plugging the above inequality into Eq. (21) and using Lemma 10 with $\delta_k = \frac{1}{k^4}$ and $L_X = \frac{H}{\sqrt{\alpha\lambda}}$, we
1692 have

$$\begin{aligned}
1693 & \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \\
1694 & = \sum_{k=1}^K \left(2\sqrt{2\pi e} \cdot \mathbb{E} \left[|\hat{\phi}_{k-1}(\pi^k)^\top \xi_k| + |\hat{\phi}_{k-1}(\pi^k)^\top \xi'_k| | F_{k-1} \right] \right. \\
1695 & \quad \left. + \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \left(\hat{\theta}_{k-1} + b_{k-1}^{pv} + \xi_k \right) - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \right) \\
1696 & = \sum_{k=1}^K \left(\left(2\sqrt{2\pi e} + 1 \right) \cdot \mathbb{E} \left[|\hat{\phi}_{k-1}(\pi^k)^\top \xi_k| | F_{k-1} \right] + 2\sqrt{2\pi e} \cdot \mathbb{E} \left[|\hat{\phi}_{k-1}(\pi^k)^\top \xi'_k| | F_{k-1} \right] \right. \\
1697 & \quad \left. + \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \left(\hat{\theta}_{k-1} + b_{k-1}^{pv} \right) - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \right) \\
1698 & \leq \sum_{k=1}^K \left(\left(4\sqrt{2\pi e} + 1 \right) \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 4\sqrt{\log(k)} \right) \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) \right\|_{\Sigma_{k-1}^{-1}} | F_{k-1} \right] \right. \\
1699 & \quad \left. + \left(4\sqrt{2\pi e} + 1 \right) \sqrt{\alpha} \cdot \nu(k-1) \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{k^2} \cdot \frac{H}{\sqrt{\alpha\lambda}} \right. \\
1700 & \quad \left. + \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \left(\hat{\theta}_{k-1} + b_{k-1}^{pv} \right) - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \right). \tag{23} \\
1701 & \\
1702 & \\
1703 & \\
1704 & \\
1705 & \\
1706 & \\
1707 & \\
1708 & \\
1709 & \\
1710 & \\
1711 & \\
1712 & \\
1713 & \\
1714 & \\
1715 &
\end{aligned}$$

1716 We have

$$\begin{aligned}
1717 & \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \left(\hat{\theta}_{k-1} + b_{k-1}^{pv} \right) - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \\
1718 & = \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top \left(\hat{\theta}_{k-1} - \theta^* \right) | F_{k-1} \right] + \mathbb{E} \left[\left(\hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right)^\top \theta^* | F_{k-1} \right] \\
1719 & \quad + \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top b_{k-1}^{pv} | F_{k-1} \right] \\
1720 & \leq \sqrt{\alpha} \cdot \nu(k-1) \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) \right\|_{\Sigma_{k-1}^{-1}} | F_{k-1} \right] + r_{\max} \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right] \\
1721 & \quad + \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top b_{k-1}^{pv} | F_{k-1} \right]. \\
1722 & \\
1723 & \\
1724 & \\
1725 & \\
1726 & \\
1727 &
\end{aligned}$$

Hence, plugging the above inequality into Eq. (23), we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^*})^\top \theta^* - (\phi^{\pi^k})^\top \theta^* | F_{k-1} \right] \\
& \leq \sum_{k=1}^K \left(\left(4\sqrt{2\pi e} + 2 \right) \sqrt{\alpha} \cdot \nu(k-1) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 4\sqrt{\log(k)} \right) \cdot \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) \right\|_{\Sigma_{k-1}^{-1}} | F_{k-1} \right] \right. \\
& \quad + \left(4\sqrt{2\pi e} + 1 \right) \cdot \nu(k-1) \frac{H}{k^2} \sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{\lambda}} + r_{\max} \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right] \\
& \quad \left. + \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top b_{k-1}^{pv} | F_{k-1} \right] \right).
\end{aligned}$$

Here we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) \right\|_{\Sigma_{k-1}^{-1}} | F_{k-1} \right] \\
& \leq \sum_{k=1}^K \mathbb{E} \left[\left\| \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} | F_{k-1} \right] + \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} | F_{k-1} \right] \\
& \leq \sum_{k=1}^K \mathbb{E} \left[\left\| \phi^{\pi^k} \right\|_{\Sigma_{k-1}^{-1}} | F_{k-1} \right] + \frac{1}{\sqrt{\alpha\lambda}} \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right] \\
& \stackrel{(a)}{\leq} 4H \sqrt{\frac{K}{\alpha\lambda} \log\left(\frac{4K}{\delta'}\right)} + \sqrt{2Km|\mathcal{S}||\mathcal{A}| \cdot \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \cdot \log\left(1 + \frac{KH^2}{\alpha\lambda|\mathcal{S}||\mathcal{A}|m}\right)} \\
& \quad + \frac{1}{\sqrt{\alpha\lambda}} \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right],
\end{aligned}$$

where inequality (a) uses Eq. (16).

In addition, we have

$$\begin{aligned}
& \sum_{k=1}^K \mathbb{E} \left[\hat{\phi}_{k-1}(\pi^k)^\top b_{k-1}^{pv} | F_{k-1} \right] \\
& \leq \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^k})^\top b_{k-1}^{pv} | F_{k-1} \right] + \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 \left\| b_{k-1}^{pv} \right\|_\infty | F_{k-1} \right] \\
& \leq \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^k})^\top b_{k-1}^{pv} | F_{k-1} \right] + Hr_{\max} \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right].
\end{aligned}$$

Therefore, plugging the above three equations into Eq. (20), we have

$$\begin{aligned}
& \mathcal{R}(K) \\
& \leq \left(4\sqrt{2\pi e} + 2 \right) \sqrt{\alpha} \cdot \nu(K) \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 4\sqrt{\log(K)} \right) \cdot \\
& \quad \left(4H \sqrt{\frac{K}{\alpha\lambda} \log\left(\frac{4K}{\delta'}\right)} + \sqrt{2Km|\mathcal{S}||\mathcal{A}| \max\left\{\frac{H^2}{m\alpha\lambda}, 1\right\} \log\left(1 + \frac{KH^2}{\alpha\lambda|\mathcal{S}||\mathcal{A}|m}\right)} \right) \\
& \quad + \left(\left(4\sqrt{2\pi e} + 2 \right) \frac{\nu(K)}{\sqrt{\lambda}} \left(\sqrt{|\mathcal{S}||\mathcal{A}|} + 4\sqrt{\log(K)} \right) + 2Hr_{\max} \right) \sum_{k=1}^K \mathbb{E} \left[\left\| \hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k} \right\|_1 | F_{k-1} \right] \\
& \quad + \sum_{k=1}^K \mathbb{E} \left[(\phi^{\pi^k})^\top b_{k-1}^{pv} | F_{k-1} \right] + 2 \left(4\sqrt{2\pi e} + 1 \right) H \cdot \nu(K) \sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{\lambda}} + 4Hr_{\max} \sqrt{K \log\left(\frac{4K}{\delta'}\right)}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(a)}{=} \tilde{O} \left(\exp \left(\frac{Hr_{\max}}{m} \right) \nu(K) \sqrt{|\mathcal{S}||\mathcal{A}|} \left(\sqrt{Km|\mathcal{S}||\mathcal{A}| \max \left\{ \frac{H^2}{m\alpha\lambda}, 1 \right\}} + H\sqrt{\frac{K}{\alpha\lambda}} \right) \right. \\
&\quad \left. + \left(\nu(K) \sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{\lambda}} + Hr_{\max} \right) |\mathcal{S}|^2 |\mathcal{A}|^{\frac{3}{2}} H^{\frac{3}{2}} \sqrt{K} \right),
\end{aligned}$$

where in equality (a), we use Lemmas 18 and 19, and the last three terms are absorbed into $\tilde{O}(\cdot)$. \square

D PROOFS FOR RL WITH SUM SEGMENT FEEDBACK

In this section, we provide the proofs for RL with sum segment feedback.

D.1 PROOF FOR THE REGRET UPPER BOUND WITH KNOWN TRANSITION

We first prove the regret upper bound (Theorem 4) of algorithm E-LinUCB for known transition.

Define event

$$\begin{aligned}
\mathcal{J} &:= \left\{ \left\| \sum_{k=1}^{K_0} \left(\sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top - \mathbb{E}_{\tau_i \sim \pi^k} \left[\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top \right] \right) \right\| \right\| \\
&\leq \frac{4H^2}{m} \sqrt{K_0 \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} + \frac{4H^2}{m} \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right). \tag{24}
\end{aligned}$$

Lemma 20 (Concentration of Initial Sampling). *It holds that*

$$\Pr[\mathcal{J}] \geq 1 - \delta'.$$

Proof. Note that π^1, \dots, π^{K_0} and K_0 are fixed before sampling, $\mathbb{E}[\sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top] = \mathbb{E}_{\tau_i \sim \pi^k} [\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top]$, and $\|\sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top\| \leq \frac{H^2}{m}$. Then, using the matrix Bernstein inequality (Theorem 6.1.1 in (Tropp et al., 2015)), we can obtain this lemma. \square

Lemma 21 (E-optimal Design). *Assume that event \mathcal{J} holds. Then, we have*

$$\left\| \left(\sum_{k=1}^{K_0} \sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top \right)^{-1} \right\| \leq \frac{1}{H^2}.$$

Proof. Using the guarantee of the rounding procedure ROUND (Theorem 1.1 in (Allen-Zhu et al., 2021)) and the fact that $K_0 \geq \frac{|\mathcal{S}||\mathcal{A}|}{\gamma^2}$, we have

$$\begin{aligned}
&\left\| \left(\sum_{k=1}^{K_0} \mathbb{E}_{\tau_i \sim \pi^k} \left[\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top \right] \right)^{-1} \right\| \\
&\leq (1 + \gamma) \left\| \left(K_0 \sum_{\pi \in \Pi} w^*(\pi) \cdot \mathbb{E}_{\tau_i \sim \pi^k} \left[\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top \right] \right)^{-1} \right\| \\
&\leq \frac{(1 + \gamma)z^*}{K_0}.
\end{aligned}$$

Let $\sigma_{\min}(\cdot)$ denote the minimum eigenvalue. Then, we have

$$\begin{aligned}
&\sigma_{\min} \left(\sum_{k=1}^{K_0} \sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top \right) \\
&= \sigma_{\min} \left(\sum_{k=1}^{K_0} \mathbb{E}_{\tau_i \sim \pi^k} \left[\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top \right] + \sum_{k=1}^{K_0} \sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top - \sum_{k=1}^{K_0} \mathbb{E}_{\tau_i \sim \pi^k} \left[\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top \right] \right)
\end{aligned}$$

$$\begin{aligned}
&\geq \sigma_{\min} \left(\sum_{k=1}^{K_0} \mathbb{E}_{\tau_i \sim \pi^k} \left[\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top \right] \right) - \left\| \sum_{k=1}^{K_0} \sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top - \sum_{k=1}^{K_0} \mathbb{E}_{\tau_i \sim \pi^k} \left[\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top \right] \right\| \\
&\geq \frac{K_0}{(1+\gamma)z^*} - \frac{4H^2}{m} \sqrt{\log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} \cdot \sqrt{K_0} - \frac{4H^2}{m} \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right). \tag{25}
\end{aligned}$$

Let $x = \sqrt{K_0}$ and

$$f(x) = \frac{1}{(1+\gamma)z^*} \cdot x^2 - \frac{4H^2}{m} \sqrt{\log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} \cdot x - \frac{4H^2}{m} \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right) - H^2.$$

According to the property of quadratic functions, when

$$x \geq \frac{\frac{4H^2}{m} \sqrt{\log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} + \sqrt{\left(\frac{4H^2}{m} \sqrt{\log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} \right)^2 + 4 \cdot \frac{1}{(1+\gamma)z^*} \left(\frac{4H^2}{m} \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right) + H^2 \right)}}{2 \cdot \frac{1}{(1+\gamma)z^*}}, \tag{26}$$

we have $f(x) \geq 0$.

To make Eq. (26) hold, it suffices to set

$$\begin{aligned}
K_0 &\geq \frac{(1+\gamma)^2(z^*)^2}{4} \cdot \left(2 \cdot \left(\frac{4H^2}{m} \sqrt{\log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} \right)^2 + 2 \cdot \left(\frac{4H^2}{m} \sqrt{\log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right)} \right) \right. \\
&\quad \left. + \frac{8}{(1+\gamma)z^*} \cdot 5H^2 \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right) \right) \\
&= \left(\frac{16H^4(1+\gamma)^2(z^*)^2}{m^2} + 10H^2(1+\gamma)z^* \right) \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right).
\end{aligned}$$

Furthermore, since $\| \sum_{\pi \in \Pi} w^*(\pi) \mathbb{E}_{\tau_i \sim \pi^k} [\sum_{i=1}^m \phi(\tau_i) \phi(\tau_i)^\top] \| \leq H^2$ and then $z^* \geq \frac{1}{H^2}$, to make the right-hand-side in Eq. (25) no smaller than H^2 , it suffices to set

$$K_0 \geq 26H^4(1+\gamma)^2(z^*)^2 \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right).$$

Therefore, combining the definition of K_0 and Eq. (25), we have

$$\sigma_{\min} \left(\sum_{k=1}^{K_0} \phi(\tau^k) \phi(\tau^k)^\top \right) \geq H^2,$$

which completes the proof. \square

Lemma 22. For any $k > 0$,

$$\sum_{k'=1}^k \log \left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^{k'}} \right\|_{(\Sigma_{k'-1})^{-1}}^2 \right) = \log \left(\frac{\det(\Sigma_k)}{\det(\lambda I)} \right) \leq |\mathcal{S}||\mathcal{A}| \log \left(1 + \frac{kH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right).$$

Proof. For any $k > 0$, it holds that

$$\begin{aligned}
\det(\Sigma_k) &= \det \left(\Sigma_{k-1} + \sum_{i=1}^m \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top \right) \\
&= \det(\Sigma_{k-1}) \det \left(I + \sum_{i=1}^m (\Sigma_{k-1})^{-\frac{1}{2}} \phi^{\tau_i^k} (\phi^{\tau_i^k})^\top (\Sigma_{k-1})^{-\frac{1}{2}} \right)
\end{aligned}$$

$$\begin{aligned}
&= \det(\Sigma_{k-1}) \left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2 \right) \\
&= \det(\lambda I) \prod_{k'=1}^k \left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^{k'}} \right\|_{(\Sigma_{k'-1})^{-1}}^2 \right).
\end{aligned}$$

Taking the logarithm on both sides, we have

$$\log \det(\Sigma_k) = \log \det(\lambda I) + \sum_{k'=1}^k \log \left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^{k'}} \right\|_{(\Sigma_{k'-1})^{-1}}^2 \right).$$

Then,

$$\begin{aligned}
\sum_{k'=1}^k \log \left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^{k'}} \right\|_{(\Sigma_{k'-1})^{-1}}^2 \right) &= \log \left(\frac{\det(\Sigma_k)}{\det(\lambda I)} \right) \\
&\stackrel{(a)}{\leq} \log \left(\frac{\left(\frac{\text{tr}(\Sigma_k)}{|\mathcal{S}||\mathcal{A}|} \right)^{|\mathcal{S}||\mathcal{A}|}}{\lambda^{|\mathcal{S}||\mathcal{A}|}} \right) \\
&= |\mathcal{S}||\mathcal{A}| \log \left(\frac{\text{tr}(\Sigma_k)}{\lambda |\mathcal{S}||\mathcal{A}|} \right) \\
&\leq |\mathcal{S}||\mathcal{A}| \log \left(\frac{\lambda |\mathcal{S}||\mathcal{A}| + km \cdot \frac{H^2}{m^2}}{\lambda |\mathcal{S}||\mathcal{A}|} \right) \\
&= |\mathcal{S}||\mathcal{A}| \log \left(1 + \frac{kH^2}{\lambda |\mathcal{S}||\mathcal{A}| m} \right),
\end{aligned}$$

where (a) uses the arithmetic mean-geometric mean inequality. \square

Lemma 23 (Elliptical Potential with Optimized Initialization). *Assume that event \mathcal{J} holds. Then, for any $k \geq K_0 + 1$,*

$$\sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2 \leq 1.$$

Furthermore, for any $K \geq K_0 + 1$,

$$\sum_{k=K_0+1}^K \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}} \leq \sqrt{2Km|\mathcal{S}||\mathcal{A}| \log \left(1 + \frac{KH^2}{\lambda |\mathcal{S}||\mathcal{A}| m} \right)}.$$

Proof. Using Lemma 21, for any $k \geq K_0 + 1$, we have

$$\begin{aligned}
\sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2 &= \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{\left(\lambda I + \sum_{k'=1}^{K_0} \phi^{\tau_{k'}} (\phi^{\tau_{k'}})^\top + \sum_{k'=K_0+1}^{k-1} \phi^{\tau_{k'}} (\phi^{\tau_{k'}})^\top \right)^{-1}}^2 \\
&\leq \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{\left(\sum_{k'=1}^{K_0} \phi^{\tau_{k'}} (\phi^{\tau_{k'}})^\top \right)^{-1}}^2 \\
&\leq m \cdot \frac{H^2}{m^2} \cdot \frac{1}{H^2} \\
&\leq 1.
\end{aligned}$$

Then, we have

$$\sum_{k=K_0+1}^K \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}} \leq \sqrt{Km \sum_{k=K_0+1}^K \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2}$$

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

$$\begin{aligned}
&\stackrel{(a)}{\leq} \sqrt{Km \cdot 2 \sum_{k=K_0+1}^K \log \left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2 \right)} \\
&\leq \sqrt{Km \cdot 2 \sum_{k=1}^K \log \left(1 + \sum_{i=1}^m \left\| \phi^{\tau_i^k} \right\|_{(\Sigma_{k-1})^{-1}}^2 \right)} \\
&\stackrel{(b)}{\leq} \sqrt{2Km|\mathcal{S}||\mathcal{A}| \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)},
\end{aligned}$$

where (a) uses the fact that $x \leq 2 \log(1+x)$ for any $x \in [0, 1]$, and (b) follows from Lemma 22. \square

Define event

$$\mathcal{K} := \left\{ \left\| \hat{\theta}_k - \theta^* \right\|_{\Sigma_k} \leq \sqrt{\frac{H|\mathcal{S}||\mathcal{A}|}{m} \log \left(1 + \frac{kH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)} + 2 \log \left(\frac{1}{\delta'} \right) + r_{\max} \sqrt{\lambda|\mathcal{S}||\mathcal{A}|}, \forall k > 0 \right\}. \quad (27)$$

Lemma 24 (Concentration of $\hat{\theta}_k$ under Sum Feedback). *It holds that*

$$\Pr[\mathcal{K}] \geq 1 - \delta'.$$

Proof. Since the sum feedback on each segment is $\frac{H}{m}$ -sub-Gaussian given the observation of transition and $\|\theta^*\| \leq r_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|}$, using Lemma 2 in (Abbasi-Yadkori et al., 2011), we can obtain this lemma. \square

Define event

$$\begin{aligned}
\mathcal{F}_{\text{opt}}^{\mathcal{S}} &:= \left\{ \left| \sum_{k'=K_0+1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}| \right] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| \right. \\
&\quad \left. \leq 4\sqrt{k \log \left(\frac{4k}{\delta'} \right)}, \forall k \geq K_0 + 1 \right\}. \quad (28)
\end{aligned}$$

Lemma 25 (Concentration of Visitation Indicators). *It holds that*

$$\Pr[\mathcal{F}_{\text{opt}}^{\mathcal{S}}] \geq 1 - \delta'.$$

Proof. According to Lemma 23, we have that for any $k' \geq K_0 + 1$, $\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \leq 1$, and then $|\mathbb{E}_{\tau \sim \pi^{k'}} [\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}|] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}}| \leq 2$.

Using the Azuma-Hoeffding inequality, we have that for any fixed $k \geq K_0 + 1$, with probability at least $1 - \frac{\delta'}{2k^2}$,

$$\left| \sum_{k'=K_0+1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}| \right] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| \leq \sqrt{2 \cdot 4(k - K_0 - 1) \log \left(\frac{4k^2}{\delta'} \right)}.$$

Since $\sum_{k=K_0+1}^{\infty} \frac{\delta'}{2k^2} \leq \delta'$, by a union bound over k , we have that with probability at least δ' , for any $k \geq K_0 + 1$,

$$\begin{aligned}
\left| \sum_{k'=K_0+1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}| \right] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| &\leq \sqrt{2 \cdot 4(k - K_0 - 1) \log \left(\frac{4k^2}{\delta'} \right)} \\
&\leq 4\sqrt{k \log \left(\frac{4k}{\delta'} \right)}.
\end{aligned}$$

\square

1998 *Proof of Theorem 4.* Let $\delta' = \frac{\delta}{3}$. We have $\Pr[\mathcal{J} \cap \mathcal{K} \cap \mathcal{F}_{\text{opt}}^{\text{S}}] \geq 1 - \delta$. To prove this theorem, it
 1999 suffices to prove the regret bound when event $\mathcal{J} \cap \mathcal{K} \cap \mathcal{F}_{\text{opt}}^{\text{S}}$ holds.
 2000

2001 Assume that event $\mathcal{J} \cap \mathcal{K} \cap \mathcal{F}_{\text{opt}}^{\text{S}}$ holds. Then, we have
 2002

$$\begin{aligned}
 2003 \quad \mathcal{R}(K) &= \sum_{k=1}^K \left((\phi^{\pi^*})^\top \theta - (\phi^{\pi^k})^\top \theta \right) \\
 2004 &\stackrel{(a)}{\leq} \sum_{k=K_0+1}^K \left((\phi^{\pi^*})^\top \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\phi^{\pi^*}\|_{(\Sigma_{k-1})^{-1}} - (\phi^{\pi^k})^\top \theta \right) + K_0 H \\
 2005 &\stackrel{(b)}{\leq} \sum_{k=K_0+1}^K \left((\phi^{\pi^k})^\top \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\phi^{\pi^k}\|_{(\Sigma_{k-1})^{-1}} - (\phi^{\pi^k})^\top \theta \right) + K_0 H \\
 2006 &\leq \sum_{k=K_0+1}^K 2\beta(k-1) \cdot \|\phi^{\pi^k}\|_{(\Sigma_{k-1})^{-1}} + K_0 H \\
 2007 &= 2\beta(K) \sum_{k=K_0+1}^K \|\mathbb{E}_{\tau \sim \pi^k} [\phi^\tau | F_{k-1}]\|_{(\Sigma_{k-1})^{-1}} + K_0 H \\
 2008 &\stackrel{(c)}{\leq} 2\beta(K) \sum_{k=K_0+1}^K \mathbb{E}_{\tau \sim \pi^k} \left[\|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} | F_{k-1} \right] + K_0 H \\
 2009 &= 2\beta(K) \sum_{k=K_0+1}^K \left(\mathbb{E}_{\tau \sim \pi^k} \left[\|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} | F_{k-1} \right] - \|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} + \|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} \right) + K_0 H \\
 2010 &\leq 2\beta(K) \sum_{k=K_0+1}^K \left(\mathbb{E}_{\tau \sim \pi^k} \left[\|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} | F_{k-1} \right] - \|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} + \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}} \right) \\
 2011 &\quad + K_0 H, \tag{29}
 \end{aligned}$$

2012 where (a) follows from Eq. (27), (b) is due to the definition of π^k , and (c) uses the Jensen inequality.

2013 Plugging Eq. (28) and Lemma 23 into Eq. (29) and using the fact that $\lambda := \frac{H}{r_{\max}^2 m}$, we have

$$\begin{aligned}
 2014 \quad \mathcal{R}(K) &\leq 2 \left(\sqrt{\frac{H|\mathcal{S}||\mathcal{A}|}{m} \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)} + 2 \log \left(\frac{1}{\delta'} \right) + r_{\max} \sqrt{\lambda|\mathcal{S}||\mathcal{A}|} \right) \\
 2015 &\quad \left(4\sqrt{K \log \left(\frac{4K}{\delta} \right)} + \sqrt{2Km|\mathcal{S}||\mathcal{A}| \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)} \right) \\
 2016 &\quad + H \left[\max \left\{ 26H^4(1+\gamma)^2(z^*)^2 \log \left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta'} \right), \frac{|\mathcal{S}||\mathcal{A}|}{\gamma^2} \right\} \right] \\
 2017 &= O \left(|\mathcal{S}||\mathcal{A}| \sqrt{HK} \log \left(\left(1 + \frac{KHr_{\max}^2}{|\mathcal{S}||\mathcal{A}|m} \right) \frac{1}{\delta} \right) + (z^*)^2 H^5 \log \left(\frac{|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}||\mathcal{A}|H \right).
 \end{aligned}$$

2018 \square

2019 D.2 PROOF FOR THE REGRET LOWER BOUND WITH KNOWN TRANSITION

2020 Now we prove the regret lower bound (Theorem 5) for RL with sum segment feedback and known
 2021 transition.

2022 *Proof of Theorem 5.* We construct a random instance \mathcal{I} as follows. As shown in Figure 5, there are
 2023 n bandit states s_1, \dots, s_n (i.e., there is an optimal action and multiple suboptimal actions), a good
 2024 absorbing state s_{n+1} and a bad absorbing state s_{n+2} . The agent starts from s_1, \dots, s_n with equal
 2025 probability $\frac{1}{n}$. For any $i \in [n]$, in state s_i , one action a_j is uniformly chosen from \mathcal{A} as the optimal

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

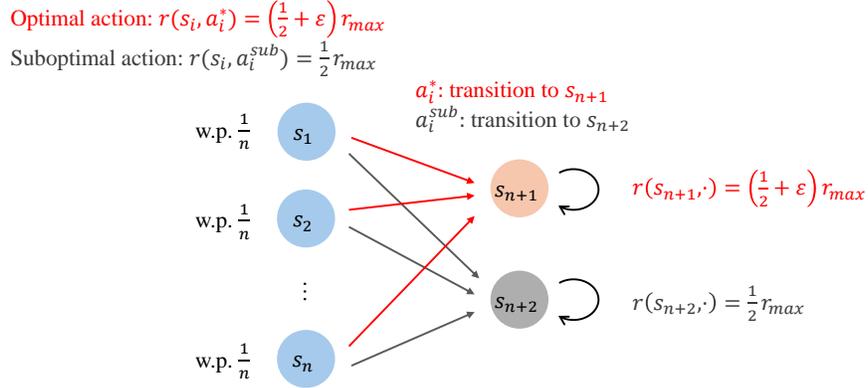


Figure 5: Instance for the lower bound under sum segment feedback and known transition.

action. In state s_i , under the optimal action a_j , the agent transitions to s_{n+1} deterministically, and $r(s_i, a_j) = \left(\frac{1}{2} + \varepsilon\right)r_{\max}$, where $\varepsilon \in (0, \frac{1}{2}]$ is a parameter specified later; Under any suboptimal action $a \in \mathcal{A} \setminus \{s_j\}$, the agent transitions to s_{n+2} deterministically, and $r(s_i, a) = \frac{1}{2}r_{\max}$. For all actions $a \in \mathcal{A}$, $r(s_{n+1}, a) = \left(\frac{1}{2} + \varepsilon\right)r_{\max}$ and $r(s_{n+2}, a) = \frac{1}{2}r_{\max}$. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, the reward distribution of (s, a) is Gaussian distribution $\mathcal{N}(r(s, a), 1)$.

In this proof, we will also use an alternative uniform instance $\mathcal{I}_{\text{unif}}$. The only difference between $\mathcal{I}_{\text{unif}}$ and \mathcal{I} is that for any $i \in [n]$, in state s_i , under all actions $a \in \mathcal{A}$, the agent transitions to s_{n+2} deterministically, and $r(s_i, a) = \frac{1}{2}r_{\max}$.

Fix an algorithm \mathbb{A} . Let $\mathbb{E}_{\text{unif}}[\cdot]$ denote the expectation with respect to $\mathcal{I}_{\text{unif}}$. Let $\mathbb{E}_*[\cdot]$ denote the expectation with respect to \mathcal{I} . For any $i \in [n]$ and $j \in [|\mathcal{A}|]$, let $\mathbb{E}_{i,j}[\cdot]$ denote the expectation with respect to the case where a_j is the optimal action in state s_i , and $N_{i,j}$ denote the number of episodes where algorithm \mathbb{A} chooses a_j in state s_i , i.e., $N_{i,j} = \sum_{k=1}^K \mathbb{1}\{\pi_1^k(s_i) = a_j\}$.

The KL divergence of the reward observations if taking a_j in s_i ($i \in [n]$) between $\mathcal{I}_{\text{unif}}$ and \mathcal{I} is

$$\begin{aligned} & \sum_{i=1}^m \text{KL} \left(\mathcal{N} \left(\frac{1}{2}r_{\max} \cdot \frac{H}{m}, \frac{H}{m} \right) \parallel \mathcal{N} \left(\left(\frac{1}{2} + \varepsilon \right) r_{\max} \cdot \frac{H}{m}, \frac{H}{m} \right) \right) \\ &= m \cdot \frac{\left(\frac{H}{m} \cdot r_{\max} \varepsilon \right)^2}{\frac{H}{m}} \\ &= Hr_{\max}^2 \varepsilon^2. \end{aligned}$$

In addition, the agent has probability only $\frac{1}{n}$ to arrive at (observe) state s_i .

Hence, using Lemma A.1 in (Auer et al., 2002), we have that for any $i \in [n]$, in state s_i ,

$$\mathbb{E}_{i,j}[N_{i,j}] \leq \mathbb{E}_{\text{unif}}[N_{i,j}] + \frac{K}{2} \sqrt{\frac{1}{n} \cdot \mathbb{E}_{\text{unif}}[N_{i,j}] \cdot Hr_{\max}^2 \varepsilon^2}.$$

Summing over $j \in [|\mathcal{A}|]$, using the Cauchy-Schwarz inequality and the fact that $\sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{\text{unif}}[N_{i,j}] = K$, we have

$$\begin{aligned} \sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{i,j}[N_{i,j}] &\leq K + \frac{K}{2} \sqrt{\frac{|\mathcal{A}|}{n} \cdot K \cdot Hr_{\max}^2 \varepsilon^2} \\ &= K + \frac{Kr_{\max} \varepsilon}{2} \sqrt{\frac{|\mathcal{A}|HK}{n}}. \end{aligned}$$

Algorithm 4: LinUCB-Tran

Input: $\delta, \delta' := \frac{\delta}{4}, \lambda := \frac{H}{m}, L := \log\left(\frac{3|\mathcal{S}||\mathcal{A}|H}{\delta'}\right) + S \log(8e(1 + KH))$. For any $k \geq 1$,

$$\beta(k) := \sqrt{\frac{H|\mathcal{S}||\mathcal{A}|}{m} \log\left(1 + \frac{kH^2}{\lambda|\mathcal{S}||\mathcal{A}|m}\right)} + 2 \log\left(\frac{1}{\delta'}\right) + r_{\max} \sqrt{\lambda|\mathcal{S}||\mathcal{A}|}.$$

1 **for** $k = 1, \dots, K$ **do**

- 2 $\hat{\theta}_{k-1} \leftarrow (\lambda I + \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau^{k'}} (\phi^{\tau^{k'}})^\top)^{-1} \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau^{k'}} R_i^{k'}$;
 3 $\Sigma_{k-1} \leftarrow \lambda I + \sum_{k'=1}^{k-1} \sum_{i=1}^m \phi^{\tau^{k'}} (\phi^{\tau^{k'}})^\top$;
 4 $\pi^k \leftarrow \operatorname{argmax}_{\pi \in \Pi} ((\hat{\phi}_{k-1}^\pi)^\top \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\hat{\phi}_{k-1}^\pi\|_{(\Sigma_{k-1})^{-1}} + \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} [B_1^{\pi; s', a'; k}(s_1)])$,
 where $B_1^{\pi; s', a'; k}(s_1)$ is defined in Eq. (31);
 5 Play episode k with policy π^k . Observe τ^k and sum segment feedback $\{R_i^k\}_{i=1}^m$;

Then, we have

$$\begin{aligned} \mathcal{R}(K) &= \sum_{k=1}^K \mathbb{E}_* [V^* - V^{\pi^k}] \\ &= \left(\frac{1}{2} + \varepsilon\right) r_{\max} HK - \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} r_{\max} HK + \varepsilon r_{\max} H \cdot \frac{1}{|\mathcal{A}|} \sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{i,j} [N_{i,j}]\right) \\ &= \varepsilon r_{\max} H \left(K - \frac{K}{|\mathcal{A}|} - \frac{K r_{\max} \varepsilon}{2} \sqrt{\frac{HK}{n|\mathcal{A}|}}\right). \end{aligned}$$

Recall that $n = |\mathcal{S}| - 2$. Let $|\mathcal{S}| \geq 3, |\mathcal{A}| \geq 2, K \geq \frac{n|\mathcal{A}|}{r_{\max}^2 H}$ and $\varepsilon = \frac{1}{2r_{\max}} \sqrt{\frac{n|\mathcal{A}|}{HK}}$. Then, we have

$$\mathcal{R}(K) = \Omega\left(\sqrt{|\mathcal{S}||\mathcal{A}|HK}\right).$$

□

D.3 PSEUDO-CODE AND DETAILED DESCRIPTION OF ALGORITHM LinUCB-Tran

Algorithm 4 presents the pseudo-code of LinUCB-Tran. In each episode k , similar to algorithm E-LinUCB, LinUCB-Tran first computes the least squares estimate of the reward parameter $\hat{\theta}_{k-1}$ and covariance matrix Σ_{k-1} with past observations (Lines 2-3).

Then, we introduce the transition estimation in LinUCB-Tran. We first define some notation which also appears in algorithm SegBiTS-Tran. For any $k > 0$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$, let $\hat{p}_k(\cdot | s, a)$ denote the empirical estimate of $p(\cdot | s, a)$, and $n_k(s, a)$ denote the number of times (s, a) was visited up to the end of episode k . In addition, for any policy π , let $\hat{\phi}_k^\pi(s, a)$ denote the expected number of times (s, a) is visited in an episode under policy π on empirical MDP \hat{p}_{k-1} (see Eq. (17) for the formal definition).

Below we establish a bound for the deviation between $\hat{\phi}_{k-1}^\pi$ and ϕ^π . For ease of analysis, we first connect ϕ^π with a newly-defined visitation value function $G_h^{\pi; s', a'}(s; p)$. For any transition model p' , policy π and $(s', a') \in \mathcal{S} \times \mathcal{A}$, if regarding hitting (s', a') as an instantaneous reward one, then we can define a visitation value function:

$$\begin{cases} G_h^{\pi; s', a'}(s; p') = \mathbb{1}\{s = s', \pi_h(s) = a'\} + p(\cdot | s, \pi_h(s))^\top G_{h+1}^{\pi; s', a'}(\cdot), & \forall s \in \mathcal{S}, \forall h \in [H], \\ G_{H+1}^{\pi; s', a'}(s; p') = 0, & \forall s \in \mathcal{S}. \end{cases} \quad (30)$$

$G_h^{\pi; s', a'}(s; p')$ denotes the expected cumulative number of times (s', a') was hit starting from s at step h under policy π on MDP p' , till the end of this episode. It holds that $\phi^\pi(s', a') = \mathbb{E}_{s_1 \sim \rho} [G_1^{\pi; s', a'}(s_1 | p)]$ and $\hat{\phi}_{k-1}^\pi(s', a') = \mathbb{E}_{s_1 \sim \rho} [G_1^{\pi; s', a'}(s_1 | \hat{p}_{k-1})]$ for any $(s', a') \in \mathcal{S} \times \mathcal{A}$.

With the definition of $G_h^{\pi; s', a'}$, bounding the deviation between $\hat{\phi}_{k-1}^\pi$ and ϕ^π is similar to bounding the gap between the estimated and true value functions. Then, we can build a Bernstein-type uncertainty bound between $\hat{\phi}_{k-1}^\pi$ and ϕ^π using the variance of $G_h^{\pi; s', a'}$. For any policy π , $(s', a') \in \mathcal{S} \times \mathcal{A}$ and $k > 0$, define

$$\begin{cases} B_h^{\pi; s', a'; k}(s) = \min \left\{ \left(4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{13H^2L}{n_{k-1}(s, \pi_h(s))} \right. \right. \\ \left. \left. + \left(1 + \frac{2}{H}\right) \hat{p}_{k-1}(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot), H \right\}, \quad \forall s \in \mathcal{S}, \forall h \in [H], \\ B_{H+1}^{\pi; s', a'; k}(s) = 0, \quad \forall s \in \mathcal{S}. \end{cases} \quad (31)$$

The construction of $B_h^{\pi; s', a'; k}(s)$ satisfies (see Lemma 29 for more details)

$$\begin{aligned} |\hat{\phi}_{k-1}^\pi(s', a') - \phi^\pi(s', a')| &\leq \mathbb{E}_{s_1 \sim \rho} [B_1^{\pi; s', a'; k}(s_1)], \quad \forall (s', a') \in \mathcal{S} \times \mathcal{A}, \\ \|\hat{\phi}_{k-1}^\pi - \phi^\pi\|_1 &\leq \sum_{(s', a')} \mathbb{E}_{s_1 \sim \rho} [B_1^{\pi; s', a'; k}(s_1)]. \end{aligned}$$

Incorporating this transition uncertainty $\mathbb{E}_{s_1 \sim \rho}[B_1^{\pi; s', a'; k}(s_1)]$ and reward uncertainty $\|\hat{\phi}_{k-1}^\pi\|_{(\Sigma_{k-1})^{-1}}$ into exploration bonuses, LinUCB-Tran computes the optimal policy π^k under optimistic estimation (Line 4). After that, LinUCB-Tran plays episode k with π^k , and collects trajectory τ^k and reward observation on each segment $\{R_i^k\}_{i=1}^m$ (Line 5).

D.4 PROOF FOR THE REGRET UPPER BOUND WITH UNKNOWN TRANSITION

In the following, we prove the regret upper bound (Theorem 6) of algorithm LinUCB-Tran for unknown transition.

Recall the definition of events \mathcal{G}_{KL} and \mathcal{H} in Eqs. (18) and (19), respectively.

For any $k > 0$, define the set of state-action pairs

$$D_k := \left\{ (s, a) \in \mathcal{S} \times \mathcal{A} : \frac{1}{4} \sum_{k'=1}^k w_{k'}(s, a) \geq H \log \left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'} \right) + H \right\}. \quad (32)$$

D_k stands for the set of state-action pairs which have sufficient visitations in expectation.

Lemma 26. *Assume that event \mathcal{H} holds. Then, if $(s, a) \in D_k$,*

$$n_{k-1}(s, a) \geq \frac{1}{4} \sum_{k'=1}^k w_{k'}(s, a).$$

Proof. We have

$$\begin{aligned} n_{k-1}(s, a) &\geq \frac{1}{2} \sum_{k'=1}^{k-1} w_{k'}(s, a) - H \log \left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'} \right) \\ &= \frac{1}{4} \sum_{k'=1}^{k-1} w_{k'}(s, a) + \frac{1}{4} \sum_{k'=1}^{k-1} w_{k'}(s, a) - H \log \left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'} \right) \\ &= \frac{1}{4} \sum_{k'=1}^k w_{k'}(s, a) + \frac{1}{4} \sum_{k'=1}^k w_{k'}(s, a) - H \log \left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'} \right) - \frac{1}{2} w_k(s, a) \\ &\stackrel{(a)}{\geq} \frac{1}{4} \sum_{k'=1}^k w_{k'}(s, a) + H - \frac{1}{2} w_k(s, a) \\ &\geq \frac{1}{4} \sum_{k'=1}^k w_{k'}(s, a), \end{aligned}$$

where (a) is due to the definition of D_k (Eq. (32)). □

Lemma 27. *It holds that*

$$\sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \notin D_k} w_{k,h}(s,a) \leq 8|\mathcal{S}||\mathcal{A}|H \log\left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'}\right).$$

Proof. If $(s,a) \notin D_k$, then

$$\frac{1}{4} \sum_{k'=1}^k w_{k'}(s,a) < H \log\left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'}\right) + H.$$

Thus, we have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \notin D_k} w_{k,h}(s,a) &= \sum_{(s,a)} \sum_{k=1}^K \sum_{h=1}^H \mathbb{1}\{(s,a) \notin D_k\} \cdot w_{k,h}(s,a) \\ &= \sum_{(s,a)} \sum_{k=1}^K \mathbb{1}\{(s,a) \notin D_k\} \cdot w_k(s,a) \\ &\leq 4|\mathcal{S}||\mathcal{A}|H \log\left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'}\right) + 4|\mathcal{S}||\mathcal{A}|H \\ &\leq 8|\mathcal{S}||\mathcal{A}|H \log\left(\frac{|\mathcal{S}||\mathcal{A}|H}{\delta'}\right). \end{aligned}$$

□

Lemma 28. *Assume that event \mathcal{H} holds. Then, we have*

$$\sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} \frac{w_{k,h}(s,a)}{n_{k-1}(s,a)} \leq 4|\mathcal{S}||\mathcal{A}| \log(2KH).$$

Proof. It holds that

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \sum_{(s,a) \in D_k} \frac{w_{k,h}(s,a)}{n_{k-1}(s,a)} &= \sum_{k=1}^K \sum_{(s,a) \in D_k} \frac{w_k(s,a)}{n_{k-1}(s,a)} \\ &= \sum_{k=1}^K \sum_{(s,a)} \frac{w_k(s,a)}{n_{k-1}(s,a)} \cdot \mathbb{1}\{(s,a) \in D_k\} \\ &\stackrel{(a)}{\leq} 4 \sum_{k=1}^K \sum_{(s,a)} \frac{w_k(s,a)}{\sum_{k'=1}^k w_k(s,a)} \cdot \mathbb{1}\{(s,a) \in D_k\} \\ &= 4 \sum_{(s,a)} \sum_{k=1}^K \frac{w_k(s,a)}{\sum_{k'=1}^k w_k(s,a)} \cdot \mathbb{1}\{(s,a) \in D_k\} \\ &\stackrel{(b)}{\leq} 4|\mathcal{S}||\mathcal{A}| \log(2KH), \end{aligned}$$

where (a) uses Lemma 26, and (b) follows from the analysis of Lemma 13 in (Zanette & Brunskill, 2019). □

Lemma 29 (Error in Visitation Vectors). *Assume that event \mathcal{G}_{KL} holds. Then, for any $k > 0$ and policy π ,*

$$\|\hat{\phi}_{k-1}(\pi) - \phi(\pi)\|_1 \leq \sum_{s',a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right].$$

2268 *Proof.* Since $\phi^\pi(s', a') = \mathbb{E}_{s_1 \sim \rho}[G_1^{\pi; s', a'}(s_1 | p)]$ and $\hat{\phi}_{k-1}^\pi(s', a') = \mathbb{E}_{s_1 \sim \rho}[G_1^{\pi; s', a'}(s_1 | \hat{p}_{k-1})]$, in
 2269 this proof, we investigate the error in $G_h^{\pi; s', a'}$ due to the estimation of the transition model.
 2270

2271 In the following, we prove by induction that for any $h \in [H]$ and $s \in \mathcal{S}$, $|G_h^{\pi; s', a'}(s | \hat{p}_{k-1}) -$
 2272 $G_h^{\pi; s', a'}(s | p)| \leq B_h^{\pi; s', a'; k}(s)$.
 2273

2274 When $h = H + 1$, by definition, we have $G_{H+1}^{\pi; s', a'}(s | \hat{p}_{k-1}) = G_{H+1}^{\pi; s', a'}(s | p) = B_{H+1}^{\pi; s', a'; k}(s) = 0$ for
 2275 any $s \in \mathcal{S}$, and then the above statement trivially holds.
 2276

2277 When $1 \leq h \leq H$, if $|G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot | p)| \leq B_{h+1}^{\pi; s', a'; k}(\cdot)$ element-wise, then for any
 2278 $s \in \mathcal{S}$, we have

$$\begin{aligned}
 & |G_h^{\pi; s', a'}(s | \hat{p}_{k-1}) - G_h^{\pi; s', a'}(s | p)| \\
 &= \left| \hat{p}_{k-1}(\cdot | s, \pi_h(s))^\top G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1}) - p(\cdot | s, \pi_h(s))^\top G_{h+1}^{\pi; s', a'}(\cdot | p) \right| \\
 &= \hat{p}_{k-1}(\cdot | s, \pi_h(s))^\top \left| G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot | p) \right| \\
 &\quad + \left| \hat{p}_{k-1}(\cdot | s, \pi_h(s)) - p(\cdot | s, \pi_h(s)) \right|^\top G_{h+1}^{\pi; s', a'}(\cdot | p) \\
 &\stackrel{(a)}{\leq} \hat{p}_{k-1}(\cdot | s, \pi_h(s))^\top \left| G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot | p) \right| + 2\sqrt{\frac{\text{Var}_{p(\cdot | s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot | p)) \cdot L}{n_{k-1}(s, \pi_h(s))}} \\
 &\quad + \frac{HL}{n_{k-1}(s, \pi_h(s))}, \tag{33}
 \end{aligned}$$

2291 where (a) is due to Lemma 37.

2292 Here, we have

$$\begin{aligned}
 & \text{Var}_{p(\cdot | s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot | p)) \\
 &\stackrel{(a)}{\leq} 2\text{Var}_{\hat{p}_{k-1}(\cdot | s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot | p)) + \frac{4H^2L}{n_{k-1}(s, \pi_h(s))} \\
 &\stackrel{(b)}{\leq} 4\text{Var}_{\hat{p}_{k-1}(\cdot | s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1})) + 4H\hat{p}_{k-1}(\cdot | s, \pi_h(s))^\top |G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot | p)| \\
 &\quad + \frac{4H^2L}{n_{k-1}(s, \pi_h(s))},
 \end{aligned}$$

2302 where (a) uses Lemma 38 and (b) comes from Lemma 39.

2303 Then,

$$\begin{aligned}
 & \sqrt{\frac{\text{Var}_{p(\cdot | s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot | p)) \cdot L}{n_{k-1}(s, \pi_h(s))}} \tag{34} \\
 &\leq \sqrt{\frac{4\text{Var}_{\hat{p}_{k-1}(\cdot | s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} \\
 &\quad + \sqrt{\frac{1}{H}\hat{p}_{k-1}(\cdot | s, \pi_h(s))^\top |G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot | p)| \cdot \frac{4H^2L}{n_{k-1}(s, \pi_h(s))} + \frac{2HL}{n_{k-1}(s, \pi_h(s))}} \\
 &\stackrel{(a)}{\leq} 2\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot | s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} \\
 &\quad + \frac{1}{H}\hat{p}_{k-1}(\cdot | s, \pi_h(s))^\top |G_{h+1}^{\pi; s', a'}(\cdot | \hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot | p)| + \frac{6H^2L}{n_{k-1}(s, \pi_h(s))}, \tag{35}
 \end{aligned}$$

2319 where (a) is due to the fact that $\sqrt{xy} \leq x + y$.

Hence, plugging Eq. (35) into Eq. (33) and using the fact that $|G_h^{\pi; s', a'}(s)| \in [0, H]$, we have

$$\begin{aligned}
& |G_h^{\pi; s', a'}(s|\hat{p}_{k-1}) - G_h^{\pi; s', a'}(s|p)| \\
& \leq \left(4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{13H^2L}{n_{k-1}(s, \pi_h(s))} \right. \\
& \quad \left. + \left(1 + \frac{2}{H}\right) \hat{p}_{k-1}(\cdot|s, \pi_h(s))^\top \left| G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot|p) \right| \right) \wedge H. \\
& \leq \left(4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{13H^2L}{n_{k-1}(s, \pi_h(s))} \right. \\
& \quad \left. + \left(1 + \frac{2}{H}\right) \hat{p}_{k-1}(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot) \right) \wedge H \\
& = B_h^{\pi; s', a'; k}(s),
\end{aligned}$$

which completes the induction proof.

Therefore,

$$\begin{aligned}
\left| \hat{\phi}_{k-1}^{\pi}(s', a') - \phi^{\pi}(s', a') \right| &= \left| \mathbb{E}_{s_1 \sim \rho} \left[G_1^{\pi; s', a'}(s_1|\hat{p}_{k-1}) \right] - \mathbb{E}_{s_1 \sim \rho} \left[G_1^{\pi; s', a'}(s_1|p) \right] \right| \\
&\leq \mathbb{E}_{s_1 \sim \rho} \left[\left| G_1^{\pi; s', a'}(s_1|\hat{p}_{k-1}) - G_1^{\pi; s', a'}(s_1|p) \right| \right] \\
&\leq \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right].
\end{aligned}$$

Summing over $(s', a') \in \mathcal{S} \times \mathcal{A}$, we obtain this lemma. \square

Lemma 30. Assume that event $\mathcal{G}_{\text{KL}} \cap \mathcal{H}$ holds. Then, for any $k > 0$ and policy π ,

$$\begin{aligned}
& \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right] \\
& \leq e^{12} \sum_{h=1}^H \sum_{s, a} w_h^{\pi}(s, a) \left(8\sqrt{\frac{\text{Var}_{p(\cdot|s, a)}(G_{h+1}^{\pi; s', a'}(\cdot|p)) \cdot L}{n_{k-1}(s, a)}} + \frac{46H^2L}{n_{k-1}(s, a)} \right) \wedge H,
\end{aligned}$$

and

$$\begin{aligned}
& \sum_{s', a'} \sum_{k=1}^K \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; k; s', a'; k}(s_1) \right] \\
& \leq 16e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H \sqrt{KL \log(2KH)} + 192e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log(2KH).
\end{aligned}$$

Proof. First, we prove the first statement.

For any policy π , $k > 0$, $(s', a') \in \mathcal{S} \times \mathcal{A}$, $h \in [H]$ and $s \in \mathcal{S}$, we have

$$\begin{aligned}
B_h^{\pi; s', a'; k}(s) &\leq 4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{13H^2L}{n_{k-1}(s, \pi_h(s))} \\
& \quad + \left(1 + \frac{2}{H}\right) \hat{p}_{k-1}(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot) \\
& = 4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{13H^2L}{n_{k-1}(s, \pi_h(s))} \\
& \quad + \left(1 + \frac{2}{H}\right) p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot) \\
& \quad + \left(1 + \frac{2}{H}\right) (\hat{p}_{k-1}(\cdot|s, \pi_h(s)) - p(\cdot|s, \pi_h(s)))^\top B_{h+1}^{\pi; s', a'; k}(\cdot)
\end{aligned}$$

$$\begin{aligned}
& \stackrel{(a)}{\leq} 4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{13H^2L}{n_{k-1}(s, \pi_h(s))} \\
& \quad + \left(1 + \frac{2}{H}\right) p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot) \\
& \quad + \left(1 + \frac{2}{H}\right) \cdot \left(2\sqrt{\frac{\text{Var}_{p(\cdot|s, \pi_h(s))}(B_{h+1}^{\pi; s', a'; k}(\cdot)) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{HL}{n_{k-1}(s, \pi_h(s))}\right) \\
& \leq 4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{13H^2L}{n_{k-1}(s, \pi_h(s))} \\
& \quad + \left(1 + \frac{2}{H}\right) p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot) \\
& \quad + \left(1 + \frac{2}{H}\right) \left(2\sqrt{\frac{1}{H}p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot) \frac{H^2L}{n_{k-1}(s, \pi_h(s))}} + \frac{HL}{n_{k-1}(s, \pi_h(s))}\right) \\
& \stackrel{(b)}{\leq} 4\sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{22H^2L}{n_{k-1}(s, \pi_h(s))} \\
& \quad + \left(1 + \frac{8}{H}\right) p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot), \tag{36}
\end{aligned}$$

where (a) uses Lemma 37, and (b) follows from the fact that $\sqrt{xy} \leq x + y$.

In addition, we have

$$\begin{aligned}
& \text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \\
& \stackrel{(a)}{=} 2\text{Var}_{p(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) + \frac{4H^2L}{n_{k-1}(s, a)} \\
& \stackrel{(b)}{\leq} 4\text{Var}_{p(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|p)) + 4Hp(\cdot|s, \pi_h(s))^\top \left| G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1}) - G_{h+1}^{\pi; s', a'}(\cdot|p) \right| \\
& \quad + \frac{4H^2L}{n_{k-1}(s, a)} \\
& \leq 4\text{Var}_{p(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|p)) + 4Hp(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1}) + \frac{4H^2L}{n_{k-1}(s, a)},
\end{aligned}$$

where (a) uses Lemma 38, and (b) comes from Lemma 39.

Then,

$$\begin{aligned}
& \sqrt{\frac{\text{Var}_{\hat{p}_{k-1}(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1})) \cdot L}{n_{k-1}(s, \pi_h(s))}} \\
& \leq \sqrt{\frac{4\text{Var}_{p(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|p)) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \sqrt{\frac{1}{H}p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1}) \cdot \frac{4H^2L}{n_{k-1}(s, \pi_h(s))}} \\
& \quad + \frac{2HL}{n_{k-1}(s, a)} \\
& \leq 2\sqrt{\frac{\text{Var}_{p(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|p)) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{1}{H}p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'}(\cdot|\hat{p}_{k-1}) \\
& \quad + \frac{6H^2L}{n_{k-1}(s, \pi_h(s))} \tag{37}
\end{aligned}$$

2430 Plugging Eq. (37) into Eq. (36) and using the clipping definition of $B_h^{\pi; s', a'; k}(s)$, we have

$$2431 B_h^{\pi; s', a'; k}(s) \leq \left(8\sqrt{\frac{\text{Var}_{p(\cdot|s, \pi_h(s))}(G_{h+1}^{\pi; s', a'}(\cdot|p)) \cdot L}{n_{k-1}(s, \pi_h(s))}} + \frac{46H^2L}{n_{k-1}(s, \pi_h(s))} \right) \wedge H$$

$$2432 + \left(1 + \frac{12}{H} \right) p(\cdot|s, \pi_h(s))^\top B_{h+1}^{\pi; s', a'; k}(\cdot)$$

2433 Using the above inequality, taking $s_1 \sim \rho$, and unfolding $B_1^{\pi; s', a'; k}(s_1)$ over h , we have

$$2434 \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right]$$

$$2435 \leq e^{12} \sum_{h=1}^H \sum_{s, a} w_h^\pi(s, a) \left(8\sqrt{\frac{\text{Var}_{p(\cdot|s, a)}(G_{h+1}^{\pi; s', a'}(\cdot|p)) \cdot L}{n_{k-1}(s, a)}} + \frac{46H^2L}{n_{k-1}(s, a)} \right) \wedge H. \quad (38)$$

2436 Next, we prove the second statement.

2437 It holds that

$$2438 \sum_{s', a'} \sum_{k=1}^K \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi^k; s', a'; k}(s_1) \right]$$

$$2439 \leq e^{12} \sum_{s', a'} \sum_{k=1}^K \sum_{h=1}^H \sum_{(s, a) \in D_k} w_{k, h}(s, a) \left(8\sqrt{\frac{\text{Var}_{p(\cdot|s, a)}(G_{h+1}^{\pi^k; s', a'}(\cdot|p)) \cdot L}{n_{k-1}(s, a)}} + \frac{46H^2L}{n_{k-1}(s, a)} \right)$$

$$2440 + e^{12} H |\mathcal{S}| |\mathcal{A}| \sum_{k=1}^K \sum_{h=1}^H \sum_{(s, a) \notin D_k} w_{k, h}(s, a)$$

$$2441 \stackrel{(a)}{\leq} 8e^{12} \sqrt{L} \sum_{s', a'} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \sum_{(s, a) \in D_k} w_{k, h}(s, a) \text{Var}_{p(\cdot|s, a)}(G_{h+1}^{\pi^k; s', a'}(\cdot|p)) \cdot}$$

$$2442 \sqrt{\sum_{k=1}^K \sum_{h=1}^H \sum_{(s, a) \in D_k} \frac{w_{k, h}(s, a)}{n_{k-1}(s, a)}} + e^{12} |\mathcal{S}| |\mathcal{A}| \cdot 46H^2L \sum_{k=1}^K \sum_{h=1}^H \sum_{(s, a) \in D_k} \frac{w_{k, h}(s, a)}{n_{k-1}(s, a)}$$

$$2443 + 8e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 \log \left(\frac{|\mathcal{S}| |\mathcal{A}| H}{\delta'} \right)$$

$$2444 \stackrel{(b)}{\leq} 8e^{12} |\mathcal{S}| |\mathcal{A}| \sqrt{L} \sqrt{KH^2} \cdot \sqrt{4|\mathcal{S}| |\mathcal{A}| \log(2KH)} + 184e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log(2KH)$$

$$2445 + 8e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 \log \left(\frac{|\mathcal{S}| |\mathcal{A}| H}{\delta'} \right)$$

$$2446 \leq 16e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H \sqrt{KL \log(2KH)} + 192e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log(2KH),$$

2447 where (a) is due to Lemma 27, and (b) follows from Lemmas 36 and 28. \square

2448 **Lemma 31** (Optimism under Sum Feedback and Unknown Transition). *Assume that event \mathcal{G}_{KL} holds. Then, for any $k > 0$ and fixed policy π ,*

$$2449 V_1^\pi(s_1) \leq \hat{\phi}_{k-1}(\pi)^\top \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\hat{\phi}_{k-1}(\pi)\|_{(\Sigma_{k-1})^{-1}} + r_{\max} \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right].$$

2450 *Proof.* It holds that

$$2451 V_1^\pi(s_1) = \phi(\pi)^\top \theta$$

$$2452 = \hat{\phi}_{k-1}(\pi)^\top \hat{\theta}_{k-1} + \phi(\pi)^\top \theta - \hat{\phi}_{k-1}(\pi)^\top \theta + \hat{\phi}_{k-1}(\pi)^\top \theta - \hat{\phi}_{k-1}(\pi)^\top \hat{\theta}_{k-1}$$

$$2453 \leq \hat{\phi}_{k-1}(\pi)^\top \hat{\theta}_{k-1} + \|\phi(\pi) - \hat{\phi}_{k-1}(\pi)\|_1 \cdot \|\theta\|_\infty + \beta(k-1) \cdot \|\hat{\phi}_{k-1}(\pi)\|_{(\Sigma_{k-1})^{-1}}$$

$$\stackrel{(a)}{\leq} \hat{\phi}_{k-1}(\pi)^\top \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\hat{\phi}_{k-1}(\pi)\|_{(\Sigma_{k-1})^{-1}} + r_{\max} \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right],$$

where (a) uses Lemma 29. \square

Lemma 32. For any $K \geq 1$, we have

$$\sum_{k=1}^K \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}} \leq H \sqrt{\frac{2K|\mathcal{S}||\mathcal{A}|}{\lambda} \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)}.$$

Proof. We have

$$\begin{aligned} \sum_{k=1}^K \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}} &\leq \sqrt{Km \sum_{k=1}^K \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}}^2} \\ &= \sqrt{Km \sum_{k=1}^K \min \left\{ \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}}^2, \frac{H^2}{m\lambda} \right\}} \\ &= \sqrt{\frac{H^2 K}{\lambda} \sum_{k=1}^K \min \left\{ \frac{m\lambda}{H^2} \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}}^2, 1 \right\}} \\ &\stackrel{(a)}{\leq} \sqrt{\frac{2H^2 K}{\lambda} \sum_{k=1}^K \log \left(1 + \min \left\{ \frac{m\lambda}{H^2} \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}}^2, 1 \right\} \right)} \\ &\stackrel{(b)}{\leq} \sqrt{\frac{2H^2 K}{\lambda} \sum_{k=1}^K \log \left(1 + \sum_{i=1}^m \|\phi^{\tau_i^k}\|_{(\Sigma_{k-1})^{-1}}^2 \right)} \\ &\stackrel{(c)}{\leq} \sqrt{\frac{2KH^2|\mathcal{S}||\mathcal{A}|}{\lambda} \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)}, \end{aligned}$$

where inequality (a) uses the fact that $x \leq 2 \log(1+x)$ for any $0 \leq x \leq 1$, inequality (b) is due to the fact that $\lambda \leq \frac{H^2}{m}$, and inequality (c) follows from Lemma 22. \square

Define event

$$\mathcal{F}_{\text{reg}}^{\mathcal{S}} := \left\{ \left| \sum_{k'=1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} | F_{k'-1} \right] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| \leq 4H \sqrt{\frac{k}{\lambda} \log \left(\frac{4k}{\delta'} \right)}, \right. \\ \left. \forall k > 0 \right\}. \quad (39)$$

Event $\mathcal{F}_{\text{reg}}^{\mathcal{S}}$ is similar to $\mathcal{F}_{\text{opt}}^{\mathcal{S}}$, except that here the universal upper bound of $\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}}$ is $\frac{H}{\sqrt{\lambda}}$ rather than 1.

Lemma 33. It holds that

$$\Pr \left[\mathcal{F}_{\text{reg}}^{\mathcal{S}} \right] \geq 1 - \delta'.$$

Proof. For any $k' \geq 1$, we have that $\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \leq \frac{H}{\sqrt{\lambda}}$, and then $|\mathbb{E}_{\tau \sim \pi^{k'}} [\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} | F_{k'-1}] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}}| \leq \frac{2H}{\sqrt{\lambda}}$.

Using the Azuma-Hoeffding inequality, we have that for any fixed $k > 0$, with probability at least $1 - \frac{\delta'}{2k^2}$,

$$\left| \sum_{k'=1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} | F_{k'-1} \right] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| \leq \sqrt{2 \cdot \frac{4H^2}{\lambda} \cdot k \log \left(\frac{4k^2}{\delta'} \right)}.$$

Since $\sum_{k=1}^{\infty} \frac{\delta'}{2k^2} \leq \delta'$, by a union bound over k , we have that with probability at least δ' , for any $k \geq 1$,

$$\begin{aligned} \left| \sum_{k'=1}^k \left(\mathbb{E}_{\tau \sim \pi^{k'}} \left[\|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} |F_{k'-1}| \right] - \|\phi^\tau\|_{(\Sigma_{k'-1})^{-1}} \right) \right| &\leq \sqrt{2 \cdot \frac{4H^2}{\lambda} \cdot k \log \left(\frac{4k^2}{\delta'} \right)} \\ &\leq 4H \sqrt{\frac{k}{\lambda} \log \left(\frac{4k}{\delta'} \right)}. \end{aligned}$$

□

Proof of Theorem 6. Let $\delta' = \frac{\delta}{4}$. Then, we have $\Pr[\mathcal{K} \cap \mathcal{F}_{\text{reg}}^S \cap \mathcal{G}_{\text{KL}} \cap \mathcal{H}] \geq 1 - \delta$. Thus, it suffices to prove the regret upper bound when event $\mathcal{K} \cap \mathcal{F}_{\text{reg}}^S \cap \mathcal{G}_{\text{KL}} \cap \mathcal{H}$ holds.

Assume that event $\mathcal{K} \cap \mathcal{F}_{\text{reg}}^S \cap \mathcal{G}_{\text{KL}} \cap \mathcal{H}$ holds. For any $k > 0$, we have

$$\begin{aligned} &\sum_{k=1}^K \left(V^*(s_1) - V^{\pi^k}(s_1) \right) \\ &\stackrel{(a)}{\leq} \sum_{k=1}^K \left(\hat{\phi}_{k-1}(\pi^*)^\top \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\hat{\phi}_{k-1}(\pi^*)\|_{(\Sigma_{k-1})^{-1}} + r_{\max} \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi^*; s', a'; k}(s_1) \right] \right. \\ &\quad \left. - V^{\pi^k} \right) \\ &\stackrel{(b)}{\leq} \sum_{k=1}^K \left(\hat{\phi}_{k-1}(\pi^k)^\top \hat{\theta}_{k-1} + \beta(k-1) \cdot \|\hat{\phi}_{k-1}(\pi^k)\|_{(\Sigma_{k-1})^{-1}} + r_{\max} \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi^k; s', a'; k}(s_1) \right] \right. \\ &\quad \left. - V^{\pi^k} \right) \\ &\leq \sum_{k=1}^K \left(\hat{\phi}_{k-1}(\pi^k)^\top \hat{\theta}_{k-1} - \hat{\phi}_{k-1}(\pi^k)^\top \theta + \hat{\phi}_{k-1}(\pi^k)^\top \theta - (\phi^{\pi^k})^\top \theta \right. \\ &\quad \left. + \beta(k-1) \cdot \|\hat{\phi}_{k-1}(\pi^k)\|_{(\Sigma_{k-1})^{-1}} + r_{\max} \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi^k; s', a'; k}(s_1) \right] \right) \\ &\stackrel{(c)}{\leq} \sum_{k=1}^K \left(2\beta(k-1) \cdot \|\hat{\phi}_{k-1}(\pi^k)\|_{(\Sigma_{k-1})^{-1}} + 2r_{\max} \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi^k; s', a'; k}(s_1) \right] \right) \\ &\leq 2\beta(K) \sum_{k=1}^K \|\hat{\phi}_{k-1}(\pi^k)\|_{(\Sigma_{k-1})^{-1}} + 2r_{\max} \sum_{k=1}^K \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi^k; s', a'; k}(s_1) \right], \end{aligned} \tag{40}$$

where (a) uses Lemma 31, (b) is due to the definition of π^k , and (c) follows from Lemma 29 and the definition of event \mathcal{K} .

Next, we first bound $\sum_{k=1}^K \|\hat{\phi}_{k-1}(\pi^k)\|_{(\Sigma_{k-1})^{-1}}$.

We have

$$\begin{aligned} \sum_{k=1}^K \|\hat{\phi}_{k-1}(\pi^k)\|_{(\Sigma_{k-1})^{-1}} &\leq \sum_{k=1}^K \left(\|\phi^{\pi^k}\|_{(\Sigma_{k-1})^{-1}} + \|\hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k}\|_{(\Sigma_{k-1})^{-1}} \right) \\ &\leq \sum_{k=1}^K \left(\|\phi^{\pi^k}\|_{(\Sigma_{k-1})^{-1}} + \frac{1}{\sqrt{\lambda}} \cdot \|\hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k}\|_2 \right) \end{aligned}$$

$$\leq \sum_{k=1}^K \left(\|\phi^{\pi^k}\|_{(\Sigma_{k-1})^{-1}} + \frac{1}{\sqrt{\lambda}} \cdot \|\hat{\phi}_{k-1}(\pi^k) - \phi^{\pi^k}\|_1 \right). \quad (41)$$

Here we have

$$\begin{aligned} & \sum_{k=1}^K \|\phi^{\pi^k}\|_{(\Sigma_{k-1})^{-1}} \\ &= \sum_{k=1}^K \|\mathbb{E}_{\tau \sim \pi^k} [\phi^\tau | F_{k-1}]\|_{(\Sigma_{k-1})^{-1}} \\ &\stackrel{(a)}{\leq} \sum_{k=1}^K \mathbb{E}_{\tau \sim \pi^k} \left[\|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} | F_{k-1} \right] \\ &= \sum_{k=1}^K \left(\mathbb{E}_{\tau \sim \pi^k} \left[\|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} | F_{k-1} \right] - \|\phi(\tau^k)\|_{(\Sigma_{k-1})^{-1}} + \|\phi(\tau^k)\|_{(\Sigma_{k-1})^{-1}} \right) \\ &\leq \sum_{k=1}^K \left(\mathbb{E}_{\tau \sim \pi^k} \left[\|\phi^\tau\|_{(\Sigma_{k-1})^{-1}} | F_{k-1} \right] - \|\phi(\tau^k)\|_{(\Sigma_{k-1})^{-1}} + \sum_{i=1}^m \|\phi^{\tau^i}\|_{(\Sigma_{k-1})^{-1}} \right) \\ &\stackrel{(b)}{\leq} 4H \sqrt{\frac{K}{\lambda} \log \left(\frac{4K}{\delta'} \right)} + H \sqrt{\frac{2K|\mathcal{S}||\mathcal{A}|}{\lambda} \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)}, \end{aligned} \quad (42)$$

where (a) uses the Jensen inequality, and (b) comes from the definition of $\mathcal{F}_{\text{reg}}^S$ and Lemma 32.

Hence, plugging Eq. (42) into Eq. (41) and using Lemma 29, we have

$$\begin{aligned} \sum_{k=1}^K \|\hat{\phi}_{k-1}(\pi^k)\|_{(\Sigma_{k-1})^{-1}} &\leq 4H \sqrt{\frac{K}{\lambda} \log \left(\frac{4K}{\delta'} \right)} + H \sqrt{\frac{2K|\mathcal{S}||\mathcal{A}|}{\lambda} \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|} \right)} \\ &\quad + \frac{1}{\sqrt{\lambda}} \sum_{k=1}^K \sum_{s', a'} \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right]. \end{aligned} \quad (43)$$

On the other hand, according to Eq. (38), we have

$$\mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi; s', a'; k}(s_1) \right] \leq e^{12} \sum_{h=1}^H \sum_{s, a} w_h^\pi(s, a) \left(8 \sqrt{\frac{\text{Var}_{p(\cdot|s, a)}(G_{h+1}^{\pi; s', a'}(\cdot|p)) \cdot L}{n_{k-1}(s, a)} + \frac{46H^2L}{n_{k-1}(s, a)}} \right) \wedge H.$$

Therefore, plugging Eqs. (43) and (38) into Eq. (40), we have

$$\begin{aligned} & \sum_{k=1}^K (V^* - V^{\pi^k}) \\ &\leq 2\beta(K) \left(4H \sqrt{\frac{K}{\lambda} \log \left(\frac{4K}{\delta'} \right)} + H \sqrt{\frac{2K|\mathcal{S}||\mathcal{A}|}{\lambda} \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)} \right) \\ &\quad + 2 \left(\frac{\beta(K)}{\sqrt{\lambda}} + r_{\max} \right) \sum_{s', a'} \sum_{k=1}^K \mathbb{E}_{s_1 \sim \rho} \left[B_1^{\pi^k; s', a'; k}(s_1) \right] \\ &\stackrel{(a)}{\leq} 2\beta(K) \left(4H \sqrt{\frac{K}{\lambda} \log \left(\frac{4K}{\delta'} \right)} + H \sqrt{\frac{2K|\mathcal{S}||\mathcal{A}|}{\lambda} \log \left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m} \right)} \right) \\ &\quad + \frac{4\beta(K)}{\sqrt{\lambda}} \left(16e^{12} |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H \sqrt{KL \log(2KH)} + 192e^{12} |\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L \log(2KH) \right) \end{aligned}$$

2645

2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

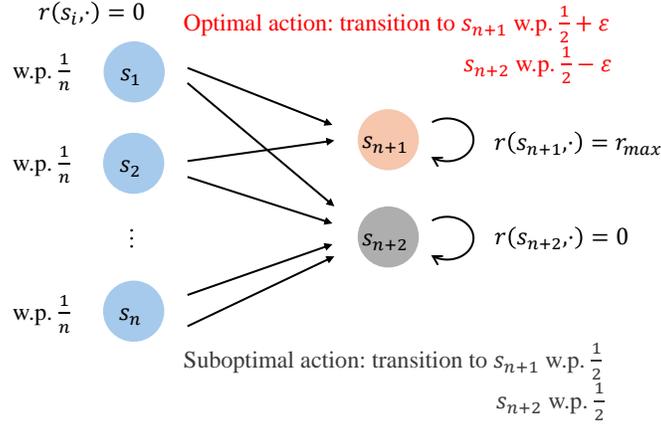


Figure 6: Instance for the lower bound under sum segment feedback and unknown transition.

$$\begin{aligned}
&= O\left(\left(\sqrt{\frac{H|\mathcal{S}||\mathcal{A}|}{m}} \log\left(\left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m}\right) \frac{1}{\delta}\right) + r_{\max} \sqrt{\lambda|\mathcal{S}||\mathcal{A}|}\right) \cdot \right. \\
&\quad \left. \left(H \sqrt{\frac{K|\mathcal{S}||\mathcal{A}|}{\lambda}} \log\left(\left(1 + \frac{KH^2}{\lambda|\mathcal{S}||\mathcal{A}|m}\right) \frac{1}{\delta}\right) + |\mathcal{S}|^{\frac{3}{2}} |\mathcal{A}|^{\frac{3}{2}} H \sqrt{\frac{KL}{\lambda}} \log(KH) \right. \right. \\
&\quad \left. \left. + \frac{|\mathcal{S}|^2 |\mathcal{A}|^2 H^2 L}{\sqrt{\lambda}} \log(KH)\right)\right) \\
&\stackrel{(b)}{=} O\left(\left(1 + r_{\max}\right) |\mathcal{S}|^2 |\mathcal{A}|^2 H \sqrt{K} \left(\log\left(\left(1 + \frac{KH}{|\mathcal{S}||\mathcal{A}|}\right) \frac{1}{\delta}\right)\right) \right. \\
&\quad \left. + \sqrt{L \log(KH)} \sqrt{\log\left(\left(1 + \frac{KH}{|\mathcal{S}||\mathcal{A}|}\right) \frac{1}{\delta}\right)} \right. \\
&\quad \left. + \left(1 + r_{\max}\right) |\mathcal{S}|^{\frac{5}{2}} |\mathcal{A}|^{\frac{5}{2}} H^2 L \log(KH) \sqrt{\log\left(\left(1 + \frac{KH}{|\mathcal{S}||\mathcal{A}|}\right) \frac{1}{\delta}\right)}\right) \\
&= \tilde{O}\left(\left(1 + r_{\max}\right) |\mathcal{S}|^{\frac{5}{2}} |\mathcal{A}|^2 H \sqrt{K} + \left(1 + r_{\max}\right) |\mathcal{S}|^{\frac{7}{2}} |\mathcal{A}|^{\frac{5}{2}} H^2\right),
\end{aligned}$$

where (a) comes from Lemma 30, and (b) uses the fact that $\lambda := \frac{H}{m}$. \square

D.5 A LOWER BOUND FOR UNKNOWN TRANSITION AND ITS PROOF

Below we provide a lower bound for RL with sum segment feedback and unknown transition with the proof.

Theorem 7. Consider the problem of RL with sum segment feedback and unknown transition. There exists a distribution of instances where the regret of any algorithm must be

$$\Omega\left(r_{\max} H \sqrt{|\mathcal{S}||\mathcal{A}|K}\right).$$

Proof of Theorem 7. We construct a random instance \mathcal{I} as follows. As shown in Figure 6, there are n bandit states s_1, \dots, s_n (i.e., there are an optimal action and multiple suboptimal actions), a good absorbing state s_{n+1} and a bad absorbing state s_{n+2} . The agent starts from s_1, \dots, s_n with equal probability $\frac{1}{n}$. For any $i \in [n]$, in state s_i , one action a_j is uniformly chosen from \mathcal{A} as the optimal action. In state s_i , under the optimal action a_j , the agent transitions to s_{n+1} and s_{n+2} with probabilities $\frac{1}{2} + \epsilon$ and $\frac{1}{2} - \epsilon$, respectively, where $\epsilon \in (0, \frac{1}{4})$ is a parameter specified later; Under any suboptimal action $a \in \mathcal{A} \setminus \{s_j\}$, the agent transitions to s_{n+1} and s_{n+2} with equal probability $\frac{1}{2}$.

The rewards are deterministic for all state-action pairs. For any $a \in \mathcal{A}$, $r(s_{n+1}, a) = r_{\max}$. For any $i \in \{1, \dots, n, n+2\}$ and $a \in \mathcal{A}$, $r(s_i, a) = 0$.

In this proof, we will also use an alternative uniform instance $\mathcal{I}_{\text{unif}}$. The only difference between $\mathcal{I}_{\text{unif}}$ and \mathcal{I} is that for any $i \in [n]$, in state s_i , under all actions $a \in \mathcal{A}$, the agent transitions to s_{n+1} and s_{n+2} with equal probability $\frac{1}{2}$.

Fix an algorithm \mathbb{A} . Let $\mathbb{E}_{\text{unif}}[\cdot]$ denote the expectation with respect to $\mathcal{I}_{\text{unif}}$. Let $\mathbb{E}_*[\cdot]$ denote the expectation with respect to \mathcal{I} . For any $i \in [n]$ and $j \in [|\mathcal{A}|]$, let $\mathbb{E}_{i,j}[\cdot]$ denote the expectation with respect to the case where a_j is the optimal action in state s_i , and $N_{i,j}$ denote the number of episodes where algorithm \mathbb{A} chooses a_j in state s_i , i.e., $N_{i,j} = \sum_{k=1}^K \mathbb{1}\{\pi_1^k(s_i) = a_j\}$.

The KL divergence of transition distribution on (s_i, a_j) ($i \in [n]$) between $\mathcal{I}_{\text{unif}}$ and \mathcal{I} is

$$\begin{aligned} \text{KL} \left(\text{Ber} \left(\frac{1}{2} \right) \parallel \text{Ber} \left(\frac{1}{2} + \varepsilon \right) \right) &= \frac{1}{2} \ln \left(\frac{\frac{1}{2}}{\frac{1}{2} - \varepsilon} \right) + \frac{1}{2} \ln \left(\frac{\frac{1}{2}}{\frac{1}{2} + \varepsilon} \right) \\ &= \frac{1}{2} \ln \left(\frac{\frac{1}{4}}{\frac{1}{4} - \varepsilon^2} \right) \\ &= -\frac{1}{2} \ln (1 - 4\varepsilon^2) \\ &\stackrel{(a)}{\leq} 4\varepsilon^2, \end{aligned}$$

where (a) uses the fact that $-\ln(1-x) \leq 2x$ when $x \in (0, \frac{1}{4})$.

In addition, the agent has probability only $\frac{1}{n}$ to arrive at (observe) state s_i .

Thus, using Lemma A.1 in (Auer et al., 2002), we have that for any $i \in [n]$, in state s_i ,

$$\begin{aligned} \mathbb{E}_{i,j}[N_{i,j}] &\leq \mathbb{E}_{\text{unif}}[N_{i,j}] + \frac{K}{2} \sqrt{\frac{1}{n} \cdot \mathbb{E}_{\text{unif}}[N_{i,j}] \cdot \text{KL} \left(\text{Ber} \left(\frac{1}{2} \right) \parallel \text{Ber} \left(\frac{1}{2} + \varepsilon \right) \right)} \\ &\leq \mathbb{E}_{\text{unif}}[N_{i,j}] + \frac{K}{2} \sqrt{\frac{1}{n} \cdot \mathbb{E}_{\text{unif}}[N_{i,j}] \cdot 4\varepsilon^2} \\ &= \mathbb{E}_{\text{unif}}[N_{i,j}] + K\varepsilon \sqrt{\frac{1}{n} \cdot \mathbb{E}_{\text{unif}}[N_{i,j}]}. \end{aligned}$$

Summing over $j \in [|\mathcal{A}|]$, using the Cauchy-Schwarz inequality and the fact that $\sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{\text{unif}}[N_{i,j}] = K$, we have

$$\sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{i,j}[N_{i,j}] \leq K + K\varepsilon \sqrt{\frac{|\mathcal{A}|}{n}} \cdot K.$$

Then, we have

$$\begin{aligned} \mathcal{R}(K) &= \sum_{k=1}^K \mathbb{E}_* \left[V^* - V^{\pi^k} \right] \\ &= \left(\frac{1}{2} + \varepsilon \right) (H-1)r_{\max}K \\ &\quad - \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} (H-1)r_{\max}K + \varepsilon(H-1)r_{\max} \cdot \frac{1}{|\mathcal{A}|} \sum_{j=1}^{|\mathcal{A}|} \mathbb{E}_{i,j}[N_{i,j}] \right) \\ &\geq \varepsilon(H-1)r_{\max} \left(K - \frac{K}{|\mathcal{A}|} - K\varepsilon \sqrt{\frac{K}{|\mathcal{A}|n}} \right). \end{aligned}$$

Recall that $n = |\mathcal{S}| - 2$. Let $|\mathcal{S}| \geq 3$, $|\mathcal{A}| \geq 2$, $H \geq 2$, $K > |\mathcal{A}|n$ and $\varepsilon = \frac{1}{4}\sqrt{\frac{|\mathcal{A}|n}{K}}$. Then, we have

$$\mathcal{R}(K) = \Omega\left(r_{\max}H\sqrt{|\mathcal{S}||\mathcal{A}|K}\right).$$

□

E TECHNICAL TOOLS

In this section, we introduce several technical tools.

Lemma 34 (Self-concordance, Lemma 9 in (Faury et al., 2020)). *For any $x_1, x_2 \in \mathbb{R}$, we have*

$$\mu'(x_1) \frac{1 - \exp(-|x_1 - x_2|)}{|x_1 - x_2|} \leq \int_{z=0}^1 \mu'((1-z)x_1 + zx_2) dz \leq \mu'(x_1) \frac{\exp(|x_1 - x_2|) - 1}{|x_1 - x_2|}.$$

Furthermore, we have

$$\int_{z=0}^1 \mu'((1-z)x_1 + zx_2) dz \geq \frac{\mu'(x_1)}{1 + |x_1 - x_2|}.$$

Lemma 35 (Value Difference Lemma, Lemma E.15 in (Dann et al., 2017)). *For any two MDPs M' and M'' with rewards r' and r'' and transition distributions p' and p'' , we have that for any $h \in [H]$ and $s \in \mathcal{S}$,*

$$V'_h(s) - V''_h(s) = \mathbb{E}_{p''} \left[\sum_{t=h}^H \left(r'(s_t, a_t) - r''(s_t, a_t) + (p'(\cdot | s_t, a_t) - p''(\cdot | s_t, a_t))^\top V'_{h+1}(\cdot) \right) \middle| s_t = s \right].$$

Lemma 36 (Law of Total Variance, Lemma 15 in (Zanette & Brunskill, 2019)). *For an MDP p and a fixed policy π , we have*

$$\mathbb{E}_{\pi, p} \left[\left(\sum_{h=1}^H r(s_h, \pi_h(s)) - V_1^\pi(s_1) \right) \middle| s_1 \right] = \mathbb{E}_{\pi, p} \left[\sum_{h=1}^H \text{Var}_{s_{h+1} \sim p(\cdot | s_h, \pi_h(s_h))} (V_{h+1}^\pi(s_{h+1})) \middle| s_1 \right].$$

The idea of Lemma 36 was also used in earlier works, e.g., (Munos & Moore, 1999; Lattimore & Hutter, 2012; Gheshlaghi Azar et al., 2013).

Lemma 37 (Lemma 10 in (Ménard et al., 2021)). *For distributions $p, q \in \Delta_{\mathcal{S}}$ and function $f : \mathcal{S} \rightarrow [0, b]$, if $\text{KL}(p, q) \leq \alpha$, then*

$$|(p(\cdot) - q(\cdot))^\top f(\cdot)| \leq \sqrt{2\text{Var}_q(f)\alpha} + \frac{2}{3}b\alpha.$$

Lemma 38 (Lemma 11 in (Ménard et al., 2021)). *For distributions $p, q \in \Delta_{\mathcal{S}}$ and function $f : \mathcal{S} \rightarrow [0, b]$, if $\text{KL}(p, q) \leq \alpha$, then*

$$\text{Var}_q(f) \leq 2\text{Var}_p(f) + 4b^2\alpha.$$

Lemma 39 (Lemma 12 in (Ménard et al., 2021)). *For distribution $p \in \Delta_{\mathcal{S}}$ and functions $f, g : \mathcal{S} \rightarrow [0, b]$, we have*

$$\text{Var}_p(f) \leq 2\text{Var}_p(g) + 2bp(\cdot)^\top |f(\cdot) - g(\cdot)|.$$