

# Variance-aware decision making with linear function approximation under heavy-tailed rewards

Anonymous authors

Paper under double-blind review

## Abstract

This paper studies how to achieve variance-aware regrets for online decision-making in the presence of heavy-tailed rewards with only finite variances. For linear stochastic bandits, we address the issue of heavy-tailed rewards by modifying the adaptive Huber regression and proposing AdaOFUL. AdaOFUL achieves a state-of-the-art regret bound of  $\tilde{O}(d(\sum_{t=1}^T \nu_t^2)^{1/2} + d)$  as if the rewards were uniformly bounded, where  $\nu_t^2$  is the conditional variance of the reward at round  $t$  and  $d$  is the feature dimension. Building upon AdaOFUL, we propose VARA for linear MDPs, which achieves a variance-aware regret bound of  $\tilde{O}(d\sqrt{HG^*K})$ . Here,  $H$  is the length of episodes,  $K$  is the number of episodes, and  $G^*$  is a smaller instance-dependent quantity that can be bounded by other instance-dependent quantities when additional structural conditions on the MDP are satisfied. Overall, our modified adaptive Huber regression algorithm may serve as a useful building block in the design of algorithms for online problems with heavy-tailed rewards.

## 1 Introduction

In many real-world scenarios, data exhibits heavy-tailed behaviors, which deviate significantly from classical assumptions in statistical analyses. Examples include stock returns in financial markets (Cont, 2001; Hull, 2012), microarray data analysis (Posekany et al., 2011), and advertiser values in online advertising (Arnosti et al., 2016). Such heavy-tailed distributions pose challenges to conventional algorithmic designs that often hinge upon uniformly bounded or sub-Gaussian reward assumptions.

A primary framework for studying decision-making under uncertainty is the multi-arm bandits problem, and its extension, the linear bandits. Regret analysis in this domain seeks to understand the suboptimality of algorithmic choices. However, traditional analyses often limit their applicability by assuming uniformly bounded or sub-Gaussian rewards. Some recent approaches address heavy-tailed behaviors by truncating rewards to achieve sub-linear worst-case regret bounds (Bubeck et al., 2013; Medina & Yang, 2016; Shao et al., 2018; Xue et al., 2021). Nevertheless, these truncation-based methods encounter estimation errors dependent on absolute moments of observations, not their central moments, suggesting suboptimality, especially in noiseless situations.

While the linear bandit framework offers a general enough setting for understanding decision-making with heavy-tailed data, reinforcement learning (RL) elevates this understanding to long-horizon decision-making processes. In RL, agents not only make decisions but also navigate through potentially infinite state and action spaces over a given horizon (Sutton & Barto, 2018). It demonstrates remarkable empirical successes in various applications, including robotics (Lillicrap et al., 2015), dialogue systems (Li et al., 2016), and Go play (Silver et al., 2016). Recent theoretical advances have expanded RL’s applicability, especially with linear function approximation in the context of linear Markov decision processes (MDPs) (Yang & Wang, 2020; Jin et al., 2020b;a; Wagenmaker et al., 2022b; Zanette et al., 2020; Ayoub et al., 2020; Zhou et al., 2021). More recently, the shift from worst-case regret analysis to variance-aware regret analysis in RL offers more nuanced insights into agent performance (Pananjady & Wainwright, 2020; Khamaru et al., 2021; Li et al., 2023; Yin & Wang, 2021; Min et al., 2021). More specifically, variance-aware regrets depend on the variances of rewards

and value functions and provide finer guarantees than worst-case bounds by characterizing problem-dependent performances across different problem instances. Yet, the challenge posed by heavy-tailed rewards remains.

This paper explores the intersection of decision-making under heavy-tailed rewards, ranging from linear bandits to RL applications. **We use the term “heavy-tailed rewards” throughout the paper to refer to the rewards that have only finite variances for simplicity.** Our aim is to achieve the variance-awareness and address the heavy-tailed issue simultaneously. A desirable algorithm should have the following two properties. First, it should possess the flexibility to function as a module, enhancing algorithms originally designed for bounded rewards to accommodate heavy-tailed rewards. Second, it should attain tight variance-aware regret bounds based on central moments, rather than absolute moments.

## 1.1 Our contributions

We provide a particular algorithm satisfying the mentioned characteristics. Our solution is motivated by adaptive Huber regression (Sun et al., 2020; Sun, 2021), which was originally proposed for analyzing offline independently and identically distributed (i.i.d.) data. It uses the (pseudo-) Huber loss to estimate the unknown coefficient with a universal robustification parameter. We adapt this method for online bandits and carefully choose different robustification parameters to handle non-i.i.d. data. The resulting algorithm, called AdaOFUL (short for Adaptive Huber regression based OFUL), achieves the state-of-the-art regret bound  $\tilde{O}\left(d\sqrt{\sum_{t \in [T]} \nu_t^2} + d\right)$  for linear bandits with heavy-tailed rewards, where  $\nu_t^2$  is the observed conditional variance of the random reward at step  $t$  and  $d$  is the feature dimension. **Here  $\tilde{O}(\cdot)$  hides constant factors and logarithmic dependence on  $T$ .** Such a variance-aware regret bound has only been obtained in the literature of linear bandits with sub-Gaussian or uniformly bounded rewards (Kirschner & Krause, 2018; Zhou & Gu, 2022). In contrast, truncation-based methods are suboptimal due to their estimation errors that depend on absolute moments instead of central moments. For example, the truncation-based algorithms from (Shao et al., 2018; Xue et al., 2021) yield regret in the form of  $\tilde{O}(v d \sqrt{T})$  where  $v^2$  is the bound for the second moment of random rewards. Our regret bound depends on the central moment instead and is thus tighter.

Using AdaOFUL as a building block, we then propose the Variance-Aware Regret via the Adaptive Huber regression (VARA) algorithm for linear MDPs with heavy-tailed rewards. In essence, VARA integrates AdaOFUL with the state-of-the-art worst-case algorithm LSVI-UCB++ from (He et al., 2022), enhancing regret performance through more careful analysis and resulting in a regret bound of  $\tilde{O}(d\sqrt{H\mathcal{G}^*K})$ . Here  $H$  is the horizon length and  $\mathcal{G}^*$  is a variance-aware quantity bounded by the sum of weighted per-step conditional variances. Our regret bound is superior to the current state-of-the-art bounds in three ways. First, it depends on a tighter instance-dependent quantity  $\mathcal{G}^*$  without knowing the value of  $\mathcal{G}^*$  in advance and has optimal dependence on  $d$  and  $H$ . Second, assuming additional structural conditions on the underlying MDP, we can obtain further instance-dependent bounds of  $\mathcal{G}^*$ , including range-dependent, first-order, and concentrability-dependent bounds. Third, our regret bound  $\tilde{O}(d\sqrt{H\mathcal{G}^*K})$  is valid even when rewards have only finite variances, which achieves a level of generality that is unmatched by previous works.

Our findings indicate that heavy-tailed rewards do not pose a limitation for developing online decision-making with linear function approximations. Our proposed modified adaptive Huber regression algorithm can be used as a general approach to adapt existing online algorithms designed for light-tailed rewards to handle heavy-tailed ones while maintaining tight dependence on variance for regret bounds.

**Overview** The rest of the paper proceeds as follows. We state our main results for heavy-tailed linear bandits in Section 2 and for linear MDPs in Section 3. We review related work in Section 4 and conclude in Section 5. Most proofs are collected in the appendix.

**Notation** We use  $\|\cdot\|$  to denote the  $\ell_2$ -norm in  $\mathbb{R}^d$ , and  $\text{Ball}_d(B)$  the  $\ell_2$ -norm ball in  $\mathbb{R}^d$  with radius  $B > 0$ . For a positive definite matrix  $\mathbf{H} \in \mathbb{R}^{d \times d}$ ,  $\|\mathbf{x}\|_{\mathbf{H}} = \sqrt{\mathbf{x}^\top \mathbf{H} \mathbf{x}}$  for a vector  $\mathbf{x} \in \mathbb{R}^d$ . For two semidefinite positive matrices  $\mathbf{H}_1, \mathbf{H}_2$ , we denote  $\mathbf{H}_1 \geq \mathbf{H}_2$  if  $\mathbf{H}_2 - \mathbf{H}_1$  is semidefinite positive. For an integer  $K \in \mathbb{N}^+$ , let  $[K] := \{1, 2, \dots, K\}$ . For a set  $\mathcal{A}$ ,  $|\mathcal{A}|$  denotes its cardinality. For real numbers  $a \leq b$  and  $x \in \mathbb{R}$ , we use  $x_{[a,b]} := \max\{a, \min\{x, b\}\}$  to denote the projection of  $x$  onto the closed interval  $[a, b]$ .

## 2 Variance-aware Regret for Heavy-tailed Linear Bandits

In this section, we first introduce the heavy-tailed linear bandit and then present the AdaOFUL algorithm, showing it achieves state-of-the-art variance-aware regret even when faced with heavy-tailed rewards.

### 2.1 Heavy-tailed Stochastic Linear Bandit

**Definition 2.1** (Heavy-tailed stochastic linear bandit). Let  $\{\mathcal{D}_t\}_{t \geq 1}$  denote a fixed sequence of decision sets and  $\{\mathcal{F}_t\}_{t \geq 1}$  a filtration. At round  $t$ , the agent chooses  $\phi_t \in \mathcal{D}_t$  and then observes the reward  $y_t$  and its conditional variance  $\nu_t^2$ . We assume  $y_t = \langle \phi_t, \theta^* \rangle + \varepsilon_t$  where  $\theta^* \in \mathbb{R}^d$  is a vector unknown to the agent and  $\varepsilon_t \in \mathbb{R}$  is a martingale difference random noise such that  $\mathbb{E}[\varepsilon_t | \mathcal{F}_{t-1}] = 0$  and  $\mathbb{E}[\varepsilon_t^2 | \mathcal{F}_{t-1}] = \nu_t^2$ . Both  $\nu_t$  and  $\phi_t$  are  $\mathcal{F}_{t-1}$ -measurable and  $\|\phi_t\| \leq L$ . We assume  $\|\theta^*\| \leq B$  with  $B$  known *a priori*. The agent aims to minimize the regret, formally defined as

$$\text{Reg}(T) := \sum_{t=1}^T \left[ \sup_{\phi \in \mathcal{D}_t} \langle \phi, \theta^* \rangle - \langle \phi_t, \theta^* \rangle \right]. \quad (2.1)$$

In heavy-tailed stochastic linear bandits, the mean-zero random noises  $\varepsilon_t$  have only bounded variances. We emphasize that in linear bandits, data are collected in an adaptive manner, and therefore, the distribution of  $\varepsilon_t$  depends on  $\phi_t$ . Moreover, the choice of  $\phi_t$  depends on all past observations  $(\phi_s, y_s, \nu_s)_{s < t}$ .

### 2.2 Algorithm Description

This section presents the AdaOFUL algorithm for heavy-tailed linear bandits. The AdaOFUL algorithm is given in Algorithm 1. AdaOFUL follows the principle of Optimism in the Face of Uncertainty (OFU) (Abbasi-Yadkori et al., 2011) to solve the heavy-tailed heterogeneous linear bandit problem. At each round  $t$ , it maintains a confidence set defined in equation 2.2 such that  $\theta^* \in \mathcal{C}_t$  uniformly for all  $t \geq 1$  with high probability when the exploration radius  $\beta_{t-1}$  is properly chosen. Unlike the standard OFUL algorithm (Abbasi-Yadkori et al., 2011) which directly selects the most optimistic estimator  $\tilde{\theta}_t$  to make an arm selection  $\phi_t$ , AdaOFUL uses adaptive Huber regression to compute a new estimator  $\theta_t$  that takes into account the heavy-tailed rewards. The agent then selects the arm  $\phi_t$  that maximizes the inner product  $\langle \phi, \theta \rangle$  over  $\theta \in \mathcal{C}_{t-1}$ . After playing the selected arm, the agent observes the reward  $y_t$  and its conditional variance  $\nu_t$ . The last step of round  $t$  updates the exploration radius  $\beta_t$  and the shape matrix  $\mathbf{H}_t$  for the confidence set construction.

**Adaptive pseudo-Huber regression** The pseudo-Huber loss (Hastie et al., 2009; Sun, 2021) is defined as

$$\ell_\tau(x) = \tau(\sqrt{\tau^2 + x^2} - \tau), \quad (2.5)$$

which is a smooth approximation to the well-known Huber loss (Huber, 1964). Similar to the Huber loss, the pseudo-Huber loss resembles a quadratic function for small values of  $|x|$  and is approximately linear when  $x$  is large in magnitude, making the loss strongly convex when close to the origin and less sensitive to changes in the tails. The parameter  $\tau$  controls the balance between the quadratic and linear regions and is referred to as the robustification parameter by Sun et al. (2020) in the case of the Huber loss. Since the value of the robustification parameter needs to be adaptive to the data for an optimal tradeoff between robustness and unbiasedness, we shall also refer to the pseudo-Huber regression with a data-adaptive  $\tau$  as adaptive pseudo-Huber regression or simply adaptive Huber regression, in line with (Sun et al., 2020).

To compute the pseudo-Huber estimator  $\theta_t$  for  $\theta^*$ , given the history  $\{(\phi_s, y_s, \nu_s)\}_{s \in [t]}$  up to time  $t$ , we solve the the convex optimization problem in equation 2.4 (Sun, 2021). Recall that  $\sigma_t$ 's are surrogate conditional variances, and  $\tau_t$ 's are the robustification parameters, given by equation 2.3, in which  $\sigma_{\min}$  is a small positive constant to avoid singularity,  $\tau_0$  is a hyper-parameter,  $w_t$ 's are importance measures,  $c_0$  and  $c_1$  are specified in Algorithm 1, and  $L$  and  $B$  are constants defined in Definition 2.1.

**Algorithm 1** Adaptive Huber regression based OFUL (AdaOFUL)

**Initialization:**  $\mathbf{H}_0 = \lambda \mathbf{I}$ ,  $\boldsymbol{\theta}_0 = \mathbf{0}$ ,  $\beta_0 = \sqrt{\lambda}B$ ,  $c_0 = \frac{1}{6\sqrt{3 \log \frac{2T^2}{\delta}}}$ ,  $c_1 = \frac{1}{42 \cdot \log \frac{2T^2}{\delta}}$ ,  $\sigma_{\min} = \frac{1}{\sqrt{T}}$ .

1 **for**  $t = 1$  **to**  $T$  **do**

2     Construct the confidence set  $\mathcal{C}_{t-1}$  as in

$$\mathcal{C}_{t-1} := \{\boldsymbol{\theta} \in \text{Ball}_d(B) : \|\boldsymbol{\theta} - \boldsymbol{\theta}_{t-1}\|_{\mathbf{H}_{t-1}} \leq \beta_{t-1}\}. \quad (2.2)$$

3     Solve  $(\boldsymbol{\phi}_t, \cdot) = \operatorname{argmax}_{\boldsymbol{\phi} \in \mathcal{D}_t, \boldsymbol{\theta} \in \mathcal{C}_{t-1}} \langle \boldsymbol{\phi}, \boldsymbol{\theta} \rangle$ .

4     Play  $\boldsymbol{\phi}_t$  and observe  $(y_t, \nu_t)$ .

5     Set  $\sigma_t, w_t$  and  $\tau_t$  according to the following equation and record  $\{\sigma_s, w_s, \tau_s : 1 \leq s \leq t\}$ .

$$\sigma_t = \max \left\{ \nu_t, \sigma_{\min}, \frac{\|\boldsymbol{\phi}_t\|_{\mathbf{H}_{t-1}^{-1}}}{c_0}, \frac{\sqrt{LB} \|\boldsymbol{\phi}_t\|_{\mathbf{H}_{t-1}^{-1}}^{\frac{1}{2}}}{c_1^{\frac{1}{4}} d^{\frac{1}{4}}} \right\}, w_t = \left\| \frac{\boldsymbol{\phi}_t}{\sigma_t} \right\|_{\mathbf{H}_{t-1}^{-1}}, \tau_t = \tau_0 \frac{\sqrt{1 + w_t^2}}{w_t}. \quad (2.3)$$

6     Compute  $\boldsymbol{\theta}_t$  by minimizing the following convex problem

$$\boldsymbol{\theta}_t := \operatorname{argmin}_{\boldsymbol{\theta} \in \text{Ball}_d(B)} L_t(\boldsymbol{\theta}) \text{ with } L_t(\boldsymbol{\theta}) := \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 + \sum_{s=1}^t \ell_{\tau_s} \left( \frac{y_s - \langle \boldsymbol{\phi}_s, \boldsymbol{\theta} \rangle}{\sigma_s} \right). \quad (2.4)$$

7     Define the confidence set radius  $\beta_t$  as in equation 2.6 and set  $\mathbf{H}_t = \mathbf{H}_{t-1} + \frac{\boldsymbol{\phi}_t \boldsymbol{\phi}_t^\top}{\sigma_t}$ .

8 **end**

**Robustification parameter** As shown in equation 2.3, the robustification parameter  $\tau_t$  is set differently for each data point  $(\boldsymbol{\phi}_t, y_t, \nu_t)$  in the pseudo-Huber regression. This is a significant departure from the case of i.i.d. data, where all robustification parameters are typically set to the same value  $\tau$ , as i.i.d. data are naturally weighted equally (Sun et al., 2020). In linear bandits, the data is generated adaptively, where the choice of  $\boldsymbol{\phi}_t$  can depend on all past observations. Since observations collected in later rounds are less important as they are based on previous observations and contribute less to the estimation accuracy, we assign greater weight to earlier observations. To measure the importance of the  $t$ -th observation, we use  $w_t = \|\boldsymbol{\phi}_t\|_{\mathbf{H}_{t-1}^{-1}} / \sigma_t$  as the importance measure for the  $t$ -th observation and set  $\tau_t = \tau_0 \sqrt{1 + w_t^2} / w_t$  as the corresponding robustification parameter.

When taking  $\tau_0 = \infty$ , the optimization problem in equation 2.4 reduces to weighted regularized least-squares, which has been proven to achieve worst-case optimality for linear bandits with uniformly bounded or sub-Gaussian rewards (Kirschner & Krause, 2018; Zhou & Gu, 2022). However, an appropriate value of  $\tau_0$  is necessary to balance robustness against heavy-tailed rewards and asymptotic unbiasedness. In Corollary 2.1, we will demonstrate that setting  $\tau_0 = \tilde{O}(\sqrt{d})$  is sufficient to achieve the state-of-the-art regret bound.

**Variance estimates** We choose  $\sigma_t \geq \sqrt{LB} \|\boldsymbol{\phi}_t\|_{\mathbf{H}_{t-1}^{-1}}^{\frac{1}{2}} / (c_1^{\frac{1}{4}} d^{\frac{1}{4}})$ , which implies  $c_1 d \geq L^2 B^2 w_t^2 / (\sigma_t^2)$ . This condition is used to lower bound the Hessian matrix  $\nabla^2 L_T(\boldsymbol{\theta})$ . For any  $\boldsymbol{\theta} \in \text{Ball}_d(B)$ , we expect  $\nabla^2 L_T(\boldsymbol{\theta}) \approx \mathbf{H}_T$  up to universal constant factors to proceed with theoretical analysis. A direct computation yields  $\nabla^2 L_T(\boldsymbol{\theta}) \leq \mathbf{H}_T$ , while for the other direction we show  $\nabla^2 L_T(\boldsymbol{\theta}) \geq \left( c - \sup_{t \in [T]} \left| \frac{\langle \boldsymbol{\phi}_t, \boldsymbol{\theta}^* - \boldsymbol{\theta} \rangle}{\tau_t \sigma_t} \right|^2 \right) \mathbf{H}_T$  for some universal constant  $c > 0$  with high probability. With the last condition on  $\sigma_t$ , for any feasible solution  $\boldsymbol{\theta} \in \text{Ball}_d(B)$ , the following quantity

$$\left| \frac{\langle \boldsymbol{\phi}_t, \boldsymbol{\theta}^* - \boldsymbol{\theta} \rangle}{\tau_t \sigma_t} \right|^2 \leq \frac{\|\boldsymbol{\phi}_t\|^2 \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|^2}{\tau_t^2 \sigma_t^2} \leq \frac{4w_t^2 L^2 B^2}{\tau_0^2 \sigma_t^2} \leq \frac{4c_1 d}{\tau_0^2}$$

can be sufficiently small provided that  $\tau_0^2 \geq c \cdot d$  for a sufficiently large constant  $c > 0$ .

### 2.3 Regret Analysis

We first validate that the optimism holds with high probability in Theorem 2.1 and then establish a high probability bound for the regret in Theorem 2.2.

**Theorem 2.1.** Let  $\kappa = d \cdot \log(1 + TL^2/(d\lambda\sigma_{\min}^2))$ . For the heavy-tailed linear bandit in Definition 2.1, if  $\tau_0\sqrt{\log(2T^2/\delta)} \geq \max\{\sqrt{2\kappa}, 2\sqrt{dLB}\}$ , then with probability at least  $1 - 4\delta$ , it holds that, for all  $0 \leq t \leq T$ ,

$$\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_{\mathbf{H}_t} \leq \beta_t,$$

where

$$\beta_t = 32 \left( \frac{\kappa}{\tau_0} + \sqrt{\kappa \log \frac{2t^2}{\delta}} + \tau_0 \log \frac{2t^2}{\delta} \right) + 5\sqrt{\lambda}B. \quad (2.6)$$

Theorem 2.1 establishes that  $\boldsymbol{\theta}^*$  is contained in the set  $\mathcal{C}_t := \{\boldsymbol{\theta} \in \text{Ball}_d(B) : \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|_{\mathbf{H}_t} \leq \beta_t\}$  for all  $t \geq 0$  with high probability. It is proved by using Bernstein-type concentration inequality for self-normalized vector-valued martingales with additional care paid to deal with heavy-tailed rewards. See the next subsection for a proof sketch.

**Theorem 2.2.** Let  $\sigma_{\min} = 1/\sqrt{T}$ . Then with probability at least  $1 - 4\delta$ , we have

$$\text{Reg}(T) \leq 2\beta_T \cdot \left[ \sqrt{2\kappa} \cdot \sqrt{\sum_{t=1}^T \nu_t^2 + 1} + \frac{2L\kappa}{c_0^2\sqrt{\lambda}} + \frac{2LB\kappa}{\sqrt{c_1d}} \right]$$

where  $\beta_T$  is defined in equation 2.6, and  $c_0, c_1 = \tilde{\mathcal{O}}(1)$  are positive constants given in Algorithm 1.

Theorem 2.2 provides a regret bound in a general form that depends on  $\beta_T$ . As shown in equation 2.6,  $\beta_t$  is a hyperbolic function of the robustification parameter  $\tau_0$ . Increasing  $\tau_0$  decreases the bias term  $\mathcal{O}(\kappa/\tau_0)$  while increasing the range term  $\mathcal{O}(\tau_0 \log(2t^2/\delta))$ . Therefore, choosing  $\tau_0$  carefully is essential to achieve the optimal trade-off between unbiasedness and robustness. Setting  $\tau_0 = \tilde{\mathcal{O}}(\sqrt{d})$  minimizes the right-hand side of equation 2.6. This, combined with Theorem 2.2, yields the simplified regret bound equation 2.7 in the following corollary.

**Corollary 2.1.** Let  $\tau_0 = \max\{\sqrt{2\kappa}, 2\sqrt{d}\}/\sqrt{\log(2T^2/\delta)}$  and  $\lambda = d/B^2$ , then  $\beta_T \leq 64 \left( 2\sqrt{\kappa \log(2T^2/\delta)} + \sqrt{d \log(2T^2/\delta)} \right) + 5\sqrt{d}$ . Consequently, the regret bound in Theorem 2.2 becomes

$$\text{Reg}(T) = \tilde{\mathcal{O}} \left( d \sqrt{\sum_{t \in [T]} \nu_t^2} + d \cdot \max\{LB, 1\} \right), \quad (2.7)$$

where  $\tilde{\mathcal{O}}(\cdot)$  hides constant factors and logarithmic dependence on  $T$ .

Corollary 2.1 demonstrates that AdaOFUL achieves state-of-the-art regret bound under heavy-tailed rewards, comparable to the case where rewards are uniformly bounded or sub-Gaussian. The regret upper bound in the noiseless case reduces to  $\tilde{\mathcal{O}}(d)$ , and in the noisy case, it reduces to  $\tilde{\mathcal{O}}(d\sqrt{\sum_{t \in [T]} \nu_t^2})$ . In the worst case scenario where  $\nu_t = \Theta(1)$  for all  $t \geq 1$ , the regret bound reduces to  $\tilde{\mathcal{O}}(d\sqrt{T})$ , which matches the worst-case minimax lower bound (Dani et al., 2008). Hence, our variance-aware regret bound equation 2.7 is tighter than the pessimistic worst-case bound  $\tilde{\mathcal{O}}(d\sqrt{T})$  when  $\sum_{t=1}^T \nu_t^2 \ll T$ . To the best of our knowledge, such a variance-aware regret bound has only been obtained in the literature for sub-Gaussian rewards (Kirschner & Krause, 2018) or uniformly bounded rewards (Zhou & Gu, 2022). We are the first to provide a variance-aware regret bound for heavy-tailed stochastic linear bandits.

## 2.4 Proof Sketch of Theorem 2.1

**Step one: Hessian approximation** Let  $z_t(\boldsymbol{\theta}) := \frac{y_t - \langle \boldsymbol{\phi}_t, \boldsymbol{\theta} \rangle}{\sigma_t}$  and  $\kappa := d \log \left( 1 + \frac{TL^2}{d\lambda\sigma_{\min}^2} \right)$  for simplicity.

**Lemma 2.1.** Assume there exists a constant  $b > 0$  such that  $\mathbb{E}[z_t^2(\boldsymbol{\theta}^*) | \mathcal{F}_{t-1}] \leq b^2$  for all  $t \geq 1$ . If  $\tau_0 \sqrt{\log \frac{2T^2}{\delta}} \geq \max\{\sqrt{2\kappa}b, 2\sqrt{d}\}$ , then with probability at least  $1 - 2\delta$ , for all  $T \geq 0$  and any  $\|\boldsymbol{\theta}\| \leq B$ ,

$$\frac{1}{4}\mathbf{H}_T \leq \nabla^2 L_T(\boldsymbol{\theta}) \leq \mathbf{H}_T.$$

Lemma 2.1 shows that with high probability and up to constant factors,  $\nabla^2 L_T(\boldsymbol{\theta})$  approximates  $\mathbf{H}_T$  well uniformly for all  $T \geq 1$  and  $\|\boldsymbol{\theta}\| \leq B$ . By contrast, in standard ridge regressions,  $\nabla^2 L_T(\boldsymbol{\theta})$  equals to  $\mathbf{H}_T$  because the corresponding loss  $L_T$  is quadratic. The proof is deferred to Section C.2.

**Step two: High probability gradient bound** In the following, we provide a high-probability bound for  $\|\nabla L_T(\boldsymbol{\theta}^*)\|_{\mathbf{H}_T^{-1}}$  in Lemma 2.2.

**Lemma 2.2.** Assume there exists a constant  $b > 0$  such that  $\mathbb{E}[z_t^2(\boldsymbol{\theta}^*) | \mathcal{F}_{t-1}] \leq b^2$  for all  $t \geq 1$ . With probability at least  $1 - 2\delta$ , for all  $T \geq 1$ , it follows that

$$\|\nabla L_T(\boldsymbol{\theta}^*)\|_{\mathbf{H}_T^{-1}} \leq 8 \left[ \underbrace{\frac{\kappa b^2}{\tau_0}}_{\text{bias term}} + \underbrace{\sqrt{\kappa b^2 \log \frac{2T^2}{\delta}}}_{\text{variance term}} + \underbrace{\tau_0 \log \frac{2T^2}{\delta}}_{\text{range term}} \right] + \underbrace{\sqrt{\lambda}B}_{\text{ridge term}}.$$

We explain briefly about each term in Lemma 2.2. Following Zhou et al. (2021), a decomposition follows that  $\|\nabla L_t(\boldsymbol{\theta}^*)\|_{\mathbf{H}_t^{-1}}^2 = \sum_{s=1}^t (X_s + Y_s)$  for two sequences of random variables  $X_t, Y_t \in \mathcal{F}_t$ . To illustrate the proof idea, we explain how to bound  $\sum_{t=1}^T X_t$ , since  $\sum_{t=1}^T Y_t$  can be bounded similarly. For the adaptive Huber regression,  $\{X_t\}_{t \in [T]}$  is not a martingale difference sequence but  $\{X_t - \mathbb{E}[X_t | \mathcal{F}_{t-1}]\}_{t \in [T]}$  is. We apply a Bernstein inequality to upper bound  $\sum_{t=1}^T (X_t - \mathbb{E}[X_t | \mathcal{F}_{t-1}])$  which contributes to the variance and range terms. Thanks to the different robustification parameters  $\tau_t$ , we can control  $\sum_{t=1}^T \mathbb{E}[X_t | \mathcal{F}_{t-1}]$  deterministically within  $\mathcal{O}(\kappa^2/\tau_0^2)$ , resulting in the bias term. Finally, the last ridge term  $\sqrt{\lambda}B$  exists because we use ridge regularization to ensure that the Hessian is always invertible. The detailed proof is in Appendix C.3.

**Step three: Combination through stationary condition** Notice that the gradient is given by

$$\nabla L_T(\boldsymbol{\theta}) := \lambda \boldsymbol{\theta} - \sum_{t=1}^T \frac{\tau_t z_t(\boldsymbol{\theta})}{\sqrt{\tau_t^2 + z_t(\boldsymbol{\theta})^2}} \frac{\boldsymbol{\phi}_t}{\sigma_t},$$

and our estimator  $\boldsymbol{\theta}_T$  is the minimizer of a constrained problem in equation 2.4. By Proposition 1.3 in (Bubeck et al., 2015), the first-order stationary condition of the constrained convex optimization equation 2.4 implies that  $\langle \nabla L_T(\boldsymbol{\theta}_T), \boldsymbol{\theta}_T - \boldsymbol{\theta}^* \rangle \leq 0$  for all  $\boldsymbol{\theta} \in \text{Ball}_d(B)$ . More specifically, due to  $\|\boldsymbol{\theta}^*\| \leq B$ , we have

$$\langle \nabla L_T(\boldsymbol{\theta}_T), \boldsymbol{\theta}_T - \boldsymbol{\theta}^* \rangle \leq 0. \quad (2.8)$$

By the mean value theorem for vector-valued functions, we have

$$\nabla L_T(\boldsymbol{\theta}_T) - \nabla L_T(\boldsymbol{\theta}^*) = \int_0^1 \nabla^2 L_T((1-\eta)\boldsymbol{\theta}^* + \eta\boldsymbol{\theta}_T) d\eta \cdot (\boldsymbol{\theta}_T - \boldsymbol{\theta}^*).$$

Using Lemma 2.1 and the fact that  $\|(1-\eta)\boldsymbol{\theta}^* + \eta\boldsymbol{\theta}_T\| \leq B$  for all  $\eta \in [0, 1]$ , we have

$$\frac{1}{4}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_{\mathbf{H}_T}^2 \leq \langle \boldsymbol{\theta}_T - \boldsymbol{\theta}^*, \nabla L_T(\boldsymbol{\theta}_T) - \nabla L_T(\boldsymbol{\theta}^*) \rangle. \quad (2.9)$$

By equation 2.9 and equation 2.8, we have

$$\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_{\mathbf{H}_T} \leq 4\|\nabla L_T(\boldsymbol{\theta}^*)\|_{\mathbf{H}_T^{-1}}. \quad (2.10)$$



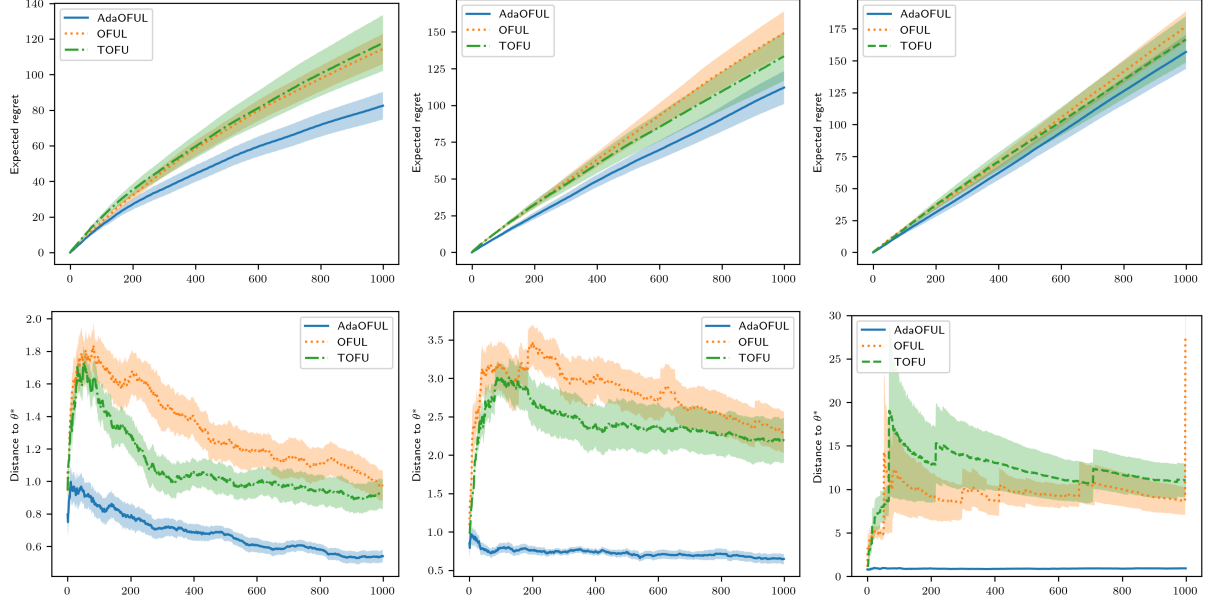


Figure 1: Regret and convergence results across three noise types: Case (a)  $\varepsilon_t \sim \mathcal{N}(0, 1)$  on the left, Case (b)  $\varepsilon_t \sim \mathbf{t}(\text{df})$  with  $\text{df} \sim \mathcal{U}(2, 3)$  in the middle, and Case (c)  $\varepsilon_t \sim \mathbf{t}(\text{df})$  with  $\text{df} \sim \mathcal{U}(1, 2)$  on the right.

Combining equation 2.10, Lemma 2.1 and Lemma 2.2, we know that if  $\tau_0 \sqrt{\log \frac{2T^2}{\delta}} = \max\{\sqrt{2kb}, 2\sqrt{d}\}$ , with probability at least  $1 - 2\delta$ , we have that  $\|\theta_t - \theta^*\|_{\mathbf{H}_t} < \beta_t$  for all  $1 \leq t \leq T$  where  $\beta_t$  is given in equation 2.6. It implies that  $\theta^*$  indeed locals in all constructed confidence regimes, i.e., for all  $1 \leq t \leq T$ ,  $\theta^* \in \mathcal{C}_t$ . Notice that by the choice of  $\beta_0$  and  $\mathbf{H}_0$ , we still have  $\theta^* \in \mathcal{C}_0$ . Finally,  $b = 1$  in our case completes the proof.

## 2.5 A Numerical Study

**Considered methods** In this subsection, we conduct a numerical comparison between AdaOFUL and two baseline algorithms: original OFUL (Abbasi-Yadkori et al., 2011) and TOFU (Shao et al., 2018). TOFU is a truncation-based variant of OFUL, designed to address the heavy-tail issue. Because these algorithms don't consider the variance information, for a fair comparison, we abstain from the variance weights and set each  $\sigma_t \equiv 1$ . Hyperparameters were chosen based on observations from the initial couple of steps so that  $\tau_0 = \sqrt{d}$  and  $c_0 = c_1 = 1$ .

**Experiment setup** We conduct an experiment with the following configuration. We set  $d = 10$  and  $|\mathcal{D}_t| \equiv 20$ . The optimal  $\theta^*$  is created by randomly sampling each coordinate from a uniform distribution  $\mathcal{U}(0, 1)$  and normalizing the resultant vector to unit length so that  $B = 1$ . To simulate varying action sets, we generate 20 distinct basic action sets,  $\{\mathcal{A}_i\}_{i \in [20]}$ , and assign  $\mathcal{D}_t = \mathcal{A}_i$  if  $t = i \bmod 20$ . For each  $\mathcal{A}_i$ , each arm vector  $\phi \in \mathcal{A}_i$  is formed in the same way as  $\theta^*$  so that  $L = 1$ . Rewards are generated by  $y_t = \langle \phi_t, \theta^* \rangle + \varepsilon_t$  with  $\varepsilon_t$  being an independent zero-mean noise. We investigate three noise types: Case (a) is Gaussian distribution  $\varepsilon_t \sim \mathcal{N}(0, 1)$ , while Case (b) and (c) correspond to Student  $t$ -distributions  $\varepsilon_t \sim \mathbf{t}(\text{df})$  with  $\text{df}$  the freedom varying. Note that if a random variable  $X$  follows a Student's  $t$ -distribution with freedom  $\text{df}$ , its mean is well defined for  $\text{df} > 1$  and its variance is well defined for  $\text{df} > 2$ , becoming infinite for  $2 \geq \text{df} > 1$ . Case (b) sets  $\text{df} \sim \mathcal{U}(2, 3)$ , while Case (c) uses  $\text{df} \sim \mathcal{U}(1, 2)$ . As we move from Case (a) to (c), the noise becomes increasingly heavy-tailed. To ensure a fair comparison, most parameters are shared, e.g.,  $\beta_t \equiv 1$  and  $\lambda = 1$ . The experiment runs for  $T = 1000$  steps and is replicated 10 times, with the outcomes averaged. Shadowing is used to depict the area within one standard deviation, calculated over these ten repetitions.

**Experiment results** Figure 1 shows the regret and convergence results across three noise cases. It is clear that for all considered algorithms, the regrets continue to grow and the  $L_2$  convergence errors  $\|\theta_t - \theta^*\|$  tend to diminish. The continuous growth in regrets is also observed in previous heavy-tailed experiments (Shao et al., 2018). A key message is that AdaOFUL consistently achieves the lowest regret, smallest convergence errors,

**and least variability.** In the context of light-tailed noise (in the left column), OFUL has a slightly smaller regret than TOFU. The result implies that the truncation technique might hurt the performance under light-tailed noises. As we transition to heavy-tailed noise with finite variance (in the middle column), TOFU outperforms OFUL instead in terms of regret and convergence, implying truncation works indeed. However, both remain suboptimal compared to AdaOFUL. In the case where the noise is predominantly heavy-tailed with only a bounded expectation (in the right column), TOFU and OFUL's convergence errors and regrets deteriorate, contrasting with the steadfast performance of AdaOFUL. These findings show AdaOFUL's empirical robustness even in the infinite variance noise regime.

### 3 An Extension to Linear MDPs

Ridge regression estimators are widely used in RL to provide confidence guarantees for bounded rewards. However, when dealing with heavy-tailed rewards, these estimators tend to degrade or even fail (as shown in Figure 1). In response, we advocate for the use of the adaptive Huber regression, or AdaOFUL, as a robust alternative to ridge regression. AdaOFUL can seamlessly enhance the original algorithm to accommodate heavy-tailed scenarios with minimal disruption to its core. In this section, we demonstrate this by integrating AdaOFUL as a foundational element to solve linear MDPs and provide variance-aware regrets. This approach can also be extended to other linear problems such as linear mixture MDPs (Zhou & Gu, 2022).

#### 3.1 Linear MDPs with Heavy-tailed Rewards

**Preliminaries about linear MDPs** An episodic finite horizon MDP is denoted by a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \{r_h\}_{h \in [H]}, \{\mathbb{P}_h\}_{h \in [H]})$  where  $\mathcal{S}$  is the state space with a possibly infinite number of states,  $\mathcal{A}$  the action space,  $H \in \mathbb{Z}^+$  the length of each episode,  $\mathbb{P}_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  the transition probability function, and  $r_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  the expected reward function. A linear MDP assumes that both the transition probability and the expected reward are linear in a known state-action feature map  $\phi(\cdot, \cdot) \in \mathbb{R}^d$  (Bradtke & Barto, 1996; Melo & Ribeiro, 2007; Yang & Wang, 2019; Jin et al., 2020b).

**Definition 3.1** (Linear MDP).  $\mathcal{M}$  is called a time-inhomogeneous linear MDP, if there exist some known feature map  $\phi(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \text{Ball}_d(1)$ , unknown signed measures  $\{\mu_h^*\}_{h \in [H]} \subseteq \mathbb{R}^{d \times |\mathcal{S}|}$ , and unknown coefficients  $\{\theta_h^*\}_{h \in [H]} \subseteq \text{Ball}_d(W)$  such that  $r_h(s, a) = \langle \phi(s, a), \theta_h^* \rangle$  and  $\mathbb{P}_h(\cdot | s, a) = \langle \phi(s, a), \mu_h^*(\cdot) \rangle$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,  $h \in [H]$ , where  $\|\mu_h^*(\mathcal{S})\| := \|\sum_{s \in \mathcal{S}} \mu_h^*(s)\| \leq \sqrt{d}$  for all  $h \in [H]$ .

For a time-inhomogeneous MDP, we denote its deterministic and time-dependent policy by  $\pi = \{\pi_h\}_{h \in [H]}$ . Let  $\{(s_h, a_h)\}_{h \in [H]}$  be state-action pairs such that  $a_h = \pi_h(s_h)$  and  $s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, a_h)$ . Define the occupancy measure for the policy  $\pi$  at the  $h$ -th round by  $d_h^\pi(s, a) = \mathbb{P}^\pi(s_h = s, a_h = a | s_1)$  where  $(a_1, s_2, a_2, \dots, s_h, a_h)$  is a trajectory starting from  $s_1$  and following the policy  $\pi$ . The state-action function  $Q_h^\pi(\cdot, \cdot)$  and value function  $V_h^\pi(\cdot)$  at the  $h$ -th round are defined as  $Q_h^\pi(\cdot, \cdot) = \mathbb{E}[\sum_{i=h}^H r_i(s_i, a_i) | (s_h, a_h) = (\cdot, \cdot)]$  and  $V_h^\pi(\cdot) = Q_h^\pi(\cdot, \pi_h(\cdot))$  respectively. **The optimal policy is denoted by  $\pi^*$  and its value function is denoted by  $V_1^*$ . One can show that  $V_1^*(s) = \sup_\pi V_1^\pi(s)$  for any  $s \in \mathcal{S}$ .** For any value function  $V$ , write  $[\mathbb{P}_h V](s, a) = \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a)} V(s')$  and  $[\mathbb{V}_h V](s, a) = [\mathbb{P}_h V^2](s, a) - [\mathbb{P}_h V]^2(s, a)$ . With a slight abuse of notation, let  $[\mathbb{P}_h R_h](s, a)$  and  $[\mathbb{V}_h R_h](s, a)$  denote the expectation and variance of the random reward  $R_h(s, a)$  at the  $h$ -th round given state-action pair  $(s, a)$ .

**Linear MDPs with heavy-tailed rewards** We consider linear MDPs with heavy-tailed random rewards that satisfy the following assumptions.

**Assumption 3.1** (Realizable reward). We assume that the following holds.

1. For all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and  $h \in [H]$ , the random reward  $R_h(s, a)$  is independent of  $s_{h+1}(s, a)$ , where  $s_{h+1}(s, a) \sim \mathbb{P}_h(\cdot | s, a)$  represents the next state transitioned from  $(s, a)$  at the  $h$ -th round.
2. There exists known feature maps  $\tilde{\phi}(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \text{Ball}_d(1)$  and unknown coefficients  $\{\psi_h^*\}_{h \in [H]} \subseteq \text{Ball}_d(W)$  so that  $[\mathbb{P}_h R_h^2](s, a) = \langle \tilde{\phi}(s, a), \psi_h^* \rangle$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and  $h \in [H]$ .

**Assumption 3.2** (Bounded variance). We assume that the following holds.



1. There exist known constants  $\sigma_R, \sigma_{R^2} > 0$  such that  $[\mathbb{V}_h R_h](s, a) \leq \sigma_R^2$  and  $[\mathbb{V}_h R_h^2](s, a) \leq \sigma_{R^2}^2$  for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A}$  and  $h \in [H]$ .
2. **There exist known upper bounds  $\mathcal{H}, \mathcal{V} > 0$  such that** for any policy  $\pi$ , we have  $0 \leq \mathbb{E} R_\pi \leq \mathcal{H}$  and  $\text{Var}(R_\pi) \leq \mathcal{V}^2$  where  $R_\pi = \sum_{h=1}^H R_h(s_h, a_h)$  denotes the sum of random rewards along the trajectory following  $\pi$ .

**Rationale behind the assumptions** Assumption 3.1 assumes that the random reward at each round is independent of future states and its second moment can be realized using a known feature map. Under this assumption, linear MDPs can recover tabular MDPs by setting the size of the state-action space as  $d = |\mathcal{S}||\mathcal{A}|$  and using the canonical basis  $\phi(s, a) = \tilde{\phi}(s, a) = \mathbf{e}_{(s,a)}$  in  $\mathbb{R}^d$ . Assumption 3.2 places upper bounds on the means and variances of every random reward and the cumulative rewards. These upper bounds are available under the classic uniformly bounded reward assumption that  $0 \leq \sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sup_{h \in [H]} R_h(s, a) \leq 1$  so that  $\sigma_R = \sigma_{R^2} = 1$  and  $\mathcal{H} = \mathcal{V} = H$ . We emphasize that almost all previous works use these "1" and "H" upper bounds implicitly in their algorithm design and regret analysis. In this way, they can't tell the effect of the expectation of cumulative rewards on the final regret from their variance. We are the first to distinguish them by separate  $\mathcal{H}$  and  $\mathcal{V}$ . Since only upper bounds for  $\mathcal{H}$  and  $\mathcal{V}$  are required, in practice one can guess them using the doubling trick.<sup>1</sup> **As we will observe, very large guessing values for  $\mathcal{H}$  and  $\mathcal{V}$  will not affect the order of the dominant (or variance-aware) term in our regret as long as  $T$  is sufficiently large.** As far as we know, Assumption 3.2 is the weakest moment condition on random rewards in the **variance-aware** RL literature.

**Learning protocol** Let  $\mathcal{F}_{h,k}$  denote the  $\sigma$ -field generated by all random variables up to, and including, the  $h$ -th round and  $k$ -th episode. At the beginning of each episode  $k$ , the environment selects the initial state  $s_{1,k}$ . The agent proposes a policy  $\pi_k = \{\pi_h^k\}_{h \in [H]}$  based on the history up to the end of episode  $k-1$ , and then executes  $\pi_k$  to generate a new trajectory  $\{(s_{h,k}, a_{h,k}, r_{h,k})\}_{h \in [H]}$ . Here  $a_{h,k} = \pi_h^k(s_{h,k})$ ,  $r_{h,k} \sim R_h(s_{h,k}, a_{h,k})$  and  $s_{h+1,k} \sim \mathbb{P}(\cdot | s_{h,k}, a_{h,k})$ . **Here  $R_h(s, a)$  denotes the distribution of the random reward conditioned on the state-action pair  $(s, a)$  at horizon  $h$ , with its expected value being  $r_h(s, a)$ .** The agent aims to minimize the cumulative regret over  $K$  episodes, given by

$$\text{Reg}(K) := \sum_{k=1}^K (V_1^* - V_1^{\pi_k})(s_{1,k}).$$

### 3.2 High-level Algorithm Description

In this subsection, we introduce VARA, an algorithm present in Algorithm 2, that extends AdaOFUL to solve linear MDPs with heavy-tailed rewards. At a high level, the VARA algorithm is built on LSVI-UCB++ (He et al., 2022), an algorithm proposed recently to achieve minimax optimality for linear MDPs. LSVI-UCB++ (He et al., 2022) uses weighted ridge regression, where the weights depend on some proper variance estimators  $\sigma_{h,k}$ 's. The variance estimation techniques in LSVI-UCB++ are important to obtain variance-aware regrets. These techniques include (i) separate variance estimation, (ii) monotonicity of value functions, and (iii) rare-switching value function update.

Due to limited space, we present the detailed and formal algorithm description in Appendix A and focus on the differences between VARA and LSVI-UCB++ (He et al., 2022) here. To obtain variance-aware regrets under heavy-tailed rewards, we made two improvements to LSVI-UCB++. First, while LSVI-UCB++ assumes a deterministic, uniformly bounded, and known reward function, we use AdaOFUL to estimate the parameters  $\theta_h^*$  and  $\psi_h^*$  for both the expected reward functions and their second-order moments. This complicates the construction of the variance estimators  $\sigma_{h,k}$  and requires a more detailed analysis of their impacts on the final regrets (see Lemma D.10). Second, previous works use the Azuma-Hoeffding inequality to analyze the concentration effect in the suboptimality gap, which leads to the regret of  $\tilde{\mathcal{O}}(\sqrt{K})$ . Instead, we use a variance-aware Bernstein inequality and produce a much tighter upper bound of  $\tilde{\mathcal{O}}(1)$  for the concentration effect (see Lemma D.8). We explain the analytical novelty in detail in Appendix D.2.

<sup>1</sup>For example, we can guess  $\mathcal{H}$  as 2, 4, 6,  $\dots$ . After a logarithmic number of guessing, we can find a true upper bound for  $\sup_\pi \mathbb{E} R_\pi$ . One can run a similar procedure for other quantities.

**Algorithm 2** The VARA algorithm (informal)

---

**Require** :  $K, H, \mathcal{H}, \mathcal{V}, W, \sigma_R, \sigma_{R^2}$ .

```

1 for episode  $k = 1$  to  $K$  do
2   for horizon  $h = H$  to  $1$  do
3     Based on all  $\{\theta_{h,k'}, \psi_{h,k'}\}_{k' \leq k}$ , estimate an optimistic  $\bar{Q}_h^k$  and a pessimistic Q-value  $\underline{Q}_h^k$  by LSVI-UCB++.
4      $\bar{V}_h^k(\cdot) = \max_a \bar{Q}_h^k(\cdot, a)$ ,  $\underline{V}_h^k(\cdot) = \max_a \underline{Q}_h^k(\cdot, a)$ ,  $\pi_h^k(\cdot) \in \operatorname{argmax}_a \bar{Q}_h^k(\cdot, a)$ .
5   end
6   for horizon  $h = 1$  to  $H$  do
7     Play  $a_{h,k} = \pi_h^k(s_{h,k})$  and observe  $r_{h,k} \sim R_h(s_{h,k}, a_{h,k})$ ,  $s_{h+1,k} \sim \mathbb{P}(\cdot | s_{h,k}, a_{h,k})$ .
8     Observe feature vectors  $\phi_{h,k} = \phi(s_{h,k}, a_{h,k})$  and  $\tilde{\phi}_{h,k} = \tilde{\phi}(s_{h,k}, a_{h,k})$ .
9     Update the estimated variance  $\sigma_{h,k}$  using observed data and estimated values  $\bar{Q}_h^k$  and  $\underline{Q}_h^k$ .
10    Update the parameters  $w_{h,k}, \tau_{h,k}, \tilde{w}_{h,k}, \tilde{\tau}_{h,k}$  following the spirit of AdaOFUL.
11    Using  $\{\phi_{h,k'}\}_{k' \leq k}$  and  $\{(\sigma_{h,k'}, w_{h,k'}, \tau_{h,k'})\}_{k' \leq k}$ , AdaOFUL produces  $\theta_{h,k}$  as the estimate for  $\theta_h^*$ .
12    Using  $\{\tilde{\phi}_{h,k'}\}_{k' \leq k}$  and  $\{(\sigma_{h,k'}, \tilde{w}_{h,k'}, \tilde{\tau}_{h,k'})\}_{k' \leq k}$ , AdaOFUL produces  $\psi_{h,k}$  as the estimate for  $\psi_h^*$ .
13  end
14 end

```

---

**3.3 Regret Analysis**

This section presents the statistical, space, and computational complexities of Algorithm 3.

**Theorem 3.1.** Consider a linear MDP satisfying Assumption 3.1 and 3.2. For any  $\delta \in (0, 1)$ , with probability at least  $1 - 2\delta$ , Algorithm 3 achieves the following regret

$$\operatorname{Reg}(K) = \tilde{\mathcal{O}} \left( d\sqrt{HK\mathcal{G}^*} + Hd\sqrt{K}\sigma_{\min} + \frac{H^{2.5}d^6\mathcal{H}^2 + Hd^2\sigma_{R^2}}{\sigma_{\min}} + H^3d^5\mathcal{H} + Hd\sigma_R + Hd^2 \right), \quad (3.1)$$

where  $\sigma_{\min}$  is a manually set and arbitrary lower bound for all variance estimators  $\sigma_{h,k}$ 's,

$$\mathcal{G}^* = \min \left\{ \sum_{h=1}^H \mathbb{E}_{(s,a) \sim \tilde{d}_h^K} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a), \mathcal{V}^2 \right\}, \quad (3.2)$$

and  $\tilde{d}_h^K(s, a) = \frac{1}{K} \sum_{k=1}^K d_h^{\pi_k}(s, a)$  with  $d_h^{\pi_k}(s, a) = \mathbb{P}^{\pi_k}(s_h = s, a_h = a | s_0 = s_{1,k})$  the probability of reaching  $(s_{h,k}, a_{h,k}) = (s, a)$  at the  $h$ -th step when the agent starts from  $s_{1,k}$  and follows the policy  $\pi_k$ .

**Trade-off by  $\sigma_{\min}$**  Theorem 3.1 reveals a trade-off arising from the choice of  $\sigma_{\min}$ . The second term in equation 3.1, stemming from the imposed lower bound on variance estimates for stability purposes, is positively dependent on  $\sigma_{\min}$ . Consequently, if  $\sigma_{\min} = 0$ , this term vanishes. The third term is negatively dependent on  $\sigma_{\min}$  due to its effect on  $\mathbf{H}_T$ . Consider an extreme case. If  $\sigma_{\min} = \infty$ ,  $\mathbf{H}_T$  is reduced to  $\lambda \mathbf{I}$ , implying the shape of the confidence region  $\mathcal{C}_t$  changes. This, then, would slightly decrease the confidence radii  $\beta_T$  and the regret. The choice of  $\sigma_{\min}$  must balance these opposing effects. Corollary 3.1 implies that choosing the optimal  $\sigma_{\min}^* = \sqrt{H^{1.5}d^5\mathcal{H}^2 + d\sigma_{R^2}} \cdot K^{-\frac{1}{4}}$  yields a regret barrier of  $\tilde{\mathcal{O}} \left( Hd \cdot \sqrt{d^5\mathcal{H}^2 + d\sigma_{R^2}} \cdot \sqrt[4]{K} \right)$ . When  $K$  is sufficiently large, the regret bound in equation 3.3 can be further simplified to  $\tilde{\mathcal{O}}(d\sqrt{H\mathcal{G}^*K})$ . To the best of our knowledge, Theorem 3.1 is the first to derive the variance-aware regret for linear MDPs, especially with heavy-tailed rewards.

**Corollary 3.1.** Under the same setting of Theorem 3.1, if we set  $\sigma_{\min} = \sqrt{H^{1.5}d^5\mathcal{H}^2 + d\sigma_{R^2}} \cdot K^{-\frac{1}{4}}$ , the regret of VARA is bounded by

$$\operatorname{Reg}(K) = \tilde{\mathcal{O}} \left( d\sqrt{H\mathcal{G}^*K} + Hd\sqrt{d^5\mathcal{H}^2 + d\sigma_{R^2}} \cdot \sqrt[4]{K} + H^3d^5\mathcal{H} + Hd\sigma_R + Hd^2 \right). \quad (3.3)$$

**Instance-dependent quantity  $\mathcal{G}^*$**  The quantity  $\mathcal{G}^*$  is given by equation 3.2. Firstly, it is bounded above by  $\mathcal{V}^2$  in Assumption 3.2, which sets an upper bound on the variance of the cumulative random rewards received when following any policy. Other upper bounds such as  $\mathcal{H}, \sigma_R$  and  $\sigma_{R^2}$  don't involve in  $\mathcal{G}^*$  and thus the regret when  $K$  is sufficiently large. Secondly, even  $\mathcal{V}$  is set to be extremely large,  $\mathcal{G}^*$  is no greater than the sum of per-round conditional variances  $[\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a)$ , weighted by an averaged occupancy measure  $\tilde{d}_h^K(s, a) := \frac{1}{K} \sum_{k=1}^K d_h^{\pi_k}(s, a)$ . The function  $\tilde{d}_h^K(\cdot, \cdot)$  introduces a probability measure on  $\mathcal{S} \times \mathcal{A}$  for any fixed  $h \in [H]$ , in accordance with the definition of  $d_h^\pi$ , which records the history of the policies taken.

Our variance-aware regret has two key features. Firstly, we do not require any prior knowledge of  $\mathcal{G}^*$  to achieve variance awareness, which is the same as (Zanette & Brunskill, 2019). Secondly, the additional conditions imposed on the MDP structure lead to other instance-dependent regrets. In the following, we also impose Assumption 3.3 for a fair comparison with related work. However, we would like to emphasize that all of our results are obtained in the presence of heavy-tailed rewards.

**Assumption 3.3.** We assume that  $0 \leq R_h(s, a) \leq 1$  for all  $h \in [H]$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

### 3.4 Other Instance-dependent Regrets

**Worst-case regret** Under Assumption 3.3,  $\mathcal{V}^2 = H^2$  according to the law of total variance (Azar et al., 2013). Consequently, we can infer that  $\mathcal{G}^* \leq H^2$ , and the regret reduces to the minimax optimal  $\tilde{\mathcal{O}}(dH\sqrt{HK})$  (He et al., 2022). The authors achieved this regret by directly setting  $\sigma_{\min} = 1/H$ , without taking into account the trade-off introduced by  $\sigma_{\min}$ . Although this was sufficient for their worst-case scenario, it was not suitable for our goal of achieving variance awareness. If we also set  $\sigma_{\min} = 1/H$ , the second term in equation 3.1 becomes  $d\sqrt{K}$ , and we cannot determine the dominant term between  $d\sqrt{HK\mathcal{G}^*}$  and  $d\sqrt{K}$ . Once we balance the trade-off of  $\sigma_{\min}$ , the second term becomes much smaller, making  $\tilde{\mathcal{O}}(d\sqrt{H\mathcal{G}^*K})$  the dominant term.

**Range-dependent regret** Let  $\mathcal{S}_{s,a}$  be the set of immediate successor states after one transition from state  $s$  upon taking action  $a$ , which is also the support set of  $\mathbb{P}(\cdot|s, a)$ . Define  $\Phi_{\text{succ}}$  as the maximum value function range when restricted to the immediate successor states:

$$\Phi_{\text{succ}} := \sup_{h \in [H]} \sup_{(s,a)} \left[ \sup_{s' \in \mathcal{S}_{s,a}} V_{h+1}^*(s') - \inf_{s' \in \mathcal{S}_{s,a}} V_{h+1}^*(s') \right].$$

Since the variance is upper bounded by one-fourth of the square range of a random variable, we have  $\sup_{h \in [H]} \sup_{(s,a)} [\mathbb{V}_h V_{h+1}^*](s, a) \leq \frac{1}{4} \Phi_{\text{succ}}^2$  and thus  $\mathcal{G}^* \leq H(\sigma_R^2 + \Phi_{\text{succ}}^2)$ . Therefore, our regret reduces to  $\tilde{\mathcal{O}}\left(dH\sqrt{(\sigma_R^2 + \Phi_{\text{succ}}^2)K}\right)$ . It is worth noting that similar range-dependent regrets have been derived for tabular MDPs with bounded rewards (Bartlett & Tewari, 2009; Fruit et al., 2018; Zanette & Brunskill, 2019), but to the best of our knowledge, we obtain the first such result for linear MDPs with heavy-tailed rewards.

**First-order regret** The first-order regret that scales proportionally to  $V_1^*$ , where  $V_1^* := V_1^*(s_1)$  is the value of the optimal value policy at the initial state  $s_1$ , has been studied for tabular MDPs (Jin et al., 2020a) and linear MDPs (Wagenmaker et al., 2022a).<sup>2</sup> However, under Assumption 3.3, the corresponding instance-dependent quantity  $H^2 V_1^*$  can be much larger than  $\mathcal{G}^*$ . This is because

$$\mathcal{G}^* \stackrel{(a)}{\leq} H \sum_{h=1}^H \mathbb{E}_{(s,a) \sim \tilde{d}_h^K} [r_h + \mathbb{P}_h V_{h+1}^*](s, a) \stackrel{(b)}{\leq} H \sum_{h=1}^H \mathbb{E}_{(s,a) \sim \tilde{d}_h^K} V_h^*(s, a) \stackrel{(c)}{\leq} H^2 V_1^*,$$

where (a) uses  $0 \leq R_h \leq 1$  and  $0 \leq V_{h+1}^* \leq H - 1$  under Assumption 3.3, (b) uses the optimality condition  $V_h^*(s) = [r_h + \mathbb{P}_h V_{h+1}^*](s, a)$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and (c) uses  $V_{h+1}^*(s) \leq V_h^*(s)$  for any  $h \in [H]$  and  $s \in \mathcal{S}$  and  $\sum_{a \in \mathcal{A}} \tilde{d}_h^K(s_1, a) = 1$  since each episode starts at a fixed state  $s_1$ . Moreover, even replacing  $\mathcal{G}^*$  with the coarse upper bound  $H^2 V_1^*$ , our regret bound becomes  $\tilde{\mathcal{O}}(\sqrt{d^2 H^3 V_1^* K})$ , which has a better dependence on  $d$  than  $\tilde{\mathcal{O}}(\sqrt{d^3 H^3 V_1^* K})$  in (Wagenmaker et al., 2022a).

<sup>2</sup>They assume all episodes start from the same initial state so that  $s_{1,k} \equiv s_1$ . However, our regret can be easily extended to the setting where initial states are different. In this case one should replace  $V_1^* K$  with  $\sum_{k=1}^K V_1^*(s_{1,k})$  in the regret bound.

**Concentrability-dependent regret** Let  $R_{\pi^*}$  denote the sum of random rewards collected in a trajectory following the optimal policy  $\pi^*$ . It is straightforward to see that  $\text{Var}(R_{\pi^*}) = \sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^{\pi^*}} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a)$ . Since  $\mathcal{G}^* \leq \sup_{\pi} \sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^{\pi}} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a)$ , we can show that  $\mathcal{G}^* \leq C^{\dagger} \cdot \text{Var}(R_{\pi^*})$  where  $C^{\dagger}$  is a data coverage measure defined as

$$C^{\dagger} := \sup_{\pi} \frac{\sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^{\pi}} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a)}{\sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^{\pi^*}} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a)}.$$

Therefore, our regret reduces to  $\tilde{\mathcal{O}}(d\sqrt{C^{\dagger} \text{Var}(R_{\pi^*})HK})$  given  $C^{\dagger} < \infty$ . The  $C^{\dagger}$  is a counterpart of the generalized concentrability coefficient which quantifies the effect of the distribution shift in offline RL (Chen & Jiang, 2019; Xie et al., 2021; Cheng et al., 2022).

### 3.5 Space and Computational Complexities

**Theorem 3.2** (Space and computational complexity). Assume the Nesterov accelerated method is used as a solver to solve the adaptive Huber regression. Solving a  $H$ -horizon finite MDP in  $K$  episodes, VARA takes  $\mathcal{O}(d^3 H^2 + d|\mathcal{A}|HK)$  space and has a running time of  $\tilde{\mathcal{O}}(d^4 |\mathcal{A}|H^3 K + HK(d + H^{-3/4} d^{-3/2} K^{3/4}))$ .

On one hand, VARA achieves the same space complexity as LSVI-UCB++ but is slightly worse than the original LSVI-UCB (Jin et al., 2020b) that needs  $\mathcal{O}(d^2 H + d|\mathcal{A}|HK)$  space. This is because the technique of monotone value function update requires remembering at most  $\tilde{\mathcal{O}}(dH)$  latest value functions, incurring a slightly worse dependence on  $d$  and  $H$ . On the other hand, the computational complexity of VARA  $\tilde{\mathcal{O}}(d^4 |\mathcal{A}|H^3 K + HK(d + H^{-3/4} d^{-3/2} K^{3/4}))$  is slightly worse than LSVI-UCB++'s  $\tilde{\mathcal{O}}(d^4 |\mathcal{A}|H^3 K)$  in terms of the dependence on  $K$ . This is because the adaptive Huber regression estimator does not have a closed-form solution. Even though the Nesterov accelerated method is used, a slightly larger computational complexity is still incurred due to the possibly large conditional number. However, VARA's computational complexity is better than LSVI-UCB's  $\tilde{\mathcal{O}}(d^2 |\mathcal{A}|HK^2)$  thanks to the rare-switching mechanism in LSVI-UCB++.

## 4 Related Work

**Heavy-tailed rewards in online decision making** The standard heavy-tailed setting assumes rewards with  $(1 + \varepsilon)$ -moments where  $\varepsilon > 0$ . There exists a large body of work considering this setting in multi-arm bandits, including deterministic (Vakili et al., 2013) and non-deterministic settings (Bubeck et al., 2013; Carpentier & Valko, 2014; Lattimore, 2017; Bhatt et al., 2022). To handle heavy-tailed rewards, robust mean estimation methods such as median of means and truncation have been applied to linear bandits (Medina & Yang, 2016; Shao et al., 2018; Lu et al., 2019; Xue et al., 2021). Given that our objective is to provide variance-aware regrets for general linear bandits, the minimal requirement is to have bounded second moments, which is the primary focus of this study. Under the assumption of rewards with bounded second moments, the minimax optimal regret is  $\tilde{\mathcal{O}}(d\sqrt{T})$ . Recently, Kang & Kim (2023) introduced the use of Huber regression to address heavy-tailed linear contextual bandits where the action set is fixed (such that  $\mathcal{D}_t \equiv \mathcal{D}$ ) and the arm  $\phi_t$  is independently and identically distributed, sampled from a fixed distribution over  $\mathcal{D}$ . In contrast, our proposed AdaOFUL is simpler and more versatile, capable of being applied to more complex scenarios where the arm  $\phi_t$  is selected adaptively based on historical observations.

On the other hand, there is a lack of RL algorithms designed to handle heavy-tailed rewards for MDPs. One exception is (Zhuang & Sui, 2021), which modifies UCRL2 and Q-Learning by using truncated rewards and achieves minimax optimal regret in tabular MDPs. However, none of these methods for linear bandits or MDPs provide variance-aware regrets, even if variance information is available. Moreover, simple truncation methods are not optimal in the noiseless setting.

**Variance-aware regrets for linear bandits** A weighted ridge regression-based algorithm proposed by Kirschner & Krause (2018) achieves the same regret in equation 2.7 by assuming each  $\varepsilon_t$  is  $\nu_t$ -sub-Gaussian. More recently, Zhou & Gu (2022) obtained the same regret assuming each  $\varepsilon_t$  is uniformly bounded and has

finite conditional variance  $\nu_t^2$ . In the case where the information of conditional variances  $\{\nu_t\}_{t \geq 0}$  is unknown, Zhang et al. (2021) and Kim et al. (2021) achieved regret bounds that involve sub-optimal dependence on  $d$ . The currently tightest variance-aware regret is achieved by Zhao et al. (2023) with an optimal dependence on  $d$ . Recently, Dai et al. (2022) explored variance-aware regrets in the context of high-dimensional and sparse linear bandits, a topic that extends beyond the scope of our paper. All of the above works consider light-tailed noises, which are either sub-Gaussian or uniformly bounded.

**Robust approach to instance-dependent bounds** Recent research explores the robust mean estimation approach to obtain instance-dependent regrets, leveraging the observation that robust estimators can achieve estimation errors that only depend on the noise scale. Such estimators often have better theoretical guarantees than non-robust ones, whose estimation errors additionally depend on the range of the problem noise. For instance, Pananjady & Wainwright (2020) use the median-of-means technique (Lecué & Lerasle, 2020) to achieve local minimax optimality that depends on the standard deviations of the optimal value function and random rewards for synchronous tabular MDPs. In linear bandits, Wagenmaker et al. (2022a) use Catoni’s estimator (Catoni, 2012) to estimate the mean of  $\mathbf{v}^\top \mathbf{H}_T^{-1} \phi_t y_t / \sigma_t^2$  for a fixed unit-norm vector  $\mathbf{v}$ . In contrast, we modify the adaptive Huber regression to estimate  $\boldsymbol{\theta}^*$  directly. This difference makes their bounds depend on the second moments of  $y_t$ ’s, while ours only relies on their variances. Moreover, all of these works, except ours, still assume light-tailed rewards.

**Variance-aware regrets for tabular and linear MDPs** In the context of online episodic MDPs, Zanette & Brunskill (2019) first derived a variance-aware regret bound in the tabular setting with uniformly bounded rewards. Their model-based algorithm, Euler, achieves a regret that can be bounded by either  $\tilde{\mathcal{O}}(\sqrt{\mathbb{Q}^S A H K})$  or  $\tilde{\mathcal{O}}(\sqrt{\mathcal{G}^2 S A K})$ , where  $\mathbb{Q}^* = \max_{(s,a,h)} (\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*)(s, a)$  is the maximum per-round conditional variance and  $\mathcal{G}$  is a deterministic upper bound on the maximum attainable reward on a single trajectory for any policy  $\pi$ , such that  $\sum_{h=1}^H R_h(s_h, \pi(s_h)) \leq \mathcal{G}$ . One can show that our instance-dependent quantity  $\mathcal{G}^*$  is smaller than  $\min\{H\mathbb{Q}^*, \mathcal{G}^2\}$  therein. Later, Jin et al. (2020a) adopted a modified analysis of Euler to obtain the regret bound  $\tilde{\mathcal{O}}(\sqrt{S A H^3 V_1^* K})$  with  $V_1^* = V_1^*(s_1)$ .<sup>3</sup> In linear MDPs, there are several recent works on obtaining regret bounds for model-free algorithms. For example, Wagenmaker et al. (2022a) proposed an optimistic algorithm with a regret bound that scales as  $\tilde{\mathcal{O}}(\sqrt{d^3 H^3 V_1^* K})$ . However, this algorithm is computationally inefficient. A computationally efficient alternative suffers from a slightly worse regret  $\tilde{\mathcal{O}}(\sqrt{d^4 H^3 V_1^* K})$ . All of these works utilize the instance-dependent quantity  $H^2 V_1^*$  (assuming  $\mathcal{H} = H$ ). However, as we argued, our proposed quantity  $\mathcal{G}^*$  is smaller than  $H^2 V_1^*$ , which implies that our algorithm may achieve better performance than these previous works. Another research direction explores variance-adaptive algorithms for linear mixture MDPs, as initially explored by (Zhou et al., 2021). Subsequent developments in this area were made by (Zhang et al., 2021; Zhou & Gu, 2022), culminating in the state-of-the-art advancements by (Zhao et al., 2023). While our study does not consider this particular setting, it is straightforward to extend our techniques and analysis to it, as linear mixture MDPs are generally considered simpler than linear MDPs.

**Other instance-dependent bounds** In the infinite-horizon setting, Pananjady & Wainwright (2020); Khamaru et al. (2021); Li et al. (2023) provided variance-aware sample complexities for Q-Learning and its variants in tabular MDPs, given a generative model that produces independent samples for all state-action pairs in every round. Variance-aware performance guarantees have also been established for offline RL optimization (Yin & Wang, 2021; Nguyen-Tang et al., 2023), off-policy evaluation (Min et al., 2021), stochastic approximation (Mou et al., 2020; 2022). Another approach to instance-dependence bounds focuses on the minimum suboptimality gap, which is the minimum gap between the best and second-best actions over all states (He et al., 2021; Wagenmaker et al., 2022c; Wagenmaker & Jamieson, 2022; Dong & Ma, 2022). However, due to the differences in the settings, we cannot make a meaningful comparison between these bounds and ours.

<sup>3</sup>Unlike our setting, they assume all initial states are the same, denoted as  $s_1$ . Furthermore, the original regret  $\tilde{\mathcal{O}}(\sqrt{S A H \cdot V_1^* K})$  by Jin et al. (2020a) was derived for an MDP where the reward function equals to one deterministically only at a single  $(h, s)$  pair. In this way, they have  $0 \leq V_1^* \leq 1$ . To convert it in the considered setting where  $0 \leq V_1^* \leq H$ , an additional factor of  $H$  should be multiplied to their regret.

## 5 Conclusion

This paper introduces two new algorithms, AdaOFUL for linear bandits and VARA for linear MDPs, both of which use modifications of the original adaptive Huber regression and are designed to handle online sequential decision-making. With only the assumption of bounded reward variances, our algorithms achieve either state-of-the-art or finer variance-aware regrets. Additionally, in linear MDPs, the instance-dependent quantity  $\mathcal{G}^*$  can be bounded by other instance-dependent quantities when additional structure assumptions are available. Our modified adaptive Huber regression can be a useful building block for algorithm design in online problems with heavy-tailed rewards.

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, volume 24, 2011.
- Nick Arnosti, Marissa Beck, and Paul Milgrom. Adverse selection and auction design for internet display advertising. *American Economic Review*, 106(10):2852–66, 2016.
- Alex Ayoub, Zeyu Jia, Csaba Szepesvári, Mengdi Wang, and Lin Yang. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pp. 463–474. PMLR, 2020.
- Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine Learning*, 91(3):325–349, 2013.
- Peter Bartlett and Ambuj Tewari. Regal: a regularization based algorithm for reinforcement learning in weakly communicating mdps. In *Conference on Uncertainty in Artificial Intelligence*, pp. 35–42. AUAI Press, 2009.
- Sujay Bhatt, Guanhua Fang, Ping Li, and Gennady Samorodnitsky. Nearly optimal Catoni’s M-estimator for infinite variance. In *International Conference on Machine Learning*, pp. 1925–1944, 2022.
- Steven J Bradtke and Andrew G Barto. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22(1):33–57, 1996.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *IEEE Transactions on Information Theory*, 59(11):7711–7717, 2013.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Alexandra Carpentier and Michal Valko. Extreme bandits. In *Advances in Neural Information Processing Systems*, 2014.
- Olivier Catoni. Challenging the empirical mean and empirical variance: A deviation study. *Annales de l’IHP Probabilités et statistiques*, 48(4):1148–1185, 2012.
- Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*, pp. 1042–1051. PMLR, 2019.
- Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Bellman-consistent pessimism for offline reinforcement learning. In *International Conference on Machine Learning*, volume 162, pp. 3852–3878, 2022.
- Rama Cont. Empirical properties of asset returns: Stylized facts and statistical issues. *Quantitative Finance*, 1(2):223, 2001.
- Yan Dai, Ruosong Wang, and Simon Shaolei Du. Variance-aware sparse linear bandits. In *International Conference on Learning Representations*, 2022.



- Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. 2008.
- Kefan Dong and Tengyu Ma. Asymptotic instance-optimal algorithms for interactive decision making. *arXiv preprint arXiv:2206.02326*, 2022.
- David A Freedman. On tail probabilities for martingales. *The Annals of Probability*, 3(1):100–118, 1975.
- Ronan Fruit, Matteo Pirotta, Alessandro Lazaric, and Ronald Ortner. Efficient bias-span-constrained exploration-exploitation in reinforcement learning. In *International Conference on Machine Learning*, pp. 1578–1586, 2018.
- Trevor Hastie, Robert Tibshirani, and Jerome H Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
- Jiafan He, Dongruo Zhou, and Quanquan Gu. Logarithmic regret for reinforcement learning with linear function approximation. In *International Conference on Machine Learning*, pp. 4171–4180. PMLR, 2021.
- Jiafan He, Heyang Zhao, Dongruo Zhou, and Quanquan Gu. Nearly minimax optimal reinforcement learning for linear markov decision processes. *arXiv preprint arXiv:2212.06132*, 2022.
- Pihe Hu, Yu Chen, and Longbo Huang. Nearly minimax optimal reinforcement learning with linear function approximation. In *International Conference on Machine Learning*, pp. 8971–9019, 2022.
- Peter J Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, pp. 73–101, 1964.
- John Hull. *Risk Management and Financial Institutions*. John Wiley & Sons, 2012.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for reinforcement learning. In *International Conference on Machine Learning*, pp. 4870–4879. PMLR, 2020a.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pp. 2137–2143, 2020b.
- Minhyun Kang and Gi-Soo Kim. Heavy-tailed linear bandit with Huber regression. In *Uncertainty in Artificial Intelligence*, pp. 1027–1036. PMLR, 2023.
- Koulik Khamaru, Ashwin Pananjady, Feng Ruan, Martin J Wainwright, and Michael I Jordan. Is temporal difference learning optimal? An instance-dependent analysis. *SIAM Journal on Mathematics of Data Science*, 3(4):1013–1040, 2021.
- Yeonung Kim, Insoon Yang, and Kwang-Sung Jun. Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture mdps. *arXiv preprint arXiv:2111.03289*, 2021.
- Johannes Kirschner and Andreas Krause. Information directed sampling and bandits with heteroscedastic noise. In *Conference on Learning Theory*, pp. 358–384, 2018.
- Tor Lattimore. A scale free algorithm for stochastic bandits with bounded kurtosis. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Guillaume Lecué and Matthieu Lerasle. Robust machine learning by median-of-means: Theory and practice. *The Annals of Statistics*, 48(2):906–931, 2020.
- Gen Li, Changxiao Cai, Yuxin Chen, Yuntao Gu, Yuting Wei, and Yuejie Chi. Is Q-learning minimax optimal? A tight sample complexity analysis. *arXiv preprint arXiv:2102.06548*, 2021.

- Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. Deep reinforcement learning for dialogue generation. In *Conference on Empirical Methods in Natural Language Processing*, pp. 1192–1202, 2016.
- Xiang Li, Wenhao Yang, Jiadong Liang, Zhihua Zhang, and Michael I Jordan. A statistical analysis of Polyak-Ruppert averaged Q-learning. In *International Conference on Artificial Intelligence and Statistics*, 2023.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Shiyin Lu, Guanghui Wang, Yao Hu, and Lijun Zhang. Optimal algorithms for Lipschitz bandits with heavy-tailed rewards. In *International Conference on Machine Learning*, pp. 4154–4163, 2019.
- Andres Munoz Medina and Scott Yang. No-regret algorithms for heavy-tailed linear bandits. In *International Conference on Machine Learning*, pp. 1642–1650, 2016.
- Francisco S Melo and M Isabel Ribeiro. Q-learning with linear function approximation. In *International Conference on Computational Learning Theory*, pp. 308–322, 2007.
- Yifei Min, Tianhao Wang, Dongruo Zhou, and Quanquan Gu. Variance-aware off-policy evaluation with linear function approximation. In *Advances in Neural Information Processing Systems*, volume 34, pp. 7598–7610, 2021.
- Wenlong Mou, Chris Junchi Li, Martin J Wainwright, Peter L Bartlett, and Michael I Jordan. On linear stochastic approximation: Fine-grained Polyak-Ruppert and non-asymptotic concentration. In *Conference on Learning Theory*, pp. 2947–2997, 2020.
- Wenlong Mou, Koulik Khamaru, Martin J Wainwright, Peter L Bartlett, and Michael I Jordan. Optimal variance-reduced stochastic approximation in banach spaces. *arXiv preprint arXiv:2201.08518*, 2022.
- Thanh Nguyen-Tang, Ming Yin, Sunil Gupta, Svetha Venkatesh, and Raman Arora. On instance-dependent bounds for offline reinforcement learning with linear function approximation. In *Association for the Advancement of Artificial Intelligence*, 2023.
- Ashwin Pananjady and Martin J Wainwright. Instance-dependent  $\ell_\infty$  bounds for policy evaluation in tabular reinforcement learning. *IEEE Transactions on Information Theory*, 67(1):566–585, 2020.
- Alexandra Posekany, Klaus Felsenstein, and Peter Sykacek. Biological assessment of robust noise models in microarray data analysis. *Bioinformatics*, 27(6):807–814, 2011.
- Han Shao, Xiaotian Yu, Irwin King, and Michael R Lyu. Almost optimal algorithms for linear stochastic bandits with heavy-tailed payoffs. In *Advances in Neural Information Processing Systems*, 2018.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Qiang Sun. Do we need to estimate the variance in robust mean estimation? *arXiv preprint arXiv:2107.00118*, 2021.
- Qiang Sun, Wen-Xin Zhou, and Jianqing Fan. Adaptive Huber regression. *Journal of the American Statistical Association*, 115(529):254–265, 2020.
- Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- Sattar Vakili, Keqin Liu, and Qing Zhao. Deterministic sequencing of exploration and exploitation for multi-armed bandit problems. *IEEE Journal of Selected Topics in Signal Processing*, 7(5):759–767, 2013.

- Andrew Wagenmaker and Kevin Jamieson. Instance-dependent near-optimal policy identification in linear MDPs via online experiment design. In *Advances in Neural Information Processing Systems*, 2022.
- Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. First-order regret in reinforcement learning with linear function approximation: A robust estimation approach. In *International Conference on Machine Learning*, pp. 22384–22429. PMLR, 2022a.
- Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. Reward-free RL is no harder than reward-aware RL in linear Markov decision processes. In *International Conference on Machine Learning*, pp. 22430–22456, 2022b.
- Andrew J Wagenmaker, Max Simchowitz, and Kevin Jamieson. Beyond no regret: Instance-dependent pac reinforcement learning. In *Conference on Learning Theory*, pp. 358–418, 2022c.
- Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent pessimism for offline reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 34, pp. 6683–6694, 2021.
- Bo Xue, Guanghui Wang, Yimu Wang, and Lijun Zhang. Nearly optimal regret for stochastic linear bandits with heavy-tailed payoffs. In *International Conference on International Joint Conferences on Artificial Intelligence*, pp. 2936–2942, 2021.
- Lin Yang and Mengdi Wang. Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning*, pp. 6995–7004, 2019.
- Lin Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pp. 10746–10756. PMLR, 2020.
- Ming Yin and Yu-Xiang Wang. Towards instance-optimal offline reinforcement learning with pessimism. In *Advances in Neural Information Processing Systems*, pp. 4065–4078, 2021.
- Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*, pp. 7304–7312. PMLR, 2019.
- Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pp. 10978–10989, 2020.
- Zihan Zhang, Jiaqi Yang, Xiangyang Ji, and Simon S Du. Improved variance-aware confidence sets for linear bandits and linear mixture MDP. In *Advances in Neural Information Processing Systems*, pp. 4342–4355, 2021.
- Heyang Zhao, Jiafan He, Dongruo Zhou, Tong Zhang, and Quanquan Gu. Variance-dependent regret bounds for linear bandits and reinforcement learning: Adaptivity and computational efficiency. *arXiv preprint arXiv:2302.10371*, 2023.
- Dongruo Zhou and Quanquan Gu. Computationally efficient horizon-free reinforcement learning for linear mixture mdps. *arXiv preprint arXiv:2205.11507*, 2022.
- Dongruo Zhou, Quanquan Gu, and Csaba Szepesvari. Nearly minimax optimal reinforcement learning for linear mixture Markov decision processes. In *Conference on Learning Theory*, pp. 4532–4576. PMLR, 2021.
- Vincent Zhuang and Yanan Sui. No-regret reinforcement learning with heavy-tailed rewards. In *International Conference on Artificial Intelligence and Statistics*, pp. 3385–3393, 2021.

# Appendix

## Overview

We describe VARA detailedly in Appendix A and explain the rationale behind its variance estimator in Appendix B. Appendix C contains proofs for Theorem 2.2 and related lemmas specifically for linear bandits. The theoretical analysis for VARA is presented in Appendix D, where we offer a proof sketch for Theorem 3.1, while all related technical lemmas are deferred to Appendices F and G. We also highlight the differences between our analysis and previous work. In Appendix E, we provide a proof for Theorem 3.2 that analyzes the space and computational complexity of VARA.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Our contributions . . . . .	2
<b>2</b>	<b>Variance-aware Regret for Heavy-tailed Linear Bandits</b>	<b>3</b>
2.1	Heavy-tailed Stochastic Linear Bandit . . . . .	3
2.2	Algorithm Description . . . . .	3
2.3	Regret Analysis . . . . .	5
2.4	Proof Sketch of Theorem 2.1 . . . . .	6
2.5	A Numerical Study . . . . .	7
<b>3</b>	<b>An Extension to Linear MDPs</b>	<b>8</b>
3.1	Linear MDPs with Heavy-tailed Rewards . . . . .	8
3.2	High-level Algorithm Description . . . . .	9
3.3	Regret Analysis . . . . .	10
3.4	Other Instance-dependent Regrets . . . . .	11
3.5	Space and Computational Complexities . . . . .	12
<b>4</b>	<b>Related Work</b>	<b>12</b>
<b>5</b>	<b>Conclusion</b>	<b>14</b>
<b>A</b>	<b>Detailed Algorithm Description for VARA</b>	<b>19</b>
<b>B</b>	<b>Variance Estimation for Value Functions</b>	<b>21</b>
<b>C</b>	<b>Proof for Section 2.3</b>	<b>23</b>
C.1	Proof of Theorem 2.2 . . . . .	23
C.2	Proof of Lemma 2.1 . . . . .	25
C.3	Proof of Lemma 2.2 . . . . .	28
C.4	Proof of Lemma C.1 . . . . .	30

C.5	Proof of Lemma C.2	31
<b>D</b>	<b>Proof of Theorem 3.1</b>	<b>32</b>
D.1	High-Probability Events	32
D.2	Regret Analysis	34
<b>E</b>	<b>Proof of Theorem 3.2</b>	<b>36</b>
<b>F</b>	<b>Omitted lemmas in Section D</b>	<b>37</b>
F.1	Proof of Lemma D.1	37
F.2	Proof of Lemma D.2	37
F.3	Proof of Lemma F.1	38
F.4	Proof of Lemma D.3	40
F.5	Proof of Lemma D.4	41
F.6	Proof of Lemma D.5	41
F.7	Proof of Lemma D.6	43
F.8	Proof of Lemma D.7	43
F.9	Proof of Lemma D.8	45
F.10	Proof of Lemma D.9	47
F.11	Proof of Lemma D.10	48
F.11.1	Proof of Lemma F.2	51
F.11.2	Proof of Lemma F.3	51
F.11.3	Proof of Lemma F.4	53
<b>G</b>	<b>Auxiliary Lemmas</b>	<b>54</b>
G.1	Concentration Inequalities	54
G.2	Elliptical Lemmas	55
G.3	Function Class and Covering Number	55

## A Detailed Algorithm Description for VARA

For each episode  $k$ , we perform optimistic value iterations (Lines 3-11), compute the greedy policy  $\pi_h^k$  with respect to the pessimistic value function  $\bar{Q}_h^k$  (Line 12), and then execute it to collect a new trajectory of data (Lines 16-17). The rest of Algorithm 3 updates maintained estimators, including the conditional variances  $\sigma_{h,k}^2$  (Line 19), the transition parameters  $\mu_{h,k}$  (Line 20), the reward parameters  $\theta_{h,k}, \psi_{h,k}$  (Lines 21-22), and the Hessian matrices  $H_{h,k}, \tilde{H}_{h,k}$  (Line 23). In what follows, we discuss in detail the key steps of Algorithm 3 in more detail.

**Reward estimation** Since rewards are collected in an adaptive manner and have only finite second moments, we use the same strategy adopted in AdaOFUL to estimate  $\theta_h^*$ :

$$\theta_{h,k} := \operatorname{argmin}_{\theta \in \text{Ball}_d(W)} \left\{ L_{h,k}^{(R)}(\theta) := \frac{\lambda}{2} \|\theta\|^2 + \sum_{j=1}^k \ell_{\tau_{h,j}} \left( \frac{r_{h,j} - \langle \phi_{h,j}, \theta \rangle}{\sigma_{h,j}} \right) \right\}. \quad (\text{A.1})$$

Following the spirit of Theorem 2.1, we set  $\tau_0 = \tilde{\mathcal{O}}(\sqrt{d})$  with its detailed expression provided in equation D.5 of the online supplement.

**Transition estimation** Let  $\delta(s) \in \mathbb{R}^{|\mathcal{S}|}$  be a one-hot vector that is zero everywhere except for the entry corresponding to the state  $s$ , which is one. We define  $\varepsilon_{h,k} = \mathbb{P}_h(\cdot | s_{h,k}, a_{h,k}) - \delta(s_{h+1,k})$ . As  $\mathbb{E}[\varepsilon_{h,k} | \mathcal{F}_{h,k}] = \mathbf{0}$ ,  $\delta(s_{h+1,k})$  is an unbiased estimator of  $\mathbb{P}_h(\cdot | s_{h,k}, a_{h,k}) = \mu_h^\top \phi(s_{h,k}, a_{h,k}) = \mu_h^\top \phi_{h,k}$ . Thus, we can learn  $\mu_h$  by regressing  $\delta(s_{h+1,k})$  on  $\phi_{h,k} := \phi(s_{h,k}, a_{h,k})$ :

$$\mu_{h,k} := \operatorname{argmin}_{\mu \in \mathbb{R}^d \times |\mathcal{S}|} \left\{ L_{h,k}^{(P)}(\mu) := \frac{\lambda}{2} \|\mu\|_F^2 + \sum_{j=1}^k \left\| \frac{\mu_h^\top \phi_{h,k} - \delta(s_{h+1,j})}{\sigma_{h,j}} \right\|^2 \right\} \quad (\text{A.2})$$

where  $\|\cdot\|_F$  denotes the Frobenius norm. This problem admits a closed-form solution given by  $\mu_{h,k} = \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} \delta(s_{h+1,j})^\top$ . We emphasize that VARA doesn't need to compute  $\mu_{h,k}$  exactly out. VARA relies on only the matrix product of  $\mu_{h,k}$  and a vectorized value function  $\mathbf{V}$  that is  $\mu_{h,k} \mathbf{V} = \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} V(s_{h+1,j})$  for any value function  $V(\cdot)$ . As Theorem 3.2 shows, both the computation and space complexity does not depend on the finite value of  $|\mathcal{S}|$ .

**Variance estimation for rewards** In linear MDPs, estimating the variance of the reward  $R_h(s, a)$  is straightforward. Since  $\mathbb{P}_h R_h^2(s, a) = \langle \tilde{\phi}(s, a), \psi_h^* \rangle$ , we estimate  $\psi_h^*$  by

$$\psi_{h,k} := \operatorname{argmin}_{\psi \in \text{Ball}_d(W)} \left\{ L_{h,k}^{(R^2)}(\psi) := \frac{\lambda}{2} \|\psi\|^2 + \sum_{j=1}^k \ell_{\tilde{\tau}_{h,j}} \left( \frac{r_{h,j}^2 - \langle \tilde{\phi}_{h,j}, \psi \rangle}{\sigma_{h,j}} \right) \right\} \quad (\text{A.3})$$

where  $\tilde{\tau}_{h,k} = \tilde{\tau}_0 \sqrt{1 + \tilde{w}_{h,k}^2} / \tilde{w}_{h,k}$  is the corresponding robustification parameter and  $\tilde{w}_{h,k} = \|\tilde{\phi}_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}$  is the importance weight. We then estimate  $[\mathbb{V}_h R_h](s_h, a_h)$  by

$$[\hat{\mathbb{V}}_h R_h](s_{h,k}, a_{h,k}) = \langle \tilde{\phi}_{h,k}, \psi_{h,k-1} \rangle - [\langle \phi_{h,k}, \theta_{h,k-1} \rangle_{[0, \mathcal{H}]}]^2. \quad (\text{A.4})$$

**Variance estimation** Inspired by Hu et al. (2022), we set the variance estimator  $\sigma_{h,k}$  to be

$$\sigma_{h,k}^2 = \max \left\{ \sigma_{\min}^2, d^3 H \cdot E_{h,k}, J_{h,k}, c_0^{-2} b_{h,k}^2, \left( \frac{W}{\sqrt{c_1 d}} + \mathcal{H} d^{2.5} H \right) b_{h,k} \right\} \quad (\text{A.5})$$

where  $\sigma_{\min}$  is a small positive constant to avoid singularity,  $b_{h,k} = \max\{\|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}\}$  is the bonus term,  $E_{h,k}$  and  $J_{h,k}$  are defined as

$$J_{h,k} = [\hat{\mathbb{V}}_{h,k} R_h + \hat{\mathbb{V}}_{h,k} \bar{V}_{h+1}^k](s_{h,k}, a_{h,k}) + R_{h,k} + U_{h,k}, \quad (\text{A.6})$$

$$E_{h,k} = \min \left\{ \mathcal{H}^2, 2\mathcal{H}\beta_0 \cdot \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \mathcal{H} \cdot [\hat{\mathbb{P}}_{h,k}(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)](s_{h,k}, a_{h,k}) \right\}, \quad (\text{A.7})$$

in which  $\beta_0 = \tilde{\mathcal{O}}(\sigma_{\min}^{-1} \mathcal{H} \sqrt{d^3 H})$  is an initial exploration radius,  $\hat{\mathbb{P}}_{h,k}(\cdot | s, a) = \mu_{h,k-1}^\top \phi(s, a)$  is the empirical transition kernel at the  $h$ -th round and  $k$ -th episode,  $\hat{\mathbb{V}}_h(\cdot)$  the empirical variance operator defined in equation A.4, and  $R_{h,k}, U_{h,k}$  are defined as

$$R_{h,k} := \beta_{R^2} \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}} + 2\mathcal{H}\beta_R \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \quad (\text{A.8})$$

$$U_{h,k} = \min \left\{ \mathcal{V}^2, 11\mathcal{H}\beta_0 \cdot \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 4\mathcal{H} \cdot \hat{\mathbb{P}}_{h,k}(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \right\}, \quad (\text{A.9})$$

with  $\beta_R = \tilde{\mathcal{O}}(\sqrt{d}), \beta_{R^2} = \tilde{\mathcal{O}}(\sqrt{d} + \sqrt{d} \frac{\sigma_{R^2}}{\sigma_{\min}})$  being two initial exploration radiuses. In Appendix B, we explain in detail why  $\sigma_{h,k}$ 's are taken in the above way.



**Algorithm 3** The VARA algorithm (formal)**Require** :  $K, H, \mathcal{H}, \mathcal{V}, W, \sigma_R, \sigma_{R^2}, \tau_0, \tilde{\tau}_0$ .**Initialization** :  $\mathbf{H}_{h,0} = \tilde{\mathbf{H}}_{h,0} = \lambda \mathbf{I}, c_0 = \frac{1}{6\sqrt{3 \log \frac{2HK^2}{\delta}}}, c_1 = \frac{1}{42 \cdot \log \frac{2HK^2}{\delta}}, \lambda = \frac{1}{\mathcal{H}^2 + W^2}, k_{\text{last}} = 1$ .

```

1 for episode  $k = 1$  to  $K$  do
2    $\bar{\mathbf{V}}_{H+1}^k(\cdot) = \underline{\mathbf{V}}_{H+1}^k(\cdot) = 0$  for round  $h = H$  to 1 do
3     if there exists a stage  $h' \in [H]$  such that  $\det(\mathbf{H}_{h',k-1}) \geq 2\det(\mathbf{H}_{h',k_{\text{last}}-1})$  then
4       Compute the products  $\boldsymbol{\mu}_{h,k-1}\bar{\mathbf{V}}_{h+1}^k$  and  $\boldsymbol{\mu}_{h,k-1}\underline{\mathbf{V}}_{h+1}^k$  with  $\boldsymbol{\mu}_{h,k}$  given in equation A.2.
5        $\hat{Q}_h^k(\cdot, \cdot) = \langle \phi(\cdot, \cdot), \boldsymbol{\theta}_{h,k-1} + \boldsymbol{\mu}_{h,k-1}\bar{\mathbf{V}}_{h+1}^k \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}}$ .
6        $\check{Q}_h^k(\cdot, \cdot) = \langle \phi(\cdot, \cdot), \boldsymbol{\theta}_{h,k-1} + \boldsymbol{\mu}_{h,k-1}\underline{\mathbf{V}}_{h+1}^k \rangle - \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}}$ .
7        $\bar{Q}_h^k(\cdot, \cdot) = \min \{ \hat{Q}_h^k(\cdot, \cdot), \bar{Q}_h^{k-1}(\cdot, \cdot), \mathcal{H} \}, \underline{Q}_h^k(\cdot, \cdot) = \max \{ \check{Q}_h^k(\cdot, \cdot), \underline{Q}_h^{k-1}(\cdot, \cdot), 0 \}$ .
8       Record the last updating episode  $k_{\text{last}} = k$ .
9     else
10       $\bar{Q}_h^k(\cdot, \cdot) = \bar{Q}_h^{k-1}(\cdot, \cdot), \underline{Q}_h^k(\cdot, \cdot) = \underline{Q}_h^{k-1}(\cdot, \cdot)$ .
11    end
12     $\bar{\mathbf{V}}_h^k(\cdot) = \max_a \bar{Q}_h^k(\cdot, a), \underline{\mathbf{V}}_h^k(\cdot) = \max_a \underline{Q}_h^k(\cdot, a)$ .
13     $\pi_h^k(\cdot) \in \operatorname{argmax}_a \bar{Q}_h^k(\cdot, a)$ .
14  end
15  Receive the initial state  $s_{1,k}$ .
16  for round  $h = 1$  to  $H$  do
17    Play  $a_{h,k} = \pi_h^k(s_{h,k})$  and observe  $r_{h,k} \sim R_h(s_{h,k}, a_{h,k}), s_{h+1,k} \sim \mathbb{P}(\cdot | s_{h,k}, a_{h,k})$ .
18    Observe feature vectors  $\phi_{h,k} = \phi(s_{h,k}, a_{h,k})$  and  $\tilde{\phi}_{h,k} = \tilde{\phi}(s_{h,k}, a_{h,k})$ .
19    Set the bonus as  $b_{h,k} = \max \{ \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}} \}$ .
20    Set the estimated variance  $\sigma_{h,k}$  as in equation A.5.
21    Compute  $\boldsymbol{\theta}_{h,k}$  via equation A.1 with  $\tau_{h,k} = \tau_0 \sqrt{1 + w_{h,k}^2/w_{h,k}}$  and  $w_{h,k} = \sigma_{h,k}^{-1} \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}$ .
22    Compute  $\boldsymbol{\psi}_{h,k}$  via equation A.3 with  $\tilde{\tau}_{h,k} = \tilde{\tau}_0 \sqrt{1 + \tilde{w}_{h,k}^2/\tilde{w}_{h,k}}$  and  $\tilde{w}_{h,k} = \sigma_{h,k}^{-1} \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}$ .
23    Update  $\mathbf{H}_{h,k} = \mathbf{H}_{h,k-1} + \sigma_{h,k}^{-2} \phi_{h,k} \phi_{h,k}^\top$  and  $\tilde{\mathbf{H}}_{h,k} = \tilde{\mathbf{H}}_{h,k-1} + \sigma_{h,k}^{-2} \tilde{\phi}_{h,k} \tilde{\phi}_{h,k}^\top$ .
24  end
25 end

```

**B Variance Estimation for Value Functions**

In order to achieve worst-case optimality, He et al. (2022) proposes two important techniques we adopt in Algorithm 3.

The first is the monotonicity of value functions. Specifically, we aim to enforce a decrease in  $k$  for the actual optimistic value function  $\bar{Q}_h^k(\cdot, \cdot)$  and an increase in  $k$  for the actual pessimistic value function  $\underline{Q}_h^k(\cdot, \cdot)$ . This concept is explained in detail below. In linear MDPs, we have  $[\mathbb{P}_h V_{h+1}](s, a) = \langle \phi(s, a), \boldsymbol{\mu}_h^* \bar{\mathbf{V}}_{h+1} \rangle$  for any value function  $V = \{V_h\}_{h \in [H]}$  and  $[\mathbb{P}_h R_h](s, a) = \langle \phi(s, a), \boldsymbol{\theta}_h^* \rangle$  for all  $h \in [H]$ . One crucial aspect of typical analysis (including ours) is demonstrating the high probability of the following event outlined in Appendix D.1. Specifically, for all  $h \in [H]$  and  $k \in [K]$ , we need to establish that the following equations hold simultaneously with high probability:

$$\begin{aligned}
\|\boldsymbol{\theta}_{h,k} - \boldsymbol{\theta}_h^*\|_{\mathbf{H}_{h,k}} &\leq \beta_R = \tilde{\mathcal{O}}(\sqrt{d}), \\
\|(\boldsymbol{\mu}_{h,k-1} - \boldsymbol{\mu}_h^*)\bar{\mathbf{V}}_{h+1}^k\|_{\mathbf{H}_{h,k-1}} &\leq \beta_V = \tilde{\mathcal{O}}(\sqrt{d}), \\
\|(\boldsymbol{\mu}_{h,k-1} - \boldsymbol{\mu}_h^*)\underline{\mathbf{V}}_{h+1}^k\|_{\mathbf{H}_{h,k-1}} &\leq \beta_V = \tilde{\mathcal{O}}(\sqrt{d}).
\end{aligned} \tag{B.1}$$

Conditional on the event that all inequalities in equation B.1 hold, we can easily verify that

$$\begin{aligned} |\langle \phi(\cdot, \cdot), \boldsymbol{\theta}_{h,k-1} \rangle - [\mathbb{P}_h R_h](\cdot, \cdot)| &\leq \beta_R \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}}, \\ |\langle \phi(\cdot, \cdot), \boldsymbol{\mu}_{h,k-1} \mathbf{V}_{h+1} \rangle - [\mathbb{P}_h V_{h+1}](\cdot, \cdot)| &\leq \beta_V \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}} \end{aligned}$$

for both  $\mathbf{V}_{h+1} \in \{\underline{\mathbf{V}}_{h+1}^k, \overline{\mathbf{V}}_{h+1}^k\}$ . Therefore, we define the temporary optimistic value function by

$$\hat{Q}_h^k(\cdot, \cdot) = \langle \phi(\cdot, \cdot), \boldsymbol{\theta}_{h,k-1} + \boldsymbol{\mu}_{h,k-1} \overline{\mathbf{V}}_{h+1}^k \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}},$$

and the temporary pessimistic value function by

$$\check{Q}_h^k(\cdot, \cdot) = \langle \phi(\cdot, \cdot), \boldsymbol{\theta}_{h,k-1} + \boldsymbol{\mu}_{h,k-1} \overline{\mathbf{V}}_{h+1}^k \rangle - \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}}$$

where  $\beta := \beta_R + \beta_V = \tilde{\mathcal{O}}(\sqrt{d})$ . The actual optimistic value function  $\overline{Q}_h^k(\cdot, \cdot)$  is the minimum function of history temporary optimistic value functions  $\hat{Q}_h^k(\cdot, \cdot)$ , and the actual pessimistic value function  $\underline{Q}_h^k(\cdot, \cdot)$  is the maximum function of history temporary pessimistic value functions  $\check{Q}_h^k(\cdot, \cdot)$  (Line 7 in Algorithm 3). In this way,  $\overline{Q}_h^k(\cdot, \cdot)$  is always non-increasing in  $k$  and  $\underline{Q}_h^k(\cdot, \cdot)$  is always non-decreasing in  $k$ .

The second is the rare-switching value function update, which updates the value function only when the determinant of the covariance matrix significantly exceeds the previous value (Line 6 in Algorithm 3). This approach allows the complexity, as measured by the metric entropy, of the function class to which  $\overline{\mathbf{V}}_h^k(\cdot)$  or  $\underline{\mathbf{V}}_h^k(\cdot)$  belongs to be independent of  $K$ . Notably, the metric entropy is linearly dependent on  $\tilde{\mathcal{O}}(dH)$ . Moreover, on the event equation B.1, we can establish optimism and pessimism in Lemma D.4, i.e., for all  $k \in [K]$  and  $h \in [H]$ ,

$$\underline{V}_{h+1}^k(\cdot) \leq V_{h+1}^*(\cdot) \leq \overline{V}_{h+1}^k(\cdot). \quad (\text{B.2})$$

Directly estimating the variance of the optimistic value function  $\overline{V}_{h+1}^k(\cdot)$  will encounter the dependence issue, which is discussed in (Jin et al., 2020b) and will introduce an additional  $\sqrt{d}$  factor in the regret due to the covering-based decoupling argument. To eliminate this factor, after noting the inequality

$$[\mathbb{V}_h \overline{V}_{h+1}^k](\cdot, \cdot) \leq 2[\mathbb{V}_h V_{h+1}^*](\cdot, \cdot) + 2[\mathbb{V}_h (\overline{V}_{h+1}^k - V_{h+1}^*)](\cdot, \cdot),$$

Hu et al. (2022) decompose the optimistic value function  $\overline{V}_{h+1}^k(\cdot)$  into the optimal value function  $V_{h+1}^*(\cdot)$  and the sub-optimality gap  $[\overline{V}_{h+1}^k - V_{h+1}^*](\cdot)$  and estimate their variances  $[\mathbb{V}_h V_{h+1}^*](\cdot, \cdot)$  and  $[\mathbb{V}_h (\overline{V}_{h+1}^k - V_{h+1}^*)](\cdot, \cdot)$  separately. The key insight is that: (i) as  $V_{h+1}^*$  is deterministic, there is no additional  $\sqrt{d}$  dependence in estimating  $[\mathbb{V}_h V_{h+1}^*](\cdot, \cdot)$ , and (ii) as  $\overline{V}_{h+1}^k$  gradually converges to  $V_{h+1}^*$ , though a uniform argument is still used, the incurred  $\sqrt{d}$  factor in the estimation of  $[\mathbb{V}_h (\overline{V}_{h+1}^k - V_{h+1}^*)]$  has ignorable effects on the final regret. We now describe the way we estimate these two variances.

- For  $[\mathbb{V}_h V_{h+1}^*](\cdot, \cdot)$ , since  $V_{h+1}^*$  is unknown, a natural choice is to estimate it by the optimistic value function  $\overline{V}_{h+1}^k$ . Hence, we estimate  $[\mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k})$  via

$$[\hat{\mathbb{V}}_{h,k} \overline{V}_{h+1}^k](s_{h,k}, a_{h,k}) := [\hat{\mathbb{P}}_{h,k}(\overline{V}_{h+1}^k)^2](s_{h,k}, a_{h,k})_{[0, \mathcal{H}^2]} - \left[ [\hat{\mathbb{P}}_{h,k} \overline{V}_{h+1}^k](s_{h,k}, a_{h,k})_{[0, \mathcal{H}^2]} \right]^2. \quad (\text{B.3})$$

To measure estimation accuracy, we introduce an error term  $U_{h,k}$  to guarantee that with high probability,  $|\hat{\mathbb{V}}_{h,k} \overline{V}_{h+1}^k(s_{h,k}, a_{h,k}) - [\mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k})| \leq U_{h,k}$  holds uniformly over all  $h, k$  where

$$U_{h,k} = \min \left\{ \mathcal{V}^2, 11\mathcal{H}\beta_0 \cdot \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 4\mathcal{H} \cdot \hat{\mathbb{P}}_{h,k}(\overline{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \right\} \quad (\text{A.9})$$

and  $\beta_0 = \tilde{\mathcal{O}}\left(\frac{\mathcal{H}}{\sigma_{\min}} \sqrt{d^3 H}\right)$  is an exploration radius.

- For  $[\mathbb{V}_h(\bar{V}_{h+1}^k - V_{h+1}^*)](\cdot, \cdot)$ , in order to meet the measurability condition of a concentration inequality (Lemma G.3), we require

$$\sigma_{h,k}^2 \geq d^3 H \cdot \sup_{k \leq j \leq K} [\mathbb{V}_h(\bar{V}_{h+1}^j - V_{h+1}^*)](s_{h,k}, a_{h,k}). \quad (\text{B.4})$$

Note that  $\sigma_{h,k}$  is  $\mathcal{F}_{h,k}$ -measurable while  $\bar{V}_{h+1}^k(\cdot)$  is  $\mathcal{F}_{H,k-1}$ -measurable. The condition equation B.4 essentially requires a  $\mathcal{F}_{h,k}$ -measurable upper bound for  $[\mathbb{V}_h(\bar{V}_{h+1}^j - V_{h+1}^*)](s_{h,k}, a_{h,k})$  even if  $j \geq k$ . Fortunately, we have for any  $k \leq j \leq K$ ,

$$\begin{aligned} [\mathbb{V}_h(\bar{V}_{h+1}^j - V_{h+1}^*)](s_{h,k}, a_{h,k}) &\leq [\mathbb{P}_h(\bar{V}_{h+1}^j - V_{h+1}^*)^2](s_{h,k}, a_{h,k}) \\ &\stackrel{(a)}{\leq} \mathcal{H}[\mathbb{P}_h(\bar{V}_{h+1}^j - V_{h+1}^*)](s_{h,k}, a_{h,k}) \\ &\stackrel{(b)}{\leq} \mathcal{H}[\mathbb{P}_h(\bar{V}_{h+1}^j - \underline{V}_{h+1}^j)](s_{h,k}, a_{h,k}) \\ &\stackrel{(c)}{\leq} \mathcal{H}[\mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)](s_{h,k}, a_{h,k}) \end{aligned} \quad (\text{B.5})$$

where (a) uses  $|\bar{V}_{h+1}^j - V_{h+1}^*|(\cdot) \leq \mathcal{H}$  and the optimism of  $\bar{V}_{h+1}^j(\cdot)$ , (b) follows from the pessimism in equation B.2, and (c) uses the monotonicity of value functions. The RHS of equation B.5 is  $\mathcal{F}_{h,k}$ -measurable but intractable due to the population expectation  $\mathbb{P}_h(\cdot)$ . By replacing  $\mathbb{P}_h(\cdot)$  with the tractable  $\hat{\mathbb{P}}_{h,k}(\cdot)$ , we introduce  $E_{h,k}$  to overestimate the RHS of equation B.5 where

$$E_{h,k} = \min \left\{ \mathcal{H}^2, 2\mathcal{H}\beta_0 \cdot \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \mathcal{H} \cdot \left[ \hat{\mathbb{P}}_{h,k}(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k) \right](s_{h,k}, a_{h,k}) \right\}. \quad (\text{A.7})$$

Hence, equation B.4 is guaranteed by  $\sigma_{h,k}^2 \geq d^3 H \cdot E_{h,k}$ . The extra  $d^3 H$  factor is introduced to offset the error caused by the covering number argument.

## C Proof for Section 2.3

### C.1 Proof of Theorem 2.2

Now, we turn to the regret equation 2.1. Recall that at iteration  $t$ , we set

$$(\phi_t, *) = \underset{\phi \in \mathcal{D}_t, \theta \in \mathcal{C}_{t-1}}{\operatorname{argmax}} \langle \phi, \theta \rangle.$$

Due to  $\sup_{\phi \in \cup_{t \geq 0} \mathcal{D}_t} |\langle \phi, \theta^* \rangle| \leq R := LB$ , it follows that

$$\begin{aligned} \operatorname{Reg}(T) &:= \sum_{t=1}^T \left[ \sup_{\phi \in \mathcal{D}_t} \langle \phi, \theta^* \rangle - \langle \phi_t, \theta^* \rangle \right] \\ &\leq \sum_{t=1}^T \left[ \sup_{\phi \in \mathcal{D}_t, \theta \in \mathcal{C}_{t-1}} \langle \phi, \theta \rangle - \langle \phi_t, \theta^* \rangle \right] \\ &= \sum_{t=1}^T \left[ \sup_{\theta \in \mathcal{C}_{t-1}} \langle \phi_t, \theta \rangle - \langle \phi_t, \theta^* \rangle \right] \\ &\leq \sum_{t=1}^T \|\phi_t\|_{\mathbf{H}_{t-1}^{-1}(\theta_{t-1})} \cdot \sup_{\theta \in \mathcal{C}_{t-1}} \|\theta - \theta^*\|_{\mathbf{H}_{t-1}} \end{aligned}$$

Notice that with probability  $1 - \delta$ ,  $\theta^* \in \mathcal{C}_t$  for all  $t \geq 1$ , i.e.,  $\|\theta_t - \theta^*\|_{\mathbf{H}_t} \leq \beta_t$ . Hence,

$$\sup_{\theta \in \mathcal{C}_t} \|\theta - \theta^*\|_{\mathbf{H}_t} \leq \sup_{\theta \in \mathcal{C}_t} \|\theta - \theta_t\|_{\mathbf{H}_t} + \|\theta_t - \theta^*\|_{\mathbf{H}_t} \leq 2\beta_t.$$

Notice that  $\beta_t$  is increasing in  $t$  and  $w_t = \left\| \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_{t-1}^{-1}}$ . Therefore,

$$\text{Reg}(T) \leq 2\beta_T \sum_{t=1}^T \|\phi_t\|_{\mathbf{H}_{t-1}^{-1}} = 2\beta_T \sum_{t=1}^T \sigma_t w_t = 2\beta_T \sum_{t=1}^T \sigma_t \min\{1, w_t\}. \quad (\text{C.1})$$

The last equality uses  $w_t \leq 1$  (which is due to  $\sigma_t \geq \|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}/c_0$  and  $c_0 \leq 1$ ). Notice that  $\|\phi_t\|/\sigma_t \leq \|\phi_t\|/\sigma_{\min} \leq L/\sigma_{\min}$ . Then by Lemma G.5,

$$\sum_{t=1}^T \min \left\{ 1, \left\| \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_{t-1}^{-1}}^2 \right\} = \sum_{t=1}^T \min \{1, w_t^2\} \leq 2d \log \left( 1 + \frac{TL^2}{d\lambda\sigma_{\min}^2} \right) = 2\kappa. \quad (\text{C.2})$$

Recall that

$$\sigma_t = \max \left\{ \nu_t, \sigma_{\min}, \frac{\|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}}{c_0}, \frac{\sqrt{LB}\|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}^{\frac{1}{2}}}{c_1^{\frac{1}{4}}d^{\frac{1}{4}}} \right\}.$$

According to what value  $\sigma_t$  takes, we decompose  $[T]$  into three sets  $[T] \subseteq \cup_{i=1}^3 \mathcal{J}_i$  where

$$\begin{aligned} \mathcal{J}_1 &= \{t \in [T] : \sigma_t \in \{\nu_t, \sigma_{\min}\}\}, \\ \mathcal{J}_2 &= \left\{ t \in [T] : \sigma_t = \frac{\|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}}{c_0} \right\}, \\ \mathcal{J}_3 &= \left\{ t \in [T] : \sigma_t = \sqrt{LB} \frac{\|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}^{\frac{1}{2}}}{c_1^{\frac{1}{4}}d^{\frac{1}{4}}} \right\}. \end{aligned}$$

First, it follows that

$$\begin{aligned} \sum_{t \in \mathcal{J}_1} \sigma_t \min\{1, w_t\} &\leq \sum_{t \in \mathcal{J}_1} \max\{\nu_t, \sigma_{\min}\} \min\{1, w_t\} \\ &\leq \sum_{t \in [T]} \max\{\nu_t, \sigma_{\min}\} \min\{1, w_t\} \\ &\stackrel{(a)}{\leq} \sqrt{\sum_{t \in [T]} (\nu_t^2 + \sigma_{\min}^2)} \sqrt{\sum_{t \in [T]} \min\{1, w_t^2\}} \\ &\stackrel{(b)}{\leq} \sqrt{2\kappa} \cdot \sqrt{\sum_{t \in [T]} \nu_t^2 + 1}. \end{aligned} \quad (\text{C.3})$$

Here (a) holds due to Cauchy-Schwarz inequality and (b) uses equation C.2 and  $\sigma_{\min} = \frac{1}{\sqrt{T}}$ .

Second, for any  $t \in \mathcal{J}_2$ , we have  $w_t = \left\| \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_{t-1}^{-1}} = c_0 \leq 1$ . Therefore,

$$\begin{aligned} \sum_{t \in \mathcal{J}_2} \sigma_t \min\{1, w_t\} &= \sum_{t \in \mathcal{J}_2} \sigma_t w_t = \frac{1}{c_0} \sum_{t \in \mathcal{J}_2} \sigma_t w_t^2 \leq \frac{\sup_{t \in \mathcal{J}_2} \sigma_t}{c_0} \sum_{t \in \mathcal{J}_2} w_t^2 \\ &\leq \frac{\sup_{t \in [T]} \|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}}{c_0^2} \cdot \sum_{t \in \mathcal{J}_2} \min\{1, w_t^2\} \\ &\leq \frac{\sup_{t \in [T]} \|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}}{c_0^2} \cdot \sum_{t \in [T]} \min\{1, w_t^2\} \leq \frac{2L\kappa}{c_0^2 \sqrt{\lambda}} \end{aligned} \quad (\text{C.4})$$

where the last inequality uses  $\|\phi_t\|_{\mathbf{H}_{t-1}^{-1}} \leq \frac{1}{\sqrt{\lambda}} \|\phi_t\| \leq \frac{L}{\sqrt{\lambda}}$  for all  $t \geq 1$  and equation C.2.

Finally, for any  $t \in \mathcal{J}_3$ , we have  $L^2 B^2 w_t^2 = c_1 d \sigma_t^2$  due to  $w_t^2 = \left\| \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_{t-1}^{-1}}^2$ . It implies  $\sigma_t = LB w_t / \sqrt{c_1 d} = LB \min\{1, w_t\} / \sqrt{c_1 d}$  with the fact that  $w_t \leq 1$ . Therefore,

$$\sum_{t \in \mathcal{J}_3} \sigma_t \min\{1, w_t\} = \frac{LB}{\sqrt{c_1 d}} \cdot \sum_{t \in \mathcal{J}_3} \min\{1, w_t^2\} \leq \frac{LB}{\sqrt{c_1 d}} \cdot \sum_{t \in [T]} \min\{1, w_t^2\} \leq \frac{2LB\kappa}{\sqrt{c_1 d}}. \quad (\text{C.5})$$

Plugging equation C.3, equation C.4 and equation C.5 into equation C.1, we have

$$\text{Reg}(T) \leq 2\beta_T \left[ \sqrt{2\kappa} \cdot \sqrt{\sum_{t \in [T]} \nu_t^2 + 1} + \frac{2L\kappa}{c_0^2 \sqrt{\lambda}} + \frac{2LB\kappa}{\sqrt{c_1 d}} \right].$$

## C.2 Proof of Lemma 2.1

Recall that  $z_t(\theta) = \frac{y_t - \langle \phi_t, \theta \rangle}{\sigma_t}$ . Direct computation yields that

$$\nabla^2 L_T(\theta) = \lambda \mathbf{I} + \sum_{t=1}^T \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta)}} \right)^3 \frac{\phi_t \phi_t^\top}{\sigma_t^2}.$$

Clearly, for any  $\theta \in \mathbb{R}^d$ ,

$$\nabla^2 L_T(\theta) \leq \lambda \mathbf{I} + \sum_{t=1}^T \frac{\phi_t \phi_t^\top}{\sigma_t^2} = \mathbf{H}_T.$$

For the other direction, we decompose it into four terms and analyze them respectively.

$$\begin{aligned} \nabla^2 L_T(\theta) &= \mathbf{H}_T - \underbrace{\sum_{t=1}^T \left[ 1 - \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta^*)}} \right)^3 \right] \frac{\phi_t \phi_t^\top}{\sigma_t^2}}_{\mathbf{H}_{1,T}} \\ &\quad + \underbrace{\sum_{t=1}^T \left[ \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta)}} \right)^3 - \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta^*)}} \right)^3 \right] \frac{\phi_t \phi_t^\top}{\sigma_t^2}}_{\mathbf{H}_{2,T}}. \end{aligned} \quad (\text{C.6})$$

where  $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | \mathcal{F}_{t-1}]$  for simplicity.

Since  $\nu_t, \theta_{t-1} \in \mathcal{F}_{t-1}$ , from Algorithm 1, we have  $\sigma_t, w_t, \tau_t \in \mathcal{F}_{t-1}$ .

**Analysis of  $\mathbf{H}_{1,T}$**  Notice that for any unit norm  $\mathbf{v} \in \mathbb{R}^d$ , it follows that

$$\begin{aligned} \mathbf{v}^\top \mathbf{H}_{1,T} \mathbf{v} &= \sum_{t=1}^T \left[ 1 - \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta^*)}} \right)^3 \right] \left\langle \frac{\phi_t}{\sigma_t}, \mathbf{v} \right\rangle^2 \\ &\leq 3 \sum_{t=1}^T \left[ 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta^*)}} \right] \left\langle \frac{\phi_t}{\sigma_t}, \mathbf{v} \right\rangle^2 \\ &\leq 3 \sum_{t=1}^T \left[ 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta^*)}} \right] \cdot \sup_{t \in [T]} \left\langle \frac{\phi_t}{\sigma_t}, \mathbf{v} \right\rangle^2 \\ &\leq 3 \sum_{t=1}^T \left[ 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\theta^*)}} \right] \cdot \sup_{t \in [T]} \left\| \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_T^{-1}}^2 \cdot \mathbf{v}^\top \mathbf{H}_T \mathbf{v} \end{aligned}$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} 3 \sum_{t=1}^T \left[ 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \right] \cdot \sup_{t \in [T]} \left\| \frac{\boldsymbol{\phi}_t}{\sigma_t} \right\|_{\mathbf{H}_t^{-1}}^2 \cdot \mathbf{v}^\top \mathbf{H}_T \mathbf{v} \\
&\stackrel{(b)}{=} 3 \sum_{t=1}^T \left[ 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \right] \cdot \sup_{t \in [T]} \frac{w_t^2}{1 + w_t^2} \cdot \mathbf{v}^\top \mathbf{H}_T \mathbf{v},
\end{aligned}$$

where (a) uses  $\mathbf{H}_T^{-1} \leq \mathbf{H}_t^{-1}$  for all  $t \in [T]$  and (b) follows from

$$\left\| \frac{\boldsymbol{\phi}_t}{\sigma_t} \right\|_{\mathbf{H}_t^{-1}}^2 = \frac{\boldsymbol{\phi}_t^\top}{\sigma_t} \left( \mathbf{H}_{t-1}^{-1} - \frac{\mathbf{H}_{t-1}^{-1} \frac{\boldsymbol{\phi}_t \boldsymbol{\phi}_t^\top}{\sigma_t} \mathbf{H}_{t-1}^{-1}}{1 + \frac{\boldsymbol{\phi}_t^\top \mathbf{H}_{t-1}^{-1} \boldsymbol{\phi}_t}{\sigma_t}} \right) \frac{\boldsymbol{\phi}_t}{\sigma_t} = w_t^2 - \frac{w_t^4}{1 + w_t^2} = \frac{w_t^2}{1 + w_t^2}.$$

By the arbitrariness of  $\mathbf{v}$ , we know that

$$\mathbf{H}_{1,T} \leq 3 \sum_{t=1}^T \left[ 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \right] \cdot \sup_{t \in [T]} \frac{w_t^2}{1 + w_t^2} \cdot \mathbf{H}_T. \quad (\text{C.7})$$

Let  $X_t = 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}}$ . It is obvious that  $0 \leq X_t \leq 1$ . We then focus on concentration of  $\sum_{t=1}^T X_t$ . To that end, we need a variance-aware Bernstein inequality Lemma G.2 for martingales. Lemma G.2 implies that with probability at least  $1 - \frac{\delta}{T^2}$ , we have

$$\sum_{t=1}^T X_t \leq \sum_{t=1}^T \mathbb{E}_t X_t + 3 \sqrt{\sum_{t=1}^T \text{Var}[X_t | \mathcal{F}_{t-1}] \cdot \log \frac{2KT^2}{\delta}} + 5 \log \frac{2KT^2}{\delta}$$

where  $K := 1 + \lceil 2 \log_2 V \rceil$  and  $V^2$  is an upper bound satisfying  $\sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] \leq V^2$ .

First notice that for any  $t \geq 1$ , we have

$$\begin{aligned}
\mathbb{E}_t X_t &= 1 - \mathbb{E}_t \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} = \mathbb{E}_t \frac{z_t^2(\boldsymbol{\theta}^*)}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}(\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)} + \tau_t)} \\
&\leq \frac{1}{2\tau_t^2} \mathbb{E}_t z_t^2(\boldsymbol{\theta}^*) \leq \frac{b^2}{2\tau_t^2} \leq \frac{b^2}{2\tau_0^2} \frac{w_t^2}{1 + w_t^2}
\end{aligned} \quad (\text{C.8})$$

which implies that

$$\sum_{t=1}^T \mathbb{E}_t X_t \leq \frac{b^2}{2\tau_0^2} \frac{w_t^2}{1 + w_t^2} \leq \frac{b^2}{2\tau_0^2} \sum_{t=1}^T \min\{1, w_t^2\} \leq \frac{\kappa b^2}{\tau_0^2}$$

where the last inequality uses Lemma G.5 and thus

$$\sum_{t=1}^T \min\{1, w_t^2\} = \sum_{t=1}^T \min \left\{ 1, \left\| \frac{\boldsymbol{\phi}_t}{\sigma_t} \right\|_{\mathbf{H}_{t-1}^{-1}}^2 \right\} \leq 2d \log \left( 1 + \frac{TL^2}{d\lambda\sigma_{\min}^2} \right) = 2\kappa.$$

Secondly, we have

$$\text{Var}[X_t | \mathcal{F}_{t-1}] \leq \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] \leq \mathbb{E}_t \left( 1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \right)^2 \stackrel{(*)}{\leq} \frac{1}{4} \frac{\mathbb{E}_t z_t^2(\boldsymbol{\theta}^*)}{\tau_t^2} \leq \frac{b^2}{4\tau_t^2}$$

where  $(*)$  uses  $1 - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \leq \frac{z_t^2(\boldsymbol{\theta}^*)}{2\tau_t \sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}}$  which is also used in equation C.8. As a result, we have

$$\sum_{t=1}^T \text{Var}[X_t | \mathcal{F}_{t-1}] \leq \sum_{t=1}^T \frac{b^2}{4\tau_t^2} \leq \frac{b^2}{4\tau_0^2} \sum_{t=1}^T \frac{w_t^2}{1 + w_t^2} \leq \frac{\kappa b^2}{2\tau_0^2}.$$



Once requiring  $\tau_0^2 \geq 2\kappa b^2$ , we have  $\sum_{t=1}^T \text{Var}[X_t | \mathcal{F}_{t-1}] \leq 1$  and thus we can set  $V = 1$  and obtain  $K = 1$ . Putting them together, if  $\tau_0^2 \geq \frac{2\kappa b^2}{\log \frac{2T^2}{\delta}}$ , with probability at least  $1 - \delta$ , we have

$$\begin{aligned} \sum_{t=1}^T X_t &\leq \frac{\kappa b^2}{\tau_0^2} + \frac{3b}{\tau_0} \sqrt{\frac{\kappa \log \frac{2T^2}{\delta}}{2}} + 5 \log \frac{2T^2}{\delta} \\ &\leq \frac{1}{2} \log \frac{2T^2}{\delta} + \frac{3}{2} \log \frac{2T^2}{\delta} + 5 \log \frac{2T^2}{\delta} \\ &\leq 9 \log \frac{2T^2}{\delta} = \frac{1}{12c_0^2} \end{aligned} \quad (\text{C.9})$$

where the last equation is due to the definition of  $c_0$ . Finally, taking a union bound for the last inequality from  $T = 1$  to  $\infty$  and using the fact that  $\sum_{t=1}^{\infty} t^{-2} < 2$ , we have  $\sum_{t=1}^T X_t \leq \frac{1}{12c_0^2}$  for all  $T \geq 1$  with probability at least  $1 - 2\delta$ .

On the other hand, by the choice of  $\sigma_t$ , we have  $\sigma_t^2 \geq \frac{1}{c_0^2} \cdot \|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}^2$ , which implies

$$\sup_{t \in [T]} \frac{w_t^2}{1 + w_t^2} \leq \sup_{t \in [T]} w_t^2 \leq c_0^2. \quad (\text{C.10})$$

Plugging equation C.9 and equation C.10 into equation C.7, we have

$$\mathbf{H}_{1,T} \leq \frac{1}{4} \mathbf{H}_T. \quad (\text{C.11})$$

**Analysis of  $\mathbf{H}_{2,T}$**  We first notice that

$$\begin{aligned} \left| \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta})}} \right)^3 - \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \right)^3 \right| &\leq 3 \left| \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta})}} - \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \right| \\ &\leq \frac{3\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta})} \sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \frac{|z_t^2(\boldsymbol{\theta}) - z_t^2(\boldsymbol{\theta}^*)|}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta})} + \sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}}. \end{aligned} \quad (\text{C.12})$$

Notice that  $z_t(\boldsymbol{\theta}) = z_t(\boldsymbol{\theta}^*) + \langle \frac{\phi_t}{\sigma_t}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle$ . It then follow that for any  $c > 0$

$$\begin{aligned} z_t^2(\boldsymbol{\theta}) &\leq \left(1 + \frac{1}{c}\right) z_t^2(\boldsymbol{\theta}^*) + (1+c) \left\langle \frac{\phi_t}{\sigma_t}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \right\rangle^2; \\ z_t^2(\boldsymbol{\theta}^*) &\leq \left(1 + \frac{1}{c}\right) z_t^2(\boldsymbol{\theta}) + (1+c) \left\langle \frac{\phi_t}{\sigma_t}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \right\rangle^2, \end{aligned}$$

By discussing which is larger between  $z_t^2(\boldsymbol{\theta})$  and  $z_t^2(\boldsymbol{\theta}^*)$ , we have

$$|z_t^2(\boldsymbol{\theta}) - z_t^2(\boldsymbol{\theta}^*)| \leq \frac{1}{c} \min \{z_t^2(\boldsymbol{\theta}), z_t^2(\boldsymbol{\theta}^*)\} + (1+c) \left\langle \frac{\phi_t}{\sigma_t}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \right\rangle^2. \quad (\text{C.13})$$

Plugging equation C.13 into equation C.12, we have that

$$\begin{aligned} \left| \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta})}} \right)^3 - \left( \frac{\tau_t}{\sqrt{\tau_t^2 + z_t^2(\boldsymbol{\theta}^*)}} \right)^3 \right| &\leq \frac{3\tau_t}{\tau_t^2 + \min \{z_t^2(\boldsymbol{\theta}), z_t^2(\boldsymbol{\theta}^*)\}} \frac{\frac{1}{c} \min \{z_t^2(\boldsymbol{\theta}), z_t^2(\boldsymbol{\theta}^*)\}}{2\sqrt{\tau_t^2 + \min \{z_t^2(\boldsymbol{\theta}), z_t^2(\boldsymbol{\theta}^*)\}}} + \frac{3(1+c)}{2\tau_t^2} \left\langle \frac{\phi_t}{\sigma_t}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \right\rangle^2 \\ &\leq \frac{3}{2c} + \frac{3(1+c)}{2\tau_t^2} \left\langle \frac{\phi_t}{\sigma_t}, \boldsymbol{\theta} - \boldsymbol{\theta}^* \right\rangle^2 \stackrel{(a)}{\leq} \frac{3}{2c} + \frac{6(1+c)}{\tau_t^2} \frac{L^2 B^2}{\sigma_t^2} \end{aligned}$$

$$\leq \frac{3}{2c} + \frac{6(1+c)}{\tau_0^2} \frac{w_t^2 L^2 B^2}{\sigma_t^2} \stackrel{(b)}{\leq} \frac{3}{2c} + \frac{6(1+c)c_1 d}{\tau_0^2} \quad (\text{C.14})$$

where (a) uses  $\left\langle \frac{\phi_t}{\sigma_t}, \theta - \theta^* \right\rangle \leq \left\| \frac{\phi_t}{\sigma_t} \right\| (\|\theta\| + \|\theta^*\|) \leq \frac{2LB}{\sigma_t}$  due to  $\|\phi_t\| \leq L$  and  $\theta, \theta^* \in \text{Ball}_d(B)$  and (b) uses the following result. By the definition of  $\sigma_t$ , we have  $\sigma_t \geq \sqrt{LB} \|\phi_t\|_{\mathbf{H}_{t-1}^{-1}}^{\frac{1}{2}} / c_1^{\frac{1}{4}} d^{\frac{1}{4}}$  which implies  $\sigma_t^2 \geq \frac{w_t^2 L^2 B^2}{c_1 d}$ . As a result of equation C.14, by definition of  $\mathbf{H}_{3,T}$ , we have

$$- \left( \frac{3}{2c} + \frac{6(1+c)c_1 d}{\tau_0^2} \right) \sum_{t=1}^T \frac{\phi_t \phi_t^\top}{\sigma_t \sigma_t} \leq \mathbf{H}_{2,T} \quad (\text{C.15})$$

**Putting pieces together** Plugging equation C.11 and equation C.15 into equation C.6, with probability at least  $1 - \delta$ , for any  $T \geq 1$  and for all  $\theta \in \text{Ball}_d(B)$ , we have

$$\begin{aligned} \nabla^2 L_T(\theta) &\geq \mathbf{H}_T - \frac{1}{4} \mathbf{H}_T - \left( \frac{3}{2c} + \frac{6(1+c)c_1 d}{\tau_0^2} \right) \sum_{t=1}^T \frac{\phi_t \phi_t^\top}{\sigma_t \sigma_t} \\ &\geq \frac{3\lambda}{4} \mathbf{I} + \left( 1 - \frac{1}{4} - \frac{3}{2c} - \frac{6(1+c)c_1 d}{\tau_0^2} \right) \sum_{t=1}^T \frac{\phi_t \phi_t^\top}{\sigma_t \sigma_t}. \end{aligned}$$

Notice that  $c_1 = \frac{1}{42 \cdot \log \frac{2T^2}{\delta}}$ . If we set  $c = 6$  and  $\tau_0 \sqrt{\frac{2T^2}{\delta}} \geq \max\{\sqrt{2\kappa}b, 2\sqrt{d}\}$ , we have

$$\max \left\{ \frac{3}{2c}, \frac{6(1+c)c_1 d}{\tau_0^2} \right\} \leq \frac{1}{4}.$$

As a result, we have

$$\nabla^2 L_T(\theta) \geq \frac{3\lambda}{4} \mathbf{I} + \frac{1}{4} \sum_{t=1}^T \frac{\phi \phi^\top}{\sigma_t \sigma_t} \geq \frac{1}{4} \mathbf{H}_T.$$

### C.3 Proof of Lemma 2.2

For simplicity, we denote  $z_t^* = z_t(\theta^*)$  for short. By triangle inequality, we have

$$\begin{aligned} \|\nabla L_T(\theta^*)\|_{\mathbf{H}_T^{-1}} &\leq \|\lambda \theta^*\|_{\mathbf{H}_T^{-1}} + \left\| \sum_{t=1}^T \frac{\tau_t z_t^*}{\sqrt{\tau_t^2 + (z_t^*)^2}} \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_T^{-1}} \\ &\leq \|\lambda \theta^*\|_{\mathbf{H}_T^{-1}} + \underbrace{\left\| \sum_{t=1}^T \frac{\tau_t z_t^*}{\sqrt{\tau_t^2 + (z_t^*)^2}} \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_T^{-1}}}_{:= \mathbf{d}_T}. \end{aligned} \quad (\text{C.16})$$

**For the residual term  $\|\lambda \theta^*\|_{\mathbf{H}_T^{-1}}$**  Notice that  $\mathbf{H}_T \geq \lambda \mathbf{I}$  and thus  $\mathbf{H}_T^{-1} \leq \lambda^{-1} \mathbf{I}_d$ . Therefore,  $\|\lambda \theta^*\|_{\mathbf{H}_T^{-1}} \leq \sqrt{\lambda} B$ .

**For the self-normalized term  $\|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}}$**  The fact that  $\mathbf{H}_T = \mathbf{H}_{T-1} + \frac{\phi_T \phi_T^\top}{\sigma_T^2}$  together with the Woodbury matrix identity implies that

$$\mathbf{H}_T^{-1} = \mathbf{H}_{T-1}^{-1} - \frac{\mathbf{H}_{T-1}^{-1} \phi_T \phi_T^\top \mathbf{H}_{T-1}^{-1}}{\sigma_T^2 (1 + w_T^2)} \quad \text{where} \quad w_T^2 := \frac{\phi_T^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T^2} = \left\| \frac{\phi_T}{\sigma_T} \right\|_{\mathbf{H}_{T-1}}^2. \quad (\text{C.17})$$

Clearly,  $w_T$  is  $\mathcal{F}_{T-1}$ -measurable and thus is predictable. By definition of  $\mathbf{d}_T$  and equation C.17,

$$\|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}}^2 = \left( \mathbf{d}_{T-1} + \frac{\tau_T z_T^*}{\sqrt{\tau_T^2 + (z_T^*)^2}} \frac{\phi_T}{\sigma_T} \right)^\top \mathbf{H}_T^{-1} \left( \mathbf{d}_{T-1} + \frac{\tau_T z_T^*}{\sqrt{\tau_T^2 + (z_T^*)^2}} \frac{\phi_T}{\sigma_T} \right)$$

$$\begin{aligned}
&= \|\mathbf{d}_{T-1}\|_{\mathbf{H}_{T-1}^{-1}}^2 - \frac{1}{1+w_T^2} \left( \frac{\mathbf{d}_{T-1}^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T} \right)^2 \\
&\quad + \frac{2\tau_T z_T^*}{\sqrt{\tau_T^2 + (z_T^*)^2}} \frac{\mathbf{d}_{T-1}^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T} + \frac{\tau_T^2 (z_T^*)^2}{\tau_T^2 + (z_T^*)^2} \frac{\phi_T^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T^2} \\
&\leq \|\mathbf{d}_{T-1}\|_{\mathbf{H}_{T-1}^{-1}}^2 + \underbrace{\frac{2\tau_T z_T^*}{\sqrt{\tau_T^2 + (z_T^*)^2}} \frac{\mathbf{d}_{T-1}^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T}}_{I_1} + \underbrace{\frac{\tau_T^2 (z_T^*)^2}{\tau_T^2 + (z_T^*)^2} \frac{\phi_T^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T^2}}_{I_2}. \tag{C.18}
\end{aligned}$$

For  $I_1$ , by equation C.17, we have

$$\begin{aligned}
I_1 &= \frac{2\tau_T z_T^*}{\sqrt{\tau_T^2 + (z_T^*)^2}} \frac{1}{\sigma_T} \mathbf{d}_{T-1}^\top \left( \mathbf{H}_{T-1}^{-1} - \frac{\mathbf{H}_{T-1}^{-1} \phi_T \phi_T^\top \mathbf{H}_{T-1}^{-1}}{\sigma_T^2 (1+w_T^2)} \right) \phi_T \\
&= \frac{2\tau_T z_T^*}{\sqrt{\tau_T^2 + (z_T^*)^2}} \frac{1}{1+w_T^2} \frac{\mathbf{d}_{T-1}^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T}.
\end{aligned}$$

For  $I_2$ , we have

$$\begin{aligned}
I_2 &= \frac{\tau_T^2 (z_T^*)^2}{\tau_T^2 + (z_T^*)^2} \frac{\phi_T^\top \mathbf{H}_{T-1}^{-1} \phi_T}{\sigma_T^2} \\
&= \frac{\tau_T^2 (z_T^*)^2}{\tau_T^2 + (z_T^*)^2} \frac{1}{\sigma_T^2} \phi_T^\top \left( \mathbf{H}_{T-1}^{-1} - \frac{\mathbf{H}_{T-1}^{-1} \phi_T \phi_T^\top \mathbf{H}_{T-1}^{-1}}{\sigma_T^2 (1+w_T^2)} \right) \phi_T \\
&= \frac{\tau_T^2 (z_T^*)^2}{\tau_T^2 + (z_T^*)^2} \left( w_T^2 - \frac{w_T^4}{1+w_T^2} \right) \\
&= \frac{\tau_T^2 (z_T^*)^2}{\tau_T^2 + (z_T^*)^2} \frac{w_T^2}{1+w_T^2}.
\end{aligned}$$

Using the equations for  $I_1, I_2$  and iterating equation C.18, we have

$$\|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}}^2 \leq \sum_{t=1}^T \frac{\tau_t z_t^*}{\sqrt{\tau_t^2 + (z_t^*)^2}} \frac{2}{1+w_t^2} \frac{\mathbf{d}_{t-1}^\top \mathbf{H}_{t-1}^{-1} \phi_t}{\sigma_t} + \sum_{t=1}^T \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \frac{w_t^2}{1+w_t^2}. \tag{C.19}$$

Recall that

$$\kappa = d \log \left( 1 + \frac{TL^2}{d\lambda\sigma_{\min}^2} \right).$$

**Lemma C.1.** Assume  $\mathbb{E}[(z_t^*)^2 | \mathcal{F}_{t-1}] \leq b^2$  for all  $t \geq 1$ . Let  $A_t$  denotes the event where  $\|\mathbf{d}_n\|_{\mathbf{H}_n^{-1}} \leq \alpha_n$  for all  $n \in [t]$ . With probability at least  $1 - \delta$ , we have for all  $T \geq 1$ ,

$$\sum_{t=1}^T \frac{2\tau_t z_t^* \mathbf{1}_{A_{t-1}}}{(\tau_t^2 + (z_t^*)^2)^{1/2}} \frac{1}{1+w_t^2} \frac{\mathbf{d}_{t-1}^\top \mathbf{H}_{t-1}^{-1} \phi_t}{\sigma_t} \leq 4 \max_{t \in [T]} \alpha_t \cdot \left[ \frac{\kappa b^2}{4\tau_0} + b \sqrt{\kappa \log \frac{2T^2}{\delta}} + \frac{2\tau_0}{3} \log \frac{2T^2}{\delta} \right].$$

**Lemma C.2.** Assume  $\mathbb{E}[(z_t^*)^2 | \mathcal{F}_{t-1}] \leq b^2$  for all  $t \geq 1$ . For a fixed  $\tau \geq 0$ , with probability at least  $1 - \delta$ , the follow inequality uniformly holds for all  $T \geq 1$ ,

$$\sum_{t=1}^T \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \frac{w_t^2}{1+w_t^2} \leq \left[ \sqrt{2\kappa b} + \tau_0 \sqrt{\log \frac{2T^2}{\delta}} \right]^2.$$

For any  $T \geq 1$ , we define

$$\alpha_T = 8 \left[ \frac{\kappa b^2}{\tau_0} + b \sqrt{\kappa \log \frac{2T^2}{\delta}} + \tau_0 \log \frac{2T^2}{\delta} \right]. \tag{C.20}$$

As a result of Lemma C.1 and Lemma C.2, with probability at least  $1 - 2\delta$ , for all  $T \geq 0$ ,<sup>4</sup>

$$\sum_{t=1}^T \frac{2\tau_t z_t^* \mathbf{1}_{A_{t-1}}}{\sqrt{\tau_t^2 + (z_t^*)^2}} \frac{1}{1+w_t^2} \frac{\mathbf{d}_{t-1}^\top \mathbf{H}_{t-1}^{-1} \phi_t}{\sigma_t} \leq \frac{\alpha_T^2}{2} \quad \text{and} \quad \sum_{t=1}^T \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \frac{w_t^2}{1+w_t^2} \leq \frac{\alpha_T^2}{2}. \quad (\text{C.21})$$

Let  $B$  denote the event that the conditions in equation C.21 hold for  $T \geq 0$ . By Lemma C.1 and Lemma C.2, we know that  $\mathbb{P}(B) \geq 1 - 2\delta$ . We now introduce a new event  $C$  that is defined by

$$C := \left\{ \|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}} \leq \alpha_T, \text{ for all } T \geq 0 \right\} = \cap_{t=0}^\infty A_t.$$

In the following, we will show that  $B \subseteq C$  by mathematical induction. As a result, it follows that

$$\mathbb{P}(C) \geq \mathbb{P}(B) \geq 1 - \delta.$$

Finally, we use mathematical induction to show that if  $B$  is true, then  $C$  must be true, i.e., all  $A_t$  is true for all  $t \geq 0$  on the condition that the last inequalities equation C.21 are valid for all  $T \geq 0$ . When  $t = 0$ ,  $A_0$  is true by definition. Suppose that at iteration  $T - 1$ , for all  $0 \leq t \leq T - 1$ , the event  $A_t$  is true, then we are going to show that  $A_T$  is also true. By comparing the definition of  $A_T$  and  $A_{T-1}$ , we only need to show that  $\|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}} \leq \alpha_T$  which is equivalent to  $\|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}}^2 \leq \alpha_T^2$ . It follows due to the following inequality

$$\begin{aligned} \|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}}^2 &\stackrel{(\text{C.19})}{\leq} \sum_{t=1}^T \frac{\tau_t z_t^*}{\sqrt{\tau_t^2 + (z_t^*)^2}} \frac{2}{1+w_t^2} \frac{\mathbf{d}_{t-1}^\top \mathbf{H}_{t-1}^{-1} \phi_t}{\sigma_t} + \sum_{t=1}^T \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \frac{w_t^2}{1+w_t^2} \\ &\stackrel{(a)}{=} \sum_{t=1}^T \frac{\tau_t z_t^* \mathbf{1}_{A_{t-1}}}{\sqrt{\tau_t^2 + (z_t^*)^2}} \frac{2}{1+w_t^2} \frac{\mathbf{d}_{t-1}^\top \mathbf{H}_{t-1}^{-1} \phi_t}{\sigma_t} + \sum_{t=1}^T \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \frac{w_t^2}{1+w_t^2} \\ &\stackrel{(b)}{\leq} \frac{\alpha_T^2}{2} + \frac{\alpha_T^2}{2} = \alpha_T^2, \end{aligned}$$

where (a) uses the condition that all  $A_t$  is true for all  $0 \leq t \leq T - 1$  and (b) uses the conditions equation C.21.

As a result, we can conclude that all  $\{A_t\}_{t \geq 0}$  is true and thus  $\|\mathbf{d}_T\|_{\mathbf{H}_T^{-1}} \leq \alpha_T$  for all  $T \geq 1$ .

#### C.4 Proof of Lemma C.1

*Proof of Lemma C.1.* We will make use of the Freedman inequality Lemma G.1 to prove our result. Recall that  $\tau_t = \tau_0 \frac{\sqrt{1+w_t^2}}{w_t}$ . Set  $Y_t = \frac{\tau_t z_t^*}{\sqrt{\tau_t^2 + (z_t^*)^2}} \frac{2}{1+w_t^2} \frac{\mathbf{d}_{t-1}^\top \mathbf{H}_{t-1}^{-1} \phi_t \mathbf{1}_{A_{t-1}}}{\sigma_t}$  with the event  $A_{t-1}$  defined in the lemma. For simplicity, we denote  $X_t = Y_t - \mathbb{E}[Y_t | \mathcal{F}_{t-1}]$ . Notice that

$$\left| \frac{\mathbf{d}_{t-1}^\top \mathbf{H}_{t-1}^{-1} \phi_t}{\sigma_t} \cdot \mathbf{1}_{A_{t-1}} \right| \leq \|\mathbf{d}_{t-1} \mathbf{1}_{A_{t-1}}\|_{\mathbf{H}_{t-1}^{-1}} \cdot \left\| \frac{\phi_t}{\sigma_t} \right\|_{\mathbf{H}_{t-1}^{-1}} \leq \alpha_{t-1} w_t.$$

As a result, we have

$$|Y_t| \leq \tau_t \alpha_{t-1} \cdot \frac{2w_t}{1+w_t^2} \leq 2\tau_0 \alpha_{t-1} \quad \text{and thus} \quad |X_t| \leq |Y_t| + |\mathbb{E}[Y_t | \mathcal{F}_{t-1}]| \leq 4\tau_0 \alpha_{t-1}.$$

We also find that

$$\begin{aligned} \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] &\stackrel{(a)}{\leq} \mathbb{E}[Y_t^2 | \mathcal{F}_{t-1}] = \mathbb{E} \left[ \left( \frac{2w_t}{1+w_t^2} \right)^2 \|\mathbf{d}_{t-1}\|_{\mathbf{H}_{t-1}^{-1}}^2 \mathbf{1}_{A_{t-1}} \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \middle| \mathcal{F}_{t-1} \right] \\ &\stackrel{(b)}{\leq} \left( \frac{2w_t}{1+w_t^2} \right)^2 \alpha_{t-1}^2 b^2 \leq \min\{1, 2w_t\}^2 \alpha_{t-1}^2 b^2 \leq 4 \min\{1, w_t^2\} \alpha_{t-1}^2 b^2 \end{aligned}$$

<sup>4</sup>Note that it's easy to verify that the following inequalities are true when  $t = 0$ .

where (a) uses  $\mathbb{E}(X - \mathbb{E}X)^2 \leq \mathbb{E}X^2$  for any random variable  $X$  and (b) uses  $\mathbb{E}[\varepsilon_t^2 | \mathcal{F}_{t-1}] \leq b^2 \sigma_t^2$  due to  $\mathbb{E}[(z_t^*)^2 | \mathcal{F}_{t-1}] \leq b^2$ .

Notice that  $\|\phi_t\|/\sigma_t \leq \|\phi_t\|/\sigma_{\min} \leq L/\sigma_{\min}$ . Then by Lemma G.5, we have

$$\sum_{t=1}^T \min\{1, w_t^2\} \leq 2d \log \left( 1 + \frac{TL^2}{d\lambda\sigma_{\min}^2} \right) := 2\kappa. \quad (\text{C.22})$$

Hence, by equation C.22,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] &\leq 4 \sum_{t=1}^T \min\{1, w_t^2\} \alpha_{t-1}^2 b^2 \leq 4 \max_{t \in [T]} \alpha_t^2 \cdot \sum_{t=1}^T \min\{1, w_t^2\} b^2 \\ &\leq \max_{t \in [T]} \alpha_t^2 \cdot 8db^2 \log \left( 1 + \frac{TL^2}{d\lambda\sigma_{\min}^2} \right) \leq 8\kappa b^2 \cdot \max_{t \in [T]} \alpha_t^2. \end{aligned}$$

On the other hand, using  $\mathbb{E}[z_t^* | \mathcal{F}_{t-1}] = 0$  we have

$$\left| \mathbb{E} \left[ \frac{\tau_t z_t^*}{\sqrt{\tau_t^2 + (z_t^*)^2}} \middle| \mathcal{F}_{t-1} \right] \right| = \left| \mathbb{E} \left[ \left( \frac{\tau_t}{\sqrt{\tau_t^2 + (z_t^*)^2}} - 1 \right) z_t^* \middle| \mathcal{F}_{t-1} \right] \right| \leq \mathbb{E} \left[ \frac{(z_t^*)^2}{2\tau_t} \middle| \mathcal{F}_{t-1} \right] \leq \frac{b^2}{2\tau_t}$$

which implies

$$\begin{aligned} \left| \sum_{t=1}^T \mathbb{E}[Y_t | \mathcal{F}_{t-1}] \right| &\leq \sum_{t=1}^T \frac{b^2}{2\tau_t} \frac{w_t}{1 + w_t^2} \alpha_{t-1} \leq \frac{b^2}{2\tau_0} \sum_{t=1}^T \frac{w_t^2}{1 + w_t^2} \alpha_{t-1} \\ &\leq \sup_{t \in [T]} \alpha_t \cdot \frac{b^2}{2\tau} \sum_{t=1}^T \min\{1, w_t^2\} \leq \sup_{t \in [T]} \alpha_t \cdot \frac{\kappa b^2}{\tau_0}. \end{aligned}$$

By Freedman inequality in Lemma G.1, it follows that for a given  $T$  and  $\tau_0$ , with probability  $1 - \frac{\delta}{2T^2}$ ,

$$\begin{aligned} \sum_{t=1}^T Y_t &\leq \left| \sum_{t=1}^T \mathbb{E}[Y_t | \mathcal{F}_{t-1}] \right| + 4 \max_{t \in [T]} \alpha_t \cdot \left[ b \sqrt{\kappa \log \frac{2T^2}{\delta}} + \frac{2\tau_0}{3} \log \frac{2T^2}{\delta} \right] \\ &\leq 4 \max_{t \in [T]} \alpha_t \cdot \left[ \frac{\kappa b^2}{4\tau_0} + b \sqrt{\kappa \log \frac{2T^2}{\delta}} + \frac{2\tau_0}{3} \log \frac{2T^2}{\delta} \right]. \end{aligned}$$

Finally, taking a union bound for the last inequality from  $T = 1$  to  $\infty$  and using the fact that  $\sum_{t=1}^{\infty} t^{-2} < 2$  complete the proof.  $\square$

## C.5 Proof of Lemma C.2

*Proof of Lemma C.2.* Set  $Y_t = \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \frac{w_t^2}{1 + w_t^2}$  and  $X_t = Y_t - \mathbb{E}[Y_t | \mathcal{F}_{t-1}]$ . Recall that  $\tau_t = \tau_0 \frac{\sqrt{1 + w_t^2}}{w_t}$ . Clearly, we have  $|Y_t| \leq \tau_t^2 \frac{w_t^2}{1 + w_t^2} \leq \tau_0^2$  and thus  $|X_t| = |Y_t - \mathbb{E}[Y_t | \mathcal{F}_{t-1}]| \leq \max\{|Y_t|, |\mathbb{E}[Y_t | \mathcal{F}_{t-1}]|\} \leq \tau_0^2$ . We also find that

$$\begin{aligned} \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] &\stackrel{(a)}{\leq} \mathbb{E}[Y_t^2 | \mathcal{F}_{t-1}] \leq \left( \frac{w_t^2}{1 + w_t^2} \right)^2 \mathbb{E} \left[ \left( \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \right)^2 \middle| \mathcal{F}_{t-1} \right] \\ &\leq \tau_t^2 \left( \frac{w_t^2}{1 + w_t^2} \right)^2 \mathbb{E}[(z_t^*)^2 | \mathcal{F}_{t-1}] \stackrel{(b)}{\leq} \tau_0^2 b^2 \frac{w_t^2}{1 + w_t^2} \end{aligned}$$

where (a) uses  $\mathbb{E}(X - \mathbb{E}X)^2 \leq \mathbb{E}X^2$  for any random variable  $X$  and (b) uses  $\mathbb{E}[(z_t^*)^2 | \mathcal{F}_{t-1}] \leq b^2$  due to  $\mathbb{E}[\varepsilon_t^2 | \mathcal{F}_{t-1}] \leq b^2 \nu_t^2$ . Hence, by equation C.22, we have

$$\sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] \leq \tau_0^2 b^2 \sum_{t=1}^T \frac{w_t^2}{1 + w_t^2} \leq \tau_0^2 b^2 \sum_{t=1}^T \min\{1, w_t^2\} \leq 2\kappa \tau_0^2 b^2.$$

On the other hand,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[Y_t | \mathcal{F}_{t-1}] &= \sum_{t=1}^T \frac{w_t^2}{1+w_t^2} \mathbb{E} \left[ \frac{\tau_t^2 (z_t^*)^2}{\tau_t^2 + (z_t^*)^2} \middle| \mathcal{F}_{t-1} \right] \\ &\leq \sum_{t=1}^T \frac{w_t^2}{1+w_t^2} \mathbb{E} [(z_t^*)^2 | \mathcal{F}_{t-1}] \leq \sum_{t=1}^T \min\{1, w_t^2\} b^2 \leq 2\kappa b^2. \end{aligned}$$

By Lemma G.1, it follows that with probability  $1 - \frac{\delta}{2T^2}$ ,

$$\sum_{t=1}^T Y_t \leq \sum_{t=1}^T \mathbb{E}[Y_t | \mathcal{F}_{t-1}] + 2\tau_0 b \sqrt{\kappa \log \frac{2T^2}{\delta}} + \frac{2\tau_0^2}{3} \log \frac{2T^2}{\delta}$$

for a given  $T$  and  $\tau_0$ . Putting all pieces together, it follows that with probability  $1 - \frac{\delta}{2T^2}$ ,

$$\sum_{t=1}^T Y_t \leq \left[ \sqrt{2\kappa} b + \tau_0 \sqrt{\log \frac{2T^2}{\delta}} \right]^2.$$

Finally, taking a union bound for the last inequality from  $T = 1$  to  $\infty$  and using the fact that  $\sum_{t=1}^{\infty} t^{-2} < 2$  complete the proof.  $\square$

## D Proof of Theorem 3.1

**Measurability** Let  $\mathcal{F}_{h,k}$  denote the  $\sigma$ -field generated by all random variables up to and including the  $h$ -th step and  $k$ -th episode. More specifically, let  $I_{h,k} = \{(i, j) : i \in [H], j \in [k-1] \text{ or } i \in [h], j = k\}$  denote the set of index pairs up to and including the  $h$ -th step and  $k$ -th episode and then  $\mathcal{F}_{h,k} = \sigma(\cup_{(i,j) \in I_{h,k}} \{s_{i,j}, a_{i,j}, r_{i,j}\})$ . We make a convention that  $\mathcal{F}_{0,k} = \mathcal{F}_{H,k-1}$ . From our algorithm, we know that (i)  $Q_h^k, V_h^k, \pi_h^k \in \mathcal{F}_{H,k-1}$  for any  $Q \in \{\bar{Q}, \hat{Q}, \check{Q}, \underline{Q}\}$  and  $V \in \{\bar{V}, \hat{V}, \check{V}, \underline{V}\}$ , and (ii)

$$\mu_{h-1,k}, \theta_{h,k}, \psi_{h,k}, \sigma_{h,k}, U_{h,k}, J_{h,k}, E_{h,k}, \phi_{h,k}, \tilde{\phi}_{h,k}, w_{h,k}, \tilde{w}_{h,k}, \tau_{h,k}, \tilde{\tau}_{h,k}, \mathbf{H}_{h,k}, \tilde{\mathbf{H}}_{h,k} \in \mathcal{F}_{h,k}.$$

### D.1 High-Probability Events

Let  $\kappa = d \log \left( 1 + \frac{K}{d\lambda\sigma_{\min}^2} \right)$ . We first introduce the following high-probability events.

1. We define  $\mathcal{B}_{R^2}$  as the event that the following inequalities hold for all  $h \in [H]$  and  $k \in [K] \cup \{0\}$ ,

$$\psi_h^* \in \tilde{\mathcal{R}}_{h,k} := \left\{ \|\psi\| \leq W : \|\psi_{h,k} - \psi\|_{\tilde{\mathbf{H}}_{h,k}^{-1}} \leq \beta_{R^2} \right\}$$

where

$$\beta_{R^2} = 128 \left( \frac{\sqrt{\kappa}\sigma_{R^2}}{\sigma_{\min}} + \sqrt{d} \right) \sqrt{\log \frac{2HK^2}{\delta}} + 5\sqrt{\lambda}W.$$

2. We define  $\mathcal{B}_0$  as the event that the following inequalities hold for all  $h \in [H]$  and  $k \in [K]$ ,

$$\begin{aligned} \max \left\{ \left\| (\mu_h^* - \mu_{h,k-1}) \bar{\mathbf{V}}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}}, \left\| (\mu_h^* - \mu_{h,k-1}) \underline{\mathbf{V}}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}} \right\} &\leq \beta_0, \\ \left\| (\mu_h^* - \mu_{h,k-1}) [\bar{\mathbf{V}}_{h+1}^k]^2 \right\|_{\mathbf{H}_{h,k-1}} &\leq \mathcal{H}\beta_0 \end{aligned}$$



where

$$\beta_0 = \frac{4\mathcal{H}}{\sigma_{\min}} \sqrt{d^3 H \iota_0^2 + \log \frac{2H}{\delta}} + 3\sqrt{d\lambda}\mathcal{H}$$

$$\iota_0 = \max \left\{ \log \left( 1 + \frac{8LK}{\lambda\mathcal{H}\sqrt{d}\sigma_{\min}^2} \right), \log \left( 1 + \frac{32B^2K^2}{\sqrt{d}\lambda^3\mathcal{H}^2\sigma_{\min}^4} \right), \log \left( 1 + \frac{K}{\lambda\sigma_{\min}^2} \right) \right\}. \quad (\text{D.1})$$

Here we choose  $B \geq 3(\beta_R + \beta_V)$  and  $L = W + \mathcal{H}\sqrt{\frac{dK}{\lambda}}$ .

3. We define  $\mathcal{B}_R$  as the event that the following inequalities hold for all  $h \in [H]$  and  $k \in [K] \cup \{0\}$ ,

$$\begin{aligned} \boldsymbol{\theta}_h^* \in \mathcal{R}_{h,k} &:= \left\{ \|\boldsymbol{\theta}\| \leq W : \|\boldsymbol{\theta}_{h,k} - \boldsymbol{\theta}_h^*\|_{\mathbf{H}_{h,k}} \leq \beta_R \right\}, \\ \left| [\hat{\mathbb{V}}_h \hat{R}_h - \mathbb{V}_h R_h](s_{h,k}, a_{h,k}) \right| &\leq R_{h,k} := \beta_{R^2} \|\tilde{\boldsymbol{\phi}}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{k-1}} + 2\mathcal{H}\beta_R \|\boldsymbol{\phi}_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \end{aligned} \quad (\text{A.8})$$

where

$$\beta_R = 128(\sqrt{\kappa} + \sqrt{d}) \sqrt{\log \frac{2HK^2}{\delta}} + 5\sqrt{\lambda}W.$$

4. We define  $\mathcal{B}_h$  as the event such that for all episode  $k \in [K]$ , all stages  $h \leq h' \leq H$ ,

$$\max \left\{ \left\| (\boldsymbol{\mu}_{h'}^* - \boldsymbol{\mu}_{h',k-1}) \bar{\mathbf{V}}_{h'+1}^k \right\|_{\mathbf{H}_{h',k-1}}, \left\| (\boldsymbol{\mu}_{h'}^* - \boldsymbol{\mu}_{h',k-1}) \underline{\mathbf{V}}_{h'+1}^k \right\|_{\mathbf{H}_{h',k-1}} \right\} \leq \beta_V \quad (\text{D.2})$$

where

$$\beta_V = \mathcal{O} \left( \sqrt{d\iota_1^2} + \sqrt{d\lambda}\mathcal{H} \right)$$

$$\iota_1 = \max \left\{ \iota_0, \log \frac{4HK^2}{\delta}, \log \left( 1 + \frac{4L\sqrt{d^3H}}{\sigma_{\min}} \right), \log \left( 1 + \frac{8\sqrt{d^7}HB^2}{\lambda\sigma_{\min}^2} \right) \right\}. \quad (\text{D.3})$$

For simplicity, we further define  $\mathcal{B}_V := \mathcal{B}_1$ .

Our ultimate goal is to show  $\mathcal{B}_V$  holds with high probability, a target used in previous work (Hu et al., 2022; He et al., 2022). More specifically, we first obtain coarse confidence sets for all parameters in the sense that the confidence radius (that is  $\beta_{R^2}$  and  $\beta_0$ ) is loose. In our analysis,  $\mathcal{B}_{R^2} \cap \mathcal{B}_0$  serves as the ‘coarse’ event where the concentration results hold with a larger confidence radius, and  $\mathcal{B}_R \cap \mathcal{B}_V$  serves as a ‘refined’ event where the confidence radius (that is  $\beta_R$  and  $\beta_V$ ) is much tighter. Our first result is that  $\mathcal{B}_{R^2} \cap \mathcal{B}_0$  holds with high probability as shown in Lemma D.1 and D.2. Their proofs are collected in Appendix F.1 and F.2.

**Lemma D.1.** If we set

$$\tilde{\tau}_0 = \max \left\{ \frac{\sqrt{2\kappa}\sigma_{R^2}}{\sigma_{\min}}, 2\sqrt{d} \right\} / \sqrt{\log \frac{2HK^2}{\delta}}, \quad (\text{D.4})$$

the event  $\mathcal{B}_{R^2}$  holds with probability at least  $1 - 4\delta$ .

**Lemma D.2.** The event  $\mathcal{B}_0$  holds with probability at least  $1 - 3\delta$ .

These coarse confidence sets are then used to estimate variance for the reward functions and value functions. A key step is to show the adapted variance  $\sigma_{h,k}$ ’s are indeed upper bounds of these variances (that is  $[\mathbb{V}_h R_h](s_{h,k}, a_{h,k}) + [\mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k})$ ) for all  $h \in [H]$ . A frequently used argument is backward induction. That is given the estimation is optimistic at the stage  $h+1$ , we then show the optimistic estimation is maintained at the stage  $h$ . Induction over the stage  $h$  would complete the proof. The following lemma provides estimation error bounds for  $[\mathbb{V}_h R_h](s_{h,k}, a_{h,k})$  and shows that the event  $\mathcal{B}_R$  holds with high probability. Its proof is deferred in Appendix F.4.

**Lemma D.3.** If we set

$$\tau_0 = \max\{\sqrt{2\kappa}, 2\sqrt{d}\} / \sqrt{\log \frac{2HK^2}{\delta}}, \quad (\text{D.5})$$

the event  $\mathcal{B}_R$  holds with probability at least  $1 - 8\delta$ .

In Lemma D.4, we show that that our constructed value functions  $\bar{V}$  and  $\underline{V}$  are indeed optimistic and pessimistic estimators of the true value functions under the event defined before. Its proof is deferred in Appendix F.5.

**Lemma D.4** (Optimism and pessimism). For any  $h \in [H]$ , if  $\mathcal{B}_R \cap \mathcal{B}_h$  holds, for any  $k \in [K] \cup \{0\}$ ,

$$\underline{V}_h^k(\cdot) \leq V_h^*(\cdot) \leq \bar{V}_h^k(\cdot).$$

With the established optimism and pessimism, we can establish upper bounds for the estimation errors of the three terms, namely  $[\mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k})$ ,  $[\mathbb{V}_h(\bar{V}_{h+1}^k - V_{h+1}^*)](s_{h,k}, a_{h,k})$ , and  $[\mathbb{V}_h(\underline{V}_{h+1}^k - V_{h+1}^*)](s_{h,k}, a_{h,k})$  in the following lemmas. Their proofs are deferred in Appendix F.6 and F.7.

**Lemma D.5.** On the event  $\mathcal{B}_0 \cap \mathcal{B}_{h+1}$ , it follows that for all  $k \in [K]$

$$\left| [\mathbb{V}_h V_{h+1}^* - \hat{\mathbb{V}}_h \bar{V}_{h+1}^k](s_{h,k}, a_{h,k}) \right| \leq U_{h,k}$$

where

$$U_{h,k} = \min \left\{ \mathcal{V}^2, 11\mathcal{H}\beta_0 \cdot \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 4\mathcal{H} \cdot \hat{\mathbb{P}}_{h,k}(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \right\} \quad (\text{A.9})$$

with  $\hat{\mathbb{P}}_{h,k}(\cdot|s, a) = \boldsymbol{\mu}_{h,k-1}^\top \phi(s, a)$ .

**Lemma D.6.** On the event  $\mathcal{B}_0 \cap \mathcal{B}_R \cap \mathcal{B}_{h+1}$ , it follows that for all  $j \leq k \leq K$

$$\max \left\{ [\mathbb{V}_h(\bar{V}_{h+1}^k - V_{h+1}^*)](s_{h,j}, a_{h,j}), [\mathbb{V}_h(\underline{V}_{h+1}^k - V_{h+1}^*)](s_{h,j}, a_{h,j}) \right\} \leq E_{h,j}$$

where

$$E_{h,j} = \min \left\{ \mathcal{H}^2, 2\mathcal{H}\beta_0 \|\phi_{h,j}\|_{\mathbf{H}_{h,j-1}^{-1}} + \mathcal{H} \cdot [\hat{\mathbb{P}}_{h,j}(\bar{V}_{h+1}^j - \underline{V}_{h+1}^j)](s_{h,j}, a_{h,j}) \right\} \quad (\text{A.7})$$

with  $\hat{\mathbb{P}}_{h,j}(\cdot|s, a) = \boldsymbol{\mu}_{h,j-1}^\top \phi(s, a)$ .

With the last four lemmas, one can easily prove  $\sigma_{h,k}$  indeed serves as an upper bound of the true variance of  $V_{h+1}^*$  at stage  $h$ . Therefore, by the backward induction, we can prove the following lemma whose proof is in Appendix F.8.

**Lemma D.7.** On the event  $\mathcal{B}_0 \cap \mathcal{B}_R$ , the event  $\mathcal{B}_V$  holds with probability at least  $1 - 2\delta$ .

## D.2 Regret Analysis

In the previous subsection, we know that with probability at least  $1 - 17\delta$ , the event  $\mathcal{B}_V \cap \mathcal{B}_R$  holds. Based on Lemma D.4, the optimism implies that

$$\text{Reg}(K) := \sum_{k=1}^K (V_1^* - V_1^{\pi_k})(s_{1,k}) \leq \sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k}).$$

We then relate the suboptimality gap  $\sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k})$  to the term  $\sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}$  in Lemma D.8. We emphasize that the bound in Lemma D.8 is much finer than previous bounds (e.g., Lemma B.1 in (He et al., 2022)) in the sense that the rest term is  $\tilde{\mathcal{O}}(H\mathcal{H})$  instead of previous  $\tilde{\mathcal{O}}(\sqrt{HK}\mathcal{H})$ . This is because we adapt a variance-aware Bernstein inequality to relate the variance of  $\sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k})$  with its expectations and use a recursion argument to simplify the final expression, while previous work directly apply Azuma-Hoeffding inequality to analyze the concentration of  $\sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k})$ , which inevitably introduces the additional  $\tilde{\mathcal{O}}(\sqrt{K})$  dependence. Its proof is collected in Appendix F.9.

**Lemma D.8** (Suboptimality gap). With probability at least  $1 - \delta$ , on the event  $\mathcal{B}_R \cap \mathcal{B}_V$ , it follows that

$$\begin{aligned} \sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k}) &\leq 6\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H\mathcal{H} \log \frac{4\lceil \log_2 HK \rceil}{\delta} \text{ and} \\ \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) &\leq 8H\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H^2\mathcal{H} \log \frac{4\lceil \log_2 HK \rceil}{\delta}. \end{aligned}$$

Using a similar argument, we provide a finer bound for the gap between optimistic and pessimistic value functions  $\sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k})$  in Lemma D.9. Its proof is provided in Appendix F.10.

**Lemma D.9** (Gap between optimistic and pessimistic value functions). With probability at least  $1 - \delta$ , on the event  $\mathcal{B}_V \cap \mathcal{B}_R$ , it follows that

$$\sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \leq 12H\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H^2\mathcal{H} \log \frac{4\lceil \log_2 HK \rceil}{\delta}.$$

The following issue is to upper bound the term  $\sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}$ . Since the estimation of reward variance concerns the other term  $\sum_{k=1}^K \sum_{h=1}^H \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}$ , we are motivated to analyze them simultaneously via  $\sum_{k=1}^K \sum_{h=1}^H b_{h,k}$  where  $b_{h,k} = \max\left\{\|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}\right\}$ . Previous works (Hu et al., 2022; He et al., 2022) mainly use Cauchy-Schwarz inequality to analyze it and obtain

$$\sum_{k=1}^K \sum_{h=1}^H b_{h,k} \leq \sqrt{\left(\sum_{k=1}^K \sum_{h=1}^H \sigma_{h,k}^2\right) \left(\sum_{k=1}^K \sum_{h=1}^H \max\{w_{h,k}^2, \tilde{w}_{h,k}^2\}\right)} = \tilde{\mathcal{O}}\left(\sqrt{dH} \cdot \sqrt{\sum_{k=1}^K \sum_{h=1}^H \sigma_{h,k}^2}\right).$$

where the last equality uses the elliptical potential lemmas in Lemma G.5. A standard analysis of the law of total variation would imply  $\sqrt{\sum_{k=1}^K \sum_{h=1}^H \sigma_{h,k}^2} = \tilde{\mathcal{O}}(\sqrt{H^2 K})$ . However, this result doesn't satisfy our target for two reasons. First, due to the use of adaptive Huber regression, our definition of  $\sigma_{h,k}$  is more complicated than previous algorithms. We need a more elaborate analysis to handle the additional terms in the definition of  $\sigma_{h,k}$ 's. Second, the previous result considers the worst-case scenario, while our target is to provide a finer variance-aware regret. Therefore, it is imperative to provide a finer bound for the sum of bonuses  $\sum_{k=1}^K \sum_{h=1}^H b_{h,k}$ . We did it in Lemma D.10.

**Lemma D.10** (Sum of bonuses). Set  $\lambda = \frac{1}{\mathcal{H}^2 + W^2}$ . Let  $\mathcal{A}_0$  denote the intersection event of Lemma D.8 and D.9. With probability at least  $1 - 2\delta$ , on the event  $\mathcal{B}_R \cap \mathcal{B}_V \cap \mathcal{B}_0 \cap \mathcal{B}_{R^2} \cap \mathcal{A}_0$ , we have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H b_{h,k} &= \tilde{\mathcal{O}}\left(\sqrt{dHK\mathcal{G}^*} + Hd^{0.5}K^{0.5}\sigma_{\min} + \frac{H^{2.5}d^{5.5}\mathcal{H}^2 + Hd^{1.5}\sigma_{R^2}}{\sigma_{\min}}\right) \\ &\quad + \tilde{\mathcal{O}}\left(H^3d^{4.5}\mathcal{H} + Hd^{0.5}\sigma_R + Hd^{1.5}\right). \end{aligned}$$

where  $\tilde{\mathcal{O}}(\cdot)$  ignores constant factors and logarithmic dependence.

We emphasize that Lemma D.10 is perhaps the most technical lemma in our paper. To address the difficulty mentioned early, we divide the full index set  $\mathcal{I} := [H] \times [K]$  into three disjoint subsets  $\mathcal{I} = \cup_{i=1,2,3} \mathcal{J}_i$  according to which value  $\sigma_{h,k}$  takes (given  $\sigma_{h,k}$  is the maximum value among five quantities). For those indexes in  $\mathcal{J}_1$  where the bonuses are small enough, we still use the Cauchy-Schwarz inequality to bound  $\sum_{(h,k) \in \mathcal{J}_1} b_{h,k} \leq \tilde{\mathcal{O}}\left(\sqrt{dH} \cdot \sqrt{\sum_{(h,k) \in \mathcal{I}} \sigma_{h,k}^2}\right)$ . This sum-of-squared-bonus quantity involves  $\sum_{(h,k) \in \mathcal{I}} E_{h,k}$  and  $\sum_{(h,k) \in \mathcal{I}} J_{h,k}$  which we then pay additional efforts to analyze. For those indexes in  $\mathcal{J}_2$  or  $\mathcal{J}_3$  where the bonuses are relatively large, we directly analyze  $\sum_{(h,k) \in \mathcal{J}_2 \cup \mathcal{J}_3} b_{h,k}$ . Thanks to the particular structure,  $\sum_{(h,k) \in \mathcal{J}_2 \cup \mathcal{J}_3} b_{h,k}$  contributes to the non-leading term in the final bound. Putting pieces together, we complete the proof. A formal proof can be found in Appendix F.11.

At the end of the subsection, we summarize the proof in a few lines.

$$\begin{aligned}
\text{Reg}(K) &= \sum_{k=1}^K (V_1^* - V_1^{\pi_k})(s_{1,k}) \stackrel{(a)}{\leq} \sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k}) \\
&\stackrel{(b)}{\leq} 3\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H\mathcal{H} \log \frac{4\lceil \log_2 HK \rceil}{\delta} \\
&\stackrel{(c)}{\leq} 3\beta \sum_{k=1}^K \sum_{h=1}^H b_{h,k} + 38H\mathcal{H} \log \frac{4\lceil \log_2 HK \rceil}{\delta} \\
&\stackrel{(d)}{=} \tilde{\mathcal{O}} \left( d\sqrt{HK\mathcal{G}^*} + HdK^{0.5}\sigma_{\min} + \frac{H^{2.5}d^6\mathcal{H}^2 + Hd^2\sigma_{R^2}}{\sigma_{\min}} + H^3d^5\mathcal{H} + Hd\sigma_R + Hd^2 \right)
\end{aligned}$$

where (a) follows from the optimism result in Lemma D.4, (b) follows from the suboptimality gap result in Lemma D.8, (c) uses  $b_{h,k} = \max \left\{ \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}} \right\}$ , and (d) follows from sum-of-bonus result in Lemma D.10 and  $\beta = \beta_R + \beta_V = \tilde{\mathcal{O}}(\sqrt{d})$ .

## E Proof of Theorem 3.2

*Proof of Theorem 3.2.* We consider the two complexities respectively.

**Space Complexity** First, in order to perform AdaOFUL, VARA needs to store all seen rewards and feature vectors (i.e.,  $\phi_{h,k}, \tilde{\phi}_{h,k}$ ), which is required by all RL/bandit algorithms robust to heavy-tailed rewards (Shao et al., 2018; Xue et al., 2021; Zhuang & Sui, 2021). AdaOFUL also keeps all robustification parameters  $\tau_{h,k}, \tilde{\tau}_{h,k}$ . It then incurs  $\mathcal{O}(HKd)$  space storage in total.

Second, due to the rare-switching technique, one can show that  $\bar{Q}_h^k$  (or  $\underline{Q}_h^k$ ) is the minimum (or maximum) of at most  $\tilde{\mathcal{O}}(dH)$  temporary optimistic (or pessimistic) functions (see Lemma G.7). It means that we need to store at most  $\tilde{\mathcal{O}}(dH)$  different versions of  $\theta_{h,k-1}, \mu_{h,k-1} \mathbf{V}_{h+1}^k, \mathbf{H}_{h,k-1}$ 's. This incurs  $\mathcal{O}(d^3H^2)$  space cost.

Last, for all  $(h, k) \in [H] \times [K]$ , we need to trace  $\{\phi(s_{h,k}, a)\}_{a \in \mathcal{A}}$  to evaluate each  $\mu_{h,k} \mathbf{V} = \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} V(s_{h+1,k})$  for  $\mathbf{V} \in \{\bar{\mathbf{V}}_{h+1}^k, [\bar{\mathbf{V}}_{h+1}^k]^2, \underline{\mathbf{V}}_{h+1}^k\}$ , which takes  $\mathcal{O}(d|\mathcal{A}|HK)$  space.

To sum up, VARA takes  $\mathcal{O}(d^3H^2 + d|\mathcal{A}|HK)$  space.

**Computational Complexity** First, we use Nesterov accelerated method to compute each  $\theta_{h,k}$ . Since the loss function in equation A.1 is  $\lambda$ -strongly convex and  $\left(\lambda + \frac{K}{\sigma_{\min}^2}\right)$ -smooth, the computational cost for each  $\theta_{h,k}$  is  $\tilde{\mathcal{O}}\left(d\sqrt{1 + \frac{K}{\lambda(\sigma_{\min}^*)^2}}\right) = \tilde{\mathcal{O}}(\max\{d, H^{-3/4}d^{-3/2}K^{3/4}\})$  and the total cost is  $\tilde{\mathcal{O}}(HK(d + H^{-3/4}d^{-3/2}K^{3/4}))$ .

We emphasize that we don't need to compute  $\theta_{h,k}$  exactly. It suffices to terminate at a solution  $\hat{\theta}_{h,k}$  once its accuracy satisfies  $\|\hat{\theta}_{h,k} - \theta_{h,k}\|_{\mathbf{H}_{h,k}} \leq \sqrt{d}$ . The iteration complexity is proportional to the root of the conditional number, i.e.,  $\tilde{\mathcal{O}}(\max\{1, d^{-7/4}K^{3/4}\})$ . Since each iteration takes  $\mathcal{O}(d)$  operation, the computation complexity is  $\tilde{\mathcal{O}}(\max\{d, d^{-3/4}K^{3/4}\})$ .

Second, each time when updating the value function, we take the minimum over at most  $\tilde{\mathcal{O}}(dH)$  quadratic functions. Moreover, the Sherman-Morrison formula computes  $\mathbf{H}_{h,k}^{-1}$  and its products with any vectors, which takes  $\mathcal{O}(d^2)$  operations. As a result, it needs  $\tilde{\mathcal{O}}(d^3H)$  to evaluate the updated  $Q_{h,k}(s, a)$  for a given pair  $(s, a)$ . Hence, computing  $Q_{h,k}(s_{h,k}, \cdot)$ , choosing  $a_{h,k} = \arg\max_{a \in \mathcal{A}} Q_{h,k}(s_{h,k}, a)$ , and estimating the variance  $\sigma_{h,k}$  lead to  $\tilde{\mathcal{O}}(d^3H^2|\mathcal{A}|)$  computational complexity for each episode.

Last, note  $\mu_{h,k} \mathbf{V} = \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} V(s_{h+1,k})$  for any value function  $V(\cdot)$ . If  $V$  remains unchanged, we only need to compute the new term  $\sigma_{h,k}^{-2} \phi_{h,k} V(s_{h+1,k})$ , which has an  $\tilde{\mathcal{O}}(d^3H|\mathcal{A}|)$  complexity each time. If  $V$

changes to  $V'$ , we need to recalculate  $\boldsymbol{\mu}_{h,k} \mathbf{V}'$ , which has an  $\tilde{O}(d^3 H |\mathcal{A}| K)$  complexity each time. Combining the computational complexity for all horizons and noticing that the number of episodes that trigger the updating criterion is at most  $\tilde{O}(dH)$ , VARA has a running time of  $\tilde{O}(d^4 |\mathcal{A}| H^3 K + HK(d + H^{-3/4} d^{-3/2} K^{3/4}))$ . In terms of the dependence on  $K$ , it is slightly worse than LSVI-UCB++'s  $\tilde{O}(d^4 |\mathcal{A}| H^3 K)$  since the adaptive Huber regression doesn't have a closed-form solution, but is better than LSVI-UCB's  $\tilde{O}(d^2 |\mathcal{A}| HK^2)$  due to the rare-switching mechanism.  $\square$

## F Omitted lemmas in Section D

### F.1 Proof of Lemma D.1

*Proof of Lemma D.1.* The proof idea of Lemma D.1 is similar to that of Theorem 2.1 except for the following changes. First,  $\tilde{\boldsymbol{\phi}}_{h,k} = \tilde{\boldsymbol{\phi}}(s_{h,k}, a_{h,k}) \in \mathbb{R}^d$  is instead the feature vector. Second, in the particular setting, we should respectively replace  $L, B, T, \delta$  therein with  $1, W, K, \delta/H$  defined here and redefine  $c_0, c_1$  as  $c_0 = \frac{1}{6\sqrt{3 \log \frac{2HK^2}{\delta}}}$ ,  $c_1 = \frac{1}{42 \cdot \frac{2HK^2}{\delta}}$  respectively. Third, by the choice of  $\sigma_{h,k}$ , we have  $\sigma_{h,k}^2 \geq \left( \frac{W}{\sqrt{c_1 d}} + \mathcal{H} d^{2.5} H \right) b_{h,k} \geq \frac{W}{\sqrt{c_1 d}} \|\tilde{\boldsymbol{\phi}}_{h,k}\|_{\tilde{\mathbf{H}}_{h-1,k}^{-1}}$ , which implies that  $\frac{W^2 \tilde{w}_{h,k}^2}{\sigma_{h,k}^2} \leq c_1 d$ . Similarly, due to  $\sigma_{h,k}^2 \geq c_0^{-2} \|\tilde{\boldsymbol{\phi}}_{h,k}\|_{\tilde{\mathbf{H}}_{h-1,k}^{-1}}^2$ , we have  $\tilde{w}_{h,k}^2 \leq c_0^2$ . Last, for simplicity, we define  $\varepsilon_{h,k} = \frac{r_{h,k}^2 - \langle \tilde{\boldsymbol{\phi}}_{h,k}, \boldsymbol{\psi}_h^* \rangle}{\sigma_{h,k}}$  and  $\mathcal{G}_{h,k} = \sigma(\mathcal{F}_{h-1,k} \cup \{s_{h,k}, a_{h,k}\})$ . Then, we have  $\varepsilon_{h,k} \in \mathcal{F}_{h,k}$ ,  $\mathbb{E}[\varepsilon_{h,k} | \mathcal{G}_{h,k}] = 0$  and  $\text{Var}[\varepsilon_{h,k} | \mathcal{G}_{h,k}] \leq \left( \frac{\sigma_{R^2}}{\sigma_{\min}} \right)^2 := b^2$ . Theorem 2.1 concerns the case where  $b = 1$ , however, its proof considers the general case where  $b$  can be arbitrary. As a result, by a similar argument in Appendix C (which is doable due to the four conditions mentioned above), once setting  $\tilde{\tau}_0 \sqrt{\log \frac{2HK^2}{\delta}} = \max \left\{ \sqrt{2\kappa b}, 2\sqrt{d} \right\}$ , with probability at least  $1 - 3\delta$ , we have for all  $h \in [H]$  and  $k \in [K]$ ,  $\|\boldsymbol{\psi}_{h,k} - \boldsymbol{\psi}_h^*\|_{\tilde{\mathbf{H}}_{h,k}} \leq \beta_{R^2}$ , that is the event  $\mathcal{B}_{R^2}$  holds.  $\square$

### F.2 Proof of Lemma D.2

We will make use of the following general result frequently. The proof is quite standard (Jin et al., 2020b; Wagenmaker et al., 2022a; Hu et al., 2022). We provide proof in Appendix F.3 for completeness.

**Lemma F.1.** Fix any  $h \in [H]$ . Consider a specific value function  $f(\cdot)$  which satisfies

- (i)  $\sup_{s \in \mathcal{S}} |f(s)| \leq C_0$ ;
- (ii)  $f \in \mathcal{V}$  where  $\mathcal{V}$  is a class of functions with  $\mathcal{N}(\mathcal{V}, \varepsilon)$  the  $\varepsilon$ -covering number of  $\mathcal{V}$  with respect to the distance  $\text{dist}(f, f') := \sup_{s \in \mathcal{S}} |f(s) - f'(s)|$ .

We assume there exists a deterministic  $C_\sigma > 0$  and  $\mathcal{A}_{h,k}$  (which is  $\mathcal{F}_{h,k}$ -measurable) such that  $\mathcal{A}_{h,k} \subseteq \left\{ \sigma_{h,k}^2 \geq (\mathbb{V}_h f)(s_{h,k}, a_{h,k}) / C_\sigma^2 \right\}$  for all  $k \in [K]$ . Let  $\boldsymbol{\mu}_{h,k}$  be defined equation A.2 and  $\sigma_{h,k}, \mathbf{H}_{h,k}$  be defined in our algorithm. Under any of the following conditions, with probability at least  $1 - \delta/H$ , it follows for all  $k \in [K] \cup \{0\}$ ,

$$\boldsymbol{\mu}_h^* \in \left\{ \boldsymbol{\mu} : \|(\boldsymbol{\mu} - \boldsymbol{\mu}_{h,k}) \mathbf{f}\|_{\mathbf{H}_{h,k}} \leq \beta \right\}. \quad (\text{F.1})$$

- (i) If  $f(\cdot)$  is a deterministic function and  $\cap_{k \in [K]} \mathcal{A}_{h,k}$  is true, equation F.1 holds with

$$\beta = 8C_\sigma \sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right) \log \frac{4HK^2}{\delta}} + \frac{8C_0}{d^{2.5} H} \log \frac{4HK^2}{\delta} + \sqrt{d \lambda} C_0.$$

- (ii) If  $f(\cdot)$  is a random function and  $\cap_{k \in [K]} \mathcal{A}_{h,k}$  is true, equation F.1 holds with

$$\beta = 8C_\sigma \sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right) \log \frac{4HK^2 N_0}{\delta}} + \frac{8C_0}{d^{2.5} H \mathcal{H}} \log \frac{4HK^2 N_0}{\delta} + 3\sqrt{d \lambda} C_0$$

where  $N_0 = |\mathcal{N}(\mathcal{V}, \varepsilon_0)|$  and  $\varepsilon_0 = \min \left\{ C_\sigma \sigma_{\min}, \frac{\lambda C_0 \sqrt{d}}{K} \sigma_{\min}^2 \right\}$ .

(iii) If  $f(\cdot)$  is a random function, equation F.1 holds with

$$\beta = \frac{2C_0}{\sigma_{\min}} \sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right) + \log \frac{N_1}{\delta}} + 3\sqrt{d\lambda} C_0.$$

where  $N_1 = |\mathcal{N}(\mathcal{V}, \varepsilon_1)|$  and  $\varepsilon_1 = \frac{\lambda C_0 \sqrt{d}}{K} \sigma_{\min}^2$ .

Using the last item suffices to prove Lemma D.2.

*Proof of Lemma D.2.* Let  $\mathcal{V}^+$  denote the class of optimistic value functions mapping from  $\mathcal{S}$  to  $\mathbb{R}$  with the parametric form given in equation G.1 and  $\mathcal{V}^-$  the class of pessimistic value functions with the parametric form given in equation G.2. By Lemma G.8 and Lemma G.7,

$$\log \mathcal{N}(\mathcal{V}^\pm, \varepsilon) \leq \left[ d \log \left( 1 + \frac{4L}{\varepsilon} \right) + d^2 \log \left( 1 + \frac{8d^{1/2} B^2}{\lambda \varepsilon^2} \right) \right] \quad (\text{F.2})$$

where  $B \geq \beta_0$  and  $L = W + \mathcal{H} \sqrt{\frac{dK}{\lambda}}$ .

- (i) Let  $\mathbf{f} = \bar{\mathbf{V}}_{h+1}^k$ . One can find that  $\mathbf{f} \in \mathcal{V}_f^+$  with parameter  $L = W + \frac{K\mathcal{H}}{\lambda \sigma_{\min}^2}$ . To plug in Lemma F.1, we first specify the parameters defined therein. We have  $\|\mathbf{f}\|_\infty \leq C_0 = \mathcal{H}$  and  $\varepsilon_1 = \frac{\lambda \mathcal{H} \sqrt{d}}{K} \sigma_{\min}^2$ . By equation F.2, it follows that

$$\begin{aligned} & \log \mathcal{N}(\mathcal{V}^+, \varepsilon_1) \\ & \leq \left[ d \log \left( 1 + \frac{4LK}{\lambda \mathcal{H} \sqrt{d} \sigma_{\min}^2} \right) + d^2 \log \left( 1 + \frac{8B^2 K^2}{\sqrt{d} \lambda^3 \mathcal{H}^2 \sigma_{\min}^4} \right) \right] \cdot dH \log_2 \left( 1 + \frac{K}{\lambda \sigma_{\min}^2} \right) \\ & \leq \frac{2}{\log 2} d^3 H \iota_0^2 \leq 3d^3 H \iota_0^2, \end{aligned}$$

By the third condition of Lemma F.1, with probability at least  $1 - \frac{\delta}{2H}$ ,  $\|(\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \hat{\mathbf{V}}_{h+1}^k\| \leq \beta_0$  for all  $k \in [K]$ . Similarly, we can also show that with probability at least  $1 - \frac{\delta}{2H}$ ,  $\|(\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \check{\mathbf{V}}_{h+1}^k\| \leq \beta_0$  for all  $k \in [K]$ . Putting them together finishes the proof.

- (ii) The analysis on  $\mathbf{V}_{h+1}^k$  is similar to (i).

- (iii) The analysis on  $[\bar{\mathbf{V}}_{h+1}^k]^2$  is similar to (i) except for the following two changes. First,  $C_0 = \mathcal{H}^2$  and  $\varepsilon'_1 = \frac{\lambda \mathcal{H}^2 \sqrt{d}}{K} \sigma_{\min}^2$ . Second, with  $[\mathcal{V}^+]^2 = \{f^2 : f \in \mathcal{V}^+\}$ , we have  $[\bar{\mathbf{V}}_{h+1}^k]^2 \in [\mathcal{V}^+]^2$  and

$$\log \mathcal{N}([\mathcal{V}^+]^2, \varepsilon'_1) \stackrel{(a)}{\leq} \log \mathcal{N}(\mathcal{V}^+, \frac{\varepsilon'_1}{2\mathcal{H}}) \leq \log \mathcal{N}(\mathcal{V}^+, \frac{\varepsilon_1}{2}) \leq 3d^3 H \iota_0^2.$$

Here (a) uses the fact that the  $\frac{\varepsilon'_1}{2\mathcal{H}}$ -cover of  $\mathcal{V}^+$  is a  $\varepsilon_1$ -cover of  $[\mathcal{V}^+]^2$  (which is also supported by Lemma G.9).

□

### F.3 Proof of Lemma F.1

*Proof of Lemma F.1.* Since the case of  $k = 0$  is trivial, we focus on  $k \in [K]$ . By definition,

$$\boldsymbol{\mu}_{h,k} = \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} \delta(s_{h+1,j})^\top = \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} (\phi_{h,j}^\top \boldsymbol{\mu}_h^* - \varepsilon_{h,j})^\top$$

$$= \boldsymbol{\mu}_h^* - \lambda \mathbf{H}_{h,k}^{-1} \boldsymbol{\mu}_h^* - \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top.$$

By the triangle inequality, it follows that

$$\begin{aligned} \|(\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k}) \mathbf{f}\|_{\mathbf{H}_{h,k}} &\leq \lambda \|\mathbf{H}_{h,k}^{-1} \boldsymbol{\mu}_h^* \mathbf{f}\|_{\mathbf{H}_{h,k}} + \left\| \mathbf{H}_{h,k}^{-1} \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top \mathbf{f} \right\|_{\mathbf{H}_{h,k}} \\ &= \lambda \|\boldsymbol{\mu}_h^* \mathbf{f}\|_{\mathbf{H}_{h,k}^{-1}} + \left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top \mathbf{f} \right\|_{\mathbf{H}_{h,k}^{-1}} \\ &\leq \sqrt{d\lambda} C_0 + \left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top \mathbf{f} \right\|_{\mathbf{H}_{h,k}^{-1}} \end{aligned}$$

where the last inequality uses  $\|\boldsymbol{\mu}_h^* \mathbf{f}\| \leq \sqrt{d} C_0$ .

- Assume  $f(\cdot)$  is a deterministic function. To evoke Lemma G.3, we set  $\mathcal{G}_j = \mathcal{F}_{h,j}$ ,  $\mathbf{x}_j = \sigma_{h,j}^{-1} \boldsymbol{\phi}_{h,j}$ ,  $\eta_j = \sigma_{h,j}^{-1} \boldsymbol{\varepsilon}_{h,j}^\top \mathbf{f} \cdot 1_{\mathcal{A}_{h,j}}$  and  $\mathbf{Z}_k = \lambda \mathbf{I} + \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\phi}_{h,j}^\top = \mathbf{H}_{h,k}$ . Here  $1_{\mathcal{A}}$  is the indicator function of the event  $\mathcal{A}$ .

Clearly  $\mathbf{x}_j \in \mathcal{G}_j$ ,  $\mathbb{E}[\eta_j | \mathcal{G}_j] = 0$  and  $\mathbb{E}[\eta_j^2 | \mathcal{G}_j] \leq C_\sigma^2$ . We also have  $\|\mathbf{x}_j\| \leq \sigma_{h,j}^{-1}$ ,  $|\eta_j| \leq 2C_0 \sigma_{h,j}^{-1}$  and  $\|\mathbf{x}_j\|_{\mathbf{Z}_{j-1}} = w_{h,j}$ . As a result,  $|\eta_j| \min\{1, \|\mathbf{x}_j\|_{\mathbf{Z}_{j-1}}\} \leq 2C_0 \frac{w_{h,j}}{\sigma_{h,j}} \leq \frac{2C_0}{\mathcal{H} d^{2.5} H}$  where the last inequality uses  $\sigma_{h,j}^2 \geq \mathcal{H} d^{2.5} H \|\boldsymbol{\phi}_{h,j}\|_{\mathbf{H}_{h,k}^{-1}}$  (which is equivalent to  $\frac{w_{h,j}}{\sigma_{h,j}} \leq (d^{2.5} H \mathcal{H})^{-1}$ ). By Lemma G.3, it follows that with probability  $1 - \frac{\delta}{H}$ , for all  $k \in [K]$ ,

$$\begin{aligned} \left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top \mathbf{f} 1_{\mathcal{A}_{h,j}} \right\|_{\mathbf{H}_{h,k}^{-1}} &= \left\| \sum_{j=1}^k \mathbf{x}_j \eta_j \right\|_{\mathbf{Z}_k^{-1}} \\ &\leq 8C_\sigma \sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right)} \log \frac{4HK^2}{\delta} + \frac{8C_0}{\mathcal{H} d^{2.5} H} \log \frac{4HK^2}{\delta}. \end{aligned}$$

Finally, on the event  $\cap_{k \in [K]} \mathcal{A}_{h,k}$ , we will have all the indicator functions equal to one.

- If  $f(\cdot)$  is a random function, we would use a covering argument to handle the possible correlation between  $f(\cdot)$  and history data, which would unfortunately enlarge  $\beta$ .

Denote the  $\varepsilon_0$ -net of  $\mathcal{V}$  by  $\mathcal{N}(\mathcal{V}, \varepsilon_0)$  where  $\varepsilon_0 = \min \left\{ C_\sigma \sigma_{\min}, \frac{\lambda C_0 \sqrt{d}}{K} \sigma_{\min}^2 \right\}$ . Hence, for any  $f \in \mathcal{V}$ , there exists  $\bar{f} \in \mathcal{N}(\mathcal{V}, \varepsilon_0)$  such that  $\|\bar{f} - \mathbf{f}\|_\infty = \sup_{s \in \mathcal{S}} |f(s) - \bar{f}(s)| \leq \varepsilon_0$ . Then,

$$\left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top \mathbf{f} \right\|_{\mathbf{H}_{h,k}^{-1}} \leq \underbrace{\left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top \bar{f} \right\|_{\mathbf{H}_{h,k}^{-1}}}_{(I)} + \underbrace{\left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top (\mathbf{f} - \bar{f}) \right\|_{\mathbf{H}_{h,k}^{-1}}}_{(II)}.$$

For the term (II), due to  $\|\boldsymbol{\phi}_{h,j}\| \leq 1$  and  $|\boldsymbol{\varepsilon}_{h,j}^\top (\mathbf{f} - \bar{f})| \leq \|\boldsymbol{\varepsilon}_{h,j}\|_1 \|\mathbf{f} - \bar{f}\|_\infty \leq 2\varepsilon_0$ , we have

$$\left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \boldsymbol{\phi}_{h,j} \boldsymbol{\varepsilon}_{h,j}^\top (\mathbf{f} - \bar{f}) \right\|_{\mathbf{H}_{h,k}^{-1}} \leq \frac{2K\varepsilon_0}{\sigma_{\min}^2 \sqrt{\lambda}} \leq 2\sqrt{d\lambda} C_0.$$

For the term (I), we define  $\mathcal{V}_{h,k} = \left\{ f' \in \mathcal{V} : 4C_\sigma^2 \sigma_{h,k}^2 \geq (\mathbb{V}_h f')(s_{h,k}, a_{h,k}) \right\}$ . Since the definition of  $\mathcal{V}_{h,k}$  involves only  $\sigma_{h,k}, s_{h,k}, a_{h,k} \in \mathcal{F}_{h,k}$ , for any fixed function  $f \in \mathcal{V}$ ,  $1_{f \in \mathcal{V}_{h,k}}$  is  $\mathcal{F}_{h,k}$ -measurable. On the event  $\mathcal{A}_{h,k}$ , by definition of  $\varepsilon_0$ ,

$$(\mathbb{V}_h \bar{f})(s_{h,k}, a_{h,k}) \leq 2(\mathbb{V}_h f)(s_{h,k}, a_{h,k}) + 2(\mathbb{V}_h (\bar{f} - f))(s_{h,k}, a_{h,k}) \leq 2C_\sigma^2 \sigma_{h,k}^2 + 2\varepsilon_0^2 \leq 4C_\sigma^2 \sigma_{h,k}^2.$$



Hence,  $\mathcal{A}_{h,k} \subseteq \left\{ \sigma_{h,k}^2 \geq (\mathbb{V}_h f)(s_{h,k}, a_{h,k}) / C_\sigma^2 \right\} \subseteq \{ \exists \bar{f} \in \mathcal{N}(\mathcal{V}, \varepsilon_0) \cap \mathcal{V}_{h,k} \}$  for all  $k \in [K]$ .

In the following, we will evoke Lemma G.3 to analyze the term (I). For any fixed  $f' \in \mathcal{V}$ , we set  $\mathcal{G}_j = \mathcal{F}_{h,j}$ ,  $\mathbf{x}_j = \sigma_{h,j}^{-1} \phi_{h,j}$ ,  $\eta_j = \sigma_{h,j}^{-1} \varepsilon_{h,j}^\top \mathbf{f}' \cdot 1_{f' \in \mathcal{V}_{h,k}}$  and  $\mathbf{Z}_k = \lambda \mathbf{I} + \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} \phi_{h,j}^\top = \mathbf{H}_{h,k}$ . Moreover, due to the choice of  $\sigma_{h,j}$ , it follows that

$$\left| \eta_j \min \left\{ 1, \|\mathbf{x}_j\|_{\mathbf{Z}_{j-1}^{-1}} \right\} \right| \leq \left| \frac{C_0}{\sigma_{h,j}} \right| \cdot \left\| \frac{\phi_{h,j}}{\sigma_{h,j}} \right\|_{\mathbf{H}_{h,j-1}^{-1}} \leq C_0 \frac{b_{h,j}}{\sigma_{h,j}^2} \leq \frac{C_0}{d^{2.5} H \mathcal{H}}.$$

By Lemma G.3 and the union bound, it follows that with probability  $1 - \frac{\delta}{H}$ , for all  $k \in [K]$ ,

$$\begin{aligned} & \sup_{f' \in \mathcal{N}(\mathcal{V}, \varepsilon_0)} \left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} \varepsilon_{h,j}^\top \mathbf{f}' 1_{f' \in \mathcal{V}_{h,k}} \right\|_{\mathbf{H}_{h,k}^{-1}} \\ & \leq 8C_\sigma \sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right)} \log \frac{4HK^2 N_0}{\delta} + \frac{8C_0}{d^{2.5} H \mathcal{H}} \log \frac{4HK^2 N_0}{\delta}. \end{aligned}$$

where  $N_0 = |\mathcal{N}(\mathcal{V}, \varepsilon_0)|$ .

As a result, we know that  $\left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} \varepsilon_{h,j}^\top \bar{\mathbf{f}} 1_{\bar{\mathbf{f}} \in \mathcal{V}_{h,k}} \right\|_{\mathbf{H}_{h,k}^{-1}}$  is no more than the RHS of the last inequality.

On the event  $\cap_{k \in [K]} \mathcal{A}_{h,k}$ , we have  $\bar{\mathbf{f}} \in \cap_{k \in [K]} \mathcal{V}_{h,k}$  and thus all the indicator functions equal to one, completing the proof.

- The proof is almost similar to the second item except that we use Lemma G.4 to analyze the term (I). Noticing we also have  $|\eta_j| = |\sigma_{h,j}^{-1} \varepsilon_{h,j}^\top \mathbf{f}'| \leq \frac{2C_0}{\sigma_{\min}}$ . By Lemma G.4 and the union bound, it follows that with probability  $1 - \frac{\delta}{H}$ , for all  $k \in [K]$ ,

$$\sup_{f' \in \mathcal{N}(\mathcal{V}_f, \varepsilon_1)} \left\| \sum_{j=1}^k \sigma_{h,j}^{-2} \phi_{h,j} \varepsilon_{h,j}^\top \mathbf{f}' 1_{f' \in \mathcal{V}_{h,k}} \right\|_{\mathbf{H}_{h,k}^{-1}} \leq \frac{2C_0}{\sigma_{\min}} \sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right)} + \log \frac{N_1}{\delta}.$$

Pay attention that here we don't utilize the variance information so that we change  $N_0 := |\mathcal{N}(\mathcal{V}, \varepsilon_0)|$  to  $N_1 := |\mathcal{N}(\mathcal{V}, \varepsilon_1)|$  and don't require  $\cap_{k \in [K]} \mathcal{A}_{h,k}$  is true.

□

#### F.4 Proof of Lemma D.3

*Proof of Lemma D.3.* The proof idea of Lemma D.3 is similar to that of Lemma D.1 except that we pay more attention on the reward variance.

Given that  $\mathcal{B}_{R^2}$  holds, we have  $\psi_h^* \in \tilde{\mathcal{R}}_{h,k}$  for all  $h \in [H]$  and  $k \in [K] \cup \{0\}$ .

We will prove the lemma by induction over  $k$ . When  $k = 0$ , we have  $\theta_{h,0} = 0$ ,  $\mathbf{H}_{h,0} = \lambda \mathbf{I}$  and  $\|\theta_{h,0} - \theta_h^*\|_{\mathbf{H}_{h,0}} = \sqrt{\lambda} \|\theta_h^*\| \leq \sqrt{\lambda} W \leq \beta_R$  for all  $h \in [H]$ . If we suppose  $\theta_h^* \in \mathcal{R}_{h,j}$  holds for all  $h \in [H]$  and  $j \in [k-1]$ , we are going to prove  $\theta_h^* \in \mathcal{R}_{h,k}$  uniformly for  $h \in [H]$ . The first thing we will show is

$$\sigma_{h,j}^2 \geq [\hat{\mathbb{V}}_h R_h](s_{h,j}, a_{h,j}) + R_{h,j} \quad \text{for all } h \in [H] \text{ and } j \in [k]. \quad (\text{F.3})$$

Notice that  $[\mathbb{V}_h R_h](s_{h,k}, a_{h,k}) = \langle \tilde{\phi}_{h,k}, \psi_h^* \rangle - \langle \phi_{h,k}, \theta_h^* \rangle^2$ . We then have for all  $h \in [H]$ ,  $j \in [k]$ ,

$$\begin{aligned} & |[\hat{\mathbb{V}}_h R_h - \mathbb{V}_h R_h](s_{h,j}, a_{h,j})| \\ & \leq \left| \langle \tilde{\phi}_{h,j}, \psi_{h,j-1} \rangle - \langle \tilde{\phi}_{h,j}, \psi_h^* \rangle \right| + \left| \langle \phi_{h,j}, \theta_h^* \rangle^2 - \langle \phi_{h,j}, \theta_{h,j-1} \rangle_{[0, \mathcal{H}]}^2 \right| \\ & \leq |\langle \tilde{\phi}_{h,k}, \psi_{h,k-1} - \psi_h^* \rangle| + 2\mathcal{H} |\langle \phi_{h,k}, \theta_{h,j-1} - \theta_h^* \rangle| \end{aligned}$$

$$\begin{aligned}
&\leq \|\tilde{\phi}_{h,j}\|_{\tilde{\mathbf{H}}_{h,j-1}^{-1}} \|\psi_{h,j-1} - \psi_h^*\|_{\tilde{\mathbf{H}}_{h,j-1}} + 2\mathcal{H}\|\phi_{h,j}\|_{\mathbf{H}_{h,j-1}^{-1}} \|\theta_{h,j-1} - \theta_h^*\|_{\mathbf{H}_{h,j-1}} \\
&\leq \beta_{R^2} \|\tilde{\phi}_{h,j}\|_{\tilde{\mathbf{H}}_{h,j-1}^{-1}} + 2\mathcal{H}\beta_R \|\phi_{h,j}\|_{\mathbf{H}_{h,j-1}^{-1}} = R_{h,j}
\end{aligned}$$

where the last inequality uses the hypothesis and the condition that  $\mathcal{B}_{R^2}$  holds. As a result, we establish equation F.3.

Let  $\mathcal{G}_{h,j} = \sigma(\mathcal{F}_{h-1,j} \cup \{s_{h,j}, a_{h,j}\})$ . One can show that both  $R_{h,j}$  and  $\sigma_{h,j}^2$  are  $\mathcal{G}_{h,j}$ -measurable. As a result, the event  $\mathcal{E}_{h,j} := \left\{ \sigma_{h,j}^2 \geq [\mathbb{V}_h R_h](s_{h,j}, a_{h,j}) \right\}$  is also  $\mathcal{G}_{h,j}$ -measurable. On the event  $\mathcal{B}_{R^2}$ , it is obvious that  $\cap_{h \in [H]} \cap_{j \in [k]} \mathcal{E}_{h,j}$  is true since equation F.3 is true.

On the other hand, we set  $\varepsilon_{h,j} = \frac{r_{h,j} - \langle \phi_{h,j}, \theta_h^* \rangle}{\sigma_{h,j}} \mathbf{1}_{\mathcal{E}_{h,j}}$  as the standardized reward. We then have  $\varepsilon_{h,j} \in \mathcal{F}_{h,j}$ ,  $\mathbb{E}[\varepsilon_{h,j} | \mathcal{G}_{h,j}] = 0$  and  $\text{Var}[\varepsilon_{h,j} | \mathcal{G}_{h,j}] \leq 1$ . We define  $\hat{\theta}_{h,k}$  as the solution of adaptive Huber regression to the response  $\{r_{h,j} \mathbf{1}_{\mathcal{E}_{h,j}}\}_{j \in [k]}$  and the feature  $\{\phi_{h,j} \mathbf{1}_{\mathcal{E}_{h,j}}\}_{j \in [k]}$ . We also define  $\hat{\mathbf{H}}_{h,k-1}$  as the counterpart matrix of  $\mathbf{H}_{h,k}$  obtained by replacing  $\phi_{h,k}$  with  $\phi_{h,k} \mathbf{1}_{\mathcal{E}_{h,k}}$ . We then apply Theorem 2.1 to analyze the concentration of  $\hat{\theta}_{h,k}$ . With probability at least  $1 - 3\delta$ , it follows that  $\|\hat{\theta}_{h,k} - \theta_h^*\|_{\hat{\mathbf{H}}_{h,k-1}} \leq \beta_R$  for all  $h \in [H]$  and  $k \in [K]$ . Because  $\mathcal{B}_{R^2}$  is true, all indicator functions equal to one. Therefore, we have  $\hat{\theta}_{h,k} = \theta_{h,k}$  and  $\hat{\mathbf{H}}_{h,k-1} = \mathbf{H}_{h,k-1}$ , implying  $\theta_h^* \in \mathcal{R}_{h,k}$  uniformly for  $h \in [H]$ .  $\square$

## F.5 Proof of Lemma D.4

*Proof of Lemma D.4.* By symmetry, we only prove the RHS inequality, or say, the optimism inequality. We prove it by induction. The statement is true for  $h = H + 1$  since both  $V_{H+1}^*(\cdot) = \bar{V}_{H+1}^k(\cdot) = 0$  for all  $k \in [K]$ . Assume the statement is also true for  $h + 1$ , implying  $V_{h+1}^*(\cdot) \leq \bar{V}_{h+1}^k(\cdot)$  for all  $k \in [K]$ . We assume there exists a sequence of updating episodes  $1 \leq k_1 < \dots < k_{N_k} \leq K$  such that

$$\bar{Q}_h^k(\cdot, \cdot) = \min_{i \in [N_k]} \left\{ \langle \phi(\cdot, \cdot), \theta_{h,k_i-1} + \mu_{h,k_i-1} \bar{V}_{h+1}^{k_i} \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k_i-1}^{-1}}, \mathcal{H} \right\}. \quad (\text{F.4})$$

Using  $Q_h^*(s, a) = \langle \phi(s, a), \theta_h^* + \mu_h^* V_{h+1}^* \rangle$ , we have for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and  $k \in [K]$ ,

$$\begin{aligned}
&\langle \phi(\cdot, \cdot), \theta_{h,k-1} + \mu_{h,k-1} \bar{V}_{h+1}^k \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}} - Q_h^*(\cdot, \cdot) \\
&= \langle \phi(\cdot, \cdot), \theta_{h,k-1} - \theta_h^* \rangle + \langle \phi(\cdot, \cdot), \mu_{h,k-1} \bar{V}_{h+1}^k - \mu_h^* V_{h+1}^* \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}} \\
&\stackrel{(a)}{\geq} \langle \phi(\cdot, \cdot), \theta_{h,k-1} - \theta_h^* \rangle + \langle \phi(\cdot, \cdot), (\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}} \\
&\stackrel{(b)}{\geq} \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k-1}^{-1}} \left[ -\|\theta_{h,k-1} - \theta_h^*\|_{\mathbf{H}_{h,k-1}} - \|(\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k\|_{\mathbf{H}_{h,k-1}} + \beta \right] \stackrel{(c)}{\geq} 0
\end{aligned}$$

where (a) uses  $\langle \phi(s, a), \mu_h^* (\bar{V}_{h+1}^k - V_{h+1}^*) \rangle = \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^*)(s, a) \geq 0$  from the hypothesis, (b) follows from Cauchy-Schwarz inequality and (c) uses  $\|\theta_{h,k-1} - \theta_h^*\|_{\mathbf{H}_{h,k-1}} + \|(\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k\|_{\mathbf{H}_{h,k-1}} \leq \beta_R + \beta_V = \beta$  on the event  $\mathcal{B}_R \cap \mathcal{B}_h$ .

As a result, by the last inequality and equation F.4, it follows that for all  $k \in [K]$ ,  $\bar{Q}_h^k(\cdot, \cdot) - Q_h^*(\cdot, \cdot) \geq 0$ . Taking maximum over actions, we have  $\bar{V}_h^k(\cdot) \geq V_h^*(\cdot)$  for all  $k \in [K]$ , which implies the case of  $h$  is also true.  $\square$

## F.6 Proof of Lemma D.5

*Proof of Lemma D.5.* The proof technique has been used in Lemma C.13 in (Hu et al., 2022) and Lemma 7.2 in (He et al., 2022). We include the proof for completeness. By definition,

$$[\mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k}) = \langle \mu_h^* [V_{h+1}^*]^2, \phi_{h,k} \rangle - \langle \mu_h^* V_{h+1}^*, \phi_{h,k} \rangle^2$$

$$[\widehat{\mathbb{V}}_h \bar{V}_{h+1}^k](s_{h,k}, a_{h,k}) = \langle \mu_{h,k-1} [\bar{V}_{h+1}^k]^2, \phi_{h,k-1} \rangle_{[0, \mathcal{H}^2]} - \langle \mu_{h,k-1} \bar{V}_{h+1}^k, \phi_{h,k} \rangle_{[0, \mathcal{H}]}^2.$$

Therefore, it follows that

$$\begin{aligned} & \left| \left[ \mathbb{V}_h V_{h+1}^* - \widehat{\mathbb{V}}_h \bar{V}_{h+1}^k \right] (s_{h,k}, a_{h,k}) \right| \\ & \leq \left| \left[ \mathbb{V}_h \bar{V}_{h+1}^k - \widehat{\mathbb{V}}_h \bar{V}_{h+1}^k \right] (s_{h,k}, a_{h,k}) \right| + \left| \left[ \mathbb{V}_h V_{h+1}^* - \mathbb{V}_h \bar{V}_{h+1}^k \right] (s_{h,k}, a_{h,k}) \right|. \end{aligned}$$

We then bound the two terms in the RHS of the last inequality as follows.

$$\begin{aligned} & \left| \left[ \mathbb{V}_h \bar{V}_{h+1}^k - \widehat{\mathbb{V}}_h \bar{V}_{h+1}^k \right] (s_{h,k}, a_{h,k}) \right| \\ & \leq \left| \langle \mu_h^* [\bar{V}_{h+1}^k]^2, \phi_{h,k} \rangle - \langle \mu_{h,k-1} [\bar{V}_{h+1}^k]^2, \phi_{h,k} \rangle_{[0, \mathcal{H}^2]} \right| \\ & \quad + \left| \langle \mu_h^* \bar{V}_{h+1}^k, \phi_{h,k} \rangle^2 - \langle \mu_{h,k-1} \bar{V}_{h+1}^k, \phi_{h,k} \rangle_{[0, \mathcal{H}]}^2 \right| \\ & \leq \left| \langle (\mu_h^* - \mu_{h,k-1}) [\bar{V}_{h+1}^k]^2, \phi_{h,k} \rangle \right| + 2\mathcal{H} \cdot \left| \langle \mu_h^* \bar{V}_{h+1}^k, \phi_{h,k} \rangle - \langle \mu_{h,k-1} \bar{V}_{h+1}^k, \phi_{h,k} \rangle_{[0, \mathcal{H}]} \right| \\ & \leq \left\| (\mu_{h,k-1} - \mu_h^*) [\bar{V}_{h+1}^k]^2 \right\|_{\mathbf{H}_{h,k-1}} \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 2\mathcal{H} \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \left\| (\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}} \end{aligned}$$

where the second inequality uses the fact that both  $\langle \mu_h^* \bar{V}_{h+1}^k, \phi_{h,k} \rangle$  and  $\langle \mu_{h,k-1} \bar{V}_{h+1}^k, \phi_{h,k} \rangle_{[0, \mathcal{H}]}$  lie between 0 and  $\mathcal{H}$ . Similarly, it follows that

$$\begin{aligned} & \left| \left[ \mathbb{V}_h V_{h+1}^* - \mathbb{V}_h \bar{V}_{h+1}^k \right] (s_{h,k}, a_{h,k}) \right| \\ & \leq \left| \mathbb{P}_h [[V_{h+1}^*]^2 - [\bar{V}_{h+1}^k]^2] (s_{h,k}, a_{h,k}) \right| + \left| [\mathbb{P}_h V_{h+1}^*]^2 (s_{h,k}, a_{h,k}) - [\mathbb{P}_h \bar{V}_{h+1}^k]^2 (s_{h,k}, a_{h,k}) \right| \\ & \leq \left| \mathbb{P}_h [(\bar{V}_{h+1}^k - V_{h+1}^*)(\bar{V}_{h+1}^k + V_{h+1}^*)] (s_{h,k}, a_{h,k}) \right| \\ & \quad + \left| [\mathbb{P}_h \bar{V}_{h+1}^k - \mathbb{P}_h V_{h+1}^*] [\mathbb{P}_h \bar{V}_{h+1}^k + \mathbb{P}_h V_{h+1}^*] (s_{h,k}, a_{h,k}) \right| \\ & \leq 4\mathcal{H} \cdot \mathbb{P}_h [\bar{V}_{h+1}^k - V_{h+1}^*] (s_{h,k}, a_{h,k}) \\ & \leq 4\mathcal{H} \cdot \mathbb{P}_h [\bar{V}_{h+1}^k - \underline{V}_{h+1}] (s_{h,k}, a_{h,k}) \\ & \leq 4\mathcal{H} \cdot \widehat{\mathbb{P}}_{h,k} [\bar{V}_{h+1}^k - \underline{V}_{h+1}] (s_{h,k}, a_{h,k}) \\ & \quad + 4\mathcal{H} \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \cdot \left[ \left\| (\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}} + \left\| (\mu_{h,k-1} - \mu_h^*) \underline{V}_{h+1} \right\|_{\mathbf{H}_{h,k-1}} \right] \end{aligned}$$

where the third and forth inequality we use the optimism and pessimism in Lemma D.4 and the last inequality uses the following result.

$$\begin{aligned} & \left| [\mathbb{P}_h \bar{V}_{h+1}^k - \widehat{\mathbb{P}}_{h,k} \bar{V}_{h+1}^k] (s_{h,k}, a_{h,k}) \right| = \left| \langle (\mu_h^* - \mu_{h,k-1}) \bar{V}_{h+1}^k, \phi_{h,k} \rangle \right| \\ & \leq \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \left\| (\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}}. \end{aligned}$$

A similar inequality can be derived for  $[\mathbb{P}_h \underline{V}_{h+1}^k - \widehat{\mathbb{P}}_{h,k} \underline{V}_{h+1}^k] (s_{h,k}, a_{h,k})$ . Finally, we have

$$\begin{aligned} & \left| \left[ \mathbb{V}_h V_{h+1}^* - \widehat{\mathbb{V}}_h \bar{V}_{h+1}^k \right] (s_{h,k}, a_{h,k}) \right| \\ & \leq \left\| (\mu_{h,k-1} - \mu_h^*) [\bar{V}_{h+1}^k]^2 \right\|_{\mathbf{H}_{h,k-1}} \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 4\mathcal{H} \widehat{\mathbb{P}}_{h,k} (\bar{V}_{h+1}^k - \underline{V}_{h+1}) (s_{h,k}, a_{h,k}) \\ & \quad + \mathcal{H} \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \cdot \left[ 6 \left\| (\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}} + 4 \left\| (\mu_{h,k-1} - \mu_h^*) \underline{V}_{h+1} \right\|_{\mathbf{H}_{h,k-1}} \right]. \end{aligned}$$

We complete the proof by noting that on the event  $\mathcal{B}_0$ , we have

$$\max \left\{ \left\| (\mu_h^* - \mu_{h,k-1}) \bar{V}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}}, \left\| (\mu_h^* - \mu_{h,k-1}) \underline{V}_{h+1} \right\|_{\mathbf{H}_{h,k-1}} \right\} \leq \beta_0,$$

$$\left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) [\bar{\mathbf{V}}_{h+1}^k]^2 \right\|_{\mathbf{H}_{h,k-1}} \leq \mathcal{H} \beta_0.$$

□

### F.7 Proof of Lemma D.6

*Proof of Lemma D.6.* For any  $j \leq k$ , we have

$$\begin{aligned} \left[ \mathbb{V}_h(\bar{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*) \right] (s_{h,j}, a_{h,j}) &\stackrel{(a)}{\leq} \left[ \mathbb{P}_h(\bar{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*)^2 \right] (s_{h,j}, a_{h,j}) \\ &\stackrel{(b)}{\leq} \mathcal{H} \left[ \mathbb{P}_h(\bar{\mathbf{V}}_{h+1}^k - \underline{\mathbf{V}}_{h+1}^k) \right] (s_{h,j}, a_{h,j}) \\ &\stackrel{(c)}{\leq} \mathcal{H} \left[ \mathbb{P}_h(\bar{\mathbf{V}}_{h+1}^j - \underline{\mathbf{V}}_{h+1}^j) \right] (s_{h,j}, a_{h,j}) \end{aligned}$$

where (a) uses the fact that  $\text{Var}(X) \leq \mathbb{E}X^2$  for any random variable  $X$ , (b) uses  $0 \leq \underline{\mathbf{V}}_{h+1}^k(\cdot) \leq \mathbf{V}_{h+1}^*(\cdot) \leq \bar{\mathbf{V}}_{h+1}^k(\cdot) \leq \mathcal{H}$  on the event  $\mathcal{B}_R \cap \mathcal{B}_{h+1}$  from Lemma D.4, and (c) uses that  $\bar{\mathbf{V}}_{h+1}^j(\cdot) \geq \bar{\mathbf{V}}_{h+1}^k(\cdot)$  and  $\underline{\mathbf{V}}_{h+1}^j(\cdot) \leq \underline{\mathbf{V}}_{h+1}^k(\cdot)$  by definition. On the other hand, on the event  $\mathcal{B}_0$ , we have

$$\max \left\{ \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,j-1}) \bar{\mathbf{V}}_{h+1}^j \right\|_{\mathbf{H}_{h,j-1}}, \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,j-1}) \underline{\mathbf{V}}_{h+1}^j \right\|_{\mathbf{H}_{h,j-1}} \right\} \leq \beta_0.$$

As a result,

$$\begin{aligned} \left[ \left( \mathbb{P}_h - \hat{\mathbb{P}}_{h,j} \right) \bar{\mathbf{V}}_{h+1}^j \right] (s_{h,j}, a_{h,j}) &= \langle \phi_{h,j}, (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,j-1}) \bar{\mathbf{V}}_{h+1}^j \rangle \\ &\leq \|\phi_{h,j}\|_{\mathbf{H}_{h,j-1}^{-1}} \|(\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,j-1}) \bar{\mathbf{V}}_{h+1}^j\|_{\mathbf{H}_{h,j-1}} \leq \beta_0 \|\phi_{h,j}\|_{\mathbf{H}_{h,j-1}^{-1}}. \end{aligned}$$

Similarly, we have  $\left[ \left( \mathbb{P}_h - \hat{\mathbb{P}}_{h,j} \right) \underline{\mathbf{V}}_{h+1}^j \right] (s_{h,j}, a_{h,j}) \leq \beta_0 \|\phi_{h,j}\|_{\mathbf{H}_{h,j-1}^{-1}}$ . Therefore,

$$\begin{aligned} \left[ \mathbb{V}_h(\bar{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*) \right] (s_{h,j}, a_{h,j}) \\ \leq \mathcal{H} \left[ 2\beta_0 \|\phi_{h,k}\|_{\mathbf{H}_{h,j-1}^{-1}} + \left[ \hat{\mathbb{P}}_{h,j}(\bar{\mathbf{V}}_{h+1}^j - \underline{\mathbf{V}}_{h+1}^j) \right] (s_{h,j}, a_{h,j}) \right] =: E_{h,j} \end{aligned}$$

Repeating the above argument, we have a similar inequality for  $\bar{\mathbf{V}}_{h+1}^k$  due to symmetry. □

### F.8 Proof of Lemma D.7

*Proof of Lemma D.7.* Due to the backward recursion structure, we will use induction (over horizon  $h$ ) to prove this lemma. First, equation D.2 is true for  $h = H$  since  $\bar{\mathbf{V}}_{H+1}^k(\cdot) = \underline{\mathbf{V}}_{H+1}^k(\cdot) = 0$  for all  $k \in [K]$ . Therefore, we have  $\mathcal{B}_H$  holds. Assume equation D.2 holds for horizons no smaller than  $h + 1$ , i.e.,  $\mathcal{B}_{h+1}$  holds with  $h + 1 \leq H$ . In the following, we will show, once  $\mathcal{B}_{h+1} \cap \mathcal{B}_0$  holds,  $\mathcal{B}_h$  holds with probability at least  $1 - \frac{2\delta}{H}$ . Repeating the argument, we have, given  $\mathcal{B}_H \cap \mathcal{B}_0$  holds, with probability at least  $1 - 2\delta$ ,  $\mathcal{B}_1 \cap \mathcal{B}_0$  holds. Hence,  $\mathbb{P}(\mathcal{B}_0 \cap \mathcal{B}_1) \geq 1 - 5\delta$ .

Note that

$$\begin{aligned} \max \left\{ \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \bar{\mathbf{V}}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}}, \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \underline{\mathbf{V}}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}} \right\} &\leq \|(\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \mathbf{V}_{h+1}^*\|_{\mathbf{H}_{h,k-1}} \\ &+ \max \left\{ \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) (\bar{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*) \right\|_{\mathbf{H}_{h,k-1}}, \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) (\underline{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*) \right\|_{\mathbf{H}_{h,k-1}} \right\}. \end{aligned}$$

we would analyze the two terms in the RHS separately to proceed with the proof.

**For the first term** Since  $\mathbf{V}_{h+1}^*$  is a deterministic function, we apply the first item in Lemma F.1 to bound it. In the following, we specify the parameters defined therein. First, we have  $C_0 = \mathcal{H}$  and  $\mathcal{A}_{h,k} = \left\{ \sigma_{h,k}^2 \geq (\mathbb{V}_h V_{h+1}^*)(s_{h,k}, a_{h,k}) \right\}$  is  $\mathcal{F}_{h,k}$ -measurable. By Lemma D.5, on the event  $\mathcal{B}_0 \cap \mathcal{B}_{h+1}$ , we have for all  $k \in [K]$ ,  $\left| \left[ \mathbb{V}_h V_{h+1}^* - \widehat{\mathbb{V}}_h \overline{V}_{h+1}^k \right] (s_{h,k}, a_{h,k}) \right| \leq U_{h,k}$  with  $U_{h,k}$  defined in equation A.9. Hence,  $\sigma_{h,k}^2 \geq [\widehat{\mathbb{V}}_h \widehat{V}_{h+1}^k](s_{h,k}, a_{h,k}) + U_{h,k} \geq [\mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k})$  for all  $k \in [K]$ , implying  $\cap_{k \in [K]} \mathcal{A}_{h,k}$  holds under  $\mathcal{B}_0 \cap \mathcal{B}_{h+1}$  and  $C_\sigma = 1$ . By Lemma F.1, with probability at least  $1 - \frac{\delta}{H}$ ,  $\|(\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \mathbf{V}_{h+1}^*\|_{\mathbf{H}_{h,k-1}} \leq \beta_1$  for all  $k \in [K]$  with  $\beta_1$  defined in the following. Finally, we simplify  $\beta_1$  as

$$\begin{aligned} \beta_1 &:= 8\sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right) \log \frac{4HK^2}{\delta}} + \frac{8}{d^{2.5}H} \log \frac{4HK^2}{\delta} + \sqrt{d\lambda\mathcal{H}} \\ &\leq 8\sqrt{d\iota_1} + \frac{8\iota_1}{d^{2.5}H} + \sqrt{d\lambda\mathcal{H}} \leq 16\sqrt{d\iota_1} + \sqrt{d\lambda\mathcal{H}}. \end{aligned}$$

**For the second term** Since both  $\overline{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*$  and  $\underline{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*$  are  $\mathcal{F}_{H,k-1}$ -measurable random functions, we apply the second item in Lemma F.1 to analyze the second term. In the following, we specify parameters defined therein. First,  $C_0 = \mathcal{H}$  and  $\mathcal{A}_{h,k} = \left\{ \sigma_{h,k}^2 \geq d^3 H \cdot E_{h,k} \right\}$  is  $\mathcal{F}_{h,k}$ -measurable. By Lemma D.6, on the event  $\mathcal{B}_0 \cap \mathcal{B}_R \cap \mathcal{B}_{h+1}$ , we have simultaneously  $\left[ \mathbb{V}_h (\overline{V}_{h+1}^k - V_{h+1}^*) \right] (s_{h,j}, a_{h,j}) \leq E_{h,j}$  and  $\left[ \mathbb{V}_h (\underline{V}_{h+1}^k - V_{h+1}^*) \right] (s_{h,j}, a_{h,j}) \leq E_{h,j}$  for all  $j \leq k \leq K$  with  $E_{h,j}$  defined in equation A.7. As a result, for all  $j \leq k$ ,

$$\sigma_{h,j}^2 \geq d^3 H \cdot E_{h,j} \geq d^3 H \cdot \max \left\{ \left[ \mathbb{V}_h (\overline{V}_{h+1}^k - V_{h+1}^*) \right] (s_{h,j}, a_{h,j}), \left[ \mathbb{V}_h (\underline{V}_{h+1}^k - V_{h+1}^*) \right] (s_{h,j}, a_{h,j}) \right\}.$$

It implies  $C_\sigma = \frac{1}{\sqrt{d^3 H}}$  and for any  $j \in [k]$ ,

$$\mathcal{A}_{h,j} \subseteq \left\{ \sigma_{h,j}^2 \geq C_\sigma^{-2} \max \left\{ \left[ \mathbb{V}_h (\overline{V}_{h+1}^k - V_{h+1}^*) \right] (s_{h,j}, a_{h,j}), \left[ \mathbb{V}_h (\underline{V}_{h+1}^k - V_{h+1}^*) \right] (s_{h,j}, a_{h,j}) \right\} \right\}.$$

Finally, with by Lemma G.8 and G.7, the covering entropy for  $\varepsilon_0 = \min \left\{ \frac{\sigma_{\min}}{\sqrt{d^3 H}}, \frac{\lambda \mathcal{H} \sqrt{d}}{K} \sigma_{\min}^2 \right\}$  and the function class to which  $\overline{V}_{h+1}^k - V_{h+1}^*$  and  $\underline{V}_{h+1}^k - V_{h+1}^*$  belong is

$$\begin{aligned} \log N_0 &= |\mathcal{N}(\mathcal{V}^\pm, \varepsilon_0)| \leq \left[ d \log \left( 1 + \frac{4L}{\varepsilon_0} \right) + d^2 \log \left( 1 + \frac{8\sqrt{d}B^2}{\lambda \varepsilon_0^2} \right) \right] \cdot dH \log_2 \left( 1 + \frac{K}{\lambda \sigma_{\min}^2} \right) \\ &= \mathcal{O}(d^3 H \iota_1^2) \end{aligned}$$

By Lemma F.1, with probability at least  $1 - \frac{\delta}{H}$ ,

$$\max \left\{ \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) (\overline{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*) \right\|_{\mathbf{H}_{h,k-1}}, \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) (\underline{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^*) \right\|_{\mathbf{H}_{h,k-1}} \right\} \leq \beta_2$$

for all  $k \in [K]$  with  $\beta_2$  defined in the following. Finally, we simplify  $\beta_2$  as

$$\begin{aligned} \beta_2 &= \frac{8}{\sqrt{d^3 H}} \sqrt{d \log \left( 1 + \frac{K}{\sigma_{\min}^2 d \lambda} \right) \log \frac{4N_0 H K^2}{\delta}} + \frac{8}{d^{2.5}H} \log \frac{4N_0 H K^2}{\delta} + \sqrt{d\lambda\mathcal{H}} \\ &\leq 8\sqrt{\frac{\iota_1}{d^2 H}} \cdot (\iota_1 + \mathcal{O}(d^3 H \iota_1^2)) + \frac{8}{d^{2.5}H} (\iota_1 + \mathcal{O}(d^3 H \iota_1^2)) + \sqrt{d\lambda\mathcal{H}} \\ &= \mathcal{O}(\sqrt{d\iota_1^{1.5}} + \iota_1 + \sqrt{d\iota_1^2} + \sqrt{d\lambda\mathcal{H}}) = \mathcal{O}(\sqrt{d\iota_1^2} + \sqrt{d\lambda\mathcal{H}}). \end{aligned}$$

**Putting pieces together** we have shown that given  $\mathcal{B}_{h+1} \cap \mathcal{B}_0$  is true, with probability at least  $1 - 2\delta$ , for all  $h \in [H]$  and  $k \in [K]$ ,

$$\max \left\{ \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \bar{\mathbf{V}}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}}, \left\| (\boldsymbol{\mu}_h^* - \boldsymbol{\mu}_{h,k-1}) \mathbf{V}_{h+1}^k \right\|_{\mathbf{H}_{h,k-1}} \right\} \leq \beta_1 + \beta_2 = \mathcal{O} \left( \sqrt{d} l_1^2 + \sqrt{d} \lambda \mathcal{H} \right).$$

Therefore,  $\mathcal{B}_V := \mathcal{B}_1$  holds.  $\square$

## F.9 Proof of Lemma D.8

*Proof of Lemma D.8.* For a given  $k$ , let  $k_{\text{last}}$  denote the latest update episode before episode  $k$ , that is  $k_{\text{last}} \leq k < k_{\text{last}} + 1$ . By Lemma G.6, due to  $\mathbf{H}_{h,k-1} \geq \mathbf{H}_{h,k_{\text{last}}-1}$  and  $\det(\mathbf{H}_{h,k-1}) \leq 2 \det(\mathbf{H}_{h,k_{\text{last}}-1})$ , it follows that for any  $\mathbf{x} \in \mathbb{R}^d$ ,

$$\|\mathbf{x}\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}} \leq 2 \|\mathbf{x}\|_{\mathbf{H}_{h,k-1}^{-1}}. \quad (\text{F.5})$$

By definition,  $\bar{Q}_h^k(\cdot, \cdot) \leq \langle \phi(\cdot, \cdot), \boldsymbol{\theta}_{h,k_{\text{last}}-1} + \boldsymbol{\mu}_{h,k_{\text{last}}-1} \bar{\mathbf{V}}_{h+1}^{k_{\text{last}}} \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}}$  and  $Q_h^{\pi_k}(s, a) = \langle \phi(s, a), \boldsymbol{\theta}_h^* + \boldsymbol{\mu}_h^* \mathbf{V}_{h+1}^{\pi_k} \rangle$ . Using  $a_{h,k} = \pi_h^k(s_{h,k}) = \arg\max_{a \in \mathcal{A}} \bar{Q}_h^k(s_{h,k}, a)$ , we then have

$$\begin{aligned} (\bar{V}_h^k - V_h^{\pi_k})(s_{h,k}) &\leq (\bar{Q}_h^k - Q_h^{\pi_k})(s_{h,k}, a_{h,k}) \\ &\leq \langle \phi_{h,k}, \boldsymbol{\theta}_{h,k_{\text{last}}-1} + \boldsymbol{\mu}_{h,k_{\text{last}}-1} \bar{\mathbf{V}}_{h+1}^{k_{\text{last}}} - (\boldsymbol{\theta}_h^* + \boldsymbol{\mu}_h^* \mathbf{V}_{h+1}^{\pi_k}) \rangle + \beta \|\phi(s_{h,k}, a_{h,k})\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}} \\ &\stackrel{(a)}{\leq} \langle \phi_{h,k}, (\boldsymbol{\theta}_{h,k_{\text{last}}-1} - \boldsymbol{\theta}_h^*) + (\boldsymbol{\mu}_{h,k_{\text{last}}-1} - \boldsymbol{\mu}_h^*) \bar{\mathbf{V}}_{h+1}^{k_{\text{last}}} \rangle + \langle \phi_{h,k}, \boldsymbol{\mu}_h^* (\bar{\mathbf{V}}_{h+1}^{k_{\text{last}}} - \mathbf{V}_{h+1}^{\pi_k}) \rangle + 2\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \\ &\stackrel{(b)}{\leq} 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \langle \phi_{h,k}, \boldsymbol{\mu}_h^* (\bar{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^{\pi_k}) \rangle \\ &\stackrel{(c)}{=} 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) \\ &\stackrel{(d)}{=} 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + (\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h+1,k}) + X_{h,k}. \end{aligned}$$

Here (a) uses equation F.5, (b) uses

$$\begin{aligned} &|\langle \phi_{h,k}, (\boldsymbol{\theta}_{h,k_{\text{last}}-1} - \boldsymbol{\theta}_h^*) + (\boldsymbol{\mu}_{h,k_{\text{last}}-1} - \boldsymbol{\mu}_h^*) \bar{\mathbf{V}}_{h+1}^{k_{\text{last}}} \rangle| \\ &\leq \|\phi_{h,k}\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}} \|(\boldsymbol{\theta}_{h,k_{\text{last}}-1} - \boldsymbol{\theta}_h^*) + (\boldsymbol{\mu}_{h,k_{\text{last}}-1} - \boldsymbol{\mu}_h^*) \bar{\mathbf{V}}_{h+1}^{k_{\text{last}}}\|_{\mathbf{H}_{h,k_{\text{last}}-1}} \\ &\leq \beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}} \leq 2\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \end{aligned}$$

on  $\mathcal{B}_R \cap \mathcal{B}_V$ , (c) uses  $\langle \phi_{h,k}, \boldsymbol{\mu}_h^* (\bar{\mathbf{V}}_{h+1}^{k_{\text{last}}} - \mathbf{V}_{h+1}^{\pi_k}) \rangle = \mathbb{P}_h(\bar{\mathbf{V}}_{h+1}^{k_{\text{last}}} - \mathbf{V}_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) = \mathbb{P}_h(\bar{\mathbf{V}}_{h+1}^k - \mathbf{V}_{h+1}^{\pi_k})(s_{h,k}, a_{h,k})$ , and (d) uses the notation

$$X_{h,k} := \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) - (\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h+1,k}).$$

The last inequality implies

$$(\bar{V}_h^k - V_h^{\pi_k})(s_{h,k}) \leq (\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h+1,k}) + X_{h,k} + 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}.$$

Iterating the above inequality over  $h$  and using  $\bar{V}_{H+1}^k(\cdot) = V_{H+1}^{\pi_k}(\cdot) = 0$ , we have

$$(\bar{V}_h^k - V_h^{\pi_k})(s_{h,k}) \leq \sum_{i=h}^H \left[ X_{i,k} + 4\beta \|\phi_{i,k}\|_{\mathbf{H}_{i,k-1}^{-1}} \right]. \quad (\text{F.6})$$

Therefore, setting  $h = 1$  and summing equation F.6 over  $k \in [K]$ , we have

$$\sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k}) \leq \sum_{k=1}^K \sum_{h=1}^H \left[ X_{h,k} + 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \right]. \quad (\text{F.7})$$

We then need to analyze  $\sum_{k=1}^K \sum_{h=1}^H X_{h,k}$ . Since  $s_{h+1,k}$  is  $\mathcal{F}_{h+1,k}$ -measurable,  $\pi_k = \{\pi_h^k\}_{h \in [H]}$ ,  $\bar{V}_{h+1}^k$  is  $\mathcal{F}_{H,k-1}$ -measurable, we have  $X_{h,k}$  is  $\mathcal{F}_{h+1,k}$ -measurable. We also have  $\mathbb{E}[X_{h,k} | \mathcal{F}_{h,k}] = 0$ ,  $|X_{h,k}| \leq 2\mathcal{H}$  and

$$\begin{aligned} \mathbb{E}[X_{h,k}^2 | \mathcal{F}_{h,k}] &\leq \mathbb{E}[(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})^2(s_{h+1,k}) | \mathcal{F}_{h,k}] \stackrel{(a)}{\leq} \mathcal{H} \mathbb{E}[|\bar{V}_{h+1}^k - V_{h+1}^{\pi_k}|(s_{h+1,k}) | \mathcal{F}_{h,k}] \\ &\stackrel{(b)}{=} \mathcal{H} \mathbb{E}[(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h+1,k}) | \mathcal{F}_{h,k}] = \mathcal{H} \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) \end{aligned}$$

where (a) uses  $|\bar{V}_{h+1}^k - V_{h+1}^{\pi_k}|(\cdot) \leq \mathcal{H}$  and (b) uses the optimism in Lemma D.4. By the variance-aware Freedman inequality in Lemma G.2, with probability at least  $1 - \frac{\delta}{2}$ , it follows that

$$\left| \sum_{k=1}^K \sum_{h=1}^H X_{h,k} \right| \leq 3\sqrt{\iota} \cdot \sqrt{\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) + 10\mathcal{H} \cdot \iota} \quad (\text{F.8})$$

where  $\iota = \log \frac{4\lceil \log_2 HK \rceil}{\delta}$ . On the other hand, it follows that

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) &= \sum_{k=1}^K \sum_{h=2}^H (\bar{V}_h^k - V_h^{\pi_k})(s_{h,k}) + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} \\ &\stackrel{(a)}{\leq} \sum_{k=1}^K \sum_{h=2}^H \sum_{i=h}^H \left[ X_{i,k} + 4\beta \|\phi_{i,k}\|_{\mathbf{H}_{i,k-1}^{-1}} \right] + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} \\ &= \sum_{k=1}^K \sum_{h=2}^H (H-h+1) \left[ X_{h,k} + 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \right] + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} \\ &\stackrel{(b)}{\leq} 4H\beta \sum_{k=1}^K \sum_{h=2}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} b_h \end{aligned}$$

where (a) uses equation F.6 and (b) uses the notation  $b_h = 1$  if  $h = 1$ ; otherwise  $= H - h + 2$  for  $2 \leq h \leq H$ . Clearly, we have  $|b_h| \leq H$  for all  $h \in [H]$ . By the variance-aware Freedman inequality in Lemma G.2, with probability at least  $1 - \frac{\delta}{2}$ , it follows that

$$\left| \sum_{k=1}^K \sum_{h=1}^H X_{h,k} b_h \right| \leq 3H\sqrt{\iota} \cdot \sqrt{\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) + 10H\mathcal{H} \cdot \iota}.$$

As a result, we have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) &\leq 3H\sqrt{\iota} \cdot \sqrt{\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k})} \\ &\quad + 4H\beta \sum_{k=1}^K \sum_{h=2}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 10H\mathcal{H}\iota. \end{aligned}$$

Using the inequality that  $x \leq 2(a^2 + b^2)$  for any  $x \leq |a|\sqrt{x} + b^2$ , we have

$$\sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) \leq 8H\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H^2\mathcal{H}\iota. \quad (\text{F.9})$$

Putting pieces together, we have

$$\sum_{k=1}^K (\bar{V}_1^k - V_1^{\pi_k})(s_{1,k}) \stackrel{\text{equation F.7}}{\leq} \sum_{k=1}^K \sum_{h=1}^H \left[ X_{h,k} + 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \right]$$



$$\begin{aligned}
& \stackrel{\text{equation F.8}}{\leq} 4\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 3\sqrt{\iota} \cdot \sqrt{\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) + 10\mathcal{H} \cdot \iota} \\
& \stackrel{\text{equation F.9}}{\leq} 4\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 3\sqrt{\iota} \cdot \sqrt{\mathcal{H} \cdot \left[ 8H\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H^2\mathcal{H}\iota \right] + 10\mathcal{H} \cdot \iota} \\
& \leq 6\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H\mathcal{H}\iota
\end{aligned}$$

where the last inequality uses  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  and  $2\sqrt{ab} \leq a+b$  for non-negative numbers  $a, b \geq 0$ .  $\square$

## F.10 Proof of Lemma D.9

*Proof of Lemma D.9.* The proof main idea is similar to that in Lemma D.8. For a given  $k$ , let  $k_{\text{last}}$  denote the latest update episode before episode  $k$ , that is  $k_{\text{last}} \leq k < k_{\text{last}} + 1$ . By definition,  $\bar{Q}_h^k(\cdot, \cdot) \leq \langle \phi(\cdot, \cdot), \theta_{h,k_{\text{last}}-1} + \mu_{h,k_{\text{last}}-1} \bar{V}_{h+1}^{k_{\text{last}}} \rangle + \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}}$  and  $\underline{Q}_h^k(\cdot, \cdot) \geq \langle \phi(\cdot, \cdot), \theta_{h,k_{\text{last}}-1} + \mu_{h,k_{\text{last}}-1} \underline{V}_{h+1}^{k_{\text{last}}} \rangle - \beta \|\phi(\cdot, \cdot)\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}}$ . Using

$$a_{h,k} = \pi_h^k(s_{h,k}) = \operatorname{argmax}_{a \in \mathcal{A}} \bar{Q}_h^k(s_{h,k}, a),$$

we then have

$$\begin{aligned}
(\bar{V}_h^k - \underline{V}_h^k)(s_{h,k}) & \leq (\bar{Q}_h^k - \underline{Q}_h^k)(s_{h,k}, a_{h,k}) \\
& \leq \langle \phi_{h,k}, \mu_{h,k_{\text{last}}-1} (\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}}) \rangle + 2\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}} \\
& \stackrel{(a)}{\leq} \langle \phi_{h,k}, (\mu_{h,k-1} - \mu_h^*)(\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}}) \rangle + \langle \phi_{h,k}, \mu_h^*(\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}}) \rangle + 4\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \\
& \stackrel{(b)}{\leq} 6\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \langle \phi_{h,k}, \mu_h^*(\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}}) \rangle \\
& \stackrel{(c)}{=} 6\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \\
& \stackrel{(d)}{=} 6\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + (\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h+1,k}) + X_{h,k}.
\end{aligned}$$

Here (a) uses equation F.5, (b) uses

$$\begin{aligned}
& |\langle \phi_{h,k}, (\mu_{h,k_{\text{last}}-1} - \mu_h^*)(\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}}) \rangle| \\
& \leq \|\phi_{h,k}\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}} \|(\mu_{h,k_{\text{last}}-1} - \mu_h^*)(\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}})\|_{\mathbf{H}_{h,k_{\text{last}}-1}} \\
& \leq 2\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k_{\text{last}}-1}^{-1}} \leq 2\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}
\end{aligned}$$

on  $\mathcal{B}_V \cap \mathcal{B}_R$ , (c) uses  $\langle \phi_{h,k}, \mu_h^*(\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}}) \rangle = \mathbb{P}_h(\bar{V}_{h+1}^{k_{\text{last}}} - \underline{V}_{h+1}^{k_{\text{last}}})(s_{h,k}, a_{h,k}) = \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k})$ , and (d) uses the notation

$$X_{h,k} := \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) - (\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h+1,k}). \quad (\text{F.10})$$

The last inequality implies

$$(\bar{V}_h^k - \underline{V}_h^k)(s_{h,k}) \leq (\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h+1,k}) + X_{h,k} + 6\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}.$$

Iterating the above inequality over  $h$  and using  $\bar{V}_{H+1}^k(\cdot) = \underline{V}_{H+1}^k(\cdot) = 0$ , we have

$$(\bar{V}_h^k - \underline{V}_h^k)(s_{h,k}) \leq \sum_{i=h}^H \left[ X_{i,k} + 6\beta \|\phi_{i,k}\|_{\mathbf{H}_{i,k-1}^{-1}} \right]. \quad (\text{F.11})$$

Using the last inequality, it follows that

$$\begin{aligned}
\sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) &= \sum_{k=1}^K \sum_{h=2}^H (\bar{V}_h^k - \underline{V}_h^k)(s_{h,k}) + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} \\
&\stackrel{(a)}{\leq} \sum_{k=1}^K \sum_{h=2}^H \sum_{i=h}^H \left[ X_{i,k} + 8\beta \|\phi_{i,k}\|_{\mathbf{H}_{i,k-1}^{-1}} \right] + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} \\
&= \sum_{k=1}^K \sum_{h=2}^H (H-h+1) \left[ X_{h,k} + 6\beta \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \right] + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} \\
&\stackrel{(b)}{\leq} 6H\beta \sum_{k=1}^K \sum_{h=2}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \sum_{k=1}^K \sum_{h=1}^H X_{h,k} b_h
\end{aligned} \tag{F.12}$$

where (a) uses equation F.6 and (b) uses the notation  $b_h = 1$  if  $h = 1$ ; otherwise  $= H - h + 2$  for  $2 \leq h \leq H$ . Clearly, we have  $|b_h| \leq H$  for all  $h \in [H]$ .

We then need to analyze  $\sum_{k=1}^K \sum_{h=1}^H X_{h,k} b_h$  with  $X_{h,k}$ 's defined in equation F.10. Since  $s_{h+1,k}$  is  $\mathcal{F}_{h+1,k}$ -measurable,  $\bar{V}_{h+1}^k, \underline{V}_{h+1}^k$  is  $\mathcal{F}_{H,k-1}$ -measurable, we have  $X_{h,k}$  is  $\mathcal{F}_{h+1,k}$ -measurable. We also have  $\mathbb{E}[X_{h,k} | \mathcal{F}_{h,k}] = 0$ ,  $|X_{h,k}| \leq 2\mathcal{H}$  and

$$\begin{aligned}
\mathbb{E}[X_{h,k}^2 | \mathcal{F}_{h,k}] &\leq \mathbb{E}[(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)^2(s_{h+1,k}) | \mathcal{F}_{h,k}] \\
&\stackrel{(a)}{\leq} \mathcal{H} \mathbb{E}[|\bar{V}_{h+1}^k - \underline{V}_{h+1}^k|(s_{h+1,k}) | \mathcal{F}_{h,k}] = \mathcal{H} \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k})
\end{aligned}$$

where (a) uses  $|\bar{V}_{h+1}^k - \underline{V}_{h+1}^k|(\cdot) \leq \mathcal{H}$ . By the variance-aware Freedman inequality in Lemma G.2, with probability at least  $1 - \delta$ , it follows that

$$\left| \sum_{k=1}^K \sum_{h=1}^H X_{h,k} b_h \right| \leq 3H\sqrt{\iota} \cdot \sqrt{\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) + 10H\mathcal{H} \cdot \iota} \tag{F.13}$$

where  $\iota = \log \frac{4[\log_2 HK]}{\delta}$ . As a result, plugging equation F.13 into equation F.12, we have

$$\begin{aligned}
\sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) &\leq 3H\sqrt{\iota} \cdot \sqrt{\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k})} \\
&\quad + 6H\beta \sum_{k=1}^K \sum_{h=2}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 10H\mathcal{H}\iota.
\end{aligned}$$

Using the inequality that  $x \leq 2(a^2 + b^2)$  for any  $x \leq |a|\sqrt{x} + b^2$ , we have

$$\sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \leq 12H\beta \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H^2\mathcal{H}\iota.$$

□

### F.11 Proof of Lemma D.10

*Proof of Lemma D.10.* Recall that  $b_{h,k} = \max\{\|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}\}$ ,  $w_{h,k} = \sigma_{h,k}^{-1} \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}$  and  $\tilde{w}_{h,k} = \sigma_{h,k}^{-1} \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}$ . As a result, we have  $\sigma_{h,k}^{-1} b_{h,k} = \max\{w_{h,k}, \tilde{w}_{h,k}\}$ . On the other hand,

$$\sigma_{h,k}^2 = \max \left\{ \sigma_{\min}^2, d^3 H \cdot E_{h,k}, J_{h,k}, c_0^{-2} b_{h,k}^2, \left( \frac{W}{\sqrt{c_1 d}} + \mathcal{H} d^{2.5} H \right) b_{h,k} \right\}. \tag{A.5}$$

Based on what value  $\sigma_{h,k}$  takes, we compose the full index set  $\mathcal{I} := [H] \times [K]$  into three disjoint sets with ties broken arbitrarily:

$$\begin{aligned}\mathcal{J}_1 &= \{(h, k) \subseteq [H] \times [K] : \sigma_{h,k}^2 \in \{\sigma_{\min}^2, d^3 H \cdot E_{h,k}, U_{h,k}\}\}, \\ \mathcal{J}_2 &= \{(h, k) \subseteq [H] \times [K] : \sigma_{h,k}^2 = c_0^{-2} b_{h,k}^2\}, \\ \mathcal{J}_3 &= \left\{ (h, k) \subseteq [H] \times [K] : \sigma_{h,k}^2 = \left( \frac{W}{\sqrt{c_1 d}} + \mathcal{H} d^{2.5} H \right) b_{h,k} \right\}.\end{aligned}$$

For simplicity, we denote  $z_{h,k} := \frac{b_{h,k}}{\sigma_{h,k}} = \max\{w_{h,k}, \tilde{w}_{h,k}\}$ . Therefore,

$$\sum_{k=1}^K \sum_{h=1}^H b_{h,k} = \sum_{(h,k) \in \mathcal{I}} \sigma_{h,k} z_{h,k} = \sum_{i=1}^3 \sum_{(h,k) \in \mathcal{J}_i} \sigma_{h,k} z_{h,k}. \quad (\text{F.14})$$

Recall that  $\kappa = d \log \left( 1 + \frac{K}{d\lambda\sigma_{\min}^2} \right)$ , we have  $\sum_{(h,k) \in \mathcal{I}} z_{h,k}^2 \leq 4H\kappa$ . This is because

$$\begin{aligned}\sum_{(h,k) \in \mathcal{I}} z_{h,k}^2 &\leq \sum_{(h,k) \in \mathcal{I}} (w_{h,k}^2 + \tilde{w}_{h,k}^2) \stackrel{(a)}{=} \sum_{k=1}^K \sum_{h=1}^H \min\{1, w_{h,k}^2\} + \sum_{k=1}^K \sum_{h=1}^H \min\{1, \tilde{w}_{h,k}^2\} \\ &\stackrel{(b)}{\leq} 4Hd \log \left( 1 + \frac{K}{d\lambda\sigma_{\min}^2} \right) = 4H\kappa.\end{aligned}$$

where (a) uses  $z_{h,k} \leq c_0 \leq 1$  due to  $\sigma_{h,k} \geq c_0^{-1} b_{h,k}$ ,  $c_0 \leq 1$  and (b) uses Lemma G.5. We will frequently use the above inequality.

Now, we are ready to analyze the three terms in the RHS of equation F.14 respectively.

- For the first term, it follows that

$$\begin{aligned}\sum_{(h,k) \in \mathcal{J}_1} \sigma_{h,k} z_{h,k} &\leq \sqrt{\sum_{(h,k) \in \mathcal{J}_1} \sigma_{h,k}^2} \sqrt{\sum_{(h,k) \in \mathcal{J}_1} z_{h,k}^2} \\ &\leq \sqrt{\sum_{(h,k) \in \mathcal{J}_1} (\sigma_{\min}^2 + d^3 H \cdot E_{h,k} + J_{h,k})} \sqrt{\sum_{(h,k) \in \mathcal{J}_1} z_{h,k}^2} \\ &\leq \sqrt{\sum_{(h,k) \in \mathcal{I}} (\sigma_{\min}^2 + d^3 H \cdot E_{h,k} + J_{h,k})} \sqrt{\sum_{(h,k) \in \mathcal{I}} z_{h,k}^2} \\ &\leq \sqrt{HK\sigma_{\min}^2 + \sum_{(h,k) \in \mathcal{I}} (d^3 H \cdot E_{h,k} + J_{h,k})} \cdot \sqrt{4H\kappa}.\end{aligned}$$

We provide an upper bound for  $\sum_{k=1}^K \sum_{h=1}^H E_{h,k}$  in Lemma F.2 whose proof is deferred in Appendix F.11.1.

**Lemma F.2** (Sum of  $E_{h,k}$ ). On the event  $\mathcal{B}_0 \cap \mathcal{A}_0$ ,

$$\sum_{k=1}^K \sum_{h=1}^H E_{h,k} = \mathcal{O} \left( (\beta_0 + H\beta) \mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + H^2 \mathcal{H}^2 \log \frac{4 \lceil \log_2 HK \rceil}{\delta} \right).$$

where  $\mathcal{O}(\cdot)$  hides universal positive constants.

We also provide an upper bound for  $\sum_{k=1}^K \sum_{h=1}^H J_{h,k}$  in Lemma F.3 whose proof is deferred in Appendix F.11.2.

**Lemma F.3** (Sum of  $J_{h,k}$ ). Recall that  $J_{h,k} = [\hat{\mathbf{V}}_h R_h + \hat{\mathbf{V}}_h \bar{\mathbf{V}}_{h+1}^k](s_{h,k}, a_{h,k}) + R_{h,k} + U_{h,k}$  with  $R_{h,k}, U_{h,k}$  defined in equation A.8 and equation A.9 respectively. On the event  $\mathcal{B}_R \cap \mathcal{B}_V \cap \mathcal{B}_0 \cap \mathcal{A}_0$ , with probability at least  $1 - 2\delta$ ,

$$\sum_{k=1}^K \sum_{h=1}^H J_{h,k} = \mathcal{O} \left( \mathcal{G}^* K + [(\beta_0 + H\beta) \mathcal{H} + \beta_{R^2}] \sum_{k=1}^K \sum_{h=1}^H b_{h,k} + H^2 \mathcal{H}^2 \log \frac{4 \lceil \log_2 HK \rceil}{\delta} + H\sigma_R^2 \log \frac{1}{\delta} \right).$$

where  $\mathcal{G}^*$  is defined in equation 3.2 and  $\mathcal{O}(\cdot)$  hides universal positive constants.

Putting pieces together and using  $\sqrt{a+b+c} \leq \sqrt{a} + \sqrt{b} + \sqrt{c}$ , we have

$$\begin{aligned} \sum_{(h,k) \in \mathcal{J}_1} b_{h,k} &= \sum_{(h,k) \in \mathcal{J}_1} \sigma_{h,k} z_{h,k} = \mathcal{O} \left( \sqrt{H\kappa} \cdot \sqrt{K(H\sigma_{\min}^2 + \mathcal{G}^*)} \right) \\ &\quad + \mathcal{O} \left( \sqrt{H\kappa} \cdot \sqrt{H^3 d^3 \mathcal{H}^2 \log \frac{4[\log_2 HK]}{\delta} + H\sigma_R^2 \log \frac{1}{\delta}} \right) \\ &\quad + \mathcal{O} \left( \sqrt{H\kappa} \cdot \sqrt{[(\beta_0 + H\beta)\mathcal{H}d^3 H + \beta_{R^2}] \sum_{(h,k) \in \mathcal{I}} b_{h,k}} \right). \end{aligned} \quad (\text{F.15})$$

- For the second term, due to  $\sigma_{h,k} = c_0^{-1} b_{h,k}$ , we have  $z_{h,k} = b_{h,k}/\sigma_{h,k} = c_0 \leq 1$  for all  $(h,k) \in \mathcal{J}_2$ . Hence,

$$\begin{aligned} \sum_{(h,k) \in \mathcal{J}_2} b_{h,k} &= \sum_{(h,k) \in \mathcal{J}_2} \sigma_{h,k} z_{h,k} = \frac{1}{c_0} \sum_{(h,k) \in \mathcal{J}_2} \sigma_{h,k} z_{h,k}^2 \leq \frac{\sup_{(h,k) \in \mathcal{I}} \sigma_{h,k}}{c_0} \sum_{(h,k) \in \mathcal{J}_2} z_{h,k}^2 \\ &\leq \sup_{(h,k) \in \mathcal{I}} \frac{\max\{\|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}\}}{c_0^2} \cdot \sum_{(h,k) \in \mathcal{I}} z_{h,k}^2 \leq \frac{4H\kappa}{c_0^2 \sqrt{\lambda}} \end{aligned} \quad (\text{F.16})$$

where the last inequality uses  $\|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} \leq \frac{1}{\sqrt{\lambda}} \|\phi_{h,k}\| \leq \frac{1}{\sqrt{\lambda}}$  and  $\|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}} \leq \frac{1}{\sqrt{\lambda}} \|\tilde{\phi}_{h,k}\| \leq \frac{1}{\sqrt{\lambda}}$  for any  $(h,k) \in \mathcal{I}$ .

- For the third term,  $\sigma_{h,k}^2 = \left(\frac{W}{\sqrt{c_1 d}} + \mathcal{H}d^{2.5}H\right) b_{h,k}$  and thus  $\sigma_{h,k} = \left(\frac{W}{\sqrt{c_1 d}} + \mathcal{H}d^{2.5}H\right) z_{h,k}$ . Hence,

$$\begin{aligned} \sum_{(h,k) \in \mathcal{J}_3} b_{h,k} &= \sum_{(h,k) \in \mathcal{J}_3} \sigma_{h,k} z_{h,k} = \left(\frac{W}{\sqrt{c_1 d}} + \mathcal{H}d^{2.5}H\right) \sum_{(h,k) \in \mathcal{J}_3} z_{h,k}^2 \\ &\leq \left(\frac{W}{\sqrt{c_1 d}} + \mathcal{H}d^{2.5}H\right) \sum_{(h,k) \in \mathcal{I}} z_{h,k}^2 \leq 4H\kappa \cdot \left(\frac{W}{\sqrt{c_1 d}} + \mathcal{H}d^{2.5}H\right). \end{aligned} \quad (\text{F.17})$$

Combing equation F.15, equation F.16 and equation F.17, we have

$$\sum_{(h,k) \in \mathcal{I}} b_{h,k} = \mathcal{O} \left( C + \sqrt{H\kappa} \sqrt{[(\beta_0 + H\beta)\mathcal{H}d^3 H + \beta_{R^2}] \cdot \sum_{(h,k) \in \mathcal{I}} b_{h,k}} \right)$$

where

$$\begin{aligned} C &= \sqrt{H\kappa} \cdot \sqrt{K(H\sigma_{\min}^2 + \mathcal{G}^*)} + H\kappa \cdot \left( \frac{W}{\sqrt{c_1 d}} + \frac{1}{c_0^2 \sqrt{\lambda}} + \mathcal{H}d^{2.5}H \right) \\ &\quad + \sqrt{H\kappa} \cdot \sqrt{H^3 d^3 \mathcal{H}^2 \log \frac{4[\log_2 HK]}{\delta} + H\sigma_R^2 \log \frac{1}{\delta}}. \end{aligned}$$

Using the inequality that  $x \leq 2(a^2 + b^2)$  for any  $x \leq |a|\sqrt{x} + b^2$ , we have

$$\sum_{(h,k) \in \mathcal{I}} b_{h,k} = \mathcal{O} \left( C + H^2 \mathcal{H} \kappa d^3 (\beta_0 + H\beta) + H\kappa \beta_{R^2} \right).$$

In the following, we are going to simplify the last inequality. We will use  $\tilde{\mathcal{O}}(\cdot)$  to hide logarithmic factors for simplicity. Notice that  $\kappa = \tilde{\mathcal{O}}(d)$ . By setting  $\lambda = \frac{1}{\mathcal{H}^2 + W^2}$ , we have  $\beta_R = \beta_V = \tilde{\mathcal{O}}(\sqrt{d})$  and thus

$\beta = \beta_V + \beta_R = \tilde{\mathcal{O}}(\sqrt{d})$ . Moreover,  $\beta_{R^2} = \tilde{\mathcal{O}}\left(\sqrt{d} + \sqrt{d} \frac{\sigma_{R^2}}{\sigma_{\min}} + \sqrt{\lambda} W\right) = \tilde{\mathcal{O}}\left(\sqrt{d} + \sqrt{d} \frac{\sigma_{R^2}}{\sigma_{\min}}\right)$  and  $\beta_0 = \tilde{\mathcal{O}}\left(\frac{\sqrt{d^3 H \mathcal{H}}}{\sigma_{\min}} + \sqrt{d \lambda} \mathcal{H}\right) = \tilde{\mathcal{O}}\left(\frac{\sqrt{d^3 H \mathcal{H}}}{\sigma_{\min}} + \sqrt{d}\right)$ . Therefore,

$$\begin{aligned} \sum_{(h,k) \in \mathcal{I}} b_{h,k} &= \mathcal{O}\left(C + H^2 \mathcal{H} \kappa d^3 (\beta_0 + H\beta) + H\kappa \beta_{R^2}\right) \\ &= \mathcal{O}(C) + \tilde{\mathcal{O}}\left(\frac{H^{2.5} d^{5.5} \mathcal{H}^2 + H d^{1.5} \sigma_{R^2}}{\sigma_{\min}} + H^3 d^{4.5} \mathcal{H} + H d^{1.5}\right). \end{aligned}$$

We then analyze  $C$ . Using  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for non-negative numbers  $a, b \geq 0$ , we have

$$C = \tilde{\mathcal{O}}\left(\sqrt{d H K \mathcal{G}^*} + H d^{0.5} K^{0.5} \sigma_{\min} + H^2 d^{3.5} \mathcal{H} + H^2 d^2 \mathcal{H} + H d^{0.5} \sigma_R + H d\right).$$

Putting the results together, we have

$$\sum_{(h,k) \in \mathcal{I}} b_{h,k} = \tilde{\mathcal{O}}\left(\sqrt{d H K \mathcal{G}^*} + H d^{0.5} K^{0.5} \sigma_{\min} + \frac{H^{2.5} d^{5.5} \mathcal{H}^2 + H d^{1.5} \sigma_{R^2}}{\sigma_{\min}} + H^3 d^{4.5} \mathcal{H} + H d^{0.5} \sigma_R + H d^{1.5}\right).$$

□

### F.11.1 Proof of Lemma F.2

*Proof of Lemma F.2.* By the definition of  $E_{h,k}$  in equation A.7, it follows that

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H E_{h,k} &\leq \sum_{k=1}^K \sum_{h=1}^H \left[ 2\mathcal{H} \beta_0 \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \mathcal{H} \cdot \left[ \hat{\mathbb{P}}_{h,k}(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k) \right](s_{h,k}, a_{h,k}) \right] \\ &\stackrel{(a)}{\leq} \sum_{k=1}^K \sum_{h=1}^H \left[ 4\mathcal{H} \beta_0 \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + \mathcal{H} \cdot \left[ \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k) \right](s_{h,k}, a_{h,k}) \right] \\ &\stackrel{(b)}{\leq} (4\beta_0 + 16H\beta) \mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 38H^2 \mathcal{H}^2 \log \frac{4[\log_2 H K]}{\delta} \\ &= \mathcal{O}\left((\beta_0 + H\beta) \mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + H^2 \mathcal{H}^2 \log \frac{4[\log_2 H K]}{\delta}\right) \end{aligned}$$

where (a) uses  $\left| \left[ (\hat{\mathbb{P}}_{h,k} - \mathbb{P}_h) \bar{V}_{h+1}^k \right](s_{h,k}, a_{h,k}) \right| = |\langle \phi_{h,k}, (\mu_{h,k-1} - \mu_h^*) \bar{V}_{h+1}^k \rangle| \leq \beta_0 \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}$  on  $\mathcal{B}_0$  and (b) follows from Lemma D.9. □

### F.11.2 Proof of Lemma F.3

*Proof of Lemma F.3.* By Lemma D.3, on the event  $\mathcal{B}_R$ , we have  $\left| [\hat{\mathbb{V}}_h \hat{R}_h - \mathbb{V}_h R_h](s_{h,k}, a_{h,k}) \right| \leq R_{h,k}$  for all  $h \in [H]$  and  $k \in [K]$ . By Lemma D.5, on the event  $\mathcal{B}_0 \cap \mathcal{B}_V$ ,  $[\hat{\mathbb{V}}_h \bar{V}_{h+1}^k](s_{h,k}, a_{h,k}) \leq [\mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k}) + U_{h,k}$  for all  $h \in [H]$  and  $k \in [K]$ . Therefore,

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H J_{h,k} &\leq \sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k}) + 2 \sum_{k=1}^K \sum_{h=1}^H R_{h,k} + 2 \sum_{k=1}^K \sum_{h=1}^H U_{h,k} \\ &:= (I) + (II) + (III). \end{aligned}$$

For the term (III), we have

$$\sum_{k=1}^K \sum_{h=1}^H U_{h,k} = \sum_{k=1}^K \sum_{h=1}^H \left[ 11\mathcal{H} \beta_0 \cdot \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 4\mathcal{H} \cdot \hat{\mathbb{P}}_{h,k}(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \right]$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} \sum_{k=1}^K \sum_{h=1}^H \left[ 19\mathcal{H}\beta_0 \cdot \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 4\mathcal{H} \cdot \mathbb{P}_h(\bar{V}_{h+1}^k - \underline{V}_{h+1}^k)(s_{h,k}, a_{h,k}) \right] \\
&\stackrel{(b)}{\leq} (19\beta_0 + 64H\beta)\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 152H^2\mathcal{H}^2 \log \frac{4\lceil \log_2 HK \rceil}{\delta}, \tag{F.18}
\end{aligned}$$

where (a) uses  $\left|[(\hat{\mathbb{P}}_{h,k} - \mathbb{P}_h)\bar{V}_{h+1}^k](s_{h,k}, a_{h,k})\right| = |\langle \phi_{h,k}, (\mu_{h,k-1} - \mu_h^*)\bar{V}_{h+1}^k \rangle| \leq \beta_0 \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}$  on  $\mathcal{B}_0$ ; and (b) follows from Lemma D.9.

For the term (II), we have

$$\sum_{k=1}^K \sum_{h=1}^H R_{h,k} = \beta_{R^2} \sum_{k=1}^K \sum_{h=1}^H \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{k-1}} + 2\mathcal{H}\beta_R \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}. \tag{F.19}$$

We provide two ways to analyze the term (I).

- On one hand, we denote  $X_k = \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k})$  for simplicity. Let  $\mathcal{G}_k := \mathcal{F}_{H,k}$  be the  $\sigma$ -field generated by all the random variables over the first  $k$  episodes. Then  $\pi_k$  is  $\mathcal{G}_{k-1}$ -measurable,  $X_k \geq 0$  is  $\mathcal{G}_k$ -measurable, and  $|X_k| \leq H(\sigma_R^2 + \mathcal{H}^2)$ . Therefore,  $|X_k - \mathbb{E}[X_k|\mathcal{G}_{k-1}]| \leq H(\sigma_R^2 + \mathcal{H}^2)$  and  $\text{Var}[X_k|\mathcal{G}_{k-1}] \leq H(\sigma_R^2 + \mathcal{H}^2) \cdot \mathbb{E}[X_k|\mathcal{G}_{k-1}]$ . By the variance-aware Freedman inequality in Lemma C.2, with probability at least  $1 - \delta$ , we have

$$\begin{aligned}
\sum_{k=1}^K X_k &\leq \sum_{k=1}^K \mathbb{E}[X_k|\mathcal{F}_{k-1}] + 3\sqrt{H(\sigma_R^2 + \mathcal{H}^2) \sum_{k=1}^K \mathbb{E}[X_k|\mathcal{G}_{k-1}] \log \frac{2\lceil \log_2 K \rceil}{\delta}} \\
&\quad + 5H(\sigma_R^2 + \mathcal{H}^2) \log \frac{2\lceil \log_2 K \rceil}{\delta} \\
&\leq 3 \sum_{k=1}^K \mathbb{E}[X_k|\mathcal{F}_{k-1}] + 7H(\sigma_R^2 + \mathcal{H}^2) \log \frac{2\lceil \log_2 K \rceil}{\delta}.
\end{aligned}$$

Notice that

$$\begin{aligned}
\mathbb{E}[X_k|\mathcal{F}_{k-1}] &= \mathbb{E} \left[ \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s_{h,k}, a_{h,k}) \middle| \mathcal{G}_{k-1} \right] \\
&= \sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^{\pi_k}} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a)
\end{aligned}$$

where  $d_h^{\pi_k}(s, a) = \mathbb{P}^{\pi_k}(s_h = s, a_h = a | s_0 = s_{1,k})$  is the probability reaching  $(s_{h,k}, a_{h,k}) = (s, a)$  at the  $h$ -th step when the agent starts from  $s_{1,k}$  and follows the policy  $\pi_k$ . Therefore, we have

$$\begin{aligned}
(I) &\leq 3 \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^{\pi_k}} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^*](s, a) + 7H(\sigma_R^2 + \mathcal{H}^2) \log \frac{2\lceil \log_2 K \rceil}{\delta} \\
&\leq 3\mathcal{G}_0^* K + 7H(\sigma_R^2 + \mathcal{H}^2) \log \frac{2\lceil \log_2 K \rceil}{\delta}
\end{aligned}$$

where

$$\mathcal{G}_0^* = \frac{1}{K} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{(s,a) \sim d_h^{\pi_k}} [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^{\pi_k}](s, a).$$

- On the other hand, we have

$$(I) = \sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h V_{h+1}^* - \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k}) + \sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k})$$

$$\begin{aligned}
& \stackrel{\text{equation F.20}}{\leq} 2\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) + \sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k}) \\
& \leq 2\mathcal{H} \cdot \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) + 2\mathcal{V}^2 K + 2H(\sigma_R^2 + \mathcal{H}^2) \log \frac{1}{\delta} \\
& \leq 2\mathcal{V}^2 K + 2H(\sigma_R^2 + \mathcal{H}^2) \log \frac{1}{\delta} + 16H\beta\mathcal{H} \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 76H^2\mathcal{H}^2 \log \frac{4[\log_2 HK]}{\delta} \\
& \leq 2\mathcal{V}^2 K + 16H\beta\mathcal{H} \sum_{k=1}^K \sum_{h=1}^H \|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}} + 78H^2\mathcal{H}^2 \log \frac{4[\log_2 HK]}{\delta} + 2H\sigma_R^2 \log \frac{1}{\delta}
\end{aligned}$$

where the first inequality uses equation F.20, the second inequality uses Lemma F.4, and the third inequality uses Lemma D.8.

$$\begin{aligned}
[\mathbb{V}_h V_{h+1}^* - \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k}) &= \mathbb{P}_h[V_{h+1}^*]^2(s_{h,k}, a_{h,k}) - [\mathbb{P}_h V_{h+1}^*(s_{h,k}, a_{h,k})]^2 \\
&\quad - (\mathbb{P}_h[V_{h+1}^{\pi_k}]^2(s_{h,k}, a_{h,k}) - [\mathbb{P}_h V_{h+1}^{\pi_k}(s_{h,k}, a_{h,k})]^2) \\
&\stackrel{(a)}{\leq} \mathbb{P}_h[V_{h+1}^*]^2(s_{h,k}, a_{h,k}) - \mathbb{P}_h[V_{h+1}^{\pi_k}]^2(s_{h,k}, a_{h,k}) \\
&\stackrel{(b)}{\leq} 2\mathcal{H} \cdot \mathbb{P}_h(V_{h+1}^* - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k}) \\
&\stackrel{(c)}{\leq} 2\mathcal{H} \cdot \mathbb{P}_h(\bar{V}_{h+1}^k - V_{h+1}^{\pi_k})(s_{h,k}, a_{h,k})
\end{aligned} \tag{F.20}$$

where (a) uses  $V_{h+1}^*(\cdot) \geq V_{h+1}^{\pi_k}(\cdot)$ , (b) uses  $V_{h+1}^{\pi_k}(\cdot) \leq V_{h+1}^*(\cdot) \leq \mathcal{H}$ , and (c) uses Lemma D.4.

Finally, we are going to put pieces together. In order to simplicity notation, we use  $b_{h,k} = \max\{\|\phi_{h,k}\|_{\mathbf{H}_{h,k-1}^{-1}}, \|\tilde{\phi}_{h,k}\|_{\tilde{\mathbf{H}}_{h,k-1}^{-1}}\}$  and  $\beta = \beta_V + \beta_R$ . From the first bullet point, we have

$$\sum_{k=1}^K \sum_{h=1}^H J_{h,k} = \mathcal{O} \left( \mathcal{G}_0^* \cdot K + [(\beta_0 + H\beta)\mathcal{H} + \beta_{R^2}] \cdot \sum_{k=1}^K \sum_{h=1}^H b_{h,k} + H^2\mathcal{H}^2 \log \frac{4[\log_2 HK]}{\delta} + H\sigma_R^2 \log \frac{1}{\delta} \right).$$

From the second bullet point, we have

$$\sum_{k=1}^K \sum_{h=1}^H J_{h,k} = \mathcal{O} \left( \mathcal{V}^2 K + [(\beta_0 + H\beta)\mathcal{H} + \beta_{R^2}] \sum_{k=1}^K \sum_{h=1}^H b_{h,k} + H^2\mathcal{H}^2 \log \frac{4[\log_2 HK]}{\delta} + H\sigma_R^2 \log \frac{1}{\delta} \right).$$

Taking minimum of the last two inequalities and using  $\min\{\mathcal{G}_0^*, \mathcal{V}^2\} \leq \mathcal{G}^*$  complete the proof.  $\square$

### F.11.3 Proof of Lemma F.4

**Lemma F.4** (Total variance lemma). With probability at least  $1 - \delta$ , we have

$$\sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k}) \leq 2\mathcal{V}^2 K + 2H(\sigma_R^2 + \mathcal{H}^2) \log \frac{1}{\delta}.$$

*Proof of Lemma F.4.* The proof uses a similar argument as Lemma C.5 in (Jin et al., 2018). Notice that the first state  $s_{1,k}$  is fixed and  $a_{h,k} = \pi_h^k(s_{h,k})$ . Therefore,  $(s_{2,k}, \dots, s_{H,k})$  is a sequence generated by following policy  $\pi_k$  starting at  $s_{1,k}$ . Let  $\mathcal{G}_k$  be the  $\sigma$ -field generated by all the random variables over the first  $k$  episodes.  $X_k = \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k})$ . We have the following properties about  $X_k$ . Clearly  $\pi_k$  is  $\mathcal{G}_{k-1}$ -measurable,  $X_k \geq 0$  is  $\mathcal{G}_k$ -measurable, and  $|X_k| \leq H(\sigma_R^2 + \mathcal{H}^2)$ .

Let  $\mathbb{E}_k(\cdot) := \mathbb{E}[\cdot | \mathcal{G}_k]$  for simplicity.

$$\mathcal{V}^2 \geq \mathbb{E}_{k-1} \left[ \sum_{h=1}^H R_h(s_{h,k}, a_{h,k}) - V_1^{\pi_k}(s_{1,k}) \right]^2$$



$$\begin{aligned}
&\stackrel{(a)}{=} \mathbb{E}_{k-1} \left[ \sum_{h=1}^H (R_h(s_{h,k}, a_{h,k}) + V_{h+1}^{\pi_k}(s_{h+1,k}) - V_h^{\pi_k}(s_{h,k})) \right]^2 \\
&\stackrel{(b)}{=} \sum_{h=1}^H \mathbb{E}_{k-1} [R_h(s_{h,k}, a_{h,k}) + V_{h+1}^{\pi_k}(s_{h+1,k}) - V_h^{\pi_k}(s_{h,k})]^2 \\
&\stackrel{(c)}{=} \sum_{h=1}^H \mathbb{E}_{k-1} \left[ [R_h - r_h]^2(s_{h,k}, a_{h,k}) + [r_h(s_{h,k}, a_{h,k}) + V_{h+1}^{\pi_k}(s_{h+1,k}) - V_h^{\pi_k}(s_{h,k})]^2 \right] \\
&\stackrel{(d)}{=} \mathbb{E}_{k-1} \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k}) = \mathbb{E}[X_k | \mathcal{F}_{k-1}]
\end{aligned}$$

where (a) uses  $V_{H+1}^{\pi_k}(\cdot) = 0$ , (b) uses the independence due to the Markov property, (c) holds since  $R_h(s_{h,k}, a_{h,k})$  is independent with  $s_{h+1,k}$  conditioning on  $(s_{h,k}, a_{h,k})$ , and (d) uses  $V_h^{\pi_k}(s_{h,k}) = r_h(s_{h,k}, a_{h,k}) + \mathbb{E}_{s_{h+1,k} \sim \mathbb{P}_h(\cdot | s_{h,k}, a_{h,k})}[V_{h+1}^{\pi_k}(s_{h+1,k})]$ . Using  $\text{Var}[X_k | \mathcal{G}_{k-1}] \leq H(\sigma_R^2 + \mathcal{H}^2) \cdot \mathbb{E}[X_k | \mathcal{G}_{k-1}]$ , we have

$$\sum_{k=1}^K \text{Var}[X_k | \mathcal{G}_{k-1}] \leq H(\sigma_R^2 + \mathcal{H}^2) \cdot \sum_{k=1}^K \mathbb{E}[X_k | \mathcal{G}_{k-1}] \leq (\sigma_R^2 + \mathcal{H}^2) \mathcal{V}^2 H K.$$

By the Freedman inequality in Lemma G.1, with probability at least  $1 - \delta$ , we have

$$\begin{aligned}
&\sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h R_h + \mathbb{V}_h V_{h+1}^{\pi_k}](s_{h,k}, a_{h,k}) \\
&= \sum_{k=1}^K X_k \leq \sum_{k=1}^K \mathbb{E}[X_k | \mathcal{F}_{k-1}] + \sqrt{2(\sigma_R^2 + \mathcal{H}^2) \mathcal{V}^2 H K \log \frac{1}{\delta}} + \frac{2}{3} H(\sigma_R^2 + \mathcal{H}^2) \log \frac{1}{\delta} \\
&\leq \mathcal{V}^2 K + 2\sqrt{\mathcal{V}^2 K \cdot H(\sigma_R^2 + \mathcal{H}^2) \log \frac{1}{\delta}} + \frac{2}{3} H(\sigma_R^2 + \mathcal{H}^2) \log \frac{1}{\delta} \\
&\leq 2\mathcal{V}^2 K + 2H(\sigma_R^2 + \mathcal{H}^2) \log \frac{1}{\delta}.
\end{aligned}$$

□

## G Auxiliary Lemmas

### G.1 Concentration Inequalities

**Lemma G.1** (Freedman inequality (Freedman, 1975)). Let  $\{X_t\}_{t \in [T]}$  be a stochastic process that adapts to the filtration  $\mathcal{F}_t$  so that  $X_t$  is  $\mathcal{F}_t$ -measurable,  $\mathbb{E}[X_t | \mathcal{F}_{t-1}] = 0$ ,  $|X_t| \leq M$  and  $\sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] \leq V$  where  $M > 0$  and  $V > 0$  are positive constants. Then with probability at least  $1 - \delta$ , we have

$$\sum_{t=1}^T X_t \leq \sqrt{2V \ln \frac{1}{\delta}} + \frac{2M}{3} \ln \frac{1}{\delta}.$$

**Lemma G.2** (Variance-aware Freedman inequality). Let  $\{X_t\}_{t \in [T]}$  be a stochastic process that adapts to the filtration  $\mathcal{F}_t$  so that  $X_t$  is  $\mathcal{F}_t$ -measurable,  $\mathbb{E}[X_t | \mathcal{F}_{t-1}] = 0$ ,  $|X_t| \leq M$  and  $\sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] \leq V^2$  where  $M > 0$  and  $V > 0$  are positive constants. Then with probability at least  $1 - \delta$ , we have

$$\left| \sum_{t=1}^T X_t \right| \leq 3\sqrt{\sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}] \cdot \log \frac{2K}{\delta}} + 5M \log \frac{2K}{\delta}$$

where  $K = 1 + \lceil 2 \log_2 \frac{V}{M} \rceil$ .

*Proof of Lemma G.2.* By Theorem 5 in (Li et al., 2021), we have for any positive integer  $K \geq 1$ ,

$$\mathbb{P} \left( \left| \sum_{t=1}^T X_t \right| \leq \sqrt{8 \max \left\{ \sum_{t=1}^T \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}], \frac{V^2}{2K} \right\} \cdot \ln \frac{2K}{\delta} + \frac{4M}{3} \ln \frac{2K}{\delta}} \right) \geq 1 - \delta.$$

By setting  $K = 1 + \lceil 2 \log_2 \frac{V}{M} \rceil$ , we have  $\frac{V^2}{2K} \leq M^2$ . Using  $\max\{a, b\} \leq a + b$ ,  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for any  $a, b \geq 0$  and  $\ln \frac{2K}{\delta} \geq 1$ , we complete the proof.  $\square$

The following two lemmas are the counterpart lemmas of Theorem 2.1 under light-tail assumption.

**Lemma G.3** (Bernstein inequality for self-normalized martingales, Lemma F.4 in (Hu et al., 2022)). Let  $\{\mathcal{G}_t\}_{t \geq 0}$  be a filtration and  $\{\mathbf{x}_t, \eta_t\}_{t \geq 0}$  be a stochastic process so that  $\mathbf{x}_t \in \mathbb{R}^d$  is  $\mathcal{G}_t$ -measurable and  $\eta_t \in \mathbb{R}$  is  $\mathcal{G}_{t+1}$ -measurable. If  $\|\mathbf{x}_t\| \leq L$  and  $\{\eta_t\}_{t \geq 1}$  satisfies that  $\mathbb{E}[\eta_t | \mathcal{G}_t] = 0$ ,  $\mathbb{E}[\eta_t^2 | \mathcal{G}_t] \leq \sigma^2$  and  $|\eta_t \min\{1, \|\mathbf{x}_t\|_{\mathbf{Z}_{t-1}^{-1}}\}| \leq M$  for all  $t \geq 1$ . Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have for all  $t \geq 1$ ,

$$\left\| \sum_{j=1}^t \mathbf{x}_j \eta_j \right\|_{\mathbf{Z}_t^{-1}} \leq 8\sigma \sqrt{d \log \left( 1 + \frac{tL^2}{d\lambda} \right) \log \frac{4t^2}{\delta}} + 4M \log \frac{4t^2}{\delta}$$

where  $\mathbf{Z}_t = \lambda \mathbf{I} + \sum_{j=1}^t \mathbf{x}_j \mathbf{x}_j^\top$  for  $t \geq 1$  and  $\mathbf{Z}_0 = \lambda \mathbf{I}$ .

**Lemma G.4** (Hoeffding inequality for self-normalized martingales, Theorem 1 in (Abbasi-Yadkori et al., 2011)). Let  $\{\mathcal{G}_t\}_{t \geq 0}$  be a filtration and  $\{\mathbf{x}_t, \eta_t\}_{t \geq 0}$  be a stochastic process so that  $\mathbf{x}_t \in \mathbb{R}^d$  is  $\mathcal{G}_t$ -measurable and  $\eta_t \in \mathbb{R}$  is  $\mathcal{G}_{t+1}$ -measurable. If  $\|\mathbf{x}_t\| \leq L$  and  $\{\eta_t\}_{t \geq 1}$  satisfies that  $\mathbb{E}[\eta_t | \mathcal{G}_t] = 0$  and  $|\eta_t| \leq M$  for all  $t \geq 1$ . Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have for all  $t \geq 1$ ,

$$\left\| \sum_{j=1}^t \mathbf{x}_j \eta_j \right\|_{\mathbf{Z}_t^{-1}} \leq M \sqrt{d \log \left( 1 + \frac{tL^2}{d\lambda} \right) + \log \frac{1}{\delta}}$$

where  $\mathbf{Z}_t = \lambda \mathbf{I} + \sum_{j=1}^t \mathbf{x}_j \mathbf{x}_j^\top$  for  $t \geq 1$  and  $\mathbf{Z}_0 = \lambda \mathbf{I}$ .

## G.2 Elliptical Lemmas

**Lemma G.5** (Lemma 11 in (Abbasi-Yadkori et al., 2011)). Let  $\{\mathbf{x}_t\}_{t \geq 1} \subset \mathbb{R}^d$  and assume  $\|\mathbf{x}_t\| \leq L$  for all  $t \geq 1$ . Set  $\mathbf{Z}_t = \sum_{s=1}^t \mathbf{x}_s \mathbf{x}_s^\top + \lambda \mathbf{I}$ . Then it follows that

$$\sum_{t=1}^T \min \left\{ 1, \|\mathbf{x}_t\|_{\mathbf{Z}_{t-1}}^2 \right\} \leq 2d \log \left( \frac{d\lambda + TL^2}{d\lambda} \right).$$

**Lemma G.6** (Lemma 12 in (Abbasi-Yadkori et al., 2011)). Suppose  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$  are two positive definite matrices satisfying that  $\mathbf{A} \succeq \mathbf{B}$ , then for any  $\mathbf{x} \in \mathbb{R}^d$ ,

$$\|\mathbf{x}\|_{\mathbf{B}^{-1}} \leq \|\mathbf{x}\|_{\mathbf{A}^{-1}} \sqrt{\frac{\det(\mathbf{A})}{\det(\mathbf{B})}}.$$

## G.3 Function Class and Covering Number

This subsection collects important lemmas in (He et al., 2022). Let  $\mathcal{K} = \{k_1, k_2, \dots\}$  denote the set of episodes where the algorithm updates the value function in Algorithm 3. For a given total number of episodes  $K$ , it definitely follows that  $|\mathcal{K}| \leq K$ . Furthermore, due to the mechanism of rare-switching value function updates,  $|\mathcal{K}|$  is actually much smaller than  $K$ .

**Lemma G.7.**

$$|\mathcal{K}| \leq dH \log_2 \left( 1 + \frac{K}{\lambda \sigma_{\min}^2} \right).$$

*Proof of Lemma G.7.* The proof is almost identical to Lemma E.1 in (He et al., 2022) except that we maintain the dependence on  $\sigma_{\min}$ . According to the determinant-based criterion, for each episode  $k_i$ , there exists a stage  $h' \in [H]$  such that  $\det(\mathbf{H}_{h',k_i-1}) \geq 2\det(\mathbf{H}_{h,k_i-1-1})$ . Since we always have  $\mathbf{H}_{h,k_i-1} \geq \mathbf{H}_{h,k_i-1-1}$  for all  $h \in [H]$ , it then follows that

$$\prod_{h \in [H]} \det(\mathbf{H}_{h,k_i-1}) \geq 2 \prod_{h \in [H]} \det(\mathbf{H}_{h,k_i-1-1}).$$

By induction, it follows that

$$\prod_{h \in [H]} \det(\mathbf{H}_{h,k_{|\mathcal{K}|}-1}) \geq 2^{|\mathcal{K}|} \prod_{h \in [H]} \det(\mathbf{H}_{h,k_1-1}) \geq 2^{|\mathcal{K}|} \prod_{h \in [H]} \det(\lambda \mathbf{I}) = 2^{|\mathcal{K}|} \lambda^{dH}$$

On the other hand, due to  $\mathbf{H}_{h,k_{|\mathcal{K}|}-1} \leq \mathbf{H}_{h,K}$  the determinant  $\det(\mathbf{H}_{h,k_{|\mathcal{K}|}-1})$  is upper bounded by

$$\prod_{h \in [H]} \det(\mathbf{H}_{h,k_{|\mathcal{K}|}-1}) \leq \prod_{h \in [H]} \det(\mathbf{H}_{h,K}) \leq \left( \lambda + \frac{K}{\sigma_{\min}^2} \right)^{dH}.$$

Combining the last two inequalities, we have

$$|\mathcal{K}| \leq dH \log_2 \left( 1 + \frac{K}{\lambda \sigma_{\min}^2} \right).$$

□

The optimistic value function  $\bar{V}_h^k(\cdot) = \min_{k_i \leq k} \max_a \bar{Q}_h^{k_i}(\cdot, a)$  belong to the function class  $\mathcal{V}^+$

$$\mathcal{V}^+ = \left\{ f | f(\cdot) = \max_{a \in \mathcal{A}} \min_{i \leq |\mathcal{K}|} \min \left\{ \mathbf{w}_i^\top \phi(\cdot, a) + \beta \|\phi(\cdot, a)\|_{\mathbf{H}_i^{-1}, \mathcal{H}} \right\}, \beta \in [0, B], \|\mathbf{w}_i\| \leq L, \mathbf{H}_i \geq \lambda \mathbf{I} \right\}. \quad (\text{G.1})$$

while the pessimistic value function  $\underline{V}_h^k(\cdot) = \max_{k_i \leq k} \max_a \underline{Q}_h^{k_i}(\cdot, a)$  belong to the function class  $\mathcal{V}^-$ ,

$$\mathcal{V}^- = \left\{ f | f(\cdot) = \max_{a \in \mathcal{A}} \max_{i \leq |\mathcal{K}|} \max \left\{ \mathbf{w}_i^\top \phi(\cdot, a) - \beta \|\phi(\cdot, a)\|_{\mathbf{H}_i^{-1}, \mathcal{H}} \right\}, \beta \in [0, B], \|\mathbf{w}_i\| \leq L, \mathbf{H}_i \geq \lambda \mathbf{I} \right\}. \quad (\text{G.2})$$

Here  $B$  upper bounds  $\beta$  and  $L = W + \mathcal{H} \sqrt{\frac{dK}{\lambda}}$  is a uniformly bound for  $\boldsymbol{\theta}_{h,k-1} + \boldsymbol{\mu}_{h,k-1} \bar{\mathbf{V}}_{h+1}^k$  because

$$\|\boldsymbol{\theta}_{h,k-1} + \boldsymbol{\mu}_{h,k-1} \bar{\mathbf{V}}_{h+1}^k\| \leq \|\boldsymbol{\theta}_{h,k-1}\| + \|\boldsymbol{\mu}_{h,k-1} \bar{\mathbf{V}}_{h+1}^k\| \leq W + \mathcal{H} \sqrt{\frac{dK}{\lambda}}$$

where the last inequality uses the boundedness of  $\boldsymbol{\theta}_{h,k-1}$ 's and the inequality  $\|\boldsymbol{\mu}_{h,k-1} \bar{\mathbf{V}}_{h+1}^k\| \leq \mathcal{H} \sqrt{\frac{dK}{\lambda}}$  (whose proof can be found in Lemma E.2 of He et al. (2022)).

**Lemma G.8** (Covering number of value functions). Let  $\mathcal{V}^\pm$  denote the class of optimistic or pessimistic value functions with definition in equation G.1 and equation G.2 respectively. Assume  $\|\phi(s, a)\| \leq 1$  for all  $(s, a)$  pairs, and let  $\mathcal{N}(\mathcal{V}, \varepsilon)$  be the  $\varepsilon$ -covering number of  $\mathcal{V}$  with respect to the distance  $\text{dist}(f, f') := \sup_{s \in \mathcal{S}} |f(s) - f'(s)|$ . Then,

$$\log \mathcal{N}(\mathcal{V}^\pm, \varepsilon) \leq \left[ d \log \left( 1 + \frac{4L}{\varepsilon} \right) + d^2 \log \left( 1 + \frac{8d^{1/2}B^2}{\lambda \varepsilon^2} \right) \right] \cdot |\mathcal{K}|.$$

*Proof of Lemma G.8.* The result about  $\mathcal{V}_f^+$  follows from Lemma E.6 in (He et al., 2022). The result about  $\mathcal{V}_f^-$  follows from Lemma E.7 in (He et al., 2022). □

**Lemma G.9** (Covering number of squared functions, Lemma E.8 in (He et al., 2022)). For the squared function class  $[\mathcal{V}^+]^2 := \{f^2 | f \in \mathcal{V}^+\}$ , let  $\mathcal{N}([\mathcal{V}^+]^2, \varepsilon)$  be the  $\varepsilon$ -covering number of  $[\mathcal{V}^+]^2$  with respect to the distance  $\text{dist}(f, f') := \sup_{s \in \mathcal{S}} |f(s) - f'(s)|$ . Then

$$\log \mathcal{N}([\mathcal{V}^+]^2, \varepsilon) \leq \left[ d \log \left( 1 + \frac{8HL}{\varepsilon} \right) + d^2 \log \left( 1 + \frac{32d^{1/2}H^2B^2}{\lambda \varepsilon^2} \right) \right] \cdot |\mathcal{K}|.$$