

---

# Improved Variance-Aware Confidence Sets for Linear Bandits and Linear Mixture MDP

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

This paper presents new *variance-aware* confidence sets for linear bandits and linear mixture Markov Decision Processes (MDPs). With the new confidence sets, we obtain the follow regret bounds:

- For linear bandits, we obtain an  $\tilde{O}(\text{poly}(d)\sqrt{1 + \sum_{k=1}^K \sigma_k^2})$  data-dependent regret bound, where  $d$  is the feature dimension,  $K$  is the number of rounds, and  $\sigma_k^2$  is the *unknown* variance of the reward at the  $k$ -th round. This is the first regret bound that only scales with the variance and the dimension but *no explicit polynomial dependency on  $K$* . When variances are small, this bound can be significantly smaller than the  $\tilde{\Theta}(d\sqrt{K})$  worst-case regret bound.
- For linear mixture MDPs, we obtain an  $\tilde{O}(\text{poly}(d, \log H)\sqrt{K})$  regret bound, where  $d$  is the number of base models,  $K$  is the number of episodes, and  $H$  is the planning horizon. This is the first regret bound that only scales *logarithmically* with  $H$  in the reinforcement learning with linear function approximation setting, thus *exponentially improving* existing results, and resolving an open problem in Zhou et al. [2020a].

We develop three technical ideas that may be of independent interest: 1) applications of the peeling technique to both the input norm and the variance magnitude, 2) a recursion-based estimator for the variance, and 3) a new convex potential lemma that generalizes the seminal elliptical potential lemma.

## 1 Introduction

In sequential decision-making problems such as bandits and reinforcement learning (RL), the agent chooses an action based on the current state, with the goal to maximize the total reward. When the state-action space is large, function approximation is often used for generalization. One of the most fundamental and widely used methods is linear function approximation.

For (infinite-actioned) linear bandits, the minimax-optimal regret bound is  $\tilde{\Theta}(d\sqrt{K})$  [Dani et al., 2008, Abbasi-Yadkori et al., 2011], where  $d$  is the feature dimension and  $K$  is the number of total rounds played by the agent.<sup>1</sup> However, oftentimes the worst-case analysis is overly pessimistic, and it is possible to obtain data-dependent bound that is substantially smaller than  $\tilde{O}(d\sqrt{K})$  in benign scenarios.

One direction to study is the variance magnitude. As a motivating example, in linear bandits, if there is no noise (variance is 0), one only needs to pay at most  $d$  regret to identify the best action because  $d$  samples are sufficient to recover the underlying linear coefficients (in general position).

---

<sup>1</sup>We follow the reinforcement learning convention to use  $K$  to denote the total number of rounds / episodes.

33 This constant-type regret bound is much smaller than the  $\sqrt{K}$ -type regret bound in the worst case  
 34 where the variance magnitude is a lower bounded constant. Therefore, a natural question is:

35 **Can we design an algorithm that adapts to the variance magnitude, and its regret degrades**  
 36 **gracefully from the benign noiseless constant-type bound to the worst-case  $\sqrt{K}$ -type bound?**

37 In RL, exploiting the variance information is also important. For tabular RL, one needs to utilize  
 38 the variance information, e.g., Bernstein-type exploration bonus to achieve the minimax optimal  
 39 regret [Azar et al., 2017, Zanette and Brunskill, 2019, Zhang et al., 2020c,a, Menard et al., 2021, Dann  
 40 et al., 2019]. For example, the recently proposed MVP algorithm [Zhang et al., 2020a], enjoys an  
 41  $\tilde{O}(\text{polylog}(H) \times (\sqrt{SAK} + S^2A))$  regret bound, where  $S$  is the number of states,  $A$  is the number  
 42 of actions,  $H$  is the planning horizon, and  $K$  is the total number of episodes.<sup>23</sup> Notably, this regret  
 43 bound only scales *logarithmically* with  $H$ . On the other hand, without using the variance information,  
 44 e.g., using Hoeffding-type bonus instead of Bernstein-type bonus, algorithms would suffer a regret  
 45 that scales *polynomially* with  $H$  [Azar et al., 2017].

46 Going beyond tabular RL, a recent line of work studied RL with linear function approximation  
 47 with different assumptions [Yang and Wang, 2019, Modi et al., 2020, Jin et al., 2020, Ayoub et al.,  
 48 2020, Zhou et al., 2020a, Modi et al., 2020]. Our paper studies the linear mixture Markov Decision  
 49 Process (MDP) setting [Modi et al., 2020, Ayoub et al., 2020, Zhou et al., 2020a], where the transition  
 50 probability can be represented by a linear function of some features or base models. This model-based  
 51 assumption is motivated by problems in robotics and queuing systems. We refer readers to Ayoub  
 52 et al. [2020] for more discussions.

53 For this linear mixture MDP setting, previous works can obtain regret bounds in the form  
 54  $\tilde{O}(\text{poly}(d, H)\sqrt{K})$ , where  $d$  is the number of base models. While these bounds do not scale  
 55 with  $SA$ , they scale *polynomially* with  $H$ , because the algorithms in previous works do not use the  
 56 variance information. In practice,  $H$  is often large, and even a polynomial dependency on  $H$  may not  
 57 be acceptable. Therefore, a natural question is

58 **Can we design an algorithm that exploits the variance information to obtain an**  
 59  **$\tilde{O}(\text{poly}(d, \log H)\sqrt{K})$  regret bound for linear mixture MDP?**

## 60 1.1 Our Contributions

61 In this paper, we develop new, *variance-aware* confidence sets for linear bandits and linear mixture  
 62 MDP and answer the above two questions affirmatively.

63 **Linear Bandits.** For linear bandits, we obtain an  $\tilde{O}(\text{poly}(d)\sqrt{1 + \sum_{k=1}^K \sigma_k^2})$  regret bound, where  
 64  $\sigma_k^2$  is the *unknown* variance at the  $k$ -th round. To our knowledge, this is the first bound that solely  
 65 depends on the variance and the feature dimension, and has no explicit polynomial dependency  
 66 on  $K$ . When the variance is very small so that  $\sigma_k^2 \ll 1$ , this bound is substantially smaller than  
 67 the worst-case  $\tilde{\Theta}(d\sqrt{K})$  bound. Furthermore, this regret bound naturally interpolates between the  
 68 worst-case  $\sqrt{K}$ -type bound and the noiseless-case constant-type bound.

69 **Linear Mixture MDP.** For linear mixture MDP, we obtain the desired  $\tilde{O}(\text{poly}(d, \log H)\sqrt{K})$   
 70 regret bound. This is the first regret bound in RL with function approximation that 1) does not scale  
 71 with the size of the state-action space, and 2) only scales *logarithmically* with the planning horizon  
 72  $H$ . Therefore, we exponentially improve existing results on RL with linear function approximation  
 73 in term of the  $H$  dependency, and resolve an open problem in [Zhou et al., 2020a]. More importantly,  
 74 our result conveys the positive conceptual message for RL: it is possible to simultaneously overcome  
 75 the two central challenges in RL, *large state-action space* and *long planning horizon*.

<sup>2</sup> $\tilde{O}(\cdot)$  hides logarithmic factors. Sometimes we write out  $\text{polylog } H$  explicitly to emphasize the logarithmic  
 dependency on  $H$ .

<sup>3</sup>This bound holds for setting where the transition is homogeneous and the total reward is bounded by 1. We  
 focus on this setting in this paper. See Section 2 and 3 for more discussions.

## 1.2 Main Difficulties and Technical Innovations

We first describe limitations of existing works why they cannot achieve the desired regret bounds described above.

**Limitations of Existing Variance-Aware Confidence Sets** Faury et al. [2020], Zhou et al. [2020a] applied Bernstein-style inequalities to construct a confidence sets of the least square estimator for linear bandits. However, their methods can not be applied directly to obtain the desired data-dependent regret bound.

We give a simple example to illustrate their limitations. Consider the case where the variance is always  $\sigma^2 \ll 1$ . Let  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{k-1}, y_{k-1})$  be the samples collected before the  $k$ -th round. Their confidence set at the  $k$ -th round is  $\Theta_k = \{\boldsymbol{\theta} \mid \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_k\|_{\Lambda_{k-1}} \leq C(\sigma\sqrt{d} + 1 + \lambda^{1/2})\}$  (See In Equation (4.3) of Zhou et al. [2020a] and Theorem 1 of Faury et al. [2020]). where  $\Lambda_{k-1} = \sum_{\tau=1}^{k-1} \mathbf{x}_\tau \mathbf{x}_\tau^\top + \lambda \mathbf{I}$  is the un-normalized covariance matrix,  $\hat{\boldsymbol{\theta}}_k = \Lambda_{k-1}^{-1} \sum_{\tau=1}^{k-1} y_\tau \mathbf{x}_\tau$  is the estimated linear coefficients by least squares,  $\lambda$  is a regularization parameter and  $C$  is a constant. Consider the case  $d = 1$  and  $\mathbf{x}_k = \sqrt{1/K}$  for  $k = 1, \dots, K$ . Their regret bound is roughly

$$\sum_{k=1}^K (\sigma\sqrt{d} + 1 + \lambda^{1/2}) \|\mathbf{x}_k\|_{\Lambda_k^{-1}} \geq (1 + \lambda^{1/2}) \sum_{i=1}^K \|\mathbf{x}_k\|_{\Lambda_k^{-1}} \geq (1 + \lambda^{1/2}) \sqrt{\frac{K}{1 + \lambda}} \geq \sqrt{K},$$

which is much larger than our bound,  $O(\sqrt{K\sigma^2 + 1})$  when  $\sigma$  is very small. For more detailed discussion, please refer to Appendix B.

Below we describe our main techniques.

**Elimination with Peeling.** Instead of using least squares and upper-confidence-bound (UCB), we use an elimination approach. More precisely, for the underlying linear coefficients  $\boldsymbol{\theta}^* \in \mathbb{R}^d$ , we build a confidence interval for  $(\boldsymbol{\theta}^*)^\top \boldsymbol{\mu}$  for every  $\boldsymbol{\mu}$  in an  $\epsilon$ -net of the  $d$ -dimensional unit ball, and we eliminate  $\boldsymbol{\theta} \in \mathbb{R}^d$  if  $\boldsymbol{\theta}^\top \boldsymbol{\mu}$  fails to fall in the confidence interval of  $\boldsymbol{\theta}^* \boldsymbol{\mu}$  for some  $\boldsymbol{\mu}$ . To build the confidence intervals, we use 1) an empirical Bernstein inequality (cf. Theorem 4) and 2) the peeling technique to both the input norm and the variance magnitude. As will be clear in the proof (cf. Section D), this peeling step is crucial to obtain a tight regret bound for the example above. The new confidence region provides a tighter estimation for  $\boldsymbol{\theta}^*$ , which helps address the drawback in least squares.

**Generalization of the Elliptical Potential Lemma.** Since we use the peeling technique which comes with a clipping operation, we cannot use the seminal elliptic potential lemma Dani et al. [2008] any more. Instead, we propose a more general lemma below, which provides a bound of potential for a general class of convex functions though with a worse dependency on  $d$  than the bound in the elliptical potential lemma. We believe this lemma can be applied to other problems as well.

**Lemma 1** (Generalized Quadratic Potential Lemma). *Let  $f(x) \geq 0$  be a convex function over  $\mathbb{R}$  such that  $\frac{f(x)}{x^2} \leq \frac{f(y)}{y^2} \leq 1$  and  $f(x) \geq f(y)$  if  $x^2 \geq y^2 > 0$ . Fix  $\ell \in (0, 1]$ . For any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t \in \mathbb{B}_2^d(1)$  and  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_t \in \mathbb{B}_2^d(1)$ , we have that*

$$\sum_{i=1}^t \min \left\{ \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^{i-1} f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2}, 1 \right\} \leq O(d^4 \log(dt/\ell)). \quad (1)$$

Note that by choosing  $f(x) = x^2$  and  $\boldsymbol{\mu}_i = \frac{\mathbf{x}_i \Lambda_i^{-1}}{\|\mathbf{x}_i \Lambda_i^{-1}\|}$  with  $\Lambda_i = \sum_{j=1}^{i-1} \mathbf{x}_j \mathbf{x}_j^\top + \ell \mathbf{I}$ , Lemma 1 reduces to the classical elliptic potential lemma [Dani et al., 2008]. Our proof consists of two major parts. We first establish a symmetric version of Equation (1) using rearrangement inequality, and then bound the number of times the energy for some  $\boldsymbol{\mu}$  (i.e.,  $\sum_{j=1}^i f(\mathbf{x}_j \boldsymbol{\mu}) + \ell^2$ ) doubles. The full proof is deferred to Appendix C.

For linear mixture MDP, we propose another technique to further reduce the dependency on  $d$ .

**Recursion-based Variance Estimation.** In linear bandits, generally it is not possible to estimate the variance because the variance at each round can arbitrarily different. On the other hand, for linear mixture MDP, the variance is a quadratic function of the underlying l coefficient  $\theta^*$ . Furthermore, the higher moments are polynomial functions of  $\theta^*$ . Utilizing this rich structure and leveraging the recursion idea in previous analyses on tabular RL [Lattimore and Hutter, 2012, Li et al., 2020, Zhang et al., 2020a], we explicitly estimate the variance and higher moments to further reduce the regret. See Section 5 for more explanations.

## 2 Related Work

**Linear Bandits.** There is a line of theoretical analyses of linear bandits problems [Auer et al., 2002, Dani et al., 2008, Chu et al., 2011, Abbasi-Yadkori et al., 2011, Li et al., 2019a,b]. For infinite-actioned linear bandits, the minimax regret bound is  $\tilde{\Theta}(d\sqrt{K})$ . and recent works tried to give fine-grained instance-dependent bounds [Katz-Samuels et al., 2020, Jedra and Proutiere, 2020]. For multi-armed bandits, Audibert et al. [2006] showed by exploiting the variance information, one can improve the regret bound. For linear bandits, only a few work studied how to use the variance information. Faury et al. [2020] studied logistic bandit problem with adaptivity to the variance of noise, where a Bernstein-style confidence set was proposed. However, they assume the variance is known and cannot attain the desired variance-dependent bound due to the example we gave above. Linear bandits can be also seen as a simplified version of RL with linear function approximation, where the planning horizon degenerates to  $H = 1$ .

**RL with Linear Function Approximation.** Recently, it is a central topic in the theoretical RL community to figure out the necessary and sufficient conditions that permit efficient learning in RL with large state-action space [Wen and Van Roy, 2013, Jiang et al., 2017, Yang and Wang, 2019, 2020, Du et al., 2019b, 2020a, 2019a, 2020b, Jiang et al., 2017, Feng et al., 2020, Sun et al., 2019, Dann et al., 2018, Krishnamurthy et al., 2016, Misra et al., 2019, Ayoub et al., 2020, Zanette et al., 2020, Wang et al., 2019, 2020c,b, Jin et al., 2020, Weisz et al., 2020, Modi et al., 2020, Shariff and Szepesvári, 2020, Jin et al., 2020, Cai et al., 2019, He et al., 2020, Zhou et al., 2020a]. However, to our knowledge, all existing regret upper bounds have a polynomial dependency on the planning horizon  $H$ , except works that assume the environment is deterministic [Wen and Van Roy, 2013, Du et al., 2020b].

This paper studies the linear mixture MDP setting [Ayoub et al., 2020, Zhou et al., 2020b,a, Modi et al., 2020], which assumes the underlying transition is a linear combination of some known base models. Ayoub et al. [2020] gave an algorithm, UCRL-VTR, with an  $\tilde{O}(dH^2\sqrt{K})$  regret in the time-inhomogeneous model.<sup>4</sup> Our algorithm improves the  $H$ -dependency from  $\text{poly}(H)$  to  $\text{polylog}(H)$ , at the cost of a worse dependency on  $d$ .

**Variance Information in Tabular MDP.** The use of the variance information in tabular MDP was first proposed by Lattimore and Hutter [2012] in the discounted MDP setting, and was later adopted in the episodic MDP setting [Azar et al., 2017, Jin et al., 2018, Zanette and Brunskill, 2019, Dann et al., 2019, Zhang et al., 2020a,b]. This technique is crucial to tighten the dependency on  $H$ .

**Concurrent Work by Zhou et al. [2020a].** While preparing this draft, we noticed a concurrent work by Zhou et al. [2020a], who also studied how to use the variance information for linear bandits and linear mixture MDPs. We first compare their results with ours. For linear bandits, they proved an  $\tilde{O}(\sqrt{dK} + d\sqrt{\sum_{i=1}^K \sigma_i^2})$  regret bound, while we prove an  $\tilde{O}(d^{4.5}\sqrt{\sum_{i=1}^K \sigma_i^2} + d^5)$  regret bound. Our bound has a worse dependency on  $d$ , but in the regime where  $K$  is very large and the sum of the variances is small, our bound is stronger. Furthermore, they assumed *the variance is known while we do not need this assumption*. For linear mixture MDP, they proved an

<sup>4</sup>The time-inhomogeneous model refers to the setting where the transition probability can vary at different levels, and the time-homogeneous model refers to the setting where the transition probability is the same at different levels. Roughly speaking, the model complexity of the time-inhomogeneous model is  $H$  times larger than that of the time-homogeneous model. In general, it is straightforward to tightly extend a result for the time-homogeneous model to the time-inhomogeneous model by extending the state-action space [Jin et al., 2018, Footnote 2], but not vice versa.

161  $\tilde{O}(\sqrt{d^2 H + d H^2} \sqrt{K} + d^2 H^2 + d^3 H)$  bound for the time-inhomogeneous model, while we prove an  
 162  $\tilde{O}(d^{4.5} \sqrt{K} + d^5) \times \text{polylog}(H)$  bound for the time-homogeneous model. Their bound has a better  
 163 dependency on  $d$  than ours and is near-optimal in the regime  $K = \Omega(\text{poly}(d, H))$  and  $H = O(d)$ .  
 164 On the other hand, we have an exponentially better dependency on  $H$  in the time-homogeneous  
 165 model. Indeed, obtaining a regret bound that is logarithmic in  $H$  (in the time-homogeneous model)  
 166 was raised as an open question in their paper [Zhou et al., 2020a, Remark 5.5].

167 Next, we compare the algorithms and the analyses. The algorithms in the two papers are very different  
 168 in nature: ours are based on elimination while theirs are based on least squares and UCB. We note  
 169 that, for linear bandits, their current analysis cannot give a  $\sqrt{K}$ -free bound because there is a term  
 170 that scales *inversely* with the variance. This can be seen by plugging the first line of their (B.25) to  
 171 their (B.23). For the same reason, they cannot give a horizon-free bound in the time-homogeneous  
 172 linear mixture MDP. In sharp contrast, our analysis does not have the term depending on the inverse  
 173 of the variance. On the other hand, their algorithms are computationally efficient (given certain  
 174 computation oracles), but our algorithms are not because ours are elimination-based. See Section 6  
 175 for more discussions.

### 176 3 Preliminaries

177 **Notations.** We use  $\mathbb{B}_p^d(r) = \{x \in \mathbb{R}^d : \|x\|_p \leq r\}$  to denote the  $d$ -dimensional  $\ell_p$ -ball of radius  $r$ .  
 178 For any set  $S \subseteq \mathbb{R}^d$ , we use  $\partial S$  to denote its boundary. For  $N \in \mathbb{N}$ , we define  $[N] = \{1, \dots, N\}$ .  
 179 One important operation used in our algorithms and analyses is clipping. Given  $\ell > 0$  and  $u \in \mathbb{R}$ , we  
 180 define

$$\text{clip}(u, \ell) = \min\{|u|, \ell\} \cdot \frac{u}{|u|}$$

181 for  $u \neq 0$  and  $\text{clip}(0, \ell) = 0$ . For any two vectors  $\mathbf{u}, \mathbf{v}$ , to save notations, we use  $\mathbf{u}\mathbf{v} = \mathbf{u}^\top \mathbf{v}$  to  
 182 denote their inner product when no ambiguity.

183 **Linear Bandits.** We use  $K$  to denote the number of rounds in the linear bandits. At each round  $k =$   
 184  $1, \dots, K$ , the algorithm is first given the context set  $\mathcal{A}_k \subseteq \mathbb{B}_2^d(1)$ , then the algorithm chooses an action  
 185  $\mathbf{x}_k \in \mathcal{A}_k$  and receives the noisy reward  $r_k = \mathbf{x}_k \boldsymbol{\theta}^* + \varepsilon_k$ , where  $\boldsymbol{\theta}^* \in \mathbb{B}_2^d(1)$  is the unknown underlying  
 186 linear coefficients and  $\varepsilon_k$  is the random noise. We define  $\mathcal{F}_k = \sigma(\mathbf{x}_1, \varepsilon_1, \dots, \mathbf{x}_k, \varepsilon_k, \mathbf{x}_{k+1})$ . We  
 187 assume that  $|r_k| \leq 1$  and that the noise  $\varepsilon_k$  satisfies  $\mathbb{E}[\varepsilon_k | \mathcal{F}_k] = 0$  and  $\mathbb{E}[\varepsilon_k^2 | \mathcal{F}_k] = \sigma_k^2$ . The goal  
 188 is to learn  $\boldsymbol{\theta}^*$  and minimize the cumulative expected regret  $\mathbb{E}[\mathfrak{R}^K]$ , where

$$\mathfrak{R}^K = \sum_{k=1}^K [\max_{\mathbf{x} \in \mathcal{A}_k} \mathbf{x} \boldsymbol{\theta}^* - \mathbf{x}_k \boldsymbol{\theta}^*].$$

189 **Remark 1.** Here we assume the reward is uniformly bounded ( $|r_k| \leq 1$ ) instead of 1-sub-Gaussian  
 190 commonly used in the literature only for the ease of presentation, because in RL, it is standard to  
 191 assume bounded reward. Note if the noise is 1-sub-Gaussian, our algorithm also applies with only  
 192 an  $O(\log T)$  overhead because a problem with 1-sub-Gaussian noise can be reduced to that with  
 193 uniformly bounded noise by clipping the noise with a threshold  $O(\log T)$ .

194 **Episodic MDP and Linear Mixture MDP.** We use a tuple  $(\mathcal{S}, \mathcal{A}, r, P, K, H)$  to define an episodic  
 195 finite-horizon MDP. Here,  $\mathcal{S}$  is its state space,  $\mathcal{A}$  is its action space,  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is its reward  
 196 function,  $P(s' | s, a)$  is the transition probability from the state-action pair  $(s, a)$  to the new state  
 197  $s'$ ,  $K$  is the number of episodes, and  $H$  is the planning horizon of each episode. Without the loss of  
 198 generality, we assume a fixed initial state  $s_1$ . A sequence of functions  $\pi = \{\pi_h : \mathcal{S} \rightarrow \Delta(\mathcal{A})\}_{h=1}^H$  is  
 199 an policy, where  $\Delta(\mathcal{A})$  denotes the set of all possible distributions over  $\mathcal{A}$ .

200 At each episode  $k = 1, \dots, K$ , the algorithm outputs a policy  $\pi^k$ , which is then executed on the  
 201 MDP by  $a_h^k \sim \pi_h^k(s_h^k)$ ,  $s_{h+1}^k \sim P(\cdot | s_h^k, a_h^k)$ . We let  $r_h^k = r(s_h^k, a_h^k)$  be the reward at time step  $h$  in  
 202 episode  $k$ . Importantly, we assume the transition model  $P(\cdot | \cdot, \cdot)$  is time-homogeneous, which is  
 203 necessary to bypass the  $\text{poly}(H)$  dependency. We assume that the reward function is known, which  
 204 is standard in the theoretical RL literature to simplify the presentation [Modi et al., 2020, Ayoub et al.,  
 205 2020]. We let  $\pi^*$  to denote the optimal policy which achieves the maximum reward in expectation.

206 We make the following regularity assumption on the rewards: the sum of reward,  $\sum_{h=1}^H r_h$ , in each  
 207 episode is bounded by 1.

---

**Algorithm 1** VOFUL: Variance-Aware Optimism in the Face of Uncertainty for Linear Bandits

---

- 1: **Initialize:**  $\ell_i = 2^{2-i}$ ,  $\iota = 16d \ln \frac{dK}{\delta}$ ,  $L_2 = \lceil \log_2 K \rceil$ ,  $\Lambda_2 = \{1, 2, \dots, L_2 + 1\}$ ,  $\Theta_1 = \mathbb{B}_2^d(1)$ ,  
Let  $\mathcal{B}$  be an  $K^{-3}$ -net of  $\mathbb{B}_2^d(2)$  with size not larger than  $(\frac{4}{K})^{3d}$
- 2: **for**  $k = 1, 2, \dots, K$  **do**
- 3:   **Optimistic Action Selection:**
- 4:   Observe context set  $\mathcal{A}_k \subseteq \mathbb{B}_2^d(1)$
- 5:   Compute  $\mathbf{x}_k \leftarrow \arg \max_{\mathbf{x} \in \mathcal{A}_k} \max_{\boldsymbol{\theta} \in \Theta_k} \mathbf{x}\boldsymbol{\theta}$ , choose action  $\mathbf{x}_k$
- 6:   Receive feedback  $y_k$
- 7:   **Construct Confidence Set:**
- 8:   For each  $\boldsymbol{\theta} \in \mathbb{B}_2^d(1)$ , define  $\epsilon_k(\boldsymbol{\theta}) = y_k - \mathbf{x}_k\boldsymbol{\theta}$ ,  $\eta_k(\boldsymbol{\theta}) = (\epsilon_k(\boldsymbol{\theta}))^2$ .
- 9:   Define confidence set  $\Theta_{k+1} = \bigcap_{j \in \Lambda_2} \Theta_{k+1}^j$ , where

$$\Theta_{k+1}^j = \left\{ \boldsymbol{\theta} \in \mathbb{B}_2^d(1) : \left| \sum_{v=1}^k \text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}) \epsilon_v(\boldsymbol{\theta}) \right| \leq \sqrt{\sum_{v=1}^k \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}) \eta_v(\boldsymbol{\theta}) \iota + \ell_j \iota}, \forall \boldsymbol{\mu} \in \mathcal{B} \right\} \quad (2)$$

and  $\text{clip}_j(\cdot) = \text{clip}(\cdot, \ell_j)$ .

10: **end for**

---

208 **Assumption 2** (Non-uniform reward).  $\sum_{h=1}^H r_h^k \leq 1$  almost surely for any policy  $\pi^k$ .

209 This assumption is much weaker than the common assumption where the reward at each time step is  
210 bounded by  $1/H$  (uniform reward) because Assumption 2 allows one spiky reward as large as  $\Omega(1)$ .  
211 See more discussions about this reward scaling in Jiang and Agarwal [2018], Wang et al. [2020a],  
212 Zhang et al. [2020a].

213 For any policy  $\pi$ , we define its  $H$ -step  $V$ -function and  $Q$ -function as

$$V_h^\pi(s) = \max_{a \in \mathcal{A}} Q_h^\pi(s, a)$$

$$\text{where } Q_h^\pi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} V_{h+1}^\pi(s') \text{ for } h = 1, \dots, H$$

214 where we set  $V_{H+1} = 0$ . For simplicity, we also denote  $V^\pi(s_1) = V_1^\pi(s_1)$  and  $V^*(s_1) = V^{\pi^*}(s_1)$ .

215 A linear mixture MDP is an episodic MDP with the extra assumption that its transition model is an  
216 unknown linear combination of a known set of models. Specifically, there is an unknown parameter  
217  $\boldsymbol{\theta}^* \in \mathbb{B}_1^d(1)$ , such that  $P = \sum_{i=1}^d \theta_i^* P_i$  where based models  $P_1, \dots, P_d$  are given. The goal is to  
218 learn  $\boldsymbol{\theta}^*$  and minimize the cumulative expected regret  $\mathbb{E}[\mathfrak{R}^K]$ , where

$$\mathfrak{R}^K = \sum_{k=1}^K [V^*(s_1) - V^k(s_1)].$$

## 219 4 Algorithm and Theory for Linear Bandits

220 In this section, we introduce our algorithm for linear bandits and analyze its regret. The pseudo-  
221 code is listed in Algorithm 1. The following theorem shows our algorithm achieves the desired  
222 variance-dependent regret bound. The proof is deferred to Section D.

223 **Theorem 3.** *The expected regret of Algorithm 1 is bounded by  $\mathbb{E}[\mathfrak{R}^K] \leq \tilde{O}(d^{4.5} \sqrt{\sum_{k=1}^K \sigma_k^2} + d^5)$ .*

224 This theorem shows our algorithm's regret has no explicit polynomial dependency on the number  
225 of rounds  $K$ . In the worst-case where the variance is  $\Omega(1)$ , our bound becomes  $\tilde{O}(d^{4.5} \sqrt{K} + d^5)$ ,  
226 which has a worse dependency on  $d$  compared with the minimax optimal algorithms [Dani et al.,  
227 2008, Abbasi-Yadkori et al., 2011]. However, in the benign case where the variance is  $o(1)$ , our  
228 bound can be much smaller. In particular, in the noiseless case, our bound is a constant-type regret  
229 bound, up to logarithmic factors. One future direction is to design an algorithm that is minimax  
230 optimal in the worst-case but also adapts to the variance magnitude like ours.

231 Now we describe our algorithm. Similar to the existing linear bandits algorithms, our algorithm  
232 maintains a confidence set for the underlying parameter  $\boldsymbol{\theta}^*$ . The confidence set  $\Theta_k$  is updated at

each round, and we choose the action greedily according to the confidence set. Note that existing confidence sets either do not exploit variance information [Dani et al., 2008, Abbasi-Yadkori et al., 2011], or require the variance to be known and do not fully exploit the variance information [Zhou et al., 2020a, Faury et al., 2020] as their regret bounds still have an  $\tilde{O}(\sqrt{K})$  term.

To relax the known variance assumption, we use the following empirical Bernstein inequality that depends on the *empirical variance*, in contrast to the Bernstein inequality that depends on the *true variance*, which was used in existing works [Zhou et al., 2020b, Faury et al., 2020].

**Theorem 4.** *Let  $\{\mathcal{F}_i\}_{i=0}^n$  be a filtration. Let  $\{X_i\}_{i=1}^n$  be a sequence of real-valued random variables such that  $X_i$  is  $\mathcal{F}_i$ -measurable. We assume that  $\mathbb{E}[X_i | \mathcal{F}_{i-1}] = 0$  and that  $|X_i| \leq b$  almost surely. For  $\delta < e^{-1}$ , we have*

$$\Pr \left[ \left| \sum_{i=1}^n X_i \right| \leq 8 \sqrt{\sum_{i=1}^n X_i^2 \ln \frac{1}{\delta}} + 16b \ln \frac{1}{\delta} \right] \geq 1 - 6\delta \log_2 n. \quad (3)$$

Importantly, this inequality controls the deviation via the empirical variance, which is  $X_i^2$  and can be computed once  $X_i$  is known. Note some previously proved inequalities require certain independence assumptions and thus cannot be directly applied to martingales [Maurer and Pontil, 2009, Peel et al., 2013], so they cannot be used for solving our linear bandits problem. The proof of the theorem is deferred to Appendix D.2.

Much more effort is devoted to designing a confidence set that fully exploits the variance information. Note Theorem 4 is for real-valued random variables, and it remains unclear how to generalize it to the linear regression setting, which is crucial for building confidence sets for linear bandits. Previous works built up their confidence sets based on analyzing the convergence of the ordinary ridged least square estimator [Dani et al., 2008, Abbasi-Yadkori et al., 2011], or the weighted one [Zhou et al., 2020a].

We drop the least square estimators and instead, we take a testing-based approach, as done in Equation (2). To illustrate the idea, we first ignore the  $\text{clip}_j(\cdot)$  operation and  $\ell_j$  terms. We define the noise function  $\epsilon_k(\theta)$  and the variance function  $\eta_k(\theta)$  (Line 8 of Algorithm 1). Note that  $\epsilon_k(\theta^*) = \varepsilon_k$  and  $\eta_k(\theta^*) = \varepsilon_k^2$ , so we have the following fact: if  $\theta = \theta^*$ , then Equation (3) would be true if we replace  $X_k = w_k(\mu)\epsilon_k(\theta)$  and  $X_k^2 = w_k^2(\mu)\eta_k(\theta)$  with high probability, where  $\{w_k(\mu)\}$  is a proper sequence of weights depending on the test direction  $\mu$ . Our approach uses the fact in the opposite direction: if weighted  $w_k(\mu)\epsilon_k(\theta)$ ,  $w_k^2(\mu)\eta_k(\theta)$  satisfies Equation (3) for all possible test directions  $\mu$ , then we put  $\theta$  into the confidence set.

Given the test direction  $\mu$ , following the least square estimation,  $w_k(\mu)$  is set to be  $x_k\mu$ . However, with  $w_k(\mu) = x_k\mu$ , the right-hand-side of Equation (3) is at least  $b \geq \max_{1 \leq k \leq n} |w_k(\mu)| = \max_{1 \leq k \leq n} |x_k\mu|$ , which might be dominant compared with  $\sum_{k=1}^n w_k^2(\mu)\eta_k(\theta)$  (See Appendix B for a toy example). To address this problem, we consider to peel  $w_k(\mu)$  for various thresholds of difference level. More precisely, we construct confidence regions respectively with  $w_k^j(\mu) = \text{clip}_j(x_k\mu)$ , where  $l_j = 2^{2-j}$  for  $j = 1, 2, \dots, \lceil \log_2 K \rceil$ . At last, we define the final confidence region as the intersections of all these confidence regions.

**Proof Sketch.** Now we explain how our confidence set enables us to obtain a variance-dependent regret bound. We define  $\theta_k = \arg \max_{\theta \in \Theta_k} x_k(\theta - \theta^*)$  and  $\mu_k = \theta_k - \theta^*$ . Then our goal is to bound the regret  $\sum_k x_k\mu_k$ . Our main idea is to consider  $\{x_k\}, \{\mu_k\}$  as two sequences of vectors. We decouple the complicated dependency between  $\{x_k\}$  and  $\{\mu_k\}$  by a union bound over the net  $\mathcal{B}$  (defined in Line 1 of Algorithm 1). To bound the regret, we implicitly divide all rounds  $k \in [K]$  into norm layers based on  $\log_2 |x_k\mu_k|$  in the analysis.<sup>5</sup> Within each layer, we apply Equation (2) to obtain the relations between  $\mu_k$  and  $\{x_1, \dots, x_{k-1}\}$ , which would self-normalize the growth of the two sequences, ensuring that their in-layer total sum is properly bounded. Since we have logarithmically many layers, the total regret is then properly bounded. We highlight that our norm peeling technique ensures that the variance-dependent term dominates the other variance-independent term in Bernstein inequalities ( $\sqrt{\sum X_i^2} \gtrsim b$  in Theorem 4), which resolves the variance-independent term in the final regret bound obtained by Zhou et al. [2020a]. See Section D for the full proof.

<sup>5</sup>This cannot be done explicitly in the algorithm, since it would re-couple the two sequences.

---

**Algorithm 2** VARLin: Variance-Aware RL with Linear Function Approximation

---

1: **Initialize:**  $\ell_i = 2^{2-i}$ ,  $\iota = 16d \ln \frac{dHK}{\delta}$ ,  $L_0 = \lceil \log_2 KH \rceil$ ,  $L_1 = L_2 = \lceil 5 \log_2(HK) + 3 \rceil$ ,  $\Lambda_0 = \{0, 1, \dots, L_0\}$ ,  $\Lambda_1 = \{1, \dots, L_1\}$ ,  $\Lambda_2 = \{1, \dots, L_2\}$ .  $\mathcal{B}$  be an  $(HK)^{-3}$ -net of  $\mathbb{B}_1^d(2)$  with size no larger than  $(\frac{4}{HK})^{3d}$ .  $\Theta_1 = \mathbb{B}_1^d(1)$ .  
2: **for**  $k = 1, 2, \dots, K$  **do**  
3:   **Optimistic Planning:**  
4:   **for**  $h = H, H-1, \dots, 1$  **do**  
5:     For each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , let  $Q_h^k(s, a) = \min\{1, r(s, a) + \max_{\theta \in \Theta_k} \sum_{i=1}^d \theta_i P_{s,a}^i V_{h+1}^k\}$ .  
6:     For each  $s \in \mathcal{S}$ , let  $V_h^k(s) = \max_{a \in \mathcal{A}} Q_h^k(s, a)$ .  
7:   **end for**  
8:   **for**  $h = 1, 2, \dots, H$  **do**  
9:     Choose action  $a_h^k \leftarrow \arg \max_{a \in \mathcal{A}} Q_h^k(s_h^k, a)$ , observe the next state  $s_{h+1}^k$ .  
10:   **end for**  
11:   **Construct Confidence Set:**  
12:   For  $m \in \Lambda_0$ ,  $h \in [H]$ , define the input  $\mathbf{x}_{k,h}^m = [P_{s_h^k, a_h^k}^1 (V_{h+1}^k)^{2^m}, \dots, P_{s_h^k, a_h^k}^d (V_{h+1}^k)^{2^m}]^\top$ .  
13:   For  $m \in \Lambda_0$ ,  $h \in [H]$ , define the variance estimate  $\eta_{k,h}^m = \max_{\theta \in \Theta_k} \{\theta \mathbf{x}_{k,h}^{m+1} - (\theta \mathbf{x}_{k,h}^m)^2\}$ .  
14:   Denote  $\epsilon_{v,u}^m(\theta) = \theta \mathbf{x}_{v,u}^m - (V_{u+1}^v(s_{u+1}^v))^{2^m}$  for  $m \in \Lambda_0$ ,  $u \in [H]$ ,  $v \in [k-1]$   
15:   Define  $\mathcal{T}_{k+1}^{m,i} = \{(v, u) \in [k] \times [H] : \eta_{v,u}^m \in (\ell_{i+1}, \ell_i]\}$ ,  $\mathcal{T}_{k+1}^{m, L_1+1} = \{(v, u) \in [k] \times [H] : \eta_{v,u}^m \leq \ell_{L_1+1}\}$ .  
16:   Define the confidence ball  $\Theta_{k+1} = \bigcap_{m,i,j} \Theta_{k+1}^{m,i,j}$ , where

$$\Theta_{k+1}^{m,i,j} = \left\{ \theta \in \mathbb{B}_1^d(1) : \left| \sum_{(v,u) \in \mathcal{T}_{k+1}^{m,i}} \text{clip}_j(\mathbf{x}_{v,u}^m \mu) \epsilon_{v,u}^m(\theta) \right| \leq 4 \sqrt{\sum_{(v,u) \in \mathcal{T}_{k+1}^{m,i}} \text{clip}_j^2(\mathbf{x}_{v,u}^m \mu) \eta_{v,u}^m} \iota + 4\ell_j \iota, \forall \mu \in \mathcal{B} \right\} \quad (4)$$

and  $\text{clip}_j(\cdot) = \text{clip}(\cdot, \ell_j)$

17: **end for**

---

## 5 Algorithm and Theory for Linear Mixture MDP

We introduce our algorithm and the regret bound for linear mixture MDP. Its pseudo-code is listed in Algorithm 2 and its regret bound is stated in the following theorem.

**Theorem 5.** *The expected regret of Algorithm 2 is bounded by  $\mathbb{E}[\mathfrak{R}^K] \leq \tilde{O}(d^{4.5} \sqrt{K} + d^9)$ .*

To our knowledge, this is the first regret bound that only scales polynomially with the dimension ( $d$ ), and does not scale polynomially with the planning horizon  $H$ . The proof of Theorem 5 is deferred to Section E.

Before describing our algorithm, we introduce some additional notations. In this section, we assume that, unless explicitly stated, the variables  $m, i, j, k, h$  iterate over the sets  $\Lambda_0, \Lambda_1, \Lambda_2, [K], [H]$ , respectively. See Line 1 of Algorithm 2 for the definitions of these sets. For example, at Line 16 of Algorithm 2, we have  $\bigcap_{m,i,j} \Theta_{k+1}^{m,i,j} = \bigcap_{m \in \Lambda_0, i \in \Lambda_1, j \in \Lambda_2} \Theta_{k+1}^{m,i,j}$ .

The starting point of our algorithm design is from Zhang et al. [2020a], in which the authors obtained a nearly horizon-free regret bound in tabular MDP. A natural idea is to combine their proof with our results for linear bandits (cf. Section 4) and obtain a nearly horizon-free regret bound for linear mixture MDP.

Note that, however, there is one caveat for such direct combination: in Section 4, the confidence set  $\Theta_k$  is updated at a per-round level, in that  $\Theta_k$  is built using all rounds prior to  $k$ ; while for the RL setting, the confidence set  $\Theta_k$  could only be updated at a per-episode level and use all time steps prior to episode  $k$ . Were it updated at a per-time-step level, severe dependency issues would prevent us from bounding the regret properly. Such discrepancy in update frequency results in a gap between

the confidence set built using data prior to episode  $k$ , and that built using data prior to time step  $(k, h)$ . Fortunately, we are able to resolve this issue. In Lemma 22, we show that we can relate these two confidence intervals, except for  $\tilde{O}(d)$  “bad” episodes. Therefore, we could adapt the analysis in Zhang et al. [2020a] only for the not “bad” episodes, and we bound the regret by 1 for the “bad” episodes. The resulting regret bound should be  $\tilde{O}(d^{6.5}\sqrt{K})$ .

To further reduce the horizon-free regret bound to  $\tilde{O}(d^{4.5}\sqrt{K})$ , we present another novel technique. We first note an important advantage of the linear mixture MDP setting over the linear bandit setting: in the latter setting, we cannot estimate the variance because there is no structure on the variance among different actions; while in the former setting, we could estimate an upper bound of the variance, because the variance is a quadratic function of  $\theta^*$ . Therefore, we can use the peeling technique on the *variance magnitude* to reduce the regret (comparing Equation (27) and Equation (40) in appendix). We note that one can also apply this step to linear bandits if the variance can be estimated.

Along the way, we also need to bound the gap between estimated variance and true variance, which can be seen as the “regret of variance predictions.” Using the same idea, we can build a confidence set using the variance sequence  $(x^2)$ , and the regret of variance predictions can be bounded by the variance of variance, namely the 4-th moment. Still, a peeling step on the 4-th moment is required to bound the regret of variance predictions, we need to bound the gap between estimated 4-th moment and true 4-th moment, which requires predicting 8-th moment. We continue to use this idea: we estimate 2-th, 4-th, 8-th,  $\dots$ ,  $O(\log KH)$ -th moments. The index  $m$  is used for moments, and  $\Lambda_0$  is the index set reserved for moments. We note that the proof in [Zhang et al., 2020a] also depends on the higher moments. The main difference is here we estimate these higher moments explicitly.

## 6 Discussions

By incorporating the variance information in the confidence set construction, we derive the first variance-dependent regret bound for linear bandits and the nearly horizon-free regret bound for linear mixture MDP. Below we discuss limitations of our work and some future directions.

One drawback of our result is that our dependency on  $d$  is large. The main reason is our bounds rely on the convex potential lemma (Lemma 17), which is  $\tilde{O}(d^4)$ . In analogous to the elliptical potential lemma in [Abbasi-Yadkori et al., 2011], we believe that this bound can be improved to  $\tilde{O}(d)$ . This improvement will directly reduce the dependencies on  $d$  in our bounds.

Another drawback is that our method is not computationally efficient. This is a common issue in elimination-based algorithms. We note that the issue of computational tractability is common in sequential decision-making problems. We list some examples. Many algorithm for tabular problems are statistically efficient but computationally inefficient [Zhang and Ji, 2019, Wang et al., 2020a, Bartlett and Tewari, 2012]. The most statistically efficient algorithm for linear MDP, ELEANOR [Zanette et al., 2020], is not computationally efficient. Algorithms for many general frameworks on RL with function approximation are elimination-based and thus not computationally efficient [Krishnamurthy et al., 2016, Jiang et al., 2017, Sun et al., 2019, Jin et al., 2021, Du et al., 2021, Dong et al., 2020]. Fortunately, later work has made progress on computational aspects for many settings [Zhang et al., 2020a, Dann et al., 2018, Du et al., 2019a, Jin et al., 2020, Fruit et al., 2018, Wang et al., 2020c, Agarwal et al., 2020]. For now, leave it as a future direction to design computationally efficient algorithms that enjoy variance-dependent bounds for linear bandits and horizon-free bounds for linear mixture MDP.

Lastly, in this paper, we only study sequential decision-making problems with linear function approximation. It would be interesting to generalize the ideas in this paper to other settings with function approximation, such as linear MDP [Yang and Wang, 2019, Jin et al., 2020], low inherent Bellman error [Zanette et al., 2020], Eluder dimension [Wang et al., 2020c, Russo and Van Roy, 2013], and various general frameworks on RL with function approximation [Jiang et al., 2017, Sun et al., 2019, Du et al., 2021, Jin et al., 2021].

## References

Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24:2312–2320, 2011.

352 Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity  
353 and representation learning of low rank mdps. *arXiv preprint arXiv:2006.10814*, 2020.

354 Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvari. Use of variance estimation in the multi-  
355 armed bandit problem. 2006.

356 Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed  
357 bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.

358 Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin F Yang. Model-based reinforcement  
359 learning with value-targeted regression. In *Proceedings of the 37th International Conference on*  
360 *Machine Learning*, 2020.

361 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforce-  
362 ment learning. In *Proceedings of the 34th International Conference on Machine Learning*, pages  
363 263–272, 2017.

364 Kazuoki Azuma. Weighted sums of certain dependent random variables. *Tohoku Mathematical*  
365 *Journal, Second Series*, 19(3):357–367, 1967.

366 Peter L Bartlett and Ambuj Tewari. Regal: A regularization based algorithm for reinforcement  
367 learning in weakly communicating mdps. *arXiv preprint arXiv:1205.2661*, 2012.

368 Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimiza-  
369 tion. *arXiv preprint arXiv:1912.05830*, 2019.

370 Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff  
371 functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and*  
372 *Statistics*, pages 208–214, 2011.

373 Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit  
374 feedback. In *Conference on Learning Theory*, 2008.

375 Christoph Dann, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E.  
376 Schapire. On oracle-efficient PAC-RL with rich observations. In *Advances in Neural Information*  
377 *Processing Systems*, 2018.

378 Christoph Dann, Lihong Li, Wei Wei, and Emma Brunskill. Policy certificates: Towards accountable  
379 reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning*,  
380 pages 1507–1516, 2019.

381 Kefan Dong, Jian Peng, Yining Wang, and Yuan Zhou. Root-n-regret for learning in Markov decision  
382 processes with function approximation and low Bellman rank. In *Conference on Learning Theory*,  
383 pages 1554–1557. PMLR, 2020.

384 Simon S Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.  
385 Provably efficient RL with rich observations via latent state decoding. In *International Conference*  
386 *on Machine Learning*, pages 1665–1674, 2019a.

387 Simon S Du, Yuping Luo, Ruosong Wang, and Hanrui Zhang. Provably efficient Q-learning with func-  
388 tion approximation via distribution shift error checking oracle. In *Advances in Neural Information*  
389 *Processing Systems*, pages 8058–8068, 2019b.

390 Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for  
391 sample efficient reinforcement learning? In *International Conference on Learning Representations*,  
392 2020a.

393 Simon S Du, Jason D Lee, Gaurav Mahajan, and Ruosong Wang. Agnostic Q-learning with func-  
394 tion approximation in deterministic systems: Tight bounds on approximation error and sample  
395 complexity. *Advances in Neural Information Processing Systems*, 2020b.

396 Simon S Du, Sham M Kakade, Jason D Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and  
397 Ruosong Wang. Bilinear classes: A structural framework for provable generalization in RL. *arXiv*  
398 *preprint arXiv:2103.10897*, 2021.

399 Louis Faury, Marc Abeille, Clément Calauzènes, and Olivier Fercoq. Improved optimistic algorithms  
 400 for logistic bandits. In *International Conference on Machine Learning*, pages 3052–3060. PMLR,  
 401 2020.

402 Fei Feng, Ruosong Wang, Wotao Yin, Simon S Du, and Lin F Yang. Provably efficient exploration  
 403 for RL with unsupervised learning. *arXiv preprint arXiv:2003.06898*, 2020.

404 Ronan Fruit, Matteo Pirodda, Alessandro Lazaric, and Ronald Ortner. Efficient bias-span-constrained  
 405 exploration-exploitation in reinforcement learning. In *International Conference on Machine  
 406 Learning*, pages 1578–1586. PMLR, 2018.

407 Jiafan He, Dongruo Zhou, and Quanquan Gu. Logarithmic regret for reinforcement learning with  
 408 linear function approximation. *arXiv preprint arXiv:2011.11566*, 2020.

409 Yassir Jedra and Alexandre Proutiere. Optimal best-arm identification in linear bandits. *Advances in  
 410 Neural Information Processing Systems*, 33, 2020.

411 Nan Jiang and Alekh Agarwal. Open problem: The dependence of sample complexity lower bounds  
 412 on planning horizon. In *Conference On Learning Theory*, pages 3395–3398, 2018.

413 Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Con-  
 414 textual decision processes with low Bellman rank are PAC-learnable. In *Proceedings of the 34th  
 415 International Conference on Machine Learning*, pages 1704–1713, 2017.

416 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably efficient?  
 417 In *Advances in Neural Information Processing Systems*, pages 4863–4873, 2018.

418 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement  
 419 learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143,  
 420 2020.

421 Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman Eluder dimension: New rich classes of RL  
 422 problems, and sample-efficient algorithms. *arXiv preprint arXiv:2102.00815*, 2021.

423 Julian Katz-Samuels, Lalit Jain, Kevin G Jamieson, et al. An empirical process approach to the union  
 424 bound: Practical algorithms for combinatorial and linear bandits. *Advances in Neural Information  
 425 Processing Systems*, 33, 2020.

426 Akshay Krishnamurthy, Alekh Agarwal, and John Langford. PAC reinforcement learning with rich  
 427 observations. In *Advances in Neural Information Processing Systems*, pages 1840–1848, 2016.

428 Tor Lattimore and Marcus Hutter. Pac bounds for discounted mdps. In *International Conference on  
 429 Algorithmic Learning Theory*, pages 320–334. Springer, 2012.

430 Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Breaking the sample size barrier in  
 431 model-based reinforcement learning with a generative model. In *Advances in Neural Information  
 432 Processing Systems*, 2020.

433 Yingkai Li, Yining Wang, and Yuan Zhou. Nearly minimax-optimal regret for linearly parameterized  
 434 bandits. In *Conference on Learning Theory*, pages 2173–2174, 2019a.

435 Yingkai Li, Yining Wang, and Yuan Zhou. Tight regret bounds for infinite-armed linear contextual  
 436 bandits. *arXiv preprint arXiv:1905.01435*, 2019b.

437 Andreas Maurer and Massimiliano Pontil. Empirical Bernstein bounds and sample variance penaliza-  
 438 tion. In *Conference on Learning Theory*, 2009.

439 Pierre Menard, Omar Darwiche Domingues, Xuedong Shang, and Michal Valko. Ucb momentum  
 440 q-learning: Correcting the bias without forgetting. *arXiv preprint arXiv:2103.01312*, 2021.

441 Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state  
 442 abstraction and provably efficient rich-observation reinforcement learning. *arXiv preprint  
 443 arXiv:1911.05815*, 2019.

Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. In *International Conference on Artificial Intelligence and Statistics*, pages 2010–2020. PMLR, 2020.

Thomas Peel, Sandrine Anthoine, and Liva Ralaivola. Empirical bernstein inequality for martingales: Application to online learning. 2013.

Dan Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, pages 2256–2264, 2013.

Roshan Shariff and Csaba Szepesvári. Efficient planning in large mdps with weak linear function approximation. *arXiv preprint arXiv:2007.06184*, 2020.

Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based RL in contextual decision processes: PAC bounds and exponential improvements over model-free approaches. In *Conference on Learning Theory*, pages 2898–2933, 2019.

Ruosong Wang, Simon S Du, Lin F Yang, and Sham M Kakade. Is long horizon reinforcement learning more difficult than short horizon reinforcement learning? In *Advances in Neural Information Processing Systems*, 2020a.

Ruosong Wang, Simon S Du, Lin F Yang, and Ruslan Salakhutdinov. On reward-free reinforcement learning with linear function approximation. In *Advances in Neural Information Processing Systems*, 2020b.

Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Provably efficient reinforcement learning with general value function approximation. *Advances in Neural Information Processing Systems*, 2020c.

Yining Wang, Ruosong Wang, Simon S Du, and Akshay Krishnamurthy. Optimism in reinforcement learning with generalized linear function approximation. *arXiv preprint arXiv:1912.04136*, 2019.

Gellert Weisz, Philip Amortila, and Csaba Szepesvári. Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions. *arXiv preprint arXiv:2010.01374*, 2020.

Zheng Wen and Benjamin Van Roy. Efficient exploration and value function generalization in deterministic systems. In *Advances in Neural Information Processing Systems*, pages 3021–3029, 2013.

Lin Yang and Mengdi Wang. Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004, 2019.

Lin F Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *International Conference on Machine Learning*, 2020.

Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*, pages 7304–7312, 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*, 2020.

Zihan Zhang and Xiangyang Ji. Regret minimization for reinforcement learning by evaluating the optimal bias function. In *Advances in Neural Information Processing Systems*, pages 2823–2832, 2019.

Zihan Zhang, Xiangyang Ji, and Simon S Du. Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. *arXiv preprint arXiv:2009.13503*, 2020a.

Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Almost optimal model-free reinforcement learning via reference-advantage decomposition. In *Advances in Neural Information Processing Systems*, 2020b.

- 491 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Model-free reinforcement learning: from clipped  
492 pseudo-regret to sample complexity. *arXiv preprint arXiv:2006.03864*, 2020c.
- 493 Dongruo Zhou, Quanquan Gu, and Csaba Szepesvari. Nearly minimax optimal reinforcement learning  
494 for linear mixture markov decision processes. *arXiv preprint arXiv:2012.08507*, 2020a.
- 495 Dongruo Zhou, Jiafan He, and Quanquan Gu. Provably efficient reinforcement learning for discounted  
496 mdps with feature mapping. *arXiv preprint arXiv:2006.13165*, 2020b.

## 497 Checklist

- 498 1. For all authors...
- 499 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
500 contributions and scope? [Yes]
- 501 (b) Did you describe the limitations of your work? [Yes] We discuss the limitations in  
502 Section 6.
- 503 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This work  
504 is theoretical so the boarder impact does not apply.
- 505 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
506 them? [Yes]
- 507 2. If you are including theoretical results...
- 508 (a) Did you state the full set of assumptions of all theoretical results? [Yes] We present the  
509 main assumption in Section 3.
- 510 (b) Did you include complete proofs of all theoretical results? [Yes] We present the proofs  
511 in Appendix.
- 512 3. If you ran experiments...
- 513 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
514 mental results (either in the supplemental material or as a URL)? [N/A] We have no  
515 experiments.
- 516 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
517 were chosen)? [N/A]
- 518 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
519 ments multiple times)? [N/A]
- 520 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
521 of GPUs, internal cluster, or cloud provider)? [N/A]
- 522 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 523 (a) If your work uses existing assets, did you cite the creators? [N/A] We do not use  
524 existing models.
- 525 (b) Did you mention the license of the assets? [N/A]
- 526 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 527
- 528 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
529 using/curating? [N/A]
- 530 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
531 information or offensive content? [N/A]
- 532 5. If you used crowdsourcing or conducted research with human subjects...
- 533 (a) Did you include the full text of instructions given to participants and screenshots, if  
534 applicable? [N/A] This work is irrelevant with human subjects.
- 535 (b) Did you describe any potential participant risks, with links to Institutional Review  
536 Board (IRB) approvals, if applicable? [N/A]
- 537 (c) Did you include the estimated hourly wage paid to participants and the total amount  
538 spent on participant compensation? [N/A]

## 539 A Technical Lemmas

540 **Lemma 6** ([Azuma, 1967]). *Let  $(M_n)_{n \geq 0}$  be a martingale such that  $M_0 = 0$  and  $|M_n - M_{n-1}| \leq b$*   
 541 *almost surely for every  $n \geq 1$ . Then we have*

$$\Pr\left[|M_n| \geq b\sqrt{2n \log(2/\delta)}\right] \leq \delta.$$

542 **Lemma 7** ([Zhang et al., 2020c], Lemma 9). *Let  $\{\mathcal{F}_i\}_{i \geq 0}$  be a filtration. Let  $\{X_i\}_{i \geq 1}$  be a real-*  
 543 *valued stochastic process adapted to  $\{\mathcal{F}_i\}_{i \geq 0}$  such that  $0 \leq X_i \leq 1$  almost surely and that  $X_i$  is*  
 544  *$\mathcal{F}_i$ -measurable. For every  $\delta \in (0, 1)$ ,  $c \geq 1$ , we have*

$$\Pr\left[\exists n \geq 1 : \sum_{i=1}^n \mathbb{E}[X_i | \mathcal{F}_{i-1}] \geq 4c \ln \frac{4}{\delta}, \sum_{i=1}^n X_i \leq c \ln \frac{4}{\delta}\right] \leq \delta.$$

545 **Lemma 8.** *Let  $\{\mathcal{F}_i\}_{i \geq 0}$  be a filtration. Let  $\{X_i\}_{i \geq 1}$  be a real-valued stochastic process adapted to*  
 546  *$\{\mathcal{F}_i\}_{i \geq 0}$  such that  $0 \leq X_i \leq 1$  almost surely and that  $X_i$  is  $\mathcal{F}_i$ -measurable. For every  $\delta \in (0, 1)$ ,  $c \geq$*   
 547 *1, we have*

$$\Pr\left[\exists n \geq 1 : \sum_{i=1}^n X_i \geq 4c \ln \frac{4}{\delta}, \sum_{i=1}^n \mathbb{E}[X_i | \mathcal{F}_{i-1}] \leq c \ln \frac{4}{\delta}\right] \leq \delta.$$

548 *Proof.* We follow the proof of Lemma 9 in [Zhang et al., 2020c]. Let  $\lambda > 0$  be a parameter,  
 549  $\mu_i = \mathbb{E}[X_i | \mathcal{F}_{i-1}]$ . Define  $Y_n = \exp(\lambda \sum_{i=1}^n X_i - (e^\lambda - 1) \sum_{i=1}^n \mu_i)$  for  $n \geq 0$ . Note that  
 550  $\mathbb{E}[e^{\lambda X}] \leq \mu e^\lambda + (1 - \mu) \leq e^{\mu(e^\lambda - 1)}$ , so  $\mathbb{E}[e^{\lambda X_i - (e^\lambda - 1)\mu_i} | \mathcal{F}_{i-1}] \leq 1$ , thus  $\{Y_n\}_{n \geq 0}$  is a  
 551 super-martingale. Let  $\tau = \min\{n : \sum_{i=1}^n X_i \geq 4c \ln(4/\delta)\}$  be a stopping time, then we have  
 552  $|Y_{\min\{\tau, n\}}| \leq e^{\lambda(4c \ln(4/\delta) + 1)} < +\infty$  almost surely for every  $n \geq 0$ . Therefore, by the optional  
 553 stopping theorem, we have  $\mathbb{E}[Y_\tau] \leq 1$ . Finally, we have

$$\begin{aligned} \Pr\left[\exists n \geq 1 : \sum_{i=1}^n X_i \geq 4c \ln \frac{4}{\delta}, \sum_{i=1}^n \mu_i \leq c \ln \frac{4}{\delta}\right] &\leq \Pr\left[\sum_{i=1}^{\tau} \mu_i \leq c \ln \frac{4}{\delta}\right] \\ &\leq \Pr\left[Y_\tau \geq \exp\left(\lambda \sum_{i=1}^{\tau} X_i - (e^\lambda - 1) c \ln \frac{2}{\delta}\right)\right] \\ &\leq \Pr\left[Y_\tau \geq \exp\left(\lambda(4c \ln \frac{4}{\delta} - 1) - (e^\lambda - 1) c \ln \frac{2}{\delta}\right)\right] \\ &\leq \exp\left(\lambda(1 - 4c \ln \frac{2}{\delta}) + (e^\lambda - 1) c \ln \frac{2}{\delta}\right) \\ &= e^\lambda e^{(e^\lambda - 1 - 4\lambda) c \ln(4/\delta)}. \end{aligned}$$

554 Choosing  $\lambda = 1$ , we have

$$e^\lambda e^{(e^\lambda - 1 - 4\lambda) c \ln(4/\delta)} \leq e \cdot e^{-2c \ln(4/\delta)} = e\left(\frac{\delta}{4}\right)^c \leq \frac{e}{4}\delta \leq \delta,$$

555 which concludes the proof.  $\square$

556 **Lemma 9.** *Let  $\{\mathcal{F}_i\}_{i \geq 0}$  be a filtration. Let  $\{X_i\}_{i=1}^n$  be a sequence of random variables such that*  
 557  *$|X_i| \leq 1$  almost surely, that  $X_i$  is  $\mathcal{F}_i$ -measurable. For every  $\delta \in (0, 1)$ , we have*

$$\Pr\left[\sum_{i=1}^n \mathbb{E}[X_i^2 | \mathcal{F}_{i-1}] \geq \sum_{i=1}^n 8X_i^2 + 4 \ln \frac{4}{\delta}\right] \leq (\lceil \log_2 n \rceil + 1)\delta.$$

558 *Proof.* Let  $Y = \sum_{i=1}^n \mathbb{E}[X_i^2 | \mathcal{F}_{i-1}]$ ,  $Z = \sum_{i=1}^n X_i^2$ . Applying Lemma 7 with the sequence  
 559  $\{X_i^2\}_{i=1}^n$ , we have for every  $c \geq 1$ ,

$$\Pr\left[Y \geq 4c \ln \frac{4}{\delta}, Z \leq c \ln \frac{4}{\delta}\right] \leq \delta.$$

560 Therefore, we have

$$\begin{aligned}
& \Pr\left[Y \geq 8Z + 4 \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 8Z + 4 \ln \frac{4}{\delta}, 2^{j-1} \ln \frac{4}{\delta} \leq Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 8Z + 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 8Z, 2^{j-1} \ln \frac{4}{\delta} \leq Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 8 \cdot 2^{j-1} \ln \frac{4}{\delta}, 2^{j-1} \ln \frac{4}{\delta} \leq Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 4 \cdot 2^j \ln \frac{4}{\delta}, Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq (\lfloor \log_2 n \rfloor + 1)\delta
\end{aligned}$$

561 as desired.  $\square$

562 **Lemma 10.** Let  $\{\mathcal{F}_i\}_{i \geq 0}$  be a filtration. Let  $\{X_i\}_{i=1}^n$  be a sequence of random variables such that  
563  $|X_i| \leq 1$  almost surely, that  $X_i$  is  $\mathcal{F}_i$ -measurable. For every  $\delta \in (0, 1)$ , we have

$$\Pr\left[\sum_{i=1}^n X_i^2 \geq \sum_{i=1}^n 8 \mathbb{E}[X_i^2 \mid \mathcal{F}_{i-1}] + 4 \ln \frac{4}{\delta}\right] \leq (\lfloor \log_2 n \rfloor + 1)\delta.$$

564 *Proof.* Let  $Y = \sum_{i=1}^n X_i^2$ ,  $Z = \sum_{i=1}^n \mathbb{E}[X_i^2 \mid \mathcal{F}_{i-1}]$ . Applying Lemma 8 with the sequence  
565  $\{X_i^2\}_{i=1}^n$ , we have for every  $c \geq 1$ ,

$$\Pr\left[Y \geq 4c \ln \frac{4}{\delta}, Z \leq c \ln \frac{4}{\delta}\right] \leq \delta.$$

566 Therefore, we have

$$\begin{aligned}
& \Pr\left[Y \geq 8Z + 4 \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 8Z + 4 \ln \frac{4}{\delta}, 2^{j-1} \ln \frac{4}{\delta} \leq Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 8Z + 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 8Z, 2^{j-1} \ln \frac{4}{\delta} \leq Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 8 \cdot 2^{j-1} \ln \frac{4}{\delta}, 2^{j-1} \ln \frac{4}{\delta} \leq Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq \sum_{j=1}^{\lfloor \log_2 n \rfloor} \Pr\left[Y \geq 4 \cdot 2^j \ln \frac{4}{\delta}, Z \leq 2^j \ln \frac{4}{\delta}\right] + \Pr\left[Y \geq 4 \ln \frac{4}{\delta}, Z \leq \ln \frac{4}{\delta}\right] \\
& \leq (\lfloor \log_2 n \rfloor + 1)\delta
\end{aligned}$$

567 as desired.  $\square$

568 **Lemma 11** ([Zhang et al., 2020c], Lemma 11). Let  $(M_n)_{n \geq 0}$  be a martingale such that  $M_0 = 0$  and  
569  $|M_n - M_{n-1}| \leq b$  almost surely for every  $n \geq 1$ . For each  $n \geq 0$ , let  $\mathcal{F}_n = \sigma(M_0, \dots, M_n)$  and  
570 let  $\text{Var}_n = \sum_{i=1}^n \mathbb{E}[(M_i - M_{i-1})^2 \mid \mathcal{F}_{i-1}]$ . Then for any  $n \geq 1$  and  $\epsilon, \delta > 0$ , we have

$$\Pr\left[|M_n| \geq 2\sqrt{2\text{Var}_n \ln(1/\delta)} + 2\sqrt{\epsilon \ln(1/\delta)} + 2b \ln(1/\delta)\right] \leq 2(\log_2(b^2 n / \epsilon) + 1)\delta.$$

**Lemma 12.** Let  $\lambda_1, \lambda_2, \lambda_4 > 0$ ,  $\lambda_3 \geq 1$  and  $\kappa = \max\{\log_2(\lambda_1), 1\}$ . Let  $a_1, a_2, \dots, a_\kappa$  be non-negative reals such that  $a_i \leq \lambda_1$  and  $a_i \leq \lambda_2 \sqrt{a_i + a_{i+1} + 2^{i+1} \lambda_3} + \lambda_4$  for any  $1 \leq i \leq \kappa$  (with  $a_{\kappa+1} = \lambda_1$ ). Then we have that

$$a_1 \leq 22\lambda_2^2 + 6\lambda_4 + 4\lambda_2 \sqrt{2\lambda_3}.$$

*Proof.* Note that

$$a_i \leq \lambda_2 \sqrt{a_i} + \lambda_2 \sqrt{a_{i+1} + 2^{i+1} \lambda_3} + \lambda_4,$$

so we have

$$a_i \leq \left( \lambda_2 + \sqrt{\lambda_2 \sqrt{a_{i+1} + 2^{i+1} \lambda_3} + \lambda_4} \right)^2 \leq 2\lambda_2^2 + 2\lambda_2 \sqrt{a_{i+1} + 2^{i+1} \lambda_3} + 2\lambda_4.$$

By Lemma 11 in [Zhang et al., 2020a], we have

$$\begin{aligned} a_1 &\leq \max\left\{ \left( 2\lambda_2 + \sqrt{(2\lambda_2)^2 + (2\lambda_2^2 + 2\lambda_4)} \right)^2, 2\lambda_2 \sqrt{8\lambda_3} + 2\lambda_2^2 + 2\lambda_4 \right\} \\ &\leq \max\{20\lambda_2^2 + 4\lambda_4, 2\lambda_2 \sqrt{8\lambda_3} + 2\lambda_2^2 + 2\lambda_4\} \leq 22\lambda_2^2 + 6\lambda_4 + 4\lambda_2 \sqrt{2\lambda_3}, \end{aligned}$$

which concludes the proof.  $\square$

## B Difficulty with Previous Approaches

In the example in Section 1, if we know  $x_i \leq \sqrt{\frac{1}{K}}$  for  $1 \leq i \leq K$ , the best confidence region for  $\theta^*$  should be  $\Theta_t = \{\theta \mid \|\theta - \hat{\theta}_t\|_{\Lambda_{t-1}} \leq C(\sigma\sqrt{d} + \lambda^{1/2})\}$ , and we can obtain a variance-aware regret bound by letting  $\lambda = \sigma^2$ . However, if we let  $x_{K+1} = 1$  and use the same concentration inequality as before, the confidence region would be  $\Theta_{K+1} = \{\theta \mid \|\theta - \hat{\theta}_t\|_{\Lambda_{t-1}} \leq C(\sigma\sqrt{d} + 1 + \lambda^{1/2})\}$ .

We present the detailed computation as below. Choose  $\theta^* = \Theta(1)$ .  $\theta^* - \hat{\theta}_{K+1} = -\frac{\sum_{i=1}^{K+1} x_i \epsilon_i}{\lambda + \sum_{i=1}^{K+1} x_i^2} + \frac{\lambda \theta^*}{\lambda + \sum_{i=1}^{K+1} x_i^2}$ . When  $\epsilon_i$  is bounded in  $[-1, 1]$  with variance  $\sigma^2$ , following Bernstein inequality, we have that  $\left| \frac{\sum_{i=1}^{K+1} x_i \epsilon_i}{\lambda + \sum_{i=1}^{K+1} x_i^2} \right| \leq \frac{\sqrt{\sigma^2 \sum_{i=1}^{K+1} x_i^2} + \max_i x_i}{\lambda + \sum_{i=1}^{K+1} x_i^2}$ . Therefore, the best confidence interval we have is

$$\|\theta^* - \hat{\theta}_{K+1}\|_{\Lambda_K} \lesssim \sqrt{\frac{\sigma^2 \sum_{i=1}^{K+1} x_i^2}{\lambda + \sum_{i=1}^{K+1} x_i^2}} + \frac{\max_i x_i}{\sqrt{\lambda + \sum_{i=1}^{K+1} x_i^2}} + \frac{\lambda \theta^*}{\sqrt{\lambda + \sum_{i=1}^{K+1} x_i^2}} = \Theta\left(\sqrt{\frac{\sigma^2}{\lambda + 1}} + \frac{1 + \lambda}{\sqrt{1 + \lambda}}\right),$$

i.e.,  $|\theta^* - \hat{\theta}_{K+1}| \lesssim \Theta(\sigma + \lambda^{1/2} + 1)$ . Therefore, to maintain a confidence region for the general case following methods in [Zhou et al., 2020a, Fauray et al., 2020], the term  $1 + \lambda^{1/2}$  is unavoidable.

**Remark 2.** We highlight that the analysis above is for the uniformly bounded noise. For sub-Gaussian noise, we can ensure that  $\left| \frac{\sum_{i=1}^{K+1} x_i \epsilon_i}{\lambda + \sum_{i=1}^{K+1} x_i^2} \right| \leq \frac{\sqrt{\sigma^2 \sum_{i=1}^{K+1} x_i^2}}{\lambda + \sum_{i=1}^{K+1} x_i^2}$ , which help to reduce the width of confidence interval. More precisely, in the way we have that  $|\theta^* - \hat{\theta}_{K+1}| \leq O(\sigma + \lambda^{1/2})$

## C Proof of Lemma 1

In this section, we present the proof of Lemma 1.

**Restatement of Lemma 1** Let  $f(x) \geq 0$  be a convex function over  $\mathbb{R}$  such that  $\frac{f(x)}{x^2} \leq \frac{f(y)}{y^2} \leq 1$  and  $f(x) \geq f(y)$  if  $x^2 \geq y^2 > 0$ . Fix  $\ell \in (0, 1]$ . For any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t \in \mathbb{B}_2^d(1)$  and  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_t \in \mathbb{B}_2^d(1)$ , we have that

$$\sum_{i=1}^t \min \left\{ \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^{i-1} f(\mathbf{x}_j \boldsymbol{\mu}_j) + \ell^2}, 1 \right\} \leq O(d^4 \log(Cdt/\ell)). \quad (5)$$

Let  $f(x)$  and  $\ell$  be fixed. To prove Lemma 1, we have the lemmas below.

597 **Lemma 13.** For any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t \in \mathbb{B}_2^d(1)$  and  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_n \in \mathbb{B}_2^d(1)$ , we have that

$$\sum_{i=1}^t \min \left\{ \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^t f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2}, 1 \right\} \leq O(d \log(Cdt/\ell)). \quad (6)$$

598 **Lemma 14.** Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t \in \mathbb{B}_2^d(1)$  be a sequence of vectors. If there exists a sequence  
599  $0 = \tau_0 < \tau_1 < \tau_2 < \dots < \tau_z = t$  such that for each  $1 \leq \zeta \leq z$ , there exists  $\boldsymbol{\mu}_\zeta \in \mathbb{B}_2^d(1)$  such that

$$\sum_{i=1}^{\tau_\zeta} f(\mathbf{x}_i \boldsymbol{\mu}_\zeta) + \ell^2 > 4(d+2)^2 \times \left( \sum_{i=1}^{\tau_{\zeta-1}} f(\mathbf{x}_i \boldsymbol{\mu}_\zeta) + \ell^2 \right), \quad (7)$$

600 then  $z \leq O(d \log^2(dt/\ell))$ .

601 We present the proofs of Lemma 13 and 14 respectively in Section C.1 and C.2. Given these two  
602 lemmas, we continue analysis as below.

603 Let  $\tau_0 = 0$  and for  $i \geq 1$ , we let

$$\tau_i = \min\{t+1\} \cup \left\{ \tau \mid \exists \tau_{i-1} \leq \tau' < \tau, \sum_{j=1}^{\tau} f(\mathbf{x}_j \boldsymbol{\mu}_{\tau'}) + \ell^2 > 4(d+2)^2 \left( \sum_{j=1}^{\tau'} f(\mathbf{x}_j \boldsymbol{\mu}_{\tau'}) + \ell^2 \right) \right\}.$$

604 Let  $k = \min\{i \mid \tau_i = t+1\}$ . Then  $k$  is well-defined and  $k \leq O(d \log^2(dt))$  by Lemma 14.

605 Furthermore, for any  $\kappa < k$  and any  $\tau_\kappa \leq i_1 < i_2 < \tau_{\kappa+1}$ , we have

$$\sum_{j=1}^{i_2} f(\mathbf{x}_j \boldsymbol{\mu}_{i_1}) + \ell^2 \leq 4(d+2)^2 \left( \sum_{j=1}^{i_1} f(\mathbf{x}_j \boldsymbol{\mu}_{i_1}) + \ell^2 \right). \quad (8)$$

606 Now we are ready to prove Lemma 1. We have

$$\begin{aligned} \sum_{i=1}^t \min \left\{ \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^{i-1} f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2}, 1 \right\} &\leq 2 \sum_{i=1}^t \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^i f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2} \\ &\leq 8(d+2)^2 \sum_{\kappa=1}^k \left( \sum_{i=\tau_{\kappa-1}}^{\tau_\kappa-1} \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^{\tau_{\kappa-1}-1} f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2} \right) \\ &\leq 8(d+2)^2 \sum_{\kappa=1}^k \left( \sum_{i=\tau_{\kappa-1}}^{\tau_\kappa-1} \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=\tau_{\kappa-1}}^{\tau_{\kappa-1}-1} f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2} \right), \\ &\leq k \times O(d^2) \times O(d \log(t/\ell)) \leq O(d^4 \log^3(dt)), \end{aligned} \quad (9) \quad (10)$$

607 where (9) uses (8) and (10) uses Lemma 13.

### 608 C.1 Proof of Lemma 13

609 **Restatement of Lemma 13** For any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t \in \mathbb{B}_2^d(1)$  and  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_n \in \mathbb{B}_2^d(1)$ , we have  
610 that

$$\sum_{i=1}^t \min \left\{ \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^t f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2}, 1 \right\} \leq O(d \log(Cdt/\ell)). \quad (11)$$

611 *Proof.* Let  $S_t$  be the permutation group over  $[t]$ . We claim that if

$$\sum_{i=1}^t \frac{f(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^t f(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2} = \max_{\xi \in S_t} \sum_{i=1}^t \frac{f(\mathbf{x}_{\xi(i)} \boldsymbol{\mu}_i)}{\sum_{j=1}^t f(\mathbf{x}_{\xi(j)} \boldsymbol{\mu}_i) + \ell^2}, \quad (12)$$

612 then there exists some  $i$  such that  $(\mathbf{x}_i \boldsymbol{\mu}_i)^2 \geq (\mathbf{x}_j \boldsymbol{\mu}_i)^2$  for any  $j \in [t]$ . Otherwise, we construct a  
613 directed graph  $G = (V, E)$  where  $V = [t]$  and edge  $(i, j)$  with  $i \neq j$  is in  $E$  if and only if  $(\mathbf{x}_j \boldsymbol{\mu}_i)^2 \geq$

614  $(x_{j'}\mu_i)^2$  for any  $j' \in [t]$ . Let  $d(i)$  be the out degree of  $i$ . By assuming  $\{(x_i\mu_i)^2 \geq (x_j^\top \mu_i)^2, \forall j \in$   
615  $[t]\}$  fails to hold, we learn that  $d(i) \geq 1$  for every  $i$ , so there exists a circle  $(i_1, i_2, \dots, i_k)$  in  $G$ .  
616 Consider the permutation  $\xi$  such that  $\xi(i_j) = i_{j+1}$  for  $j \in [k]$  (with  $i_{k+1} := i_1$ ) and  $\xi(i) = i$  for  
617  $i \notin \{i_1, \dots, i_k\}$ . By definition, we have  $(\mu_{i_j} x_{\xi(i_j)})^2 > (\mu_{i_j} x_{i_j})^2$  for  $j \in [k]$ , which implies that  
618  $f(\mu_{i_j} x_{\xi(i_j)}) > f(\mu_{i_j} x_{i_j})$  for  $j \in [k]$ . Therefore

$$\sum_{i=1}^t \frac{f(x_i \mu_i)}{\sum_{j=1}^t f(x_j \mu_i) + \ell^2} < \sum_{i=1}^t \frac{f(x_{\xi(i)} \mu_i)}{\sum_{j=1}^t f(x_j \mu_i) + \ell^2} = \sum_{i=1}^t \frac{f(x_{\xi(i)} \mu_i)}{\sum_{j=1}^t f(x_{\xi(j)} \mu_i) + \ell^2},$$

619 which leads to contradiction.

620 We assume that (12) holds, otherwise we can bound an upper bound of the original quantity. Therefore,  
621 we can find an index  $i$  such that  $(x_i \mu_i)^2 \geq (x_j^\top \mu_i)^2$  for any  $j \in [t]$ . Without loss of generality, we  
622 assume  $i = 1$ . Because  $\frac{f(x)}{x^2}$  is decreasing in  $x$ , so we have

$$\frac{f(x_1 \mu_1)}{(x_1 \mu_1)^2} \leq \frac{f(x_j \mu_1)}{(x_j \mu_1)^2}$$

623 for any  $j \in [t]$ , which implies

$$\frac{f(x_1 \mu_1)}{\sum_{j=1}^t f(x_j \mu_1) + \ell^2} = \frac{(x_1 \mu_1)^2}{\left(\sum_{j=1}^t f(x_j \mu_1) + \ell^2\right) \cdot \frac{(x_1 \mu_1)^2}{f(x_1 \mu_1)}} \leq \frac{(x_1 \mu_1)^2}{\sum_{j=1}^t (x_j \mu_1)^2 + \ell^2}. \quad (13)$$

624 Therefore, we have

$$\begin{aligned} \sum_{i=1}^t \frac{f(x_i \mu_i)}{\sum_{j=1}^t f(x_i \mu_i) + \ell^2} &\leq \frac{(x_1 \mu_1)^2}{\sum_{j=1}^t (x_j \mu_1)^2 + \ell^2} + \sum_{i=2}^t \frac{f(x_i \mu_i)}{\sum_{j=1}^t f(x_i \mu_i) + \ell^2} \\ &\leq \frac{(x_1 \mu_1)^2}{\sum_{j=1}^t (x_j \mu_1)^2 + \ell^2} + \sum_{i=2}^t \frac{f(x_i \mu_i)}{\sum_{j=2}^t f(x_i \mu_i) + \ell^2}. \end{aligned} \quad (14)$$

625 Similarly, we can show that there exists a permutation  $\xi^* \in S_t$  such that

$$\sum_{i=1}^t \frac{f(x_i \mu_i)}{\sum_{j=1}^t f(x_j \mu_i) + \ell^2} \leq \sum_{i=1}^t \frac{(x_{\xi^*(i)}^\top \mu_i)^2}{\sum_{j=i}^t (x_{\xi^*(j)}^\top \mu_i)^2 + \ell^2}. \quad (15)$$

626 Finally, by Lemma 15, we have that

$$\sum_{i=1}^t \frac{(x_{\xi^*(i)}^\top \mu_i)^2}{\sum_{j=i}^t (x_{\xi^*(j)}^\top \mu_i)^2 + \ell^2} = \sum_{i=1}^t \min \left\{ \frac{(x_{\xi^*(i)}^\top \mu_i)^2}{\sum_{j=i}^t (x_{\xi^*(j)}^\top \mu_i)^2 + \ell^2}, 1 \right\} \leq O(d \log(t/\ell)).$$

627 □

## 628 C.2 Proof of Lemma 14

629 **Restatement of Lemma 14** Let  $x_1, x_2, \dots, x_t \in \mathbb{B}_2^d(1)$  be a sequence of vectors. If there exists a  
630 sequence  $0 = \tau_0 < \tau_1 < \tau_2 < \dots < \tau_z = t$  such that for each  $1 \leq \zeta \leq z$ , there exists  $\mu_\zeta \in \mathbb{B}_2^d(1)$   
631 such that

$$\sum_{i=1}^{\tau_\zeta} f(x_i \mu_\zeta) + \ell^2 > 4(d+2)^2 \times \left( \sum_{i=1}^{\tau_{\zeta-1}} f(x_i \mu_\zeta) + \ell^2 \right), \quad (16)$$

632 then  $z \leq O(d \log^2(dt/\ell))$ .

633 *Proof.* If  $f(1) \leq \ell^2/t$ , then the conclusion holds trivially because  $0 \leq f(x) \leq f(1) \leq \ell^2/t$  for  
634 all  $x \in [-1, 1]$ . Suppose  $f(1) > \ell^2/t$ . Since  $\frac{f(x)}{x^2} \leq \frac{f(y)}{y^2} \leq 1$  for all  $x^2 \geq y^2$ , we have that for  
635  $0 < \lambda \leq 1$  and any  $x \in \mathbb{R}$ ,  $f(\lambda x) \geq \lambda^2 f(x)$ .

Let  $e_i = [0, \dots, 1, \dots, 0]$  be the one-hot vector whose only 1 entry is at its  $i$ -th coordinate. Noting that  $f(x) \leq x^2$ ,  $|\mathbf{x}_i \boldsymbol{\mu}_\zeta| \leq \|\boldsymbol{\mu}_\zeta\|_2$  and

$$\sum_{i=1}^{\tau_\zeta} f(\mathbf{x}_i \boldsymbol{\mu}_\zeta) > 4(d+2)^2 \times \left( \sum_{i=1}^{\tau_\zeta-1} f(\mathbf{x}_i \boldsymbol{\mu}_\zeta) + \ell^2 \right) - \ell^2 \geq 4d^2 \ell^2$$

we have that  $|\boldsymbol{\mu}_\zeta|_2 \geq \sqrt{\frac{4d^2 \ell^2}{t}}$ . Define  $E_\tau(\boldsymbol{\mu}) = \sum_{i=1}^t f(\mathbf{x}_i \boldsymbol{\mu}) + \frac{\ell^2}{d} \sum_{i=1}^d f(e_i \boldsymbol{\mu})$ . Then  $E_\tau(\boldsymbol{\mu})$  is convex in  $\boldsymbol{\mu}$  because  $f(x)$  is convex in  $x$ . By definition, we have that

$$E_\tau(\boldsymbol{\mu}) \leq \sum_{i=1}^{\tau} f(\mathbf{x}_i \boldsymbol{\mu}) + \ell^2.$$

By (16), we have that

$$E_{\tau_\zeta}(\boldsymbol{\mu}_\zeta) \geq \sum_{i=1}^{\tau_\zeta} f(\mathbf{x}_i \boldsymbol{\mu}_\zeta) \geq 4d^2 \left( \sum_{i=1}^{\tau_\zeta-1} f(\mathbf{x}_i \boldsymbol{\mu}_\zeta) + \ell^2 \right) \geq 4d^2 E_{\tau_\zeta-1}(\boldsymbol{\mu}_\zeta). \quad (17)$$

Define

$$\Lambda = \{i \in \mathbb{Z} : \lfloor \log_2(d\ell^4/t^2) + 2 \rfloor \leq i \leq 2 \lfloor \log_2 t + 2 \rfloor\}.$$

We consider the convex set  $D_{\tau,i} = \{\boldsymbol{\mu} : E_\tau(\boldsymbol{\mu}) \leq 2^i\}$  for  $i \in \Lambda$ . Let  $\zeta$  be fixed. Because  $\|\boldsymbol{\mu}_\zeta\| \geq \sqrt{\frac{4d^2 \ell^2}{t}}$  and  $\sup_i f(e_i \boldsymbol{\mu}) \geq \frac{4d\ell^2}{t} \cdot f(1) \geq \frac{4d\ell^4}{t^2}$ , we have that  $\frac{4d\ell^4}{t^2} \leq E_\tau(\boldsymbol{\mu}_\zeta) \leq t + \ell^2 \leq t+1$  for any  $1 \leq \tau \leq t$ . Then we can find  $i_\zeta \in \Lambda$  such that  $E_{\tau_\zeta-1}(\boldsymbol{\mu}_\zeta) \in (2^{i_\zeta-1}, 2^{i_\zeta}]$ , which means that  $\boldsymbol{\mu}_\zeta \in D_{\tau_\zeta-1, i_\zeta}$ . Note that for  $0 \leq \lambda \leq 1$ ,  $f(\lambda x) \geq \lambda^2 f(x)$  for any  $x$ , it then follows that  $E_t(\lambda \boldsymbol{\mu}) \geq \lambda^2 E_t(\boldsymbol{\mu})$  for any  $t, \boldsymbol{\mu}$ . Choosing  $\lambda = \frac{1}{d}$ , we have that  $E_{\tau_\zeta}(\frac{\boldsymbol{\mu}_\zeta}{d}) \geq \frac{1}{d^2} E_{\tau_\zeta}(\boldsymbol{\mu}_\zeta) \geq 4E_{\tau_\zeta-1}(\boldsymbol{\mu}_\zeta) \geq 2^{i_\zeta}$ . Therefore,  $\frac{\boldsymbol{\mu}_\zeta}{d} \notin D_{\tau_\zeta, i_\zeta}$ . In words, the intercept of  $D_{\tau_\zeta, i_\zeta}$  in the direction  $\boldsymbol{\mu}_\zeta$  is at most  $1/d$  times of that of  $D_{\tau_\zeta-1, i_\zeta}$ .

Note that  $D_{t,i}$  is decreasing in  $t$  for any  $i$ , so by Lemma 16, we have

$$\text{Volume}(D_{\tau_\zeta, i_\zeta}) \leq \frac{6}{7} \text{Volume}(D_{\tau_\zeta-1, i_\zeta}).$$

Also note that  $\text{Volume}(D_{0,i}) \leq (\frac{2t}{\ell})^d$  and  $\text{Volume}(D_{t,i}) \geq (\frac{1}{dt^3})^d$ , so we conclude that  $z \leq d|\Lambda| \log_{7/6}(2dt^4/\ell) \leq O(d \log^2(td/\ell))$ .

□

### C.3 Other Lemmas and Proofs

**Lemma 15.** Fix  $\ell \in (0, 1]$ . Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t \in \mathbb{B}_2^d(1)$  and  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_t \in \mathbb{B}_2^d(1)$  be two sequences of vectors. Then we have

$$\sum_{i=1}^t \mathbb{I} \left\{ (\mathbf{x}_i \boldsymbol{\mu}_i)^2 > \sum_{j=1}^{i-1} (\mathbf{x}_j \boldsymbol{\mu}_i)^2 + \ell^2 \right\} \leq \sum_{i=1}^t \min \left\{ \frac{(\mathbf{x}_i \boldsymbol{\mu}_i)^2}{\sum_{j=1}^{i-1} (\mathbf{x}_j \boldsymbol{\mu}_i)^2 + \ell^2}, 1 \right\} \leq O(d \log \frac{t}{\ell}). \quad (18)$$

*Proof.* The first inequality in (18) holds clearly. To prove the second inequality, we define  $\mathbf{U}_0 = \ell^2 \mathbf{I}$  and  $\mathbf{U}_i = \ell^2 \mathbf{I} + \sum_{j=1}^i \mathbf{x}_j \mathbf{x}_j^\top$  for  $i \geq 1$ . Note that

$$\frac{(\mathbf{x}_i \boldsymbol{\mu}_i)^2}{\sum_{j=1}^{i-1} (\mathbf{x}_j \boldsymbol{\mu}_i)^2 + \ell^2} \leq \frac{(\mathbf{x}_i \boldsymbol{\mu}_i)^2}{\boldsymbol{\mu}_i^\top \mathbf{U}_{i-1} \boldsymbol{\mu}_i} \leq \mathbf{x}_i^\top \mathbf{U}_{i-1}^{-1} \mathbf{x}_i,$$

where the first inequality is because  $\|\boldsymbol{\mu}_i\|_2 \leq 1$  and the second inequality uses the Cauchy's inequality, so we have

$$\sum_{i=1}^t \min \left\{ \frac{(\mathbf{x}_i \boldsymbol{\mu}_i)^2}{\sum_{j=1}^{i-1} (\mathbf{x}_j \boldsymbol{\mu}_i)^2 + \ell^2}, 1 \right\} \leq \sum_{i=1}^t \min \{ \mathbf{x}_i^\top \mathbf{U}_{i-1}^{-1} \mathbf{x}_i, 1 \} \leq 2d \ln(t/\ell^2) \leq 4d \ln(t/\ell),$$

where the second-to-third inequality uses the elliptical potential lemma. □

**Lemma 16.** Given  $x \in \mathbb{R}^d$ , we use  $(u(x), l(x))$  to denote the polar coordinate of  $x$  where  $\|\mu(x)\|_2 = \frac{x}{\|x\|_2}$  is the direction and  $l(x) = \|x\|_2$ . We also use  $(u, \ell)$  to denote the unique element  $x$  in  $\mathbb{R}^d$  such that  $(u(x), l(x)) = (u, \ell)$ . Let  $D$  be a bounded symmetric convex subset of  $\mathbb{R}^d$  with  $d \geq 2$ . Given any direction  $\mu \in \partial\mathbb{B}_d$ , there exists a unique  $l(u) \in \mathbb{R}$  such that  $(u, l(u)), (-u, l(u)) \in \partial D$  are on its boundary. Let  $D'$  be a bounded symmetric convex subset of  $\mathbb{R}^d$  containing  $D \subseteq D'$  such that  $(u, d \cdot l(u)) \in D'$  for some direction  $u \in \partial\mathbb{B}_d$ . Then we have that

$$\text{Volume}(D') \geq \frac{7}{6} \text{Volume}(D).$$

*Proof.* Let  $A = (u, l(u))$  and  $B = (u, d \cdot l(u))$ . Since  $A$  is on the boundary of  $D$ , we can find a hyperplane  $h_1$  such that  $A \in h_1$  and  $h_1$  is tangent to  $D$ . Let  $h_2$  be the parallel hyperplane of  $h_1$  containing the origin  $O \in h_2$ . Define

$$H = \left\{ x \in \mathbb{R}^d \mid d(x, h_1) + d(x, h_2) = d(h_1, h_2), \exists y \in D, \lambda \in \mathbb{R}, (B - y) = \lambda(B - x) \right\}$$

It is obvious that  $\text{Volume}(H) \geq \frac{1}{2} \text{Volume}(D)$  since for each  $x \in D$  lying between  $h_1$  and  $h_2$ ,  $x \in H$ . Define

$$U = \left\{ x \in \mathbb{R}^d \mid d(x, h_2) = d(x, h_1) + d(h_1, h_2), \exists y \in H, \lambda \in [0, 1], x = \lambda y + (1 - \lambda)B \right\}.$$

We claim that

$$\text{Volume}(U) = \left(1 - \frac{1}{d}\right)^d \text{Volume}(U \cup H) = \left(1 - \frac{1}{d}\right)^d (\text{Volume}(U) + \text{Volume}(H)). \quad (19)$$

To see the first equality, we note that  $U$  and  $U \cup H$  are both  $d$ -dimensional pyramids. It then follows from the volume formula and the relation  $d(B, O) = d \times d(A, O)$ . The second equality is because by their definitions,  $U, H$  are separated by the hyperplane  $h_1$ , and thus they are disjoint. Finally, by (19), we have

$$\begin{aligned} \text{Volume}(D') &\geq \text{Volume}(U) + \text{Volume}(H) = \left(1 + \frac{1}{1 - (1 - 1/d)^d}\right) \text{Volume}(H) \\ &\geq \frac{1}{2} \left(1 + \frac{1}{(1 - (1 - 1/d)^d)}\right) \text{Volume}(D) \geq \frac{7}{6} \text{Volume}(D). \end{aligned}$$

□

## D Missing Proofs in Section 4

### D.1 Application of the General Potential Lemma

As an application of Lemma 1 on linear bandit and linear RL, we have the lemma as below

**Lemma 17.** Fix  $\ell \in (0, 1]$ . Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t \in \mathbb{B}_2^d(1)$  be a sequence of vectors, and  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_t \in \mathbb{B}_2^d(1)$  be another sequence of vectors. Then we have

$$\sum_{i=1}^t \frac{\text{clip}^2(\mathbf{x}_i \boldsymbol{\mu}_i, \ell)}{\sum_{j=1}^{i-1} \text{clip}(\mathbf{x}_j \boldsymbol{\mu}_i, \ell) \mathbf{x}_j^\top \boldsymbol{\mu}_i + \ell^2} \leq O(d^4 \log^3(dt)). \quad (20)$$

*Proof.* Let

$$f_\ell(x) = \begin{cases} x^2, & |x| \leq \ell, \\ 2\ell x - \ell^2, & x > \ell, \\ -2\ell x - \ell^2, & x < -\ell \end{cases}$$

be a convex relaxation of the function  $x \mapsto \text{clip}(x, \ell)x$ . It is easy to see that  $f_\ell(x)$  is convex in  $x$  and for any  $x \in \mathbb{R}, \ell > 0$ ,

$$\text{clip}(x, \ell)x \leq f_\ell(x) \leq 2\text{clip}(x, \ell)x \leq 2x^2. \quad (21)$$

686 Let  $h(x) = \frac{f_\ell(x)}{2}$ . It is easy to see that if  $x^2 \geq y^2$ ,  $\frac{h(x)}{x^2} = \frac{\text{clip}(x, \ell)}{2x} \leq \frac{\text{clip}(y, \ell)}{2y} = \frac{h(y)}{y^2} \leq 1$ . By  
 687 Lemma 1 with  $f(x) = h(x)$ , we have that

$$\sum_{i=1}^t \frac{h(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^{i-1} h(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2} \leq O(d^4 \log^3(dt)).$$

688 By (21), we obtain that

$$\begin{aligned} \sum_{i=1}^t \frac{\text{clip}^2(\mathbf{x}_i \boldsymbol{\mu}_i, \ell)}{\sum_{j=1}^{i-1} \text{clip}(\mathbf{x}_j \boldsymbol{\mu}_i, \ell) \mathbf{x}_j^\top \boldsymbol{\mu}_i + \ell^2} &\leq \sum_{i=1}^t \frac{\text{clip}(\mathbf{x}_i \boldsymbol{\mu}_i, \ell) \mathbf{x}_i \boldsymbol{\mu}_I}{\sum_{j=1}^{i-1} \text{clip}(\mathbf{x}_j \boldsymbol{\mu}_i, \ell) \mathbf{x}_j^\top \boldsymbol{\mu}_i + \ell^2} \\ &\leq 4 \sum_{i=1}^t \frac{h(\mathbf{x}_i \boldsymbol{\mu}_i)}{\sum_{j=1}^{i-1} h(\mathbf{x}_j \boldsymbol{\mu}_i) + \ell^2} \\ &\leq O(d^4 \log^3(dt)). \end{aligned}$$

689 The proof is completed. □

## 690 D.2 Proof of Theorem 4

### 691 D.2.1 Optimism

692 The equation (2) accounts for the main novelty of our algorithm. We note that our confidence set is  
 693 different from all previous ones [Dani et al., 2008, Abbasi-Yadkori et al., 2011]. Our confidence set  
 694 is built based on the following new inequality, which may be of independent interest.

695 With Lemma 4 in hand, we can easily prove that the optimal  $\boldsymbol{\theta}^*$  is always in our confidence set with  
 696 high probability. The proof details can be found in Appendix D.3.

697 **Lemma 18.** *With probability at least  $1 - O(\delta \log K)$ , we have  $\boldsymbol{\theta}^* \in \Theta_k$  for all  $k \in [K]$ .*

### 698 D.2.2 Bounding the Regret

699 We bound the regret under the event specified in Lemma 18. We have

$$\begin{aligned} \mathfrak{R}^K &= \sum_{k=1}^K (\max_{\mathbf{x} \in \mathcal{A}_k} \mathbf{x} \boldsymbol{\theta}^* - \mathbf{x}_k \boldsymbol{\theta}^*) \\ &\leq \sum_{k=1}^K \left( \max_{\mathbf{x} \in \mathcal{A}_k, \boldsymbol{\theta} \in \Theta_k} \mathbf{x} \boldsymbol{\theta} - \mathbf{x}_k \boldsymbol{\theta}^* \right) \leq \sum_{k=1}^K \mathbf{x}_k (\boldsymbol{\theta}_k - \boldsymbol{\theta}^*) = \sum_k \mathbf{x}_k \boldsymbol{\mu}_k, \end{aligned}$$

700 where second inequality follows from Lemma 18. Therefore, it suffices to bound  $\sum_k \mathbf{x}_k \boldsymbol{\mu}_k$ , for  
 701 which we have the following lemma.

702 **Lemma 19.** *With probability  $1 - O(\delta \log K)$ , we have*

$$\sum_k \mathbf{x}_k \boldsymbol{\mu}_k \leq O \left( d^{4.5} (\log^4 dK) \left( \log \frac{dK}{\delta} \right) \left( \sqrt{d} + \sqrt{\sum_{k=1}^K \sigma_k^2} \right) \right).$$

703 Since this lemma is one of our main technical contribution, we provide more proof details.

704 *Proof.* First, we define the desired event  $\mathcal{E} = \mathcal{E}_1 \cap \mathcal{E}_2$ , where

$$\mathcal{E}_1 = \{\forall k \in [K] : \boldsymbol{\theta}^* \in \Theta_k\}, \quad \mathcal{E}_2 = \left\{ \sum_{k=1}^K \eta_k(\boldsymbol{\theta}^*) \leq \sum_{k=1}^K 8\sigma_k^2 + 4 \ln \frac{4}{\delta} \right\}.$$

705 By Lemma 18, we have  $\Pr[\mathcal{E}_1] \geq 1 - O(\delta)$ . By Lemma 10, we have  $\Pr[\mathcal{E}_2] \geq 1 - O(\delta \log K)$ .  
 706 Therefore, by union bound, we have  $\Pr[\mathcal{E}] \geq 1 - O(\delta \log K)$ .

Now we bound  $\sum_k \mathbf{x}_k \boldsymbol{\mu}_k$  under the event  $\mathcal{E}$  to prove the lemma. Before presenting the proof, we define

$$\Phi_k^j(\boldsymbol{\mu}) = \sum_{v=1}^{k-1} \text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}) \mathbf{x}_v \boldsymbol{\mu} + \ell_j^2, \quad \Psi_k^j(\boldsymbol{\mu}) = \sum_{v=1}^{k-1} \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}) \eta_v(\boldsymbol{\theta}^*). \quad (22)$$

Also for each  $k \in [K]$ , we define  $j_k \in \Lambda_2$  which satisfies  $\mathbf{x}_k \boldsymbol{\mu}_k \in (\ell_{j_k}/2, \ell_{j_k}]$ . If such  $j_k$  does not exist (because  $\mathbf{x}_k \boldsymbol{\mu}_k \leq \ell_{L_2+1}/2$ ), we assign  $j_k = L_2 + 1$ .

To proceed, we need the following claim.

**Claim 20.** *We have*

$$\begin{aligned} \sum_k \mathbf{x}_k \boldsymbol{\mu}_k &= \sum_{k: j_k = L_2+1} \mathbf{x}_k \boldsymbol{\mu}_k + \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k \\ &\leq 1 + \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k \times \frac{3\sqrt{\Psi_k^{j_k}(\boldsymbol{\mu}_k)\ell} + \sqrt{\sum_{v=1}^{k-1} 2\text{clip}_{j_k}^2(\mathbf{x}_v \boldsymbol{\mu}_k)(\mathbf{x}_v \boldsymbol{\mu}_k)^2\ell} + 3\ell_{j_k}\ell}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)}. \end{aligned} \quad (23)$$

We defer the proof of the claim to Appendix D.4 and continue to bound the three terms in (23). For the second term, we have

$$\begin{aligned} &\sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k \frac{\sqrt{\sum_{v=1}^{k-1} 2\text{clip}_{j_k}^2(\mathbf{x}_v \boldsymbol{\mu}_k)(\mathbf{x}_v \boldsymbol{\mu}_k)^2\ell}}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} \\ &\leq \frac{1}{2} \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k + \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k \mathbb{I} \left\{ \frac{\sqrt{\sum_{v=1}^{k-1} 2\text{clip}_{j_k}^2(\mathbf{x}_v \boldsymbol{\mu}_k)(\mathbf{x}_v \boldsymbol{\mu}_k)^2\ell}}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} > \frac{1}{2} \right\}. \end{aligned} \quad (24)$$

We note that

$$\begin{aligned} \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k \mathbb{I} \left\{ \frac{\sqrt{\sum_{v=1}^{k-1} 2\text{clip}_{j_k}^2(\mathbf{x}_v \boldsymbol{\mu}_k)(\mathbf{x}_v \boldsymbol{\mu}_k)^2\ell}}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} > \frac{1}{2} \right\} &\leq \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k \mathbb{I} \left\{ \Phi_k^{j_k}(\boldsymbol{\mu}_k) \leq 4\ell_{j_k}\ell \right\} \\ &\leq \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k \frac{4\ell_{j_k}\ell}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} \\ &\leq \sum_{k: j_k \leq L_2} \frac{4\text{clip}_{j_k}^2(\mathbf{x}_k \boldsymbol{\mu}_k)\ell}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} \\ &\leq O(d^4 |\Lambda_2| \ell \log^3(dK)), \end{aligned} \quad (25)$$

where the last inequality uses Lemma 17. Collecting (23), (24) and (25), we have

$$\sum_k \mathbf{x}_k \boldsymbol{\mu}_k \leq 1 + \sum_{k: j_k \leq L_2} 3\mathbf{x}_k \boldsymbol{\mu}_k \times \frac{\sqrt{\Psi_k^{j_k}(\boldsymbol{\mu}_k)\ell} + \ell_{j_k}\ell}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} + \frac{1}{2} \sum_{k: j_k \leq L_2} \mathbf{x}_k \boldsymbol{\mu}_k + O(d^4 |\Lambda_2| \ell \log^3(dK)).$$

Solving  $\sum_k \mathbf{x}_k \boldsymbol{\mu}_k$ , we obtain

$$\begin{aligned} \sum_k \mathbf{x}_k \boldsymbol{\mu}_k &\leq O(d^4 |\Lambda_2| \ell \log^3(dK)) + \sum_{k: j_k \leq L_2} 6\mathbf{x}_k \boldsymbol{\mu}_k \times \frac{\sqrt{\Psi_k^{j_k}(\boldsymbol{\mu}_k)\ell} + \ell_{j_k}\ell}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} \\ &\leq O(d^4 |\Lambda_2| \ell \log^3(dK)) + \sum_{k: j_k \leq L_2} 6\mathbf{x}_k \boldsymbol{\mu}_k \times \frac{\sqrt{\Psi_k^{j_k}(\boldsymbol{\mu}_k)\ell}}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)}, \end{aligned} \quad (26)$$

718 where (26) uses the last two steps in (25). The remaining term in (26) can be bounded as

$$\sum_{k:j_k \leq L_2} 6\mathbf{x}_k \boldsymbol{\mu}_k \times \frac{\sqrt{\Psi_k^{j_k}(\boldsymbol{\mu}_k)^\iota}}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} \leq \sum_{k:j_k \leq L_2} 6\mathbf{x}_k \boldsymbol{\mu}_k \ell_{j_k} \frac{\sqrt{\sum_{v=1}^{k-1} \eta_v(\boldsymbol{\theta}^*)^\iota}}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} \quad (27)$$

$$\leq \sum_{k:j_k \leq L_2} \frac{6\mathbf{x}_k \boldsymbol{\mu}_k \ell_{j_k}}{\Phi_k^{j_k}(\boldsymbol{\mu}_k)} \sqrt{\sum_{k=1}^K \eta_k(\boldsymbol{\theta}^*)^\iota} \leq O(d^4 |\Lambda_2| \log^3(dK)) \times \sqrt{\sum_{k=1}^K \eta_k(\boldsymbol{\theta}^*)^\iota} \quad (28)$$

$$\leq O(d^4 |\Lambda_2| \log^3(dK)) \times \sqrt{\left( \ln \frac{1}{\delta} + \sum_{k=1}^K \sigma_k^2 \right)^\iota}, \quad (29)$$

719 where (27) uses the definition of  $\Psi_k^j(\cdot)$ , (28) again uses the last two steps in (25), and (29) uses the  
720 event  $\mathcal{E}_2$ .  $\square$

721 Now we can finish the proof of Theorem 3. We choose  $\delta = O((K \log K)^{-1})$ . Since on the event  $\mathcal{E}^C$ ,  
722 we have  $\mathfrak{R}^K \leq K$ . Therefore, together with the bound on  $\mathcal{E}$  from Lemma 19, we conclude that the  
723 expected regret is bounded by  $\mathbb{E}[\mathfrak{R}^K] \leq \tilde{O}(d^{4.5} \sqrt{\sum_{k=1}^K \sigma_k^2} + d^5)$ .

724 *Proof.* It suffices to prove the theorem for  $b = 1$ , because otherwise we can apply  $\{X_i/b\}_{i=1}^n$  to the  
725  $b = 1$  case. By Lemma 11 with  $\epsilon = 1$  and  $\delta < 1/e$ , we have

$$\Pr \left[ \left| \sum_{i=1}^n X_i \right| \geq 2 \sqrt{\sum_{i=1}^n 2 \mathbb{E}[X_i^2 \mid \mathcal{F}_{i-1}] \ln \frac{1}{\delta}} + 4 \ln \frac{1}{\delta} \right] \leq 4\delta \log_2 n. \quad (30)$$

726 By Lemma 9, we have

$$\Pr \left[ \sum_{i=1}^n \mathbb{E}[X_i^2 \mid \mathcal{F}_{i-1}] \geq \sum_{i=1}^n 8X_i^2 + 4 \ln \frac{4}{\delta} \right] \leq ([\log_2 n] + 1)\delta. \quad (31)$$

727 Therefore, by a union bound over (30) and (31), we have with probability at least  $1 - 6\delta \log_2 n$ ,

$$\begin{aligned} \left| \sum_{i=1}^n X_i \right| &\leq \sqrt{\sum_{i=1}^n 8 \mathbb{E}[X_i^2 \mid \mathcal{F}_{i-1}] \ln \frac{1}{\delta}} + 4 \ln \frac{1}{\delta} \\ &\leq \sqrt{8 \left( \sum_{i=1}^n 8X_i^2 + 4 \ln \frac{4}{\delta} \right) \ln \frac{1}{\delta}} + 4 \ln \frac{1}{\delta} \leq 8 \sqrt{\sum_{i=1}^n X_i^2 \ln \frac{1}{\delta}} + 16 \ln \frac{1}{\delta}, \end{aligned}$$

728 which concludes the proof.  $\square$

### 729 D.3 Proof of Lemma 18

730 *Proof.* Let  $\delta' = e^{-\iota}$ . We define the desired event  $\mathcal{E} = \bigcap_{k \in [K], j \in \Lambda_2} \mathcal{E}_k^j$ , where

$$\mathcal{E}_k^j = \left\{ \left| \sum_{v=1}^k \text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}) \epsilon_v(\boldsymbol{\theta}^*) \right| \leq \sqrt{\sum_{v=1}^k \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}) \eta_v(\boldsymbol{\theta}^*)^\iota} + \ell_j \iota, \forall \boldsymbol{\mu} \in \mathcal{B} \right\}.$$

731 Note that for each  $v$ , we have that  $|\text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}) \epsilon_v(\boldsymbol{\theta}^*)| \leq \ell_j$  and that  $(\text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}) \epsilon_v(\boldsymbol{\theta}^*))^2 =$   
 732  $\text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}) \eta_v(\boldsymbol{\theta}^*)$ , so by Theorem 4, we have

$$\begin{aligned} \Pr \left[ \left| \sum_{v=1}^k \text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}) \epsilon_v \right| \leq \sqrt{\sum_{v=1}^k \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}) \text{Var}(\epsilon_v \mid \mathcal{F}_v) \ell} + \ell_j \ell \right] \\ \geq 1 - O \left( e^{-\frac{\ell}{\log_2 \log_2 K}} \right) \\ \geq 1 - O \left( \frac{\delta}{K |\mathcal{B}| |\Lambda_2|} \log K \right), \end{aligned}$$

733 where  $\mathcal{F}_v$  is as defined in Section 3. Finally, using a union bound over  $(\boldsymbol{\mu}, j, k) \in \mathcal{B} \times \Lambda_2 \times [K]$ , we  
 734 have  $\Pr[\mathcal{E}] \geq 1 - O(\delta \log K)$ .  $\square$

#### 735 D.4 Proof of Claim 20

736 *Proof.* We elaborate on (23). We will prove it by showing that the numerator is always greater than  
 737 the denominator in the fraction in (23), so each term  $\mathbf{x}_k \boldsymbol{\mu}_k$  is multiplied by a number greater than 1.  
 738 We have for every  $j \in \Lambda_2$ ,

$$\begin{aligned} \Phi_k^j(\boldsymbol{\mu}_k) &= \sum_{v=1}^{k-1} \text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}_k) \mathbf{x}_v \boldsymbol{\mu}_k + \ell_j^2 \\ &\leq \left| \sum_{v=1}^{k-1} \text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}_k) \epsilon_v(\boldsymbol{\theta}^*) \right| + \left| \sum_{v=1}^{k-1} \text{clip}_j(\mathbf{x}_v \boldsymbol{\mu}_k) \epsilon_v(\boldsymbol{\theta}_k) \right| + \ell_j^2 \end{aligned} \quad (32)$$

$$\leq \sqrt{\Psi_k^j(\boldsymbol{\mu}_k) \ell} + \sqrt{\sum_{v=1}^{k-1} \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) \eta_v(\boldsymbol{\theta}_k) \ell} + 3\ell_j \ell \quad (33)$$

$$\leq \sqrt{\Psi_k^j(\boldsymbol{\mu}_k) \ell} + \sqrt{\sum_{v=1}^{k-1} \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) \eta_v(\boldsymbol{\theta}^*) \ell} + \sqrt{\sum_{v=1}^{k-1} \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) |\eta_v(\boldsymbol{\theta}_k) - \eta_v(\boldsymbol{\theta}^*)| \ell} + 3\ell_j \ell$$

$$\begin{aligned} &= 2\sqrt{\Psi_k^j(\boldsymbol{\mu}_k) \ell} + \sqrt{\sum_{v=1}^{k-1} \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) |\eta_v(\boldsymbol{\theta}_k) - \eta_v(\boldsymbol{\theta}^*)| \ell} + 3\ell_j \ell \\ &\leq 2\sqrt{\Psi_k^j(\boldsymbol{\mu}_k) \ell} + \sqrt{\sum_{v=1}^{k-1} \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) (\eta_v(\boldsymbol{\theta}^*) + 2(\mathbf{x}_v \boldsymbol{\mu}_k)^2) \ell} + 3\ell_j \ell \end{aligned} \quad (34)$$

$$\begin{aligned} &\leq 2\sqrt{\Psi_k^j(\boldsymbol{\mu}_k) \ell} + \sqrt{\sum_{v=1}^{k-1} \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) \eta_v(\boldsymbol{\theta}^*) \ell} + \sqrt{\sum_{v=1}^{k-1} 2 \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) (\mathbf{x}_v \boldsymbol{\mu}_k)^2 \ell} + 3\ell_j \ell, \\ &= 3\sqrt{\Psi_k^j(\boldsymbol{\mu}_k) \ell} + \sqrt{\sum_{v=1}^{k-1} 2 \text{clip}_j^2(\mathbf{x}_v \boldsymbol{\mu}_k) (\mathbf{x}_v \boldsymbol{\mu}_k)^2 \ell} + 3\ell_j \ell, \end{aligned} \quad (35)$$

739 where (32) uses  $\epsilon_v(\boldsymbol{\theta}_k) - \epsilon_v(\boldsymbol{\theta}^*) = \mathbf{x}_v(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*) = \mathbf{x}_v \boldsymbol{\mu}_k$ , (33) uses that  $\boldsymbol{\theta}^*, \boldsymbol{\theta}_k \in \Theta_k$  and the  
 740 definition of  $\Theta_k$  in (2), and (34) uses

$$\begin{aligned} |\eta_v(\boldsymbol{\theta}_k) - \eta_v(\boldsymbol{\theta}^*)| &= |(\epsilon_v(\boldsymbol{\theta}^*) - \mathbf{x}_v \boldsymbol{\mu}_k)^2 - (\epsilon_v(\boldsymbol{\theta}^*))^2| \\ &\leq 2|\epsilon_v(\boldsymbol{\theta}^*)| \mathbf{x}_v \boldsymbol{\mu}_k + (\epsilon_v(\boldsymbol{\theta}^*))^2 \leq (\mathbf{x}_v \boldsymbol{\mu}_k)^2 + 2(\epsilon_v(\boldsymbol{\theta}^*))^2. \end{aligned}$$

741 Since (35) holds for every  $j \in \Lambda_2$ , it holds for  $j = j_k$ , and thus (23) follows.  $\square$

## E Missing Proofs in Section 5

### E.1 Proof of Theorem 5

Before introducing our proof, we make some definitions. We let  $\theta_{k,h}^m = \arg \max_{\theta \in \Theta_k} x_{k,h}^m(\theta - \theta^*)$  and  $\mu_{k,h}^m = \theta_{k,h}^m - \theta^*$ . Recall that  $\mathcal{T}_k^{m,i}$  is defined in Algorithm 2. We define

$$\Phi_k^{m,i,j}(\mu) = \sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j(x_{v,u}^m \mu) x_{v,u}^m \mu + \ell_j^2, \quad \Psi_k^{m,i,j}(\mu) = \sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j^2(x_{v,u}^m \mu) \eta_{v,u}^m. \quad (36)$$

Note that our definitions in (36) are similar to those for linear bandits in (22). The main differences are: 1) we define  $\Phi(\cdot), \Psi(\cdot)$  also for higher moments, as indicated by the index  $m$  in their superscripts; 2) we add the variance layer, so that we only use samples from  $\mathcal{T}^{m,i}$ ; 3) since we can now estimate variance, we use the upper bound of estimated variance in lieu of the empirical variance. For  $h \in [H+1]$ , we further define

$$I_h^k = \mathbb{I}\{\forall u \leq h, m, i, j : \Phi_{k,u}^{m,i,j}(\mu_{k,u}^m) \leq 4(d+2)^2 \Phi_k^{m,i,j}(\mu_{k,u}^m)\}, \quad (37)$$

where  $I_h^k = 1$  indicates that for every  $u \leq h$ , the confidence set using data prior to the time step  $(k, u)$  can be properly approximated by the confidence set with data prior to the episode  $k$ . We define  $I_h^k$  in this way to ensure that it is  $\mathcal{F}_h^k$ -measurable. The following lemma ensures that  $Q_h^k$  is optimistic with high probability. Its proof is deferred to Appendix E.3.

**Lemma 21.**  $\Pr[\forall k, h, s, a : Q_h^k(s, a) \geq Q_h^*(s, a)] \geq \Pr[\forall k \in [K] : \theta^* \in \Theta_k] \geq 1 - O(\delta)$ .

When the event specified in Lemma 21 holds, the regret can be decomposed as

$$\mathfrak{R}^K = \sum_{k=1}^K (V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k)) \leq \sum_{k=1}^K (V_1^k(s_1^k) - V_1^{\pi_k}(s_1^k)) \leq \check{\mathfrak{R}}_1 + \check{\mathfrak{R}}_2 + \mathfrak{R}_3 + \sum_{k,h} (I_h^k - I_{h+1}^k),$$

where

$$\check{\mathfrak{R}}_1 = \sum_{k,h} (P_{s_h^k, a_h^k} V_{h+1}^k - V_{h+1}^k(s_{h+1}^k)) I_h^k, \quad \check{\mathfrak{R}}_2 = \sum_{k,h} (V_h^k(s_h^k) - r_h^k - P_{s_h^k, a_h^k} V_{h+1}^k) I_h^k,$$

$$\mathfrak{R}_3 = \sum_{k=1}^K \left( \sum_{h=1}^H r_h^k - V_1^{\pi_k}(s_1^k) \right).$$

Next we analyze these terms. First, we observe that  $\mathfrak{R}_3$  is a sum of a martingale difference sequence, so by Lemma 6, we have  $\mathfrak{R}_3 \leq O(\sqrt{K \log(1/\delta)})$  with probability at least  $1 - \delta$ . Next, we use the following lemma to bound  $\sum_{k,h} (I_h^k - I_{h+1}^k)$ . We defer its proof to Appendix E.4.

**Lemma 22.**  $\sum_{k,h} (I_h^k - I_{h+1}^k) \leq O(d \log^5(dHK))$ .

To bound  $\check{\mathfrak{R}}_1$  and  $\check{\mathfrak{R}}_2$ , we need to define the following quantities. First, we denote  $\check{x}_{k,h} = x_{k,h} I_h^k$  and  $\check{\eta}_{k,h}^m = \eta_{k,h}^m I_h^k$ . Next, for  $m \in \Lambda_0$ , we define

$$\check{R}_m = \sum_{k,h} \check{x}_{k,h}^m \mu_{k,h}^m, \quad \check{M}_m = \sum_{k,h} \left( P_{s_h^k, a_h^k} (V_{h+1}^k)^{2^m} - (V_{h+1}^k(s_{h+1}^k))^{2^m} \right) I_h^k.$$

Intuitively,  $\check{R}_m$  represents the “regret” of  $2^m$ -th moment prediction and  $\check{M}_m$  represents the total variance of  $2^m$ -th order value function. We have  $\check{\mathfrak{R}}_1 = \check{M}_0$  by definition and using that

$$Q_h^k(s, a) - r(s, a) - P_{s,a} V_{h+1}^k \leq \max_{\theta \in \Theta_k} x_{k,h}^0(\theta - \theta^*),$$

we have  $\check{\mathfrak{R}}_2 \leq \check{R}_0$ . So it suffices to bound  $\check{R}_0 + \check{M}_0$ , which is done by the following lemma.

**Lemma 23.** *With probability at least  $1 - \delta$ , we have*

$$\check{R}_0 + |\check{M}_0| \leq O\left(d^{4.5} \sqrt{K \log^5(dHK) \log(1/\delta)} + d^9 \log^6(dHK) \log(1/\delta)\right).$$

Lemma 23 is the main technical part of our result in Section 5, so we sketch its proof in the next subsection. With the lemma in hand, we have with probability  $1 - \delta$  that  $\mathfrak{R}^K \leq \tilde{O}(d^{4.5} \sqrt{K} + d^9)$ . Finally, We conclude the proof to Theorem 5 by choosing  $\delta = 1/K$  and noting that  $\mathfrak{R}^K \leq K$ .

## 771 E.2 Bounding $\check{R}$ and $\check{M}$

772 We sketch the proof for Lemma 23. The first step to bound  $\check{R}_m$  is to relate it to the variance  $\check{\eta}^m$ .

773 **Lemma 24.** *With probability at least  $1 - \delta$ , we have  $\check{R}_m \leq O(d^4 \sqrt{\sum_{k,h} \check{\eta}_{k,h}^m} \iota \log^7(dHK)) +$*   
 774  *$d^6 \iota \log^5(dHK)$ .*

775 We defer the proof to Appendix E.5. The proof is spiritually similar to proof of Lemma 19. The main  
 776 difference is that we use the peeling technique to the magnitude of the variance.

777 Based on Lemma 24, we use the following recursion lemma to relate  $\check{R}_m, \check{M}_m$  to  $\check{R}_{m+1}, \check{M}_{m+1}$ . We  
 778 defer the proof to Appendix E.6. It mainly uses similar ideas in Zhang et al. [2020a].

779 **Lemma 25 (Recursions).** *With probability at least  $1 - \delta$ , we have*

$$\begin{aligned} \check{R}_m &\leq O\left(d^4 \sqrt{(\check{M}_{m+1} + 2^{m+1}(K + \check{R}_0) + \check{R}_{m+1} + \check{R}_m) \iota \log^7(dHK)} + d^6 \iota \log^5(dHK)\right), \\ |\check{M}_m| &\leq O\left(\sqrt{(\check{M}_{m+1} + O(d \log^5(dHK)) + 2^{m+1}(K + \check{R}_0)) \log(1/\delta)} + \log(1/\delta)\right). \end{aligned}$$

780 Finally, we can prove Lemma 23 by collecting Lemma 24,25 and using a technical lemma about  
 781 recursion (Lemma 12). The details are in Appendix E.7.

## 782 E.3 Proof of Lemma 21

783 *Proof.* The lemma consists of two inequalities. The first inequality is proved using backward  
 784 induction, where the induction step is given as

$$\begin{aligned} Q_h^k(s, a) &= \min\{1, r(s, a) + \max_{\theta \in \Theta_k} \sum_{i=1}^d \theta_i P_{s,a}^i V_{h+1}^k\} \\ &\geq \min\{1, r(s, a) + \sum_{i=1}^d \theta_i^* P_{s,a}^i V_{h+1}^k\} \geq \min\{1, r(s, a) + \sum_{i=1}^d \theta_i^* P_{s,a}^i V_{h+1}^*\} = Q_h^*(s, a), \\ V_h^k(s) &= \max_a Q_h^k(s, a) \geq \max_a Q_h^*(s, a) = V_h^*(s). \end{aligned}$$

785 We now prove the second inequality. Let  $\delta' = e^{-\iota}$ . We define the desired event  $\mathcal{E} = \bigcap_{k,m,i,j} \mathcal{E}_k^{m,i,j}$ ,  
 786 where

$$\mathcal{E}_k^{m,i,j} = \left\{ \left| \sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}) \varepsilon_{\kappa,h}^m \right| \leq 4 \sqrt{\sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j^2(\mathbf{x}_{v,u}^m \boldsymbol{\mu}) \text{Var}(\varepsilon_{v,u}^m \mid \mathcal{F}_u^v) \ln \frac{1}{\delta'}} + 4\ell_j \ln \frac{1}{\delta'} \right\}, \forall \boldsymbol{\mu} \in \mathcal{B}.$$

787 Note that for a fixed  $k$ , we have that  $|\text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}) \varepsilon_{v,u}^m| \leq \ell_j \leq 1$  and that

$$\text{Var}\left(\text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}) \varepsilon_{v,u}^m \mathbb{I}\{(v, u) \in \mathcal{T}_k^{m,i}\} \mid \mathcal{F}_u^v\right) = \text{clip}_j^2(\mathbf{x}_{k,h}^m \boldsymbol{\mu})^2 \mathbb{I}\{(v, u) \in \mathcal{T}_k^{m,i}\} \text{Var}(\varepsilon_{v,u}^m \mid \mathcal{F}_u^v),$$

788 so by Lemma 11 with  $b = \ell_j, \epsilon = 1$ , we have

$$\begin{aligned} \Pr\left[\left| \sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}) \varepsilon_{v,u}^m \right| \geq 4 \sqrt{\sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j^2(\mathbf{x}_{v,u}^m \boldsymbol{\mu}) \text{Var}(\varepsilon_{v,u}^m \mid \mathcal{F}_u^v) \ln \frac{1}{\delta'}} + 4\ell_j \ln \frac{1}{\delta'}\right] \\ \leq 4\delta' \log_2(HK). \end{aligned}$$

789 Using a union bound over  $(\boldsymbol{\mu}, m, i, j, k) \in \mathcal{B} \times \Lambda_0 \times \Lambda_1 \times \Lambda_2 \times [K]$ , we have  $\Pr[\mathcal{E}] \geq 1 -$   
 790  $O(\delta' K |\mathcal{B}| \log^4(HK)) \geq 1 - O(\delta)$ .

791 Next we show that the event  $\mathcal{E}$  implies that  $\boldsymbol{\theta}^* \in \Theta_k$  for every  $k \in [K]$ . We show by induction  
 792 over  $k$ . For  $k = 1$  it is clear. For  $k \geq 1$ , since  $\boldsymbol{\theta}^* \in \Theta_k$ , for every  $h \in [H]$ , we have  $\eta_{k,h}^m =$   
 793  $\max_{\boldsymbol{\theta} \in \Theta_k} \{\boldsymbol{\theta} \mathbf{x}_{k,h}^{m+1} - (\boldsymbol{\theta} \mathbf{x}_{k,h}^m)^2\} \geq \boldsymbol{\theta}^* \mathbf{x}_{k,h}^{m+1} - (\boldsymbol{\theta}^* \mathbf{x}_{k,h}^m)^2 \geq \text{Var}(\varepsilon_{k,h}^m \mid \mathcal{F}_h^k)$ , which, together with  
 794 the event  $\bigcap_{m,i,j} \mathcal{E}_{k+1}^{m,i,j}$ , implies that  $\boldsymbol{\theta}^* \in \Theta_{k+1}$ .  $\square$

#### 795 E.4 Proof of Lemma 22

796 *Proof.* We define

$$I_{k,h}^{m,i,j} = \mathbb{I}\{\forall u \leq h : \Phi_{k,u}^{m,i,j}(\boldsymbol{\mu}_{k,u}^m) \leq 4(d+2)^2 \Phi_k^{m,i,j}(\boldsymbol{\mu}_{k,u}^m)\}.$$

797 Then we have  $I_h^k = \prod_{m,i,j} I_{k,h}^{m,i,j}$ . Also we have

$$\sum_h (I_h^k - I_{h+1}^k) \leq \sum_{m,i,j} \sum_h (I_{k,h}^{m,i,j} - I_{k,h+1}^{m,i,j}).$$

798 Note that  $I_h^k \geq I_{h+1}^k$  and  $I_{k,h}^{m,i,j} \geq I_{k,h+1}^{m,i,j}$ . For each fixed  $m, i, j$ , if  $\sum_h (I_{k,h}^{m,i,j} - I_{k,h+1}^{m,i,j}) = 1$ , then  
 799 there exists  $h \in [H]$ , such that for the time step  $(k, h)$ , we have  $\Phi_{k,h}^{m,i,j}(\boldsymbol{\mu}) > 4(d+2)^2 \Phi_k^{m,i,j}(\boldsymbol{\mu})$   
 800 for some  $\boldsymbol{\mu}$ . By Lemma 14 with  $f(x) = \text{clip}(x, \ell_j)x$  and  $\ell = \ell_j$ , there are at most  $O(d \log^2(dHK))$   
 801 such time steps. We conclude by noting that we have  $|\Lambda_0 \times \Lambda_1 \times \Lambda_2| \leq O(\log^3(dHK))$  possible  
 802  $m, i, j$  pairs.  $\square$

#### 803 E.5 Proof of Lemma 24

804 To prove this lemma, we define the index sets to help us apply the peeling technique. We denote

$$\begin{aligned} \mathcal{T}_k^{m,i,j} &= \{(v, u) \in \mathcal{T}_k^{m,i} : |\mathbf{x}_{v,u}^m \boldsymbol{\mu}_{v,u}^m| \in (\ell_{j+1}, \ell_j]\}, \\ \mathcal{T}_k^{m,i,L_2+1} &= \{(v, u) \in \mathcal{T}_k^{m,i} : |\mathbf{x}_{v,u}^m \boldsymbol{\mu}_{v,u}^m| \in [0, \ell_{L_2+1}]\}, \end{aligned}$$

805 and  $\tilde{\mathcal{T}}_k^{m,i,j} = \{(v, u) \in \mathcal{T}_k^{m,i,j} : I_u^v = 1\}$ . We also denote  $\mathcal{T}^{m,i,j} = \mathcal{T}_{K+1}^{m,i,j}$ ,  $\tilde{\mathcal{T}}^{m,i,j} = \tilde{\mathcal{T}}_{K+1}^{m,i,j}$ .

806 *Proof.* Since  $\boldsymbol{\theta}_{k,h}^m \in \Theta_k \subseteq \Theta_k^{m,i,j}$ , choosing  $\boldsymbol{\mu} = \boldsymbol{\mu}_{k,h}^m$  in the confidence set definition and using that  
 807  $\mathbf{x}_{v,u}^m \boldsymbol{\mu}_{k,h}^m = \epsilon_{v,u}^m(\boldsymbol{\theta}^*) - \epsilon_{v,u}^m(\boldsymbol{\theta}_{k,h}^m)$ , we have

$$\begin{aligned} \Phi_k^{m,i,j}(\boldsymbol{\mu}_{k,h}^m) &= \sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}_{k,h}^m) \mathbf{x}_{v,u}^m \boldsymbol{\mu}_{k,h}^m + \ell_j^2 \\ &\leq \left| \sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}_{k,h}^m) \epsilon_{v,u}^m(\boldsymbol{\theta}^*) \right| + \left| \sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}_{k,h}^m) \epsilon_{v,u}^m(\boldsymbol{\theta}_{k,h}^m) \right| + \ell_j^2 \\ &\leq 8 \sqrt{\sum_{(v,u) \in \mathcal{T}_k^{m,i}} \text{clip}_j(\mathbf{x}_{v,u}^m \boldsymbol{\mu}_{k,h}^m) \eta_{v,u}^m \iota} + 8 \ell_j \iota + \ell_j^2 \\ &\leq 8 \sqrt{\Psi_k^{m,i,j}(\boldsymbol{\mu}_{k,h}^m) \iota} + 16 \ell_j \iota. \end{aligned} \tag{38}$$

808 Therefore, when  $I_k^h = 0$ , we have

$$\frac{\Phi_{k,h}^{m,i,j}(\boldsymbol{\mu}_{k,h}^m)}{4(d+2)^2} \leq \Phi_k^{m,i,j}(\boldsymbol{\mu}_{k,h}^m) \leq 16(\sqrt{\Psi_k^{m,i,j}(\boldsymbol{\mu}_{k,h}^m) \iota} + \ell_j \iota).$$

809 Next we analyze the sum. Using the fact that

$$\frac{64(d+2)^2 \left( \sqrt{\Psi_k^{m,i,j}(\boldsymbol{\mu}_{k,h}^m) \iota} + \ell_j \iota \right)}{\Phi_{k,h}^{m,i,j}(\boldsymbol{\mu}_{k,h}^m)} \geq 1,$$

810 we obtain

$$\sum_{(k,h) \in \tilde{\mathcal{T}}^{m,i,j}} \mathbf{x}_{k,h}^m \boldsymbol{\mu}_{k,h}^m \leq \sum_{(k,h) \in \tilde{\mathcal{T}}^{m,i,j}} \mathbf{x}_{k,h}^m \boldsymbol{\mu}_{k,h}^m \frac{64(d+2)^2 \left( \sqrt{\Psi_k^{m,i,j}(\boldsymbol{\mu}_{k,h}^m) \iota} + \ell_j \iota \right)}{\Phi_{k,h}^{m,i,j}(\boldsymbol{\mu}_{k,h}^m)} \tag{39}$$

$$\leq 64(d+2)^2 \sum_{(k,h) \in \tilde{\mathcal{T}}^{m,i,j}} \left( \frac{\mathbf{x}_{k,h}^m \boldsymbol{\mu}_{k,h}^m \sqrt{\ell_i \iota}}{\sqrt{\Phi_{k,h}^{m,i,j}(\boldsymbol{\mu}_{k,h}^m)}} + \frac{\mathbf{x}_{k,h}^m \boldsymbol{\mu}_{k,h}^m \ell_j \iota}{\Phi_{k,h}^{m,i,j}(\boldsymbol{\mu}_{k,h}^m)} \right), \tag{40}$$

811 where the last inequality uses that for every  $\mu$ , we have

$$\Psi_{k,h}^{m,i,j}(\mu) = \sum_{(v,u) \in \mathcal{T}_{k,h}^{m,i}} \text{clip}_j^2(\mathbf{x}_{v,u}^m \mu) \eta_{v,u}^m \leq \ell_i \sum_{(v,u) \in \mathcal{T}_{k,h}^{m,i}} \text{clip}_j(\mathbf{x}_{k,h}^m \mu) \mathbf{x}_{k,h}^m \mu \leq \ell_i \Phi_{k,h}^{m,i,j}(\mu). \quad (41)$$

812 In (41), the first inequality uses that  $\eta_{v,u}^m \leq \ell_i$  for  $(v,u) \in \mathcal{T}_{k,h}^{m,i}$  and that  $\text{clip}_j^2(\alpha) \leq \text{clip}_j(\alpha)\alpha$  for  
 813  $\alpha \in \mathbb{R}$ , and the second inequality uses the definition of  $\Phi_{k,h}^{m,i,j}(\mu)$ . Next we bound the two terms in  
 814 (40). To bound the first term, we note that

$$\sum_{(k,h) \in \mathcal{T}^{m,i,j}} \frac{\mathbf{x}_{k,h}^m \mu_{k,h}^m}{\sqrt{\Phi_{k,h}^{m,i,j}(\mu_{k,h}^m)}} \leq \sqrt{|\mathcal{T}^{m,i,j}|} \sqrt{\sum_{(k,h) \in \mathcal{T}^{m,i,j}} \frac{(\mathbf{x}_{k,h}^m \mu_{k,h}^m)^2}{\Phi_{k,h}^{m,i,j}(\mu_{k,h}^m)}} \quad (42)$$

$$\leq \sqrt{|\mathcal{T}^{m,i,j}|} \sqrt{\sum_{(k,h) \in \mathcal{T}^{m,i,j}} \frac{\text{clip}_j^2(\mathbf{x}_{k,h}^m \mu_{k,h}^m)}{\Phi_{k,h}^{m,i,j}(\mu_{k,h}^m)}} \quad (43)$$

$$\leq \sqrt{|\mathcal{T}^{m,i,j}|} \sqrt{\sum_{(k,h) \in \mathcal{T}^{m,i,j}} \frac{\text{clip}_j^2(\mathbf{x}_{k,h}^m \mu_{k,h}^m)}{\sum_{(v,u) \in \mathcal{T}_{k,h}^{m,i,j}} \text{clip}_j(\mathbf{x}_{v,u}^m \mu_{k,h}^m) \mathbf{x}_{v,u}^m \mu_{k,h}^m + \ell_j^2}} \quad (44)$$

$$\leq \sqrt{|\mathcal{T}^{m,i,j}|} \times O(\sqrt{d^4 \log^3(dHK)}), \quad (45)$$

815 where (42) uses Cauchy's inequality, (43) uses that  $\mathbf{x}_{k,h}^m \mu_{k,h}^m \leq \ell_j$  for  $(k,h) \in \mathcal{T}^{m,i,j}$ , (44) uses the  
 816 definition of  $\Phi_{k,h}^{m,i,j}(\mu)$ , and (45) uses Lemma 17. To bound the second term in (40), we have

$$\sum_{(k,h) \in \mathcal{T}^{m,i,j}} \frac{\mathbf{x}_{k,h}^m \mu_{k,h}^m \ell_j}{\Phi_{k,h}^{m,i,j}(\mu_{k,h}^m)} \leq \sum_{(k,h) \in \mathcal{T}^{m,i,j}} \frac{2 \text{clip}_j^2(\mathbf{x}_{k,h}^m \mu_{k,h}^m)}{\Phi_{k,h}^{m,i,j}(\mu_{k,h}^m)} \leq O(d^4 \log^3(dHK)), \quad (46)$$

817 where the first inequality uses that  $\mathbf{x}_{k,h}^m \mu_{k,h}^m \geq \ell_j/2$  for  $(k,h) \in \mathcal{T}^{m,i,j}$  and the second inequality is  
 818 the same as what we have shown from (43) to (45). As a result, combining (40), (45) and (46), we  
 819 have

$$\sum_{(k,h) \in \mathcal{T}^{m,i,j}} \mathbf{x}_{k,h}^m \mu_{k,h}^m \leq 64(d+2)^2 \times O\left(\sqrt{d^4 \ell_i |\mathcal{T}^{m,i,j}| \log^3(dHK)} + d^4 \ell \log^3(dHK)\right) \quad (47)$$

$$\leq O\left(d^4 \sqrt{\ell_i |\mathcal{T}^{m,i,j}| \log^3(dHK)} + d^6 \ell \log^3(dHK)\right). \quad (48)$$

820 Recall that (48) requires  $\mathbf{x}_{k,h}^m \mu_{k,h}^m \in [\ell_j/2, \ell_j]$ , which would be false for  $j = L_2 + 1$ . In this corner  
 821 case,  $j = L_2 + 1$ , we have

$$\sum_i \sum_{(k,h) \in \mathcal{T}^{m,i,j}} \mathbf{x}_{k,h}^m \mu_{k,h}^m \leq KH \ell_j \leq O(1). \quad (49)$$

822 Finally, combining (48) and (49), we have

$$\begin{aligned} \sum_{k,h} \tilde{\mathbf{x}}_{k,h}^m \mu_{k,h}^m &= \sum_{i,j} \sum_{(k,h) \in \mathcal{T}^{m,i,j}} \mathbf{x}_{k,h}^m \mu_{k,h}^m \\ &\leq O(1) + \sum_{i,j} O\left(d^4 \sqrt{\ell_i |\mathcal{T}^{m,i,j}| \log^3(dHK)} + L_2 d^6 \ell \log^3(dHK)\right) \\ &\leq O\left(d^4 \sqrt{\sum_{k,h} \tilde{\eta}_{k,h}^m \log^7(dHK)} + d^6 \ell \log^5(dHK)\right), \end{aligned} \quad (50)$$

823 where (50) uses that  $\ell_i |\mathcal{T}^{m,i,j}| \leq O(1 + \sum_{k,h} \tilde{\eta}_{k,h}^m)$ , which can be proved as follows: for  $i \leq L_1$ , it  
 824 is due to  $\eta_{k,h}^m \geq \ell_i/2$ ; for  $i = L_1 + 1$ , it is due to  $1/\ell_i \geq KH \geq |\mathcal{T}^{m,i,j}|$ .  $\square$

825 **E.6 Proof of Lemma 25**

826 *Proof.* Define

$$\check{\zeta}_{k,h}^m = (P_{s_h^k, a_h^k}(V_{h+1}^k)^{2^{m+1}} - (P_{s_h^k, a_h^k}(V_{h+1}^k)^{2^m})^2) I_h^k.$$

827 We note that  $\check{M}_m$  is a martingale, so by Lemma 11 with a union bound over  $m$ , we have

$$\Pr \left[ \forall m \in \Lambda_0 : |\check{M}_m| \leq 2 \sqrt{2 \sum_{k,h} \check{\zeta}_{k,h}^m \ln \frac{1}{\delta}} + 4 \ln \frac{1}{\delta} \right] \geq 1 - O(\delta \log^2(dKH)). \quad (51)$$

828 By the definition of  $\check{\eta}_{k,h}^m$ , we have

$$\sum_{k,h} \check{\eta}_{k,h}^m \leq \sum_{k,h} \left( \check{\zeta}_{k,h}^m + \max_{\theta \in \Theta_k} \check{x}_{k,h}^{m+1}(\theta - \theta^*) + 2 \max_{\theta \in \Theta_k} \check{x}_{k,h}^m(\theta - \theta^*) \right) \quad (52)$$

$$= \sum_{k,h} \check{\zeta}_{k,h}^m + \check{R}_{m+1} + 2\check{R}_m, \quad (53)$$

829 We have that

$$\begin{aligned} \sum_{k,h} \check{\zeta}_{k,h}^m &= \sum_{k,h} \left( P_{s_h^k, a_h^k}(V_{h+1}^k)^{2^{m+1}} - (P_{s_h^k, a_h^k}(V_{h+1}^k)^{2^m})^2 \right) I_h^k \\ &\leq \sum_{k,h} \left( P_{s_h^k, a_h^k}(V_{h+1}^k)^{2^{m+1}} - (V_{h+1}^k(s_{h+1}^k))^{2^{m+1}} \right) I_h^k + \sum_{k,h} (V_h^k(s_h^k))^{2^{m+1}} (I_h^k - I_{h+1}^k) \\ &\quad + \sum_{k,h} \left( (V_h^k(s_h^k))^{2^{m+1}} - (P_{s_h^k, a_h^k}(V_{h+1}^k)^{2^m})^2 \right) I_h^k \\ &\leq \check{M}_{m+1} + O(d \log^5(dHK)) + \sum_{k,h} \left( (V_h^k(s_h^k))^{2^{m+1}} - (P_{s_h^k, a_h^k}(V_{h+1}^k)^{2^m})^2 \right) I_h^k \\ &\leq \check{M}_{m+1} + O(d \log^5(dHK)) + \sum_{k,h} \left( (V_h^k(s_h^k))^{2^{m+1}} - (P_{s_h^k, a_h^k} V_{h+1}^k)^{2^{m+1}} \right) \\ &\leq \check{M}_{m+1} + O(d \log^5(dHK)) + 2^{m+1} \sum_{k,h} I_h^k \cdot \max\{V_h^k(s_h^k) - P_{s_h^k, a_h^k} V_{h+1}^k, 0\} \\ &\leq \check{M}_{m+1} + O(d \log^5(dHK)) + 2^{m+1} \sum_{k,h} I_h^k \left( r(s_h^k, a_h^k) + \max_{\theta \in \Theta_k} \mathbf{x}_{k,h}^0(\theta - \theta^*) \right) \\ &\leq \check{M}_{m+1} + O(d \log^5(dHK)) + 2^{m+1} (K + \check{R}_0). \end{aligned} \quad (54)$$

830 Finally, by (53), (54) and Lemma 24, we have

$$\begin{aligned} \check{R}_m &\leq O \left( d^4 \sqrt{(\check{M}_{m+1} + O(d \log^5(dHK)) + 2^{m+1}(K + \check{R}_0) + \check{R}_{m+1} + 2\check{R}_m) \iota \log^7(dHK)} + d^6 \iota \log^5(dHK) \right) \\ &\leq O \left( d^4 \sqrt{(\check{M}_{m+1} + 2^{m+1}(K + \check{R}_0) + \check{R}_{m+1} + \check{R}_m) \iota \log^7(dHK)} + d^6 \iota \log^5(dHK) \right), \end{aligned} \quad (55)$$

831 which proves the first part of the lemma. By (51) and (54), we have

$$|\check{M}_m| \leq O \left( \sqrt{(\check{M}_{m+1} + O(d \log^5(dHK)) + 2^{m+1}(K + \check{R}_0)) \log(1/\delta)} + \log(1/\delta) \right), \quad (56)$$

832 which proves the second part of the lemma.  $\square$

833 **E.7 Proof of Lemma 23**

834 *Proof.* Let  $b_m = \check{R}_m + |\check{M}_m|$ . By (55) and (56), we can bound  $b_m$  recursively as

$$b_m \leq O \left( \sqrt{d^9 \log^5(Td) \log \frac{1}{\delta}} \sqrt{b_m + b_{m+1} + 2^{m+1}(K + \check{R}_0) + d^7 \log^6(Td) \log \frac{1}{\delta}} \right). \quad (57)$$

835 Note that  $b_m \leq 2KH$  for  $m \in \Lambda_1$ . By Lemma 12 with parameters

$$\lambda_1 = 2KH, \quad \lambda_2 = \sqrt{d^9 \log^5(Td) \log(1/\delta)}, \quad \lambda_3 = K + \check{R}_0, \quad \lambda_4 = d^7 \log^6(Td) \log(1/\delta),$$

836 we obtain that

$$\check{R}_0 \leq b_0 \leq O\left(\sqrt{d^9(K + \check{R}_0) \log^5(Td) \log(1/\delta)} + d^9 \log^6(Td) \log(1/\delta)\right),$$

837 which implies

$$b_0 \leq O\left(d^{4.5} \sqrt{K \log^5(Td) \log(1/\delta)} + d^9 \log^6(Td) \log(1/\delta)\right)$$

838 and completes the proof. □