

---

# Bypassing the Simulator: Near-Optimal Adversarial Linear Contextual Bandits

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We consider the adversarial linear contextual bandit problem, where the loss vectors  
2 are selected fully adversarially and the per-round action set (i.e. the context) is  
3 drawn from a fixed distribution. Existing methods for this problem either require  
4 access to a simulator to generate free i.i.d. contexts, achieve a sub-optimal regret no  
5 better than  $\tilde{O}(T^{5/6})$ , or are computationally inefficient. We greatly improve these  
6 results by achieving a regret of  $\tilde{O}(\sqrt{T})$  without a simulator, while maintaining  
7 computational efficiency when the action set in each round is small. In the special  
8 case of sleeping bandits with adversarial loss and stochastic arm availability, our  
9 result answers affirmatively the open question by [SGV20] on whether there exists  
10 a polynomial-time algorithm with  $\text{poly}(d)\sqrt{T}$  regret. Our approach naturally  
11 handles the case where the loss is linear up to an additive misspecification error,  
12 and our regret shows near-optimal dependence on the magnitude of the error.

## 13 1 Introduction

14 Contextual bandit is a widely used model for sequential decision making. The interaction between the  
15 learner and the environment proceeds in rounds: in each round, the environment provides a context;  
16 based on it, the learner chooses an action and receive a reward. The goal is to maximize the total  
17 reward across multiple rounds. This model has found extensive applications in fields such as medical  
18 treatment [TM17], personalized recommendations [BLL<sup>+</sup>11], and online advertising [CLRS11].

19 Algorithms for contextual bandits with provable guarantees have been developed under various  
20 assumptions. In the linear regime, the most extensively studied model is the *stochastic linear*  
21 *contextual bandit*, in which the context can be arbitrarily distributed in each round, while the reward  
22 is determined by a fixed linear function of the context-action pair. Near-optimal algorithms for  
23 this setting have been established in, e.g., [CLRS11, AYPS11, LWZ19, FGMZ20]. Another model,  
24 which is the focus of this paper, is the *adversarial linear contextual bandit*, in which the context is  
25 drawn from a fixed distribution, while the reward is determined by a time-varying linear function of  
26 the context-action pair. <sup>1</sup> A computationally efficient algorithm for this setting is first proposed by  
27 [NO20]. However, existing research for this setting still faces challenges in achieving near-optimal  
28 regret and sample complexity when the context distribution is unknown.

29 The algorithm by [NO20] requires the learner to have *full knowledge* on the context distribution, and  
30 access to an *exploratory policy* that induces a feature covariance matrix with a smallest eigenvalue  
31 at least  $\lambda$ . Under these assumptions, their algorithm provides a regret guarantee of  $\tilde{O}(\sqrt{dT}/\lambda)$ ,

---

<sup>1</sup>Apparently, the stochastic and adversarial linear contextual bandits defined here are incomparable, and their names do not fully capture their underlying assumptions. However, these are the terms commonly used in the literature (e.g., [AYPS11, NO20]).

Table 1: Related works in the “S-A” category. CB stands for contextual bandits and SB stands for semi-bandits. The relations among settings are as follows: Sleeping Bandit  $\subset$  Contextual SB  $\subset$  Linear CB, Linear CB  $\subset$  Linear MDP, and Linear CB  $\subset$  General CB. The table compares our results with the Pareto frontier of the literature. For algorithms dealing more general settings, we have carefully translated their techniques to Linear CB and reported the resulting bounds.  $\Sigma_\pi$  denotes the feature covariance matrix induced by policy  $\pi$ .  $|\mathcal{A}|$  and  $|\Pi|$  are sizes of the action set and the policy set.

Target Setting	Algorithm	Regret	Simulator	Computation	Assumption
General CB	[SLKS16]	$(\log  \Pi )^{1/3} ( \mathcal{A} T)^{2/3}$	✓	$\text{poly}( \mathcal{A} , \log  \Pi , T)$	ERM Oracle
Linear MDP	[DLWZ23]	$\sqrt{dT \log  \mathcal{A} }$	✓	$\text{poly}( \mathcal{A} , d, T)$	
	[DLWZ23, SKM23]	$d(\log  \mathcal{A} )^{1/6} T^{5/6}$		$\text{poly}( \mathcal{A} , d, T)$	
	[KZWL23]	$(dT^4)^{1/5} + \text{poly}(\frac{1}{\lambda})$		$T^d$	$\exists \pi, \Sigma_\pi \succeq \lambda I$
Linear CB	Algorithm 1	$d^2 \sqrt{T}$		$\text{poly}( \mathcal{A} , d, T)$	
	Algorithm 2	$d\sqrt{T}$		$T^d$	
Contextual SB	[NV14]	$(dT)^{2/3}$		$\text{poly}(d, T)$	
Sleeping Bandit	[SGV20]	$\sqrt{2dT}$		$\text{poly}(d, T) ( \mathcal{A}  \leq d)$	

where  $d$  is the feature dimension and  $T$  is the number of rounds. These assumptions are relaxed in the work of [LWL21], who studied a more general linear MDP setting. When specialized to linear contextual bandits, [LWL21] only requires access to a *simulator* from which the learner can draw free i.i.d. contexts. Their algorithm achieves a  $\tilde{O}((dT)^{2/3})$  regret. The regret is further improved to the near-optimal one  $\tilde{O}(d\sqrt{T})$  by [DLWZ23] through refined loss estimator construction.

All results that attain  $\tilde{O}(T^{2/3})$  or  $\tilde{O}(\sqrt{T})$  regret bound discussed above rely on access to the simulator. In their algorithms, the number of calls to the simulator significantly exceeds the number of interactions between the environment and the learner, but this is concealed from the regret bound. Therefore, their regret bounds do not accurately reflect the sample complexity of their algorithms. Another set of results for linear MDPs [LWL21, DLWZ23, SKM23, KZWL23] also consider the simulator-free scenario, essentially using interactions with the environment to fulfill the original purpose of the simulator. When applying their techniques to linear contextual bandits, their algorithms only achieve a regret bound of  $\tilde{O}(T^{5/6})$  at best (see detailed analysis and comparison in Appendix G).

Our result significantly improves the previous ones: without simulators, we develop an algorithm that ensures a regret bound of order  $\tilde{O}(d^2 \sqrt{T})$ , and it is computationally efficient as long as the size of the action set is small in each round (similar to all previous work). Unlike previous algorithms which always collect new contexts (through simulators or interactions with the environment) to estimate the feature covariance matrix, we leverage the context samples the learner received in the past to do this. Although natural, establishing a near-tight regret requires highly efficient use of context samples, necessitating a novel way to construct the estimator of feature covariance matrix and a tighter concentration bound for it. Additionally, to address the potentially large magnitude and the bias of the loss estimator, we turn to the use of log-determinant (logdet) barrier in the follow-the-regularized-leader (FTRL) framework. Logdet accommodates larger loss estimators and induces a larger bonus term to cancel the bias of the loss estimator, both of which are crucial for our result.

Our setting subsumes sleeping bandits with stochastic arm availability [KMB09, SGV20] and combinatorial semi-bandits with stochastic action sets [NV14]. Our result answers affirmatively the main open question left by [SGV20] on whether there exists a polynomial-time algorithm with  $\text{poly}(d)\sqrt{T}$  regret for sleeping bandits with adversarial loss and stochastic availability.

As a side result, we give a computationally inefficient algorithm that achieves an improved  $\tilde{O}(d\sqrt{T})$  regret without a simulator. While this is a direct extension from the EXP4 algorithm [ACBFS02], such a result has not been established to our knowledge, so we include it for completeness.

## 1.1 Related work

We review the literature of various contextual bandit problems, classifying them based on the nature of the context and the reward function, specifically whether they are stochastic/fixed or adversarial.

66 **Contextual bandits with i.i.d. contexts and fixed reward functions (S-S)** Significant progress has  
 67 been made in contextual bandits with i.i.d. contexts and fixed reward functions, under general reward  
 68 function classes or policy classes [LZ07, DHK<sup>+</sup>11, ADK<sup>+</sup>12, AHK<sup>+</sup>14, SLX22]. In [DHK<sup>+</sup>11,  
 69 ADK<sup>+</sup>12, AHK<sup>+</sup>14], the algorithms also use previously collected contexts to estimate the inverse  
 70 probability of selecting actions under the current policy. However, these results only obtain regret  
 71 bounds that polynomially depend on the number of actions. Furthermore, these results rely on having  
 72 a fixed reward function, making their techniques not directly applicable to our case even if we allow  
 73 poly-action dependence. For the linear case, [HYF22] provides a reduction from the original problem  
 74 to one with a fixed action set and fixed reward function. Our work can be viewed as a generalization  
 75 of their result to the adversarial reward setting.

76 **Contextual bandits with adversarial contexts and fixed reward functions (A-S)** In this category,  
 77 the most well-known results are in the linear setting [CLRS11, AYPS11, ZHZ<sup>+</sup>23]. Besides the linear  
 78 case, previous work has investigated specific reward function classes [RVR13, LKFS22, FAD<sup>+</sup>18].  
 79 Recently, [FR20] introduced a general approach to deal with general function classes with a finite  
 80 number of actions, which has since been improved or extended by [FK21, FRSLX21, Zha22]. This  
 81 category of problems is not directly comparable to the setting studied in this paper, but both capture a  
 82 certain degree of non-stationarity of the environment.

83 **Contextual bandits with i.i.d. contexts and adversarial reward functions (S-A)** This is the  
 84 category which our work falls into. Several oracle efficient algorithms that require simulators have  
 85 been proposed for general policy classes [RS16, SLKS16]. The oracle they use (i.e., the empirical risk  
 86 minimization, or ERM oracle), however, is not generally implementable in an efficient manner. For  
 87 the linear case, the first computationally efficient algorithm is by [NO20], under the assumption that  
 88 the context distribution is known. This is followed by [OMvE<sup>+</sup>23] to obtain refined data-dependent  
 89 bounds. A series of works [NO21, LWL21, DLWZ23, SKM23] apply similar techniques to linear  
 90 MDPs, but when specialized to linear contextual bandits, they all assume known context distribution,  
 91 or access to a simulator, or only achieves a regret no better than  $\tilde{O}(T^{5/6})$ . The work of [KZWL23]  
 92 also studies linear MDPs; when specialized to contextual bandits, they obtain a regret bound of  
 93  $\tilde{O}(T^{4/5} + \text{poly}(\frac{1}{\lambda}))$  without a simulator but with a computationally inefficient algorithm and an  
 94 undesired inverse dependence on the smallest eigenvalue of the covariance matrix. Related but  
 95 simpler settings have also been studied. The sleeping bandit problem with stochastic arm availability  
 96 and adversarial reward [KNMS10, KMB09, SGV20] is a special case of our problem where the  
 97 context is always a subset of standard unit vectors. Another special case is the combinatorial semi-  
 98 bandit problem with stochastic action sets and adversarial reward [NV14]. While these are special  
 99 cases, the regret bounds in these works are all worse than  $\tilde{O}(\text{poly}(d)\sqrt{T})$ . Therefore, our result also  
 100 improves upon theirs.<sup>2</sup>

101 **Contextual bandits with adversarial contexts and adversarial reward functions (A-A)** When  
 102 both contexts and reward functions are adversarial, there are computational [KS14] and oracle-call  
 103 [HK16] lower bounds showing that no sublinear regret is achievable unless the computational cost  
 104 scales polynomially with the size of the policy set. Even for the linear case, [NO20] argued that  
 105 the problem is at least as hard as online learning a one-dimensional threshold function, for which  
 106 sublinear regret is impossible. For this challenging category, besides using the inefficient EXP4  
 107 algorithm, previous work makes stronger assumptions on the contexts [SKS16] or resorts to alternative  
 108 benchmarks such as dynamic regret [LWAL18, CLLW19] and approximate regret [EZWLK21].

109 **Lifting and exploration bonus for high-probability adversarial linear bandits** Our technique  
 110 is related to those obtaining high-probability bounds for linear bandits. Early development in this  
 111 line of research only achieves computational efficiency when the action set size is small [BDH<sup>+</sup>08]  
 112 or only applies to special action sets such as two-norm balls [AR09]. Recently, near-optimal high-  
 113 probability bounds for general convex action sets have been obtained by lifting the problem to  
 114 a higher dimensional one, which allows for a computationally efficient way to impose bonuses  
 115 [LLWZ20, ZL22]. The lifting and the bonus ideas we use are inspired by them, though for different  
 116 purposes. However, due to the extra difficulty arising in the contextual case, currently we only obtain  
 117 a computationally efficient algorithm when the action set size is small.

<sup>2</sup>For combinatorial semi-bandit problems, our algorithm is not as computationally efficient as [NV14], which  
 can handle exponentially large action sets.

## 118 1.2 Computational Complexity

119 Our main algorithm is based on log-determinant barrier optimization similar to [FGMZ20, ZL22].  
 120 Computing its action distribution is closely related to computing the D-optimal experimental design  
 121 [KT90]. Per step, this is shown to require  $\tilde{O}(|\mathcal{A}_t|\text{poly}(d))$  computational and  $\tilde{O}(\log(|\mathcal{A}_t|)\text{poly}(d))$   
 122 memory complexity [FGMZ20, Prop 1], where  $|\mathcal{A}_t|$  is the action set size at round  $t$ . The computa-  
 123 tional bottleneck comes from (approximately) maximizing a quadratic function over the action set. It  
 124 is an open question whether linear optimization oracles or other type of oracles can lead to efficient  
 125 implementation of our algorithm for continuous action sets.

126 On the other hand, we are unaware of *any* linear context bandit algorithm that provably avoids  $|\mathcal{A}|$   
 127 computation per round while maintaining a  $o(|\mathcal{A}|)$  regret dependence in the frequentist setting. The  
 128 LinUCB algorithm [CLRS11, AYPS11] suffers from the same quadratic function maximization issue,  
 129 and therefore is computationally comparable to our algorithm. The SquareCB.Lin algorithm by  
 130 [FGMZ20] is based on the same log-determinant barrier optimization. Another recent algorithm by  
 131 [Zha22] only admits an efficient implementation for continuous action sets in the Bayesian setting  
 132 but not in the frequentist setting (though they provided an efficient heuristic implementation in their  
 133 experiments). The Thompson sampling algorithm by [AG13], which has efficient implementation,  
 134 also relies on well-specified Gaussian prior.

## 135 2 Preliminaries

136 We study the adversarial linear contextual bandit problem where the loss vectors are selected fully  
 137 adversarially and the per-round action set (i.e. the context) is drawn from a fixed distribution. The  
 138 learner and the environment interact in the following way. Let  $\mathbb{B}_2^d$  be the L2-norm unit ball in  $\mathbb{R}^d$ .

139 For  $t = 1, \dots, T$ ,

- 140 1. The environment decides an adversarial loss vector  $y_t \in \mathbb{B}_2^d$ , and generates a random action  
 141 set (i.e., context)  $\mathcal{A}_t \subset \mathbb{B}_2^d$  from a fixed distribution  $D$  independent from anything else.
- 142 2. The learner observes  $\mathcal{A}_t$ , and (randomly) chooses an action  $a_t \in \mathcal{A}_t$ .
- 143 3. The learner receives the loss  $\ell_t \in [-1, 1]$  with  $\mathbb{E}[\ell_t] = \langle a_t, y_t \rangle$ .

144 A policy  $\pi$  is a mapping which, given any action set  $\mathcal{A} \subset \mathbb{R}^d$ , maps it to an element in the convex hull  
 145 of  $\mathcal{A}$  (denoted as  $\text{conv}(\mathcal{A})$ ). We use  $\pi(\mathcal{A}) \in \text{conv}(\mathcal{A})$  to refer to the element that it maps  $\mathcal{A}$  to. The  
 146 learner's *regret with respect to policy*  $\pi$  is defined as the expected performance difference between  
 147 the learner and policy  $\pi$ :

$$\text{Reg}(\pi) = \mathbb{E} \left[ \sum_{t=1}^T \langle a_t, y_t \rangle - \sum_{t=1}^T \langle \pi(\mathcal{A}_t), y_t \rangle \right]$$

148 where the expectation is taken over all randomness from the environment ( $y_t$  and  $\mathcal{A}_t$ ) and from  
 149 the learner ( $a_t$ ). The *pseudo-regret* (or just *regret*) is defined as  $\text{Reg} = \max_{\pi} \text{Reg}(\pi)$ , where the  
 150 maximization is taken over all possible policies.

151 **Notations** For any matrix  $A$ , we use  $\lambda_{\max}(A)$  and  $\lambda_{\min}(A)$  to denote the maximum and minimum  
 152 eigenvalues of  $A$ , respectively. We use  $\text{Tr}(A)$  to denote the trace of matrix  $A$ . For any action set  $\mathcal{A}$ ,  
 153 let  $\Delta(\mathcal{A})$  be the space of probability measures on  $\mathcal{A}$ . Let  $\mathcal{F}_t = \sigma(\mathcal{A}_s, a_s, \forall s \leq t)$  be the  $\sigma$ -algebra at  
 154 round  $t$ . Define  $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | \mathcal{F}_{t-1}]$ . Given a differentiable convex function  $F: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ , the  
 155 Bregman divergence with respect to  $F$  is defined as  $D_F(x, y) = F(x) - F(y) - \langle \nabla F(y), x - y \rangle$ .  
 156 Given a positive semi-definite (PSD) matrix  $A$ , for any vector  $x$ , define the norm generated by  
 157  $A$  as  $\|x\|_A = \sqrt{x^\top A x}$ . For any context  $\mathcal{A} \subset \mathbb{R}^d$  and  $p \in \Delta(\mathcal{A})$ , define  $\mu(p) = \mathbb{E}_{a \sim p}[a]$  and  
 158  $\text{Cov}(p) = \mathbb{E}_{a \sim p}[(a - \mu(p))(a - \mu(p))^\top]$ . For any  $a$ , define the lifted action  $\mathbf{a} = (a, 1)^\top$  and the lifted

159 covariance matrix  $\widehat{\text{Cov}}(p) = \mathbb{E}_{a \sim p}[\mathbf{a}\mathbf{a}^\top] = \mathbb{E}_{a \sim p} \begin{bmatrix} aa^\top & a \\ a^\top & 1 \end{bmatrix} = \begin{bmatrix} \text{Cov}(p) + \mu(p)\mu(p)^\top & \mu(p) \\ \mu(p)^\top & 1 \end{bmatrix}$ .

160 We use **bold** matrices to denote matrices in the lifted space (e.g., in Algorithm 1 and Definition 1).

---

**Algorithm 1** Logdet-FTRL for linear contextual bandits
 

---

**Definitions:**  $F(\mathbf{H}) = -\log \det(\mathbf{H})$ ,  $\eta_t = \frac{1}{64d\sqrt{t}}$ ,  $\alpha_t = \frac{d}{\sqrt{t}}$ ,  $\beta_t = \frac{100(d+1)^3 \log(3T)}{t-1}$ .

- 1 **for**  $t = 1, 2, \dots$  **do**
  - 2     For all  $\mathcal{A}$ , define  $\mathbf{H}_t^{\mathcal{A}} = \operatorname{argmin}_{\mathbf{H} \in \mathcal{H}^{\mathcal{A}}} \sum_{s=1}^{t-1} \langle \mathbf{H}, \hat{\gamma}_s - \alpha_s \hat{\Sigma}_s^{-1} \rangle + \frac{F(\mathbf{H})}{\eta_t}$ .
  - 3     For all  $\mathcal{A}$ , define  $p_t^{\mathcal{A}} \in \Delta(\mathcal{A})$  such that  $\mathbf{H}_t^{\mathcal{A}} = \widehat{\operatorname{Cov}}(p_t^{\mathcal{A}})$ .
  - 4     Receive  $\mathcal{A}_t$  and sample  $a_t \sim p_t^{\mathcal{A}_t}$ .
  - 5     Observe  $\ell_t \in [-1, 1]$  with  $\mathbb{E}[\ell_t] = a_t^\top y_t$  and construct  $\hat{y}_t = \hat{\Sigma}_t^{-1}(a_t - \hat{x}_t)\ell_t$ , where
 
$$\hat{x}_t = \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbb{E}_{a \sim p_t^{\mathcal{A}_\tau}}[a], \quad \hat{H}_t = \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbb{E}_{a \sim p_t^{\mathcal{A}_\tau}}[(a - \hat{x}_t)(a - \hat{x}_t)^\top], \quad \hat{\Sigma}_t = \hat{H}_t + \beta_t \mathbf{I}.$$
  - 6     Define  $\hat{\mathbf{H}}_t = \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbf{H}_t^{\mathcal{A}_\tau}$  and  $\hat{\Sigma}_t = \hat{\mathbf{H}}_t + \beta_t \mathbf{I}$ .  
 (If  $t = 1$ , define  $\hat{\Sigma}_t^{-1}$  and  $\hat{\Sigma}_t^{-1}$  as zeros).
- 

### 161 3 Follow-the-Regularized-Leader with the Log-Determinant Barrier

162 In this section, we present our main algorithm, [Algorithm 1](#). This algorithm can be viewed as  
 163 instantiating an individual Follow-The-Regularized-Leader (FTRL) algorithm on each action set  
 164 ([Line 2](#)), with all FTRLs sharing the same loss vectors. This perspective has been taking by previous  
 165 works [[NO20](#), [OMvE<sup>+</sup>23](#)] and simplifies the understanding of the problem. The rationale comes  
 166 from the following calculation due to [[NO20](#)]: for any policy  $\pi$  that may depend on  $\mathcal{F}_{t-1}$ ,

$$\mathbb{E}_t[\langle \pi(\mathcal{A}_t), y_t \rangle] = \mathbb{E}_{\mathcal{A}_t}[\mathbb{E}_{y_t}[\langle \pi(\mathcal{A}_t), y_t \rangle \mid \mathcal{F}_{t-1}]] = \mathbb{E}_{\mathcal{A}_0}[\mathbb{E}_{y_t}[\langle \pi(\mathcal{A}_0), y_t \rangle \mid \mathcal{F}_{t-1}]] = \mathbb{E}_t[\langle \pi(\mathcal{A}_0), y_t \rangle]$$

167 where  $\mathcal{A}_0$  is a sample drawn from  $D$  independent of all interaction history. This allows us to calculate  
 168 the regret as

$$\mathbb{E} \left[ \sum_{t=1}^T \langle \pi_t(\mathcal{A}_t) - \pi(\mathcal{A}_t), y_t \rangle \right] = \mathbb{E} \left[ \sum_{t=1}^T \langle \pi_t(\mathcal{A}_0) - \pi(\mathcal{A}_0), y_t \rangle \right] \quad (1)$$

169 where  $\pi_t$  is the policy used by the learner at time  $t$ . Note that this view does not require the learner  
 170 to simultaneously “run” an algorithm on every action set since the learner only needs to calculate  
 171 the policy on  $\mathcal{A}$  whenever  $\mathcal{A}_t = \mathcal{A}$ . In the regret analysis, in view of [Eq. \(1\)](#), it suffices to consider  
 172 a single fixed action set  $\mathcal{A}_0$  drawn from  $D$  and bound the regret on it, even though the learner may  
 173 never execute the policy on it. This  $\mathcal{A}_0$  is called a “ghost sample” in [[NO20](#)].

#### 174 3.1 The lifting idea and the execution of [Algorithm 1](#)

175 Our algorithm is built on the logdet-FTRL algorithm developed by [[ZL22](#)] for high-probability  
 176 adversarial linear bandits, which lifts the original  $d$ -dimensional problem over the feature space to  
 177 a  $(d+1) \times (d+1)$  one over the covariance matrix space, with the regularizer being the negative  
 178 log-determinant function. In our case, we instantiate an individual logdet-FTRL on each action set.  
 179 The motivation behind [[ZL22](#)] to lift the problem to the space of covariance matrix is that it casts the  
 180 problem to one in the positive orthant, which allows for an easier way to construct the *bonus* term that  
 181 is crucial to compensate the variance of the losses, enabling a high-probability bound in their case. In  
 182 our case, we use the same technique to introduce the bonus term, but the goal is to compensate the  
 183 *bias* resulting from the estimation error in the covariance matrix (see [Section 3.4](#)). This bias only  
 184 appears in our contextual case but not in the linear bandit problem originally considered in [[ZL22](#)].

185 As argued previously, we can focus on the learning problem over a fixed action set  $\mathcal{A}$ , and our  
 186 algorithm operates in the lifted space of covariance matrices  $\mathcal{H}^{\mathcal{A}} = \{\widehat{\operatorname{Cov}}(p) : p \in \Delta(\mathcal{A})\} \subset$

187  $\mathbb{R}^{(d+1) \times (d+1)}$ . For this space, we define the lifted loss  $\gamma_t = \begin{bmatrix} 0 & \frac{1}{2}y_t \\ \frac{1}{2}y_t^\top & 0 \end{bmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}$  so that

188  $\langle \widehat{\operatorname{Cov}}(p), \gamma_t \rangle = \mathbb{E}_{a \sim p}[a^\top y_t] = \langle \mu(p), y_t \rangle$  and thus the loss value in the lifted space is the same as  
 189 that in the original space.

190 In each round  $t$ , the FTRL on  $\mathcal{A}$  outputs a lifted covariance matrix  $\mathbf{H}_t^A \in \mathcal{H}^A$  that corresponds to a  
 191 probability distribution  $p_t^A \in \Delta(\mathcal{A})$  such that  $\widehat{\text{Cov}}(p_t^A) = \mathbf{H}_t^A$  (Line 2 and Line 3). Upon receiving  
 192  $\mathcal{A}_t$ , the learner samples an action from  $p_t^A$  and the agent constructs the loss estimator  $\hat{y}_t$  (Line 5).

193 Similarly to the construction of  $\gamma_t$ , we define the lifted loss estimator  $\hat{\gamma}_t = \begin{bmatrix} 0 & \frac{1}{2}\hat{y}_t \\ \frac{1}{2}\hat{y}_t^\top & 0 \end{bmatrix}$  which

194 makes  $\langle \widehat{\text{Cov}}(p), \hat{\gamma}_t \rangle = \mathbb{E}_{a \sim p}[a^\top \hat{y}_t] = \langle \mu(p), \hat{y}_t \rangle$ . The lifted loss estimator is then fed to the FTRL  
 195 on all  $\mathcal{A}$ 's.

196 In the rest of this section, we use the following notation in addition to those defined in Algorithm 1.

197 **Definition 1.** Define  $x_t^A = \mathbb{E}_{a \sim p_t^A}[a]$ ,  $x_t = \mathbb{E}_{\mathcal{A} \sim D}[x_t^A]$ ,  $H_t^A = \mathbb{E}_{a \sim p_t^A}[(a - \hat{x}_t)(a - \hat{x}_t)^\top]$ ,  $H_t =$   
 198  $\mathbb{E}_{\mathcal{A} \sim D}[H_t^A]$ ,  $\mathbf{H}_t = \mathbb{E}_{\mathcal{A} \sim D}[\mathbf{H}_t^A]$ . Let the regret comparator on  $\mathcal{A}$  be  $p_\star^A \in \Delta(\mathcal{A})$ , and define  
 199  $u^A = \mathbb{E}_{a \sim p_\star^A}[a]$ ,  $u = \mathbb{E}_{\mathcal{A} \sim D}[u^A]$ ,  $\mathbf{U}^A = \mathbb{E}_{a \sim p_\star^A}[\mathbf{a}\mathbf{a}^\top]$ ,  $\mathbf{U} = \mathbb{E}_{\mathcal{A} \sim D}[\mathbf{U}^A]$ . Notice that the  $x_t^A$  and  
 200  $u^A$  defined here is equivalent to the  $\pi_t(\mathcal{A})$  and  $\pi(\mathcal{A})$  in Eq. (1), respectively.

### 201 3.2 The construction of loss estimators and feature covariance matrix estimators

202 Our goal is to make  $\hat{y}_t$  in Line 5 an estimator of  $y_t$  with controllable bias and variance. If the context  
 203 distribution is known (as in [NO20]), then a standard unbiased estimator of  $y_t$  is

$$\hat{y}_t = \hat{\Sigma}_t^{-1} a_t \ell_t, \quad \text{where} \quad \hat{\Sigma}_t = \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p_t^A} [a a^\top]. \quad (2)$$

204 To see its unbiasedness, notice that  $\mathbb{E}[a_t \ell_t] = \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p_t^A} [a a^\top y_t]$  and thus  $\mathbb{E}[\hat{y}_t] = y_t$ . This  $\hat{y}_t$ ,  
 205 however, can have a variance that is inversely related to the smallest eigenvalue of the covariance  
 206 matrix  $\hat{\Sigma}_t$ , which can be unbounded in the worst case. This is the main reason why [NO20] does  
 207 not achieve the optimal bound, and requires the bias-variance-tradeoff techniques in [DLWZ23] to  
 208 close the gap. When the context distribution is unknown but the learner has access to a simulator  
 209 [LWL21, DLWZ23, SKM23, KZWL23], the learner can draw free contexts to estimate the covariance  
 210 matrix  $\hat{\Sigma}_t$  up to a very high accuracy without interacting with the environment, making the problem  
 211 close to the case of known context distribution.

212 Challenges arise when the learner has no knowledge about the context distribution and there is no  
 213 simulator. In this case, there are two natural ways to estimate the covariance matrix under the current  
 214 policy. One is to draw new samples from the environment, treating the environment like a simulator.  
 215 This approach is essentially taken by all previous work studying linear models in the ‘‘S-A’’ category.  
 216 However, this is very expensive, and it causes the simulator-equipped bound  $\sqrt{T}$  in [DLWZ23] to  
 217 deteriorate to the simulator-free bound  $T^{5/6}$  at best (see Appendix G for details). The other is to use  
 218 the contexts received in time 1 to  $t$  to estimate the covariance matrix under the policy at time  $t$ . This  
 219 demands a very high efficiency in reusing the contexts samples, and existing ways of constructing the  
 220 covariance matrix and the accompanied analysis by [DLWZ23, SKM23] are insufficient to achieve  
 221 the near-optimal bound even with context reuse. This necessitates our tighter construction of the  
 222 covariance matrix estimator and tighter concentration bounds for it.

223 Our construction of the loss estimator (Line 5) is

$$\hat{y}_t = \hat{\Sigma}_t^{-1} (a_t - \hat{x}_t) \ell_t \quad \text{where} \quad \hat{\Sigma}_t = \mathbb{E}_{\mathcal{A} \sim \hat{D}_t} \mathbb{E}_{a \sim p_t^A} [(a - \hat{x}_t)(a - \hat{x}_t)^\top] + \beta_t I \quad (3)$$

224 where  $\hat{D}_t = \text{Uniform}\{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_{t-1}\}$ ,  $\hat{x}_t = \mathbb{E}_{\mathcal{A} \sim \hat{D}_t} \mathbb{E}_{a \sim p_t^A}[a]$ , and  $\beta_t = \tilde{\mathcal{O}}(d^3/t)$ . Comparing  
 225 Eq. (3) with Eq. (2), we see that besides using the empirical context distribution  $\hat{D}_t$  in place of the  
 226 ground truth  $D$  and adding a small term  $\beta_t I$  to control the smallest eigenvalue of the covariance  
 227 matrix, we also centralize the features by  $\hat{x}_t$ , an estimation of the mean features under the current  
 228 policy. The centralization is important in making the bias  $y_t - \hat{y}_t$  appear in a nice form that can  
 229 be compensated by a bonus term. The estimator might seem problematic on first sight, because  $p_t^A$   
 230 is strongly dependent on  $\hat{D}_t$ , which rules out canonical concentration bounds. We circumvent this  
 231 issue by leveraging the special structure of  $p_t$  in Algorithm 1, which allows for a union bound over  
 232 a sufficient covering of all potential policies (Appendix C.3). The analysis on the bias of this loss  
 233 estimator is also non-standard, which is the key to achieve the near-optimal bound. In the next two  
 234 subsections, we explain how to bound the bias of this loss estimator (Section 3.3), and how the bonus  
 235 term can be used to compensate the bias (Section 3.4).

236 **3.3 The bias of the loss estimator**

237 Since the true loss vector is  $y_t$  and we use the loss estimator  $\hat{y}_t$  in the update, there is a bias term  
 238 emerging in the regret bound at time  $t$ :

$$\mathbb{E}_t \left[ \langle x_t^{A_0} - u^{A_0}, y_t - \hat{y}_t \rangle \right] = \mathbb{E}_t \left[ \langle x_t - u, y_t - \hat{y}_t \rangle \right] = \mathbb{E}_t \left[ (x_t - u)^\top \left( I - \hat{\Sigma}_t^{-1} (a_t - \hat{x}_t) a_t^\top \right) y_t \right]$$

239 where definitions of  $x_t^A, u^A, x_t, u$  can be found in [Definition 1](#), and we use the definition of  $\hat{y}_t$  in  
 240 [Eq. \(3\)](#) in the last equality. Now taking expectation over  $\mathcal{A}_t$  and  $a_t$  conditioned on  $\mathcal{F}_{t-1}$ , we can  
 241 further bound the expectation in the last expression by

$$\begin{aligned} & (x_t - u)^\top \left( I - \hat{\Sigma}_t^{-1} H_t \right) y_t - (x_t - u)^\top \hat{\Sigma}_t^{-1} (x_t - \hat{x}_t) \hat{x}_t^\top y_t \\ & \leq \|x_t - u\|_{\hat{\Sigma}_t^{-1}} \|(\hat{\Sigma}_t - H_t) y_t\|_{\hat{\Sigma}_t^{-1}} + \|x_t - u\|_{\hat{\Sigma}_t^{-1}} \|x_t - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}} \end{aligned} \quad (4)$$

242 (see [Definition 1](#) for the definition of  $H_t$ ). The two terms  $\|(\hat{\Sigma}_t - H_t) y_t\|_{\hat{\Sigma}_t^{-1}}$  and  $\|x_t -$   
 243  $\hat{x}_t\|_{\hat{\Sigma}_t^{-1}}$  in [Eq. \(4\)](#) are related to the error between the empirical context distribution  $\hat{D}_t =$   
 244  $\text{Uniform}\{\mathcal{A}_1, \dots, \mathcal{A}_{t-1}\}$  and the true distribution  $D$ . We handle them through novel analysis and  
 245 bound both of them by  $\tilde{\mathcal{O}}(\sqrt{d^3/t})$ . See [Lemma 13](#) and [Lemma 14](#) for details. The techniques we  
 246 use in these two lemmas surpass those in [\[DLWZ23, SKM23\]](#). As a comparison, a similar term as  
 247  $\|(\hat{\Sigma}_t - H_t) y_t\|_{\hat{\Sigma}_t^{-1}}$  is also presented in [Eq. \(16\)](#) of [\[DLWZ23\]](#) and [Lemma B.5](#) of [\[SKM23\]](#) when  
 248 bounding the bias. While they ensure that this term can be bounded by  $\mathcal{O}(\sqrt{\beta})$  after collecting  
 249  $\mathcal{O}(\beta^{-2})$  new samples ([Lemma 5.1](#) of [\[DLWZ23\]](#) and [Lemma B.1](#) of [\[SKM23\]](#)), we are able to bound  
 250 it by  $\mathcal{O}(1/\sqrt{t})$  only using  $t$  samples that the learner received up to time  $t$ . This essentially improves  
 251 their  $\mathcal{O}(\beta^{-2})$  sample complexity bound to  $\mathcal{O}(\beta^{-1})$ , and can be directly used to obtain an improved  
 252 result for their linear MDP problem. See [Appendix G](#) for detailed comparison.

253 Now we have bounded the regret due to bias of  $\hat{y}_t$  by the order of  $\sqrt{d^3/t} \|x_t - u\|_{\hat{\Sigma}_t^{-1}}$ . The next  
 254 problem is how to mitigate this term. This is also a problem in previous work [\[LWL21, DLWZ23,](#)  
 255 [SKM23\]](#), and it has become clear that this can be handled by incorporating *bonus* in the algorithm.

256 **3.4 The bonus term**

257 To handle a bias term in the form of  $\|x_t - u\|_{\hat{\Sigma}_t^{-1}}$ , we resort to the idea of *bonus*. To illustrate this,  
 258 suppose that instead of feeding  $\hat{y}_t$  to the FTRLs, we feed  $\hat{y}_t - b_t$  for some  $b_t$ . Then this would give  
 259 us a regret bound of the following form:

$$\begin{aligned} \text{Reg} &= \mathbb{E} \left[ \sum_{t=1}^T \langle x_t - u, \hat{y}_t - b_t \rangle \right] + \mathbb{E} \left[ \sum_{t=1}^T \langle x_t - u, y_t - \hat{y}_t \rangle \right] + \mathbb{E} \left[ \sum_{t=1}^T \langle x_t - u, b_t \rangle \right] \\ &\lesssim \tilde{\mathcal{O}}(d^2 \sqrt{T}) + \mathbb{E} \left[ \sum_{t=1}^T \sqrt{\frac{d^3}{t}} \|x_t - u\|_{\hat{\Sigma}_t^{-1}} \right] + \mathbb{E} \left[ \sum_{t=1}^T \langle x_t - u, b_t \rangle \right] \end{aligned} \quad (5)$$

260 where we assume that FTRL can give us  $\tilde{\mathcal{O}}(d^2 \sqrt{T})$  bound for the loss sequence  $\hat{y}_t - b_t$ . Our hope  
 261 here is to design a  $b_t$  such that  $\langle x_t - u, b_t \rangle$  provides a negative term that can be used to cancel the  
 262 bias term  $\sqrt{d^3/t} \|x_t - u\|_{\hat{\Sigma}_t^{-1}}$  in the following manner:

$$\text{bias} + \text{bonus} = \sum_{t=1}^T \left( \sqrt{\frac{d^3}{t}} \|x_t - u\|_{\hat{\Sigma}_t^{-1}} + \langle x_t - u, b_t \rangle \right) \lesssim \tilde{\mathcal{O}}(d^2 \sqrt{T}). \quad (6)$$

263 which gives us a  $\tilde{\mathcal{O}}(d^2 \sqrt{T})$  overall regret by [Eq. \(5\)](#). This approach relies on two conditions to be  
 264 satisfied. First, we have to find a  $b_t$  that makes [Eq. \(6\)](#) hold. Second, we have to ensure that the FTRL  
 265 algorithm achieves a  $\tilde{\mathcal{O}}(d^2 \sqrt{T})$  bound under the loss sequence  $\hat{y}_t - b_t$ .

266 To meet the first condition, we take inspiration from [\[ZL22\]](#) and lift the problem to the space of  
 267 covariance matrix in  $\mathbb{R}^{(d+1) \times (d+1)}$ . Considering the bonus term  $\alpha_t \hat{\Sigma}_t^{-1}$  in the lifted space, we have

$$\langle \mathbf{H}_t - \mathbf{U}, \alpha_t \hat{\Sigma}_t^{-1} \rangle = \alpha_t \text{Tr}(\mathbf{H}_t \hat{\Sigma}_t^{-1}) - \alpha_t \text{Tr}(\mathbf{U} \hat{\Sigma}_t^{-1}) \quad (7)$$

268 Using [Lemma 15](#) and [Corollary 20](#), we can upper bound [Eq. \(7\)](#) by  $\mathcal{O}(d\alpha_t) - \frac{\alpha_t}{4}\|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2$ .  
 269 Though the negative part does not match the bias  $\sqrt{\frac{d^3}{t}}\|x_t - u\|_{\hat{\Sigma}_t^{-1}}$ , cancellation still happens since  
 bias + bonus  $\leq \sum_{t=1}^T \left( \sqrt{\frac{d^3}{t}}\|x_t - u\|_{\hat{\Sigma}_t^{-1}} + d\alpha_t - \frac{\alpha_t}{4}\|\hat{x}_t - u\|_{\hat{\Sigma}_t^{-1}}^2 \right)$   
 $\leq \tilde{\mathcal{O}}(d^2\sqrt{T}) + \sum_{t=1}^T \sqrt{\frac{d^3}{t}}\|x_t - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}} + \sum_{t=1}^T \left( \sqrt{\frac{d^3}{t}}\|\hat{x}_t - u\|_{\hat{\Sigma}_t^{-1}} - \frac{\alpha_t}{4}\|\hat{x}_t - u\|_{\hat{\Sigma}_t^{-1}}^2 \right)$ .

270 Using [Lemma 16](#) to bound the second term above by  $\tilde{\mathcal{O}}(d^3)$ , and AM-GM to bound the third term by  
 271  $\tilde{\mathcal{O}}(\sum_t d^3/(t\alpha_t)) = \tilde{\mathcal{O}}(d^2\sqrt{T})$ , we get [Eq. \(6\)](#), through the help of lifting.

272 To meet the second condition, we have to analyze the regret of FTRL under the loss  $\hat{y}_t - b_t$ . The key  
 273 is to show that the bonus  $\alpha_t\hat{\Sigma}_t^{-1}$  introduces small *stability term* overhead. Thanks to the use of the  
 274 logdet regularizer and its self-concordance property, the extra stability term introduced by the bonus  
 275 can indeed be controlled by the order  $\sqrt{T}$ . The key analysis is in [Lemma 25](#).

276 Previous works rely on exponential weights [[LWL21](#), [DLWZ23](#), [SKM23](#)] rather than logdet-FTRL,  
 277 which comes with the following drawbacks. 1) In [[LWL21](#), [SKM23](#)] where exponential weights is  
 278 combined with standard loss estimators, the bonus introduces large stability term overhead. Therefore,  
 279 their bound can only be  $T^{2/3}$  at best even with simulators. 2) In [[DLWZ23](#)] where exponential weights  
 280 is combined with magnitude-reduced loss estimators, the loss estimator for action  $a$  can no longer  
 281 be represented as a simple linear function  $a^\top \hat{y}_t$ . Instead, it becomes a complex non-linear function.  
 282 This restricts the algorithm's potential to leverage linear optimization oracle over the action set and  
 283 achieve computational efficiency.

### 284 3.5 Overall regret analysis

285 With all the algorithmic elements discussed above, now we give a formal statement for our regret  
 286 guarantee and perform a complete regret analysis. Our main theorem is the following.

287 **Theorem 2.** *Algorithm 1 ensures  $\text{Reg} \leq \mathcal{O}(d^2\sqrt{T} \log T)$ .*

288 *Proof sketch.* Let  $\mathcal{A}_0$  be drawn from  $D$  independently from all the interaction history between the  
 289 learner and the environment. Recalling the definitions in [Definition 1](#), we have

$$\begin{aligned} \text{Reg} &= \mathbb{E} \left[ \sum_{t=1}^T \langle a_t - u^{\mathcal{A}_t}, y_t \rangle \right] = \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_t} - \mathbf{U}^{\mathcal{A}_t}, \gamma_t \rangle \right] = \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \gamma_t \rangle \right] \\ &\leq \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \gamma_t - \hat{\gamma}_t \rangle \right]}_{\text{Bias}} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \rangle \right]}_{\text{Bonus}} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right]}_{\text{FTRL-Reg}} \end{aligned}$$

290 Each term can be bounded as follows:

- 291 • **Bias**  $\leq \mathcal{O}(d^2\sqrt{T} \log T) + \frac{1}{4} \sum_{t=1}^T \alpha_t \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2$  (discussed in [Section 3.3](#)).
- 292 • **Bonus**  $\leq \mathcal{O}(d^2\sqrt{T} \log T) - \frac{1}{4} \sum_{t=1}^T \alpha_t \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2$  (discussed in [Section 3.4](#)).
- 293 • **FTRL-Reg**  $\leq \mathcal{O}(d^2\sqrt{T} \log T)$ .

294 Combining all terms gives the desired bound. The complete proof is provided in [Appendix D](#).  $\square$

### 295 3.6 Handling Misspecification

296 In this subsection, we show how our approach naturally handles the case when the expectation of the  
 297 loss cannot be exactly realized by a linear function but with a misspecification error. In this case, we  
 298 assume that the expectation of the loss is given by  $\mathbb{E}[\ell_t|a_t = a] = f_t(a)$  for some  $f_t: \mathbb{R}^d \rightarrow [-1, 1]$ .  
 299 We define the following notion of misspecification (slightly more refined than that in [\[NO20\]](#)):

300 **Assumption 1** (misspecification).  $\sqrt{\frac{1}{T} \sum_{t=1}^T} \inf_{y \in \mathbb{B}_2^d} \sup_{\mathcal{A} \in \text{supp}(D)} \sup_{a \in \mathcal{A}} (f_t(a) - \langle a, y \rangle)^2 \leq \varepsilon$ .

301 Based on previous discussions, the design idea of [Algorithm 1](#) is to 1) identify the bias of the loss  
 302 estimator, and 2) add necessary bonus to compensate the bias. When there is misspecification, this  
 303 design idea still applies. The difference is that now the loss estimator  $\hat{y}_t$  potentially has more bias due  
 304 to misspecification. Therefore, the bias becomes larger by an amount related to  $\varepsilon$ . Consequently, we  
 305 need to enlarge bonus (raising  $\alpha_t$ ) to compensate it. Due to the larger bonus, we further need to tune  
 306 down the learning rate  $\eta_t$  to make the algorithm stable. Overall, to handle misspecification, when  $\varepsilon$  is  
 307 known, it boils down to using the same algorithm ([Algorithm 1](#)) with adjusted  $\alpha_t$  and  $\eta_t$ . The case  
 308 of unknown  $\varepsilon$  can be handled by the standard meta-learning technique *Corral* [[ALNS17](#), [FGMZ20](#)].  
 309 We defer all details to [Appendix E](#) and only state the final bound here.

310 **Theorem 3.** *Under misspecification, there is an algorithm ensuring  $\text{Reg} \leq \tilde{\mathcal{O}}(d^2\sqrt{T} + \sqrt{d\varepsilon T})$ .*

## 311 4 Linear EXP4

312 To tighten the  $d$ -dependence in the regret bound, we can use the computationally inefficient algorithm  
 313 EXP4 [[ACBFS02](#)]. The original regret bound for EXP4 has a polynomial dependence on the number  
 314 of actions, but here we take the advantage of the linear structure to show a bound that only depends  
 315 on the feature dimension  $d$ . The algorithm is presented in [Algorithm 2](#).

---

### Algorithm 2 Linear EXP4

---

**input:**  $\Pi, \eta, \gamma$ .

**for**  $t = 1, 2, \dots$  **do**

    Receive  $\mathcal{A}_t \subset \mathbb{R}^d$ .

    Construct  $\nu_t \in \Delta(\mathcal{A}_t)$  such that  $\max_{a \in \mathcal{A}_t} \|a\|_{G_t^{-1}}^2 \leq d$ , where  $G_t = \mathbb{E}_{a \sim \nu_t}[aa^\top]$ . Set

$$P_{t,\pi} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_{s,\pi}\right)}{\sum_{\pi' \in \Pi} \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_{s,\pi'}\right)}$$

    and define  $p_{t,a} = \sum_{\pi \in \Pi} P_{t,\pi} \mathbb{I}\{\pi(\mathcal{A}_t) = a\}$ .

    Sample  $a_t \sim \tilde{p}_t = (1 - \gamma)p_t + \gamma\nu_t$  and receive  $\ell_t \in [-1, 1]$  with  $\mathbb{E}[\ell_t] = \langle a_t, y_t \rangle$ .

    Construct  $\forall \pi \in \Pi$ :  $\hat{\ell}_{t,\pi} = \langle \pi(\mathcal{A}_t), \tilde{H}_t^{-1} a_t \ell_t \rangle$ , where  $\tilde{H}_t = \mathbb{E}_{a \sim \tilde{p}_t}[aa^\top]$ .

---

316 To run [Algorithm 2](#), we restrict ourselves to a finite policy class. The policy class we use in the  
 317 algorithm is the set of linear policies defined as

$$\Pi = \left\{ \pi_\theta : \theta \in \Theta, \pi_\theta(\mathcal{A}) = \underset{a \in \mathcal{A}}{\operatorname{argmin}} a^\top \theta \right\} \quad (8)$$

318 where  $\Theta$  is an 1-net of  $[-T, T]^d$ . The next theorem shows that this suffices to give us near-optimal  
 319 bounds for our problem. The proof is given in [Appendix F](#).

320 **Theorem 4.** *With  $\gamma = 2d\sqrt{(\log T)/T}$  and  $\eta = \sqrt{(\log T)/T}$ , [Algorithm 2](#) with the policy class  
 321 defined in [Eq. \(8\)](#) guarantees  $\text{Reg} = \mathcal{O}(d\sqrt{T \log T})$ .*

322 Note that this result technically also holds in the ‘‘A-A’’ category with respect to the policy class  
 323 defined in [Eq. \(8\)](#). However, this policy class is *not* necessarily a sufficient cover of all policies of  
 324 interest when the contexts and losses are adversarial.

## 325 5 Conclusions

326 We derived the first algorithm that obtains  $\sqrt{T}$  regret in contextual linear bandits with stochastic  
 327 action sets in the absence of a simulator or strong assumptions on the distribution. As a side result,  
 328 we obtained the first computationally efficient  $\text{poly}(d)\sqrt{T}$  algorithm for adversarial sleeping bandits  
 329 with general stochastic arm availabilities. We believe the techniques in this paper will be useful for  
 330 improving results for simulator-free linear MDPs as well.

331 **References**

- 332 [ACBFS02] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The non-  
333 stochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77,  
334 2002.
- 335 [ADK<sup>+</sup>12] Alekh Agarwal, Miroslav Dudík, Satyen Kale, John Langford, and Robert Schapire.  
336 Contextual bandit learning with predictable rewards. In *Artificial Intelligence and*  
337 *Statistics*, pages 19–26. PMLR, 2012.
- 338 [AG13] Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with  
339 linear payoffs. In *International conference on machine learning*, pages 127–135.  
340 PMLR, 2013.
- 341 [AHK<sup>+</sup>14] Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert  
342 Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In  
343 *International Conference on Machine Learning*, pages 1638–1646. PMLR, 2014.
- 344 [AHR09] Jacob D Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An  
345 efficient algorithm for bandit linear optimization. 2009.
- 346 [ALNS17] Alekh Agarwal, Haipeng Luo, Behnam Neyshabur, and Robert E Schapire. Corraling  
347 a band of bandit algorithms. In *Conference on Learning Theory*, pages 12–38. PMLR,  
348 2017.
- 349 [Ano23] Anonymous, 2023.
- 350 [AR09] Jacob Abernethy and Alexander Rakhlin. Beating the adaptive bandit with high  
351 probability. In *2009 Information Theory and Applications Workshop*, pages 280–289.  
352 IEEE, 2009.
- 353 [AYPS11] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for  
354 linear stochastic bandits. *Advances in neural information processing systems*, 24:2312–  
355 2320, 2011.
- 356 [BDH<sup>+</sup>08] Peter L Bartlett, Varsha Dani, Thomas P Hayes, Sham M Kakade, Alexander Rakhlin,  
357 and Ambuj Tewari. High-probability regret bounds for bandit online linear optimization.  
358 In *COLT*, pages 335–342, 2008.
- 359 [BLL<sup>+</sup>11] Alina Beygelzimer, John Langford, Lihong Li, Lev Reyzin, and Robert Schapire.  
360 Contextual bandit algorithms with supervised learning guarantees. In *Proceedings of*  
361 *the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages  
362 19–26. JMLR Workshop and Conference Proceedings, 2011.
- 363 [CLLW19] Yifang Chen, Chung-Wei Lee, Haipeng Luo, and Chen-Yu Wei. A new algorithm for  
364 non-stationary contextual bandits: Efficient, optimal and parameter-free. In *Conference*  
365 *on Learning Theory*, pages 696–726. PMLR, 2019.
- 366 [CLRS11] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with  
367 linear payoff functions. In *Proceedings of the Fourteenth International Conference on*  
368 *Artificial Intelligence and Statistics*, pages 208–214. JMLR Workshop and Conference  
369 Proceedings, 2011.
- 370 [DHK<sup>+</sup>11] M Dudik, D Hsu, S Kale, N Karampatziakis, J Langford, L Reyzin, and T Zhang.  
371 Efficient optimal learning for contextual bandits. In *Proceedings of the 27th Conference*  
372 *on Uncertainty in Artificial Intelligence, UAI 2011*, page 169, 2011.
- 373 [DLWZ23] Yan Dai, Haipeng Luo, Chen-Yu Wei, and Julian Zimmert. Refined regret for adver-  
374 sarial mdps with linear function approximation. *arXiv preprint arXiv:2301.12942*,  
375 2023.
- 376 [DWZ23a] Christoph Dann, Chen-Yu Wei, and Julian Zimmert. Best of both worlds policy  
377 optimization. *arXiv preprint arXiv:2302.09408*, 2023.

- 378 [DWZ23b] Christoph Dann, Chen-Yu Wei, and Julian Zimmert. A blackbox approach to best of  
379 both worlds in bandits and beyond. *arXiv preprint arXiv:2302.09739*, 2023.
- 380 [EZWLK21] Ehsan Emamjomeh-Zadeh, Chen-Yu Wei, Haipeng Luo, and David Kempe. Adversarial  
381 online learning with changing action sets: Efficient algorithms with approximate regret  
382 bounds. In *Algorithmic Learning Theory*, pages 599–618. PMLR, 2021.
- 383 [FAD<sup>+</sup>18] Dylan Foster, Alekh Agarwal, Miroslav Dudík, Haipeng Luo, and Robert Schapire.  
384 Practical contextual bandits with regression oracles. In *International Conference on*  
385 *Machine Learning*, pages 1539–1548. PMLR, 2018.
- 386 [FGMZ20] Dylan J Foster, Claudio Gentile, Mehryar Mohri, and Julian Zimmert. Adapting to  
387 misspecification in contextual bandits. *Advances in Neural Information Processing*  
388 *Systems*, 33:11478–11489, 2020.
- 389 [FK21] Dylan J Foster and Akshay Krishnamurthy. Efficient first-order contextual bandits:  
390 Prediction, allocation, and triangular discrimination. *arXiv preprint arXiv:2107.02237*,  
391 2021.
- 392 [FR20] Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual  
393 bandits with regression oracles. In *International Conference on Machine Learning*,  
394 pages 3199–3210. PMLR, 2020.
- 395 [FRSLX21] Dylan Foster, Alexander Rakhlin, David Simchi-Levi, and Yunzong Xu. Instance-  
396 dependent complexity of contextual bandits and reinforcement learning: A  
397 disagreement-based perspective. In *Conference on Learning Theory*, pages 2059–  
398 2059. PMLR, 2021.
- 399 [HK16] Elad Hazan and Tomer Koren. The computational power of optimization in online  
400 learning. In *Proceedings of the forty-eighth annual ACM symposium on Theory of*  
401 *Computing*, pages 128–141, 2016.
- 402 [HYF22] Osama A Hanna, Lin F Yang, and Christina Fragouli. Contexts can be cheap:  
403 Solving stochastic contextual bandits with linear bandit algorithms. *arXiv preprint*  
404 *arXiv:2211.05632*, 2022.
- 405 [KMB09] Varun Kanade, H Brendan McMahan, and Brent Bryan. Sleeping experts and bandits  
406 with stochastic action availability and adversarial rewards. In *Artificial Intelligence*  
407 *and Statistics*, pages 272–279. PMLR, 2009.
- 408 [KNMS10] Robert Kleinberg, Alexandru Niculescu-Mizil, and Yogeshwer Sharma. Regret bounds  
409 for sleeping experts and bandits. *Machine learning*, 80(2-3):245–272, 2010.
- 410 [KS14] Varun Kanade and Thomas Steinke. Learning hurdles for sleeping experts. *ACM*  
411 *Transactions on Computation Theory (TOCT)*, 6(3):1–16, 2014.
- 412 [KT90] Leonid G Khachiyan and Michael J Todd. On the complexity of approximating the  
413 maximal inscribed ellipsoid for a polytope. Technical report, Cornell University  
414 Operations Research and Industrial Engineering, 1990.
- 415 [KZWL23] Fang Kong, Xiangcheng Zhang, Baoxiang Wang, and Shuai Li. Improved regret bounds  
416 for linear adversarial mdps via linear optimization. *arXiv preprint arXiv:2302.06834*,  
417 2023.
- 418 [LKFS22] Gene Li, Pritish Kamath, Dylan J Foster, and Nati Srebro. Understanding the eluder  
419 dimension. *Advances in Neural Information Processing Systems*, 35:23737–23750,  
420 2022.
- 421 [LLWZ20] Chung-Wei Lee, Haipeng Luo, Chen-Yu Wei, and Mengxiao Zhang. Bias no more:  
422 high-probability data-dependent regret bounds for adversarial bandits and mdps. *arXiv*  
423 *preprint arXiv:2006.08040*, 2020.
- 424 [LS02] Tzon-Tzer Lu and Sheng-Hua Shiou. Inverses of  $2 \times 2$  block matrices. *Computers &*  
425 *Mathematics with Applications*, 43(1-2):119–129, 2002.

- 426 [LS20] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press,  
427 2020.
- 428 [LWAL18] Haipeng Luo, Chen-Yu Wei, Alekh Agarwal, and John Langford. Efficient contextual  
429 bandits in non-stationary worlds. In *Conference On Learning Theory*, pages 1739–1776.  
430 PMLR, 2018.
- 431 [LWL21] Haipeng Luo, Chen-Yu Wei, and Chung-Wei Lee. Policy optimization in adversarial  
432 mdps: Improved exploration via dilated bonuses. *Advances in Neural Information  
433 Processing Systems*, 34:22931–22942, 2021.
- 434 [LWZ19] Yingkai Li, Yining Wang, and Yuan Zhou. Nearly minimax-optimal regret for linearly  
435 parameterized bandits. In *Conference on Learning Theory*, pages 2173–2174. PMLR,  
436 2019.
- 437 [LZ07] John Langford and Tong Zhang. The epoch-greedy algorithm for contextual multi-  
438 armed bandits. *Advances in neural information processing systems*, 20(1):96–1, 2007.
- 439 [LZZZ22] Haipeng Luo, Mengxiao Zhang, Peng Zhao, and Zhi-Hua Zhou. Corraling a larger  
440 band of bandits: A case study on switching regret for linear bandits. In *Conference on  
441 Learning Theory*, pages 3635–3684. PMLR, 2022.
- 442 [Nem04] Arkadi Nemirovski. Interior point polynomial time methods in convex programming.  
443 *Lecture notes*, 42(16):3215–3224, 2004.
- 444 [NO20] Gergely Neu and Julia Olkhovskaya. Efficient and robust algorithms for adversarial  
445 linear contextual bandits. In *Conference on Learning Theory*, pages 3049–3068. PMLR,  
446 2020.
- 447 [NO21] Gergely Neu and Julia Olkhovskaya. Online learning in mdps with linear function  
448 approximation and bandit feedback. *Advances in Neural Information Processing  
449 Systems*, 34:10407–10417, 2021.
- 450 [NV14] Gergely Neu and Michal Valko. Online combinatorial optimization with stochastic  
451 decision sets and adversarial losses. *Advances in Neural Information Processing  
452 Systems*, 27, 2014.
- 453 [OMvE<sup>+</sup>23] Julia Olkhovskaya, Jack Mayo, Tim van Erven, Gergely Neu, and Chen-Yu Wei. First-  
454 and second-order bounds for adversarial linear contextual bandits. *arXiv preprint  
455 arXiv:2305.00832*, 2023.
- 456 [RS16] Alexander Rakhlin and Karthik Sridharan. Bistro: An efficient relaxation-based  
457 method for contextual bandits. In *International Conference on Machine Learning*,  
458 pages 1977–1985. PMLR, 2016.
- 459 [RVR13] Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of  
460 optimistic exploration. *Advances in Neural Information Processing Systems*, 26, 2013.
- 461 [SGV20] Aadirupa Saha, Pierre Gaillard, and Michal Valko. Improved sleeping bandits with  
462 stochastic action sets and adversarial rewards. In *International Conference on Machine  
463 Learning*, pages 8357–8366. PMLR, 2020.
- 464 [SKM23] Uri Sherman, Tomer Koren, and Yishay Mansour. Improved regret for efficient  
465 online reinforcement learning with linear function approximation. *arXiv preprint  
466 arXiv:2301.13087*, 2023.
- 467 [SKS16] Vasilis Syrgkanis, Akshay Krishnamurthy, and Robert Schapire. Efficient algorithms  
468 for adversarial contextual learning. In *International Conference on Machine Learning*,  
469 pages 2159–2168. PMLR, 2016.
- 470 [SLKS16] Vasilis Syrgkanis, Haipeng Luo, Akshay Krishnamurthy, and Robert E Schapire.  
471 Improved regret bounds for oracle-based adversarial contextual bandits. *Advances in  
472 Neural Information Processing Systems*, 29, 2016.

- 473 [SLX22] David Simchi-Levi and Yunzong Xu. Bypassing the monster: A faster and simpler  
474 optimal algorithm for contextual bandits under realizability. *Mathematics of Operations*  
475 *Research*, 47(3):1904–1931, 2022.
- 476 [TM17] Ambuj Tewari and Susan A Murphy. From ads to interventions: Contextual bandits in  
477 mobile health. *Mobile Health: Sensors, Analytic Methods, and Applications*, pages  
478 495–517, 2017.
- 479 [WDZ22] Chen-Yu Wei, Christoph Dann, and Julian Zimmert. A model selection approach for  
480 corruption robust reinforcement learning. In *International Conference on Algorithmic*  
481 *Learning Theory*, pages 1043–1096. PMLR, 2022.
- 482 [WL18] Chen-Yu Wei and Haipeng Luo. More adaptive algorithms for adversarial bandits. In  
483 *Conference On Learning Theory*, pages 1263–1291. PMLR, 2018.
- 484 [ZAK22] Julian Zimmert, Naman Agarwal, and Satyen Kale. Pushing the efficiency-regret  
485 pareto frontier for online learning of portfolios and quantum states. In *Conference on*  
486 *Learning Theory*, pages 182–226. PMLR, 2022.
- 487 [Zha22] Tong Zhang. Feel-good thompson sampling for contextual bandits and reinforcement  
488 learning. *SIAM Journal on Mathematics of Data Science*, 4(2):834–857, 2022.
- 489 [ZHZ<sup>+</sup>23] Heyang Zhao, Jiafan He, Dongruo Zhou, Tong Zhang, and Quanquan Gu. Variance-  
490 dependent regret bounds for linear bandits and reinforcement learning: Adaptivity and  
491 computational efficiency. *arXiv preprint arXiv:2302.10371*, 2023.
- 492 [ZL22] Julian Zimmert and Tor Lattimore. Return of the bias: Almost minimax optimal high  
493 probability bounds for adversarial linear bandits. In *Conference on Learning Theory*,  
494 pages 3285–3312. PMLR, 2022.

# 495 Appendices

496	<b>A Summary of Notation</b>	<b>15</b>
497	<b>B Auxiliary Lemmas</b>	<b>15</b>
498	<b>C Concentration Inequalities</b>	<b>16</b>
499	C.1 General Concentration Inequalities . . . . .	16
500	C.2 Concentration Inequalities under a Fixed Policy $p$ . . . . .	17
501	C.3 Union Bound over Policies . . . . .	21
502	<b>D Regret Analysis</b>	<b>26</b>
503	D.1 Bounding the Bias term . . . . .	26
504	D.2 Bounding the Bonus term . . . . .	27
505	D.3 Bounding the Penalty term . . . . .	30
506	D.4 Bounding the Stability-1 term . . . . .	30
507	D.5 Bounding the Stability-2 term . . . . .	32
508	D.6 Bounding the Error term . . . . .	33
509	D.7 Finishing up . . . . .	33
510	<b>E Handling Misspecification</b>	<b>34</b>
511	E.1 Known misspecification . . . . .	34
512	E.2 Unknown misspecification . . . . .	36
513	<b>F Analysis for Linear EXP4</b>	<b>42</b>
514	<b>G Comparison with [DLWZ23, SKM23]</b>	<b>43</b>
515	G.1 Regret Analysis Sketch . . . . .	44

516 **A Summary of Notation**

517 We summarize the notations that have been defined in [Algorithm 1](#) and [Definition 1](#).

$$\begin{aligned}
\beta_t &= \Theta\left(\frac{(d+1)^3 \log(T/\delta)}{t-1}\right) \\
\hat{x}_t &= \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbb{E}_{a \sim p_t^{A_\tau}} [a] \\
\hat{H}_t &= \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbb{E}_{a \sim p_t^{A_\tau}} [(a - \hat{x}_t)(a - \hat{x}_t)^\top] \\
\hat{\mathbf{H}}_t &= \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbb{E}_{a \sim p_t^{A_\tau}} \begin{bmatrix} aa^\top & a \\ a^\top & 1 \end{bmatrix} = \begin{bmatrix} \hat{H}_t + \hat{x}_t \hat{x}_t^\top & \hat{x}_t \\ \hat{x}_t^\top & 1 \end{bmatrix} \\
\hat{\Sigma}_t &= \hat{H}_t + \beta_t I \\
\hat{\mathbf{\Sigma}}_t &= \hat{\mathbf{H}}_t + \beta_t \mathbf{I} = \begin{bmatrix} \hat{\Sigma}_t + \hat{x}_t \hat{x}_t^\top & \hat{x}_t \\ \hat{x}_t^\top & 1 + \beta_t \end{bmatrix} \\
x_t &= \mathbb{E}_{\mathcal{A} \sim \mathcal{D}} \mathbb{E}_{a \sim p_t^A} [a] \\
H_t &= \mathbb{E}_{\mathcal{A} \sim \mathcal{D}} \mathbb{E}_{a \sim p_t^A} [(a - \hat{x}_t)(a - \hat{x}_t)^\top] \\
\mathbf{H}_t &= \mathbb{E}_{\mathcal{A} \sim \mathcal{D}} \mathbb{E}_{a \sim p_t^A} \begin{bmatrix} aa^\top & a \\ a^\top & 1 \end{bmatrix}
\end{aligned}$$

518 **B Auxiliary Lemmas**

519 **Lemma 5** (FTRL regret bound, Lemma 18 of [\[DWZ23a\]](#)). *Let  $\Omega \subset \mathbb{R}^d$  be a convex set,  $g_1, \dots, g_T \in$*   
520  *$\mathbb{R}^d$ , and  $\eta_1, \dots, \eta_T > 0$ . Then the FTRL update*

$$w_t = \operatorname{argmin}_{w \in \Omega} \left\{ \left\langle w, \sum_{\tau=1}^{t-1} g_\tau \right\rangle + \frac{1}{\eta_t} \psi(w) \right\}$$

521 *ensures for any  $u \in \Omega$  and  $\eta_0 > 0$ ,*

$$\begin{aligned}
& \sum_{t=1}^T \langle w_t - u, g_t \rangle \\
& \leq \underbrace{\frac{\psi(u) - \min_{w \in \Omega} \psi(w)}{\eta_0} + \sum_{t=1}^T (\psi(u) - \psi(w_t)) \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right)}_{\text{Penalty}} + \underbrace{\sum_{t=1}^T \left( \max_{w \in \Omega} \langle w_t - w, g_t \rangle - \frac{D_\psi(w, w_t)}{\eta_t} \right)}_{\text{Stability}}.
\end{aligned}$$

522 *When  $\eta_0, \eta_1, \dots, \eta_T$  is non-increasing, the penalty term can further be upper bounded by*

$$\text{Penalty} \leq \frac{\psi(u) - \min_{w \in \Omega} \psi(w)}{\eta_T}.$$

523 **Lemma 6** (Bernstein's inequality). *Let  $X_1, \dots, X_n$  be iid random variables; let  $\mathbb{E}[X]$  be the*  
524 *expectation and  $\text{Var}(X)$  be the variance of these random variables. If for any  $i$ ,  $|X_i - \mathbb{E}[X_i]| \leq R$ ,*  
525 *then with probability of at least  $1 - \delta$ ,*

$$\left| \frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X] \right| \leq \sqrt{\frac{4 \text{Var}(X) \log \frac{2}{\delta}}{n}} + \frac{4R \log \frac{2}{\delta}}{3n}.$$

526 **Lemma 7** (Hoeffding's inequality). Let  $X_1, \dots, X_n$  be iid random variables; let  $a \leq X_i \leq b$  and  
 527 let  $\mathbb{E}[X]$  be the expectatio. Then with probability of at least  $1 - \delta$ ,

$$\left| \frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X] \right| \leq (b - a) \sqrt{\frac{1}{2n} \log\left(\frac{2}{\delta}\right)}$$

528 Given  $F(X) = -\log \det(X)$ ,  $D^2F(X) = X^{-1} \otimes X^{-1}$  where  $\otimes$  is the Kronecker prod-

529 uct. For any matrix  $A = \begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix}$ , let  $\text{vec}(A) = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}$  which vectorizes matrix

530  $A$  to a column vector by stacking the columns  $A$ . The second order directional derivative  
 531 for  $F$  is  $D^2F(X)[A, A] = \text{vec}(A)^T (X^{-1} \otimes X^{-1}) \text{vec}(A) = \text{Tr}(A^T X^{-1} A X^{-1})$ . We define  
 532  $\|A\|_{\nabla^2 F(X)} = \sqrt{\text{Tr}(A^T X^{-1} A X^{-1})}$  and  $\|A\|_{\nabla^{-2} F(X)} = \sqrt{\text{Tr}(A^T X A X)}$ . It is a pseudo-norm,  
 533 and more discussion can be found in Appendix D of [ZAK22]. In the following analysis, we will  
 534 only use one property of this pseudo-norm which is similar to the Holder inequality.

535 **Lemma 8.** For any two symmetric matrices  $A, B$  and positive definite matrix  $X$ ,

$$\langle A, B \rangle \leq \|A\|_{\nabla^2 F(X)} \|B\|_{\nabla^{-2} F(X)}$$

536 *Proof.* Since  $(X \otimes X)^{-1} = X^{-1} \otimes X^{-1}$ , from Holder inequality, we have

$$\langle A, B \rangle = \langle \text{vec}(A), \text{vec}(B) \rangle \leq \|\text{vec}(A)\|_{X^{-1} \otimes X^{-1}} \|\text{vec}(B)\|_{(X^{-1} \otimes X^{-1})^{-1}} = \|A\|_{\nabla^2 F(X)} \|B\|_{\nabla^{-2} F(X)}$$

537

□

## 538 C Concentration Inequalities

539 The goal of this section is to show [Lemma 16](#) and [Lemma 17](#), which are key to bound the bias  
 540 term. We first introduce a useful lemma from [DLWZ23], which will be used later to prove our  
 541 concentration bounds.

### 542 C.1 General Concentration Inequalities

543 **Lemma 9** (Lemma A.4 in [DLWZ23]). Let  $H_1, H_2, \dots, H_n$  be i.i.d. PSD matrices such that  
 544  $\mathbb{E}[H_i] = H$ ,  $H_i \preceq I$  almost surely and  $H \succeq \frac{1}{dn} \log \frac{d}{\delta} I$ . Then with probability  $1 - \delta$ ,

$$\frac{1}{n} \sum_{i=1}^n H_i - H \succeq -\sqrt{\frac{d}{n} \log \frac{d}{\delta}} H^{1/2}$$

545 **Corollary 10.** Let  $H_1, H_2, \dots, H_n$  be i.i.d. PSD matrices such that  $\mathbb{E}[H_i] = H$  and  $H_i \preceq cI$  almost  
 546 surely for some positive constant  $c$ . Let  $\hat{H} = \frac{1}{n} \sum_{i=1}^n H_i$ , then with probability  $1 - \delta$ ,

$$\hat{H} + \frac{3c}{2} \cdot \frac{d}{n} \log \left( \frac{d}{\delta} \right) I \succeq \frac{1}{2} H \tag{9}$$

547 *Proof.* A simple corollary of [Lemma 9](#) under the condition of [Lemma 9](#) is that

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n H_i - H \succeq -\sqrt{\frac{d}{n} \log \frac{d}{\delta}} H^{1/2} \succeq -\frac{1}{2} H - \frac{d}{2n} \log \left( \frac{d}{\delta} \right) I \\ \Rightarrow & \frac{1}{n} \sum_{i=1}^n H_i + \frac{d}{2n} \log \left( \frac{d}{\delta} \right) I \succeq \frac{1}{2} H, \end{aligned} \tag{10}$$

548 where we use that  $H^{\frac{1}{2}} \preceq \frac{k}{2} H + \frac{1}{2k} I$  for any  $k > 0$ .

549 Now consider the condition of this corollary. We first consider the case where  $\frac{d}{n} \log\left(\frac{d}{\delta}\right) \leq 1$ . In this  
 550 case, we apply Eq. (10) with  $H'_i = \frac{1}{2c}H_i + \frac{d}{2n} \log\left(\frac{d}{\delta}\right)I$ , which satisfies the condition for Eq. (10) to  
 551 hold. This gives

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{2c}H_i + \frac{d}{2n} \log\left(\frac{d}{\delta}\right)I \right) + \frac{d}{2n} \log\left(\frac{d}{\delta}\right)I \succeq \frac{1}{2} \left( \frac{1}{2c}H + \frac{d}{2n} \log\left(\frac{d}{\delta}\right)I \right) \\ \Rightarrow & \hat{H} + \frac{3c}{2} \cdot \frac{d}{n} \log\left(\frac{d}{\delta}\right)I \succeq \frac{1}{2}H \end{aligned}$$

552 with probability at least  $1 - \delta$ . When  $\frac{d}{n} \log\left(\frac{d}{\delta}\right) > 1$ , Eq. (9) is trivial because  $\frac{1}{2}H \preceq \frac{c}{2}I \preceq$   
 553  $\frac{c}{2} \cdot \frac{d}{n} \log\left(\frac{d}{\delta}\right)I$ .

554 □

## 555 C.2 Concentration Inequalities under a Fixed Policy $p$

556 In this subsection, we establish concentration bounds for a *fixed* policy  $p$  (with  $p^A \in \Delta(\mathcal{A})$  denoting  
 557 the action distribution it uses over  $\mathcal{A}$ ) over i.i.d. contexts. The results in this subsection are preparation  
 558 for Appendix C.3 where we take union bounds over policies.

559 The setting and notation to be used in this subsection are defined in Definition 11.

560 **Definition 11.** Let  $\{\mathcal{A}_1, \dots, \mathcal{A}_n\}$  be i.i.d. context samples drawn from  $D$ . Let  $\hat{D}$  be the uniform  
 561 distribution over  $\{\mathcal{A}_1, \dots, \mathcal{A}_n\}$ .

562 Over this set of context samples, define for any policy  $p$ ,

$$\begin{aligned} x(p) &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^A} [a], \\ \hat{x}(p) &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^A} [a], \\ H(p) &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^A} [(a - \hat{x}(p))(a - \hat{x}(p))^\top], \\ \hat{H}(p) &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^A} [(a - \hat{x}(p))(a - \hat{x}(p))^\top], \\ \mathbf{H}(p) &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^A} [\mathbf{a}\mathbf{a}^\top], \\ \hat{\mathbf{H}}(p) &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^A} [\mathbf{a}\mathbf{a}^\top], \\ \hat{\Sigma}(p) &= \hat{H}(p) + \beta I, \\ \hat{\mathbf{\Sigma}}(p) &= \hat{\mathbf{H}}(p) + \beta I, \end{aligned}$$

563 where  $\beta = \frac{5d \log(6d/\delta)}{n}$ .

564 **Lemma 12.** Under the setting of Definition 11, for any fixed  $p$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} \hat{H}(p) + \frac{4d \log(6d/\delta)}{n} I &\succeq \frac{1}{2} H(p), \\ \hat{\mathbf{H}}(p) + \frac{3d \log(d/\delta)}{n} \mathbf{I} &\succeq \frac{1}{2} \mathbf{H}(p). \end{aligned}$$

565 *Proof.* In this proof, we use  $\hat{x}, x, \hat{H}, H, \hat{\mathbf{H}}, \mathbf{H}$  to denote  $\hat{x}(p), x(p), \hat{H}(p), H(p), \hat{\mathbf{H}}(p), \mathbf{H}(p)$  since  
 566  $p$  is fixed throughout the proof.

567 Since  $\|a\| \leq 1$ ,  $\mathbf{H} \preceq 2I$  and  $\hat{\mathbf{H}} \preceq 2I$ . Thus, we can directly apply Corollary 10 with  $c = 2$  to get  
 568 with probability  $1 - \frac{\delta}{3}$

$$\hat{\mathbf{H}} + \frac{3d \log(3d/\delta)}{n} \mathbf{I} \succeq \frac{1}{2} \mathbf{H}.$$

569 To prove the first inequality, we first decompose  $H$  and  $\hat{H}$

$$\begin{aligned} H &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - \hat{x})(a - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x + x - \hat{x})(a - x + x - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] + (x - \hat{x})(x - \hat{x})^\top \quad (\text{because } \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}}(a - x) = 0) \end{aligned} \quad (11)$$

$$\begin{aligned} \hat{H} &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - \hat{x})(a - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x + x - \hat{x})(a - x + x - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] - (x - \hat{x})(x - \hat{x})^\top \\ &\quad (\text{because } \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}}(a - x) = \hat{x} - x) \end{aligned} \quad (12)$$

570 From Hoeffding inequality (Lemma 7) and union bound, with probability  $1 - \frac{\delta}{3}$ , for all  $k \in [d]$ , we  
571 have

$$|e_k^\top x - e_k^\top \hat{x}| \leq \sqrt{\frac{1}{2n} \log \left( \frac{6d}{\delta} \right)},$$

572 which implies that  $e_k^\top (x - \hat{x})(x - \hat{x})^\top e_k \leq \frac{1}{2n} \log \left( \frac{6d}{\delta} \right)$  for all  $k$ , and thus

$$(x - \hat{x})(x - \hat{x})^\top \preceq \frac{1}{2n} \log \left( \frac{6d}{\delta} \right) I. \quad (13)$$

573 By directly applying Corollary 10 with  $c = 2$ , we get with probability at least  $1 - \frac{\delta}{3}$ ,

$$\mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] + \frac{3d \log(3d/\delta)}{n} I \succeq \frac{1}{2} \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top]$$

574 Further using Eq. (11), Eq. (12) and Eq. (13), we get with probability at least  $1 - \frac{2\delta}{3}$ ,

$$\hat{H} + \frac{4d \log(6d/\delta)}{n} I \succeq \frac{1}{2} H$$

575 Taking union bound for both inequality finishes the proof.  $\square$

576 **Lemma 13.** Under the setting of Definition 11, for any fixed policy  $p$ , with probability at least  
577  $1 - \mathcal{O}(\delta)$ ,

$$\|x(p) - \hat{x}(p)\|_{\hat{\Sigma}(p)^{-1}}^2 \leq \mathcal{O} \left( \frac{d \log(d/\delta)}{n} \right)$$

578 *Proof.* In this proof, we use  $\hat{x}, x, \hat{H}, H, \hat{\mathbf{H}}, \mathbf{H}, \hat{\Sigma}, \hat{\mathbf{\Sigma}}$  to denote  $\hat{x}(p), x(p), \hat{H}(p), H(p), \hat{\mathbf{H}}(p), \mathbf{H}(p),$   
579  $\hat{\Sigma}(p), \hat{\mathbf{\Sigma}}(p)$  since  $p$  is fixed throughout the proof.

580 We first rewrite  $H$ .

$$\begin{aligned} H &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - \hat{x})(a - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x + x - \hat{x})(a - x + x - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] + (x - \hat{x})(x - \hat{x})^\top \quad (\text{because } \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}}(a - x) = 0) \end{aligned} \quad (14)$$

581 To simplify analysis, we perform diagonalization. Suppose that  $\mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top]$   
582 admits the following eigen-decomposition:

$$\mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] = V \Lambda V^\top$$

583 where  $V$  is an orthogonal matrix and  $\Lambda$  is a diagonal matrix. By Lemma 12 and the definition of  $\beta$  in  
584 Definition 11, we have with probability  $1 - \delta$ ,

$$\hat{\Sigma} \succeq \frac{1}{2} H + \rho I \succeq \frac{1}{2} V \Lambda V^\top + \rho I$$

585 with some  $\rho = \Theta\left(\frac{d \log(d/\delta)}{n}\right)$ , where the second inequality is by Eq. (14). Thus,

$$\begin{aligned} \|x - \hat{x}\|_{\hat{\Sigma}^{-1}}^2 &= (x - \hat{x})^\top \hat{\Sigma}^{-1} (x - \hat{x}) \\ &\leq (x - \hat{x})^\top \left( \frac{1}{2} V \Lambda V^\top + \rho I \right)^{-1} (x - \hat{x}) \\ &= (\hat{x} - x)^\top V \left( \frac{1}{2} \Lambda + \rho I \right)^{-1} V^\top (\hat{x} - x). \end{aligned}$$

586 Define

$$\Delta_k = e_k^\top V^\top (\hat{x} - x) = \frac{1}{n} \sum_{i=1}^n \underbrace{e_k^\top V^\top \mathbb{E}_{a \sim p^{\mathcal{A}_i}} [a]}_{\text{Define as } Z_k^{(i)}} - \underbrace{e_k^\top V^\top \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [a]}_{\text{Define as } Z_k}$$

587 Since  $\mathbb{E}_{\mathcal{A}_i \sim D} [Z_k^{(i)}] = Z_k$ , by Bernstein's inequality, with probability at least  $1 - \delta$ , we have

$$|\Delta_k| \leq \mathcal{O} \left( \sqrt{\frac{\text{Var}(Z_k^{(i)}) \log(d/\delta)}{n}} + \frac{\log(d/\delta)}{n} \right) \quad (15)$$

588 for all  $k$ , where

$$\text{Var}(Z_k^{(i)}) = \mathbb{E}_{\mathcal{A} \sim D} \left[ \left( e_k^\top V^\top \mathbb{E}_{a \sim p^{\mathcal{A}}} [a] - e_k^\top V^\top x \right)^2 \right].$$

589 On the other hand,

$$\begin{aligned} \Lambda_{kk} &= e_k^\top \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [V^\top (a - x)(a - x)^\top V] e_k \\ &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} \left[ \left( e_k^\top V^\top a - e_k^\top V^\top x \right)^2 \right]. \end{aligned}$$

590 From Jensen's inequality,

$$\Lambda_{kk} = \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} \left[ \left( e_k^\top V^\top a - e_k^\top V^\top x \right)^2 \right] \geq \mathbb{E}_{\mathcal{A} \sim D} \left[ \left( e_k^\top V^\top \mathbb{E}_{a \sim p^{\mathcal{A}}} [a] - e_k^\top V^\top x \right)^2 \right] = \text{Var}(Z_k^{(i)})$$

591 Thus,

$$\begin{aligned} \|x - \hat{x}\|_{\hat{\Sigma}^{-1}}^2 &\leq (\hat{x} - x)^\top V \left( \frac{1}{2} \Lambda + \rho I \right)^{-1} V^\top (\hat{x} - x) \\ &= \sum_{k=1}^d \frac{(\Delta_k)^2}{\frac{1}{2} \Lambda_{kk} + \rho} \\ &\leq \mathcal{O} \left( \frac{\log(d/\delta)}{n} \sum_{k=1}^d \frac{\text{Var}(Z_k^{(i)}) + \frac{\log(d/\delta)}{n}}{\Lambda_{kk} + \rho} \right) \quad (\text{by Eq. (15)}) \\ &\leq \mathcal{O} \left( \frac{d \log(d/\delta)}{n} \right). \quad (\Lambda_{kk} \geq \text{Var}(Z_k^{(i)}) \text{ and } \rho = \Theta\left(\frac{d \log(d/\delta)}{n}\right)) \end{aligned}$$

592  $\square$

593 **Lemma 14.** Under the setting of Definition 11, for any fixed policy  $p$ , with probability at least  
594  $1 - \mathcal{O}(\delta)$ ,

$$\|(\hat{\Sigma}(p) - H(p))y\|_{\hat{\Sigma}(p)^{-1}}^2 \leq \mathcal{O} \left( \frac{d \log(d/\delta)}{n} \right)$$

595 for any  $y \in \mathbb{B}_2^d$ .

596 *Proof.* In this proof, we use  $\hat{x}, x, \hat{H}, H, \hat{\mathbf{H}}, \mathbf{H}, \hat{\Sigma}, \hat{\Sigma}$  to denote  $\hat{x}(p), x(p), \hat{H}(p), H(p), \hat{\mathbf{H}}(p), \mathbf{H}(p),$   
597  $\hat{\Sigma}(p), \hat{\Sigma}(p)$  since  $p$  is fixed throughout the proof.

598 First, we re-write  $H$  and  $\hat{H}$ :

$$\begin{aligned} H &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - \hat{x})(a - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x + x - \hat{x})(a - x + x - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] + (x - \hat{x})(x - \hat{x})^\top \quad (\text{because } \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}}(a - x) = 0) \end{aligned} \quad (16)$$

$$\begin{aligned} \hat{H} &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - \hat{x})(a - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x + x - \hat{x})(a - x + x - \hat{x})^\top] \\ &= \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] - (x - \hat{x})(x - \hat{x})^\top \\ &\quad (\text{because } \mathbb{E}_{\mathcal{A} \sim \hat{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}}(a - x) = \hat{x} - x) \end{aligned}$$

599 Then, by definition (in [Definition 11](#)) and the calculation above,

$$\begin{aligned} \hat{\Sigma} - H &= \hat{H} - H + \beta I \\ &= \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{E}_{a \sim p^{\mathcal{A}_i}} [(a - x)(a - x)^\top] - \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top]}_{\text{define this as } \Gamma} - 2(x - \hat{x})(x - \hat{x})^\top + \beta I. \end{aligned}$$

600 Using  $\|a + b + c\|^2 \leq 3\|a\|^2 + 3\|b\|^2 + 3\|c\|^2$ , we have

$$\begin{aligned} \|(\hat{\Sigma} - H)y\|_{\hat{\Sigma}^{-1}}^2 &\leq 3\|\Gamma y\|_{\hat{\Sigma}^{-1}}^2 + 12\|(x - \hat{x})(x - \hat{x})^\top y\|_{\hat{\Sigma}^{-1}}^2 + \beta^2\|y\|_{\hat{\Sigma}^{-1}}^2 \\ &\leq 3\|\Gamma y\|_{\hat{\Sigma}^{-1}}^2 + 12\|x - \hat{x}\|_{\hat{\Sigma}^{-1}}^2 + \mathcal{O}(\beta). \end{aligned} \quad (17)$$

601 The second and third term are bounded by  $\mathcal{O}\left(\frac{d \log(d/\delta)}{n}\right)$  using [Lemma 13](#) and the definition of  $\beta$ ,  
602 with probability at least  $1 - \mathcal{O}(\delta)$ . Below, we further deal with the first term. To simplify analysis,  
603 we perform diagonalization. Suppose that  $\mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top]$  admits the following  
604 eigen-decomposition:

$$\mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - x)(a - x)^\top] = V \Lambda V^\top$$

605 where  $V$  is an orthogonal matrix and  $\Lambda$  is a diagonal matrix. Then

$$\|\Gamma y\|_{\hat{\Sigma}^{-1}}^2 = y^\top \Gamma \hat{\Sigma}^{-1} \Gamma y = (V^\top y)^\top (V^\top \Gamma V) (V^\top \hat{\Sigma} V)^{-1} (V^\top \Gamma V) (V^\top y). \quad (18)$$

606 Below, we further deal with the  $V^\top \Gamma V$  and  $V^\top \Lambda V$  terms in [Eq. \(18\)](#). By [Lemma 12](#), with probability  
607 at least  $1 - \delta$ ,

$$\hat{\Sigma} \succeq \frac{1}{2} H + \rho I \succeq \frac{1}{2} V \Lambda V^\top + \rho I,$$

608 for some  $\rho = \Theta\left(\frac{d \log(d/\delta)}{n}\right)$ , where we use [Eq. \(16\)](#) in the second inequality. Therefore,

$$V^\top \hat{\Sigma} V \succeq \frac{1}{2} \Lambda + \rho I. \quad (19)$$

609 Next, denote  $\Delta = V^\top \Gamma V$ . By definition, it can be written as the following:

$$\Delta = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{E}_{a \sim p^{\mathcal{A}_i}} [V^\top (a - x)(a - x)^\top V]}_{\text{defining this as } \Lambda^{(i)}} - \underbrace{\mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [V^\top (a - x)(a - x)^\top V]}_{=\Lambda}$$

610 with  $\Lambda^{(i)}$  being i.i.d. samples with mean  $\mathbb{E}[\Lambda^{(i)}] = \Lambda$ . While these are  $d \times d$  matrices, we will apply  
611 concentration inequalities to individual entries.

612 Let  $\lambda_{ikh} = e_k^\top \Lambda^{(i)} e_h$  be the  $(k, h)$ -th entry of  $\Lambda^{(i)}$ . Notice that  $\mathbb{E}[\lambda_{ikh}] = e_k^\top \Lambda e_h = \Lambda_{kh}$ , the  
613  $(k, h)$ -th entry of  $\Lambda$ .

614 By Bernstein's inequality, with probability at least  $1 - \delta$ , we have

$$|\Delta_{kh}| = \left| \frac{1}{n} \sum_{i=1}^n (\lambda_{ikh} - \Lambda_{kh}) \right| \leq \mathcal{O} \left( \sqrt{\frac{\text{Var}(\lambda_{ikh}) \log(d/\delta)}{n}} + \frac{\log(d/\delta)}{n} \right). \quad (20)$$

615 With the manipulations and notations above, we continue to bound Eq. (18) by

$$\begin{aligned} \|\Gamma y\|_{\hat{\Sigma}^{-1}}^2 &= y'^{\top} \Delta (V^{\top} \hat{\Sigma} V)^{-1} \Delta y' && \text{(let } y' = V^{\top} y) \\ &\leq 2y'^{\top} \Delta (\Lambda + \rho I)^{-1} \Delta y' && \text{(by Eq. (19))} \\ &\leq 2 \text{Tr} \left( \Delta (\Lambda + \rho I)^{-1} \Delta \right) \end{aligned}$$

616 By direct expansion and the fact that  $\Lambda$  is diagonal,

$$\begin{aligned} \text{Tr} \left( \Delta (\Lambda + \rho I)^{-1} \Delta \right) &= \sum_{k=1}^d \left( \Delta (\Lambda + \rho I)^{-1} \Delta \right)_{kk} \\ &= \sum_{k=1}^d \sum_{h=1}^d \frac{\Delta_{kh} \Delta_{hk}}{\Lambda_{hh} + \rho} \\ &\leq \mathcal{O} \left( \sum_{k=1}^d \sum_{h=1}^d \frac{1}{\Lambda_{hh} + \rho} \left( \frac{\text{Var}(\lambda_{ikh}) \log(d/\delta)}{n} + \frac{\log^2(d/\delta)}{n^2} \right) \right) \\ &\hspace{15em} \text{(by Eq. (20))} \\ &\leq \mathcal{O} \left( \sum_{k=1}^d \sum_{h=1}^d \frac{1}{\Lambda_{hh} + \rho} \frac{\mathbb{E}[\lambda_{ikh}^2] \log(d/\delta)}{n} + \frac{d^2 \log^2(d/\delta)}{\rho n^2} \right) \quad (21) \end{aligned}$$

617 By definition,

$$\lambda_{ikh} = \mathbb{E}_{a \sim p^{\mathcal{A}_i}} [\mathbf{e}_k V^{\top} (a - x)(a - x)^{\top} V \mathbf{e}_h]$$

618 and thus

$$\begin{aligned} \sum_{k=1}^d \lambda_{ikh}^2 &\leq \mathbb{E}_{a \sim p^{\mathcal{A}_i}} \left[ \sum_{k=1}^d (\mathbf{e}_k V^{\top} (a - x)(a - x)^{\top} V \mathbf{e}_h)^2 \right] \\ &= \mathbb{E}_{a \sim p^{\mathcal{A}_i}} \left[ \sum_{k=1}^d \mathbf{e}_h^{\top} V^{\top} (a - x)(a - x)^{\top} V \mathbf{e}_k \mathbf{e}_k^{\top} V^{\top} (a - x)(a - x)^{\top} V \mathbf{e}_h \right] \\ &= \mathbb{E}_{a \sim p^{\mathcal{A}_i}} [\mathbf{e}_h^{\top} V^{\top} (a - x)(a - x)^{\top} (a - x)(a - x)^{\top} V \mathbf{e}_h] \\ &\leq \mathbb{E}_{a \sim p^{\mathcal{A}_i}} [\mathbf{e}_h^{\top} V^{\top} (a - x)(a - x)^{\top} V \mathbf{e}_h] \\ &= \lambda_{ihh} \end{aligned}$$

619 and  $\sum_{k=1}^d \mathbb{E}[\lambda_{ikh}^2] \leq \mathbb{E}[\lambda_{ihh}] = \Lambda_{hh}$ . Continuing from Eq. (21) and using that  $\rho = \Theta\left(\frac{d \log(d/\delta)}{n}\right)$ ,

$$\text{Tr} \left( \Delta (\Lambda + \rho I)^{-1} \Delta \right) \leq \mathcal{O} \left( \sum_{h=1}^d \frac{\Lambda_{hh} \log(d/\delta)}{(\Lambda_{hh} + \rho)n} + \frac{d^2 \log^2(d/\delta)}{n^2} \right) \leq \mathcal{O} \left( \frac{d \log(d/\delta)}{n} \right).$$

620 This gives a bound on  $\|\Gamma y\|_{\hat{\Sigma}^{-1}}^2$  and finishes the proof after combining Eq. (17).

621 □

### 622 C.3 Union Bound over Policies

623 In Lemma 12, Lemma 13, and Lemma 14, we have obtained the desired concentration inequalities  
624 *under a fixed policy*  $p$ . In this subsection, we proceed to take union bound over *all policies* that are  
625 possibly used by Algorithm 1.

626 The set of policies that could be generated by Algorithm 1 is the following:

$$\mathbf{P} = \left\{ p : \widehat{\text{Cov}}(p^{\mathcal{A}}) = \underset{\mathbf{H} \in \mathcal{H}^{\mathcal{A}}}{\text{argmin}} \{ \langle \mathbf{H}, \mathbf{Z} \rangle + F(\mathbf{H}) \}, \text{ for } \mathbf{Z} \in \mathcal{Z} \right\}$$

627 where  $\mathcal{Z} = [-T^2, T^2]^{(d+1) \times (d+1)} \cap \mathbb{S}$  with  $\mathbb{S}$  denoting the set of symmetric matrices. To see this,  
 628 notice that Algorithm 1 at round  $t$  corresponds to the policy defined above with  $\mathbf{Z} = \eta_t \sum_{s=1}^{t-1} (\hat{\gamma}_s -$   
 629  $\alpha_s \hat{\Sigma}_s^{-1})$ .

630 Our goal is to construct a  $\epsilon$ -cover  $\mathbf{P}'$  so that every policy  $p \in \mathbf{P}$  can find a policy  $p' \in \mathbf{P}'$  making  
 631  $-\epsilon I \preceq \widehat{\text{Cov}}(p^{\mathcal{A}}) - \widehat{\text{Cov}}(p'^{\mathcal{A}}) \preceq \epsilon I$  on every action set  $\mathcal{A}$ . The size of such a cover is bounded in the  
 632 Proposition below.

633 **Proposition 1.** *There exists an  $\epsilon$ -cover  $\mathbf{P}'$  of  $\mathbf{P}$  with size  $\log |\mathbf{P}'| = \mathcal{O}(d^2 \log \frac{d}{\epsilon})$  such that for any  
 634  $p \in \mathbf{P}$ , there exists an  $p' \in \mathbf{P}'$  satisfying*

$$\left\| \widehat{\text{Cov}}(p^{\mathcal{A}}) - \widehat{\text{Cov}}(p'^{\mathcal{A}}) \right\|_F \leq \epsilon$$

635 for all  $\mathcal{A}$ .

636 *Proof.* It is straightforward to construct an  $\frac{\epsilon}{4}$ -cover  $\mathcal{C}$  for  $\mathcal{Z} = [-T^2, T^2]^{(d+1) \times (d+1)} \cap \mathbb{S}$  in Frobenius  
 637 norm with size  $|\mathcal{C}| = \left(\frac{2A(d+1)^2}{\epsilon}\right)^{(d+1)^2}$  (Exercise 27.6 of [LS20]). Now define  $\mathbf{P}'$  as

$$\mathbf{P}' = \left\{ p : \widehat{\text{Cov}}(p^{\mathcal{A}}) = \underset{\mathbf{H} \in \mathcal{H}^{\mathcal{A}}}{\text{argmin}} \{ \langle \mathbf{H}, \mathbf{Z} \rangle + F(\mathbf{H}) \}, \text{ for } \mathbf{Z} \in \mathcal{C} \right\} \quad (22)$$

638 Below, we show that this is a  $\epsilon$ -cover for  $\mathbf{P}$ .

639 Consider two policies  $p_1$  and  $p_2$  defined as the following:

$$\begin{aligned} \widehat{\text{Cov}}(p_1^{\mathcal{A}}) &= \underset{\mathbf{H} \in \mathcal{H}^{\mathcal{A}}}{\text{argmin}} \{ \langle \mathbf{H}, \mathbf{Z}_1 \rangle + F(\mathbf{H}) \} \\ \widehat{\text{Cov}}(p_2^{\mathcal{A}}) &= \underset{\mathbf{H} \in \mathcal{H}^{\mathcal{A}}}{\text{argmin}} \{ \langle \mathbf{H}, \mathbf{Z}_2 \rangle + F(\mathbf{H}) \} \end{aligned}$$

640 with  $\|\mathbf{Z}_1 - \mathbf{Z}_2\|_F \leq \frac{\epsilon}{4}$ . Consider an arbitrary  $\mathcal{A}$  and define  $\mathbf{H}_1 = \widehat{\text{Cov}}(p_1^{\mathcal{A}})$ ,  $\mathbf{H}_2 = \widehat{\text{Cov}}(p_2^{\mathcal{A}})$ . Below  
 641 we show  $\|\mathbf{H}_1 - \mathbf{H}_2\|_F \leq \epsilon$ .

642 Since  $F(\mathbf{H})$  is convex for  $\mathbf{H}$ , from the first-order optimality condition for convex function, we have

$$\begin{aligned} \langle \mathbf{H}_1, \mathbf{Z}_1 \rangle + F(\mathbf{H}_1) &\leq \langle \mathbf{H}_2, \mathbf{Z}_1 \rangle + F(\mathbf{H}_2) - D_F(\mathbf{H}_2, \mathbf{H}_1) \\ &= \langle \mathbf{H}_2, \mathbf{Z}_2 \rangle + \langle \mathbf{H}_2, \mathbf{Z}_1 - \mathbf{Z}_2 \rangle + F(\mathbf{H}_2) - D_F(\mathbf{H}_2, \mathbf{H}_1) \\ \langle \mathbf{H}_2, \mathbf{Z}_2 \rangle + F(\mathbf{H}_2) &\leq \langle \mathbf{H}_1, \mathbf{Z}_2 \rangle + F(\mathbf{H}_1) - D_F(\mathbf{H}_1, \mathbf{H}_2) \\ &= \langle \mathbf{H}_1, \mathbf{Z}_1 \rangle + \langle \mathbf{H}_1, \mathbf{Z}_2 - \mathbf{Z}_1 \rangle + F(\mathbf{H}_1) - D_F(\mathbf{H}_1, \mathbf{H}_2) \end{aligned}$$

643 Adding up these the two inequalities, we get

$$2 \min\{D_F(\mathbf{H}_1, \mathbf{H}_2), D_F(\mathbf{H}_2, \mathbf{H}_1)\} \leq D_F(\mathbf{H}_1, \mathbf{H}_2) + D_F(\mathbf{H}_2, \mathbf{H}_1) \leq \langle \mathbf{Z}_1 - \mathbf{Z}_2, \mathbf{H}_2 - \mathbf{H}_1 \rangle$$

644 Since the second order directional derivative for  $F$  is  $D^2 F(\mathbf{H})[\mathbf{X}, \mathbf{X}] = \text{Tr}(\mathbf{X} \mathbf{H}^{-1} \mathbf{X} \mathbf{H}^{-1})$  for any  
 645 symmetric matrix  $\mathbf{X}$ , from the Taylor series, there exists  $\mathbf{H}'$  that is a line segment between  $\mathbf{H}_1$  and  
 646  $\mathbf{H}_2$  such that

$$\begin{aligned} \|\mathbf{H}_1 - \mathbf{H}_2\|_{\nabla^2 F(\mathbf{H}')}^2 &= 2 \min\{D_F(\mathbf{H}_1, \mathbf{H}_2), D_F(\mathbf{H}_2, \mathbf{H}_1)\} \leq \langle \mathbf{Z}_1 - \mathbf{Z}_2, \mathbf{H}_2 - \mathbf{H}_1 \rangle \\ &\leq \|\mathbf{Z}_1 - \mathbf{Z}_2\|_{\nabla^{-2} F(\mathbf{H}')} \|\mathbf{H}_1 - \mathbf{H}_2\|_{\nabla^2 F(\mathbf{H}')} \quad (\text{Lemma 8}) \end{aligned}$$

647 Thus we have  $\|\mathbf{H}_1 - \mathbf{H}_2\|_{\nabla^2 F(\mathbf{H}')} \leq \|\mathbf{Z}_1 - \mathbf{Z}_2\|_{\nabla^{-2} F(\mathbf{H}')}$ . Since  $\|a\|_2 \leq 1$ ,  $\mathbf{H}' \preceq 2\mathbf{I}$ . The left-hand  
 648 side and right-hand side can be bounded as follows,

$$\begin{aligned} \|\mathbf{H}_1 - \mathbf{H}_2\|_{\nabla^2 F(\mathbf{H}')} &= \sqrt{\text{Tr}((\mathbf{H}_1 - \mathbf{H}_2)(\mathbf{H}')^{-1}(\mathbf{H}_1 - \mathbf{H}_2)(\mathbf{H}')^{-1})} \geq \frac{1}{2} \|\mathbf{H}_1 - \mathbf{H}_2\|_F \\ \|\mathbf{Z}_1 - \mathbf{Z}_2\|_{\nabla^{-2} F(\mathbf{H}')} &= \sqrt{\text{Tr}((\mathbf{Z}_1 - \mathbf{Z}_2)\mathbf{H}'(\mathbf{Z}_1 - \mathbf{Z}_2)\mathbf{H}')} \leq 2 \|\mathbf{Z}_1 - \mathbf{Z}_2\|_F \leq \frac{\epsilon}{2} \end{aligned}$$

649 Combining the three inequalities above, we conclude that

$$650 \quad \|\mathbf{H}_1 - \mathbf{H}_2\|_F \leq 2\|\mathbf{H}_1 - \mathbf{H}_2\|_{\nabla^2 F(\mathbf{H}')} \leq 2\|\mathbf{Z}_1 - \mathbf{Z}_2\|_{\nabla^{-2} F(\mathbf{H}')} \leq 4\|\mathbf{Z}_1 - \mathbf{Z}_2\|_F \leq \epsilon.$$

$$651 \quad -\epsilon \mathbf{I} \preceq \mathbf{H}_1 - \mathbf{H}_2 \preceq \epsilon \mathbf{I}. \quad \square$$

652 **Proposition 2.** Suppose that  $p, p'$  are two policies such that for all action set  $\mathcal{A}$ ,

$$\left\| \widehat{\text{Cov}}(p^{\mathcal{A}}) - \widehat{\text{Cov}}(p'^{\mathcal{A}}) \right\|_F \leq \epsilon \quad (23)$$

653 Then all quantities defined in [Definition 11](#) under  $p$  and  $p'$  are close. That is,

$$\|x(p) - x(p')\| \leq \epsilon \quad (24)$$

$$\|\hat{x}(p) - \hat{x}(p')\| \leq \epsilon \quad (25)$$

$$\|H(p) - H(p')\|_F \leq 7\epsilon \quad (26)$$

$$\|\hat{H}(p) - \hat{H}(p')\|_F \leq 7\epsilon \quad (27)$$

$$\|\mathbf{H}(p) - \mathbf{H}(p')\|_F \leq \epsilon \quad (28)$$

$$\|\hat{\mathbf{H}}(p) - \hat{\mathbf{H}}(p')\|_F \leq \epsilon \quad (29)$$

$$\|\hat{\Sigma}(p) - \hat{\Sigma}(p')\|_F \leq 7\epsilon \quad (30)$$

$$\|\hat{\Sigma}(p) - \hat{\Sigma}(p')\|_F \leq \epsilon \quad (31)$$

654 *Proof.* [Eq. \(28\)](#) and [Eq. \(29\)](#) are direct consequences of [Eq. \(23\)](#) since  $\mathbf{H}(p)$  and  $\hat{\mathbf{H}}(p)$  are expect-  
655 tations of  $\widehat{\text{Cov}}(p^{\mathcal{A}})$  over distributions over  $\mathcal{A}$ . [Eq. \(31\)](#) is directly implied by [Eq. \(29\)](#) because  
656  $\hat{\Sigma}(p) = \hat{\mathbf{H}}(p) + \beta \mathbf{I}$ .

657 To show [Eq. \(24\)](#) and [Eq. \(25\)](#), observe that by the definition of  $x(p)$  and  $\mathbf{H}(p)$ ,

$$\begin{aligned} \mathbf{H}(p) &= \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} \begin{bmatrix} aa^\top & a \\ a^\top & 1 \end{bmatrix} = \begin{bmatrix} \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [aa^\top] & \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [a] \\ \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [a^\top] & 1 \end{bmatrix} \\ &= \begin{bmatrix} \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [aa^\top] & x(p) \\ x(p)^\top & 1 \end{bmatrix} \end{aligned}$$

658 Therefore,  $\|x(p) - x(p')\| \leq \|\mathbf{H}(p) - \mathbf{H}(p')\|_F \leq \epsilon$ . Similarly,  $\|\hat{x}(p) - \hat{x}(p')\| \leq \|\hat{\mathbf{H}}(p) -$   
659  $\hat{\mathbf{H}}(p')\|_F \leq \epsilon$ .

660 It remains to show [Eq. \(26\)](#), [Eq. \(27\)](#) and [Eq. \(30\)](#). Next, we show [Eq. \(26\)](#):

$$\begin{aligned} &H(p) - H(p') \\ &= \mathbb{E}_{\mathcal{A} \sim D} \left[ \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - \hat{x}(p))(a - \hat{x}(p))^\top] - \mathbb{E}_{a \sim p'^{\mathcal{A}}} [(a - \hat{x}(p'))(a - \hat{x}(p'))^\top] \right] \\ &= \mathbb{E}_{\mathcal{A} \sim D} \left[ \mathbb{E}_{a \sim p^{\mathcal{A}}} [aa^\top] - \mathbb{E}_{a \sim p'^{\mathcal{A}}} [aa^\top] \right] \\ &\quad - x(p)\hat{x}(p)^\top - \hat{x}(p)x(p)^\top + x(p')\hat{x}(p')^\top + \hat{x}(p')x(p')^\top \quad (\text{using } \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p^{\mathcal{A}}} [a] = x(p)) \\ &\quad + \hat{x}(p)\hat{x}(p)^\top - \hat{x}(p')\hat{x}(p')^\top \quad (32) \end{aligned}$$

661 Using the property

$$\|ab^\top - cd^\top\|_F \leq \|ab^\top - cb^\top\|_F + \|cb^\top - cd^\top\|_F \leq \|a - c\| \|b\| + \|c\| \|b - d\|$$

662 we continue from [Eq. \(32\)](#) and bound

$$\begin{aligned} &\|H(p) - H(p')\|_F \\ &\leq \|\mathbf{H}(p) - \mathbf{H}(p')\|_F + 2(\|\hat{x}(p) - \hat{x}(p')\| + \|x(p) - x(p')\|) + \|\hat{x}(p) - \hat{x}(p')\| + \|\hat{x}(p) - \hat{x}(p')\| \\ &\leq 7\epsilon. \end{aligned}$$

663 [Eq. \(27\)](#) can be shown in the same manner, which further implies [Eq. \(30\)](#) by the definition of  $\hat{\Sigma}(p)$ .

664  $\square$

665 **Lemma 15.** *With probability  $1 - \delta$ , for all  $t = 1, \dots, T$ ,*

$$\begin{aligned} \hat{H}_t + \frac{50(d+1)^3 \log(3T/\delta)}{t-1} \mathbf{I} &\succeq \frac{1}{2} H_t, \\ \hat{\mathbf{H}}_t + \frac{50(d+1)^3 \log(3T/\delta)}{t-1} \mathbf{I} &\succeq \frac{1}{2} \mathbf{H}_t. \end{aligned}$$

666 *Proof.* Notice that  $\hat{H}_t, \hat{\mathbf{H}}_t, H_t, \mathbf{H}_t$  corresponds to  $\hat{H}(p_t), \hat{\mathbf{H}}(p_t), H(p_t), \mathbf{H}(p_t)$  defined in [Definition 11](#) with  $n = t - 1$ . To show the lemma, our strategy is to argue the following two facts: 1) the  
667 two desired inequalities hold for all policies in the cover  $\mathbf{P}'$  (defined in [Eq. \(22\)](#)) with high probability.  
668 This is simply by applying [Lemma 12](#) with an union bound over policies in  $\mathbf{P}'$ . 2)  $p_t$  is sufficiently  
669 close to the nearest element in  $\mathbf{P}'$  so the desired inequalities still approximately hold.  
670

671 By [Proposition 1](#), we can find  $p' \in \mathbf{P}'$  such that for all  $\mathcal{A}$ ,

$$\left\| \widehat{\text{Cov}}(p_t^{\mathcal{A}}) - \widehat{\text{Cov}}(p'^{\mathcal{A}}) \right\|_F \leq \epsilon.$$

672 By [Proposition 2](#), it holds that

$$\|H(p_t) - H(p')\|_F \leq 7\epsilon, \quad \|\hat{H}(p_t) - \hat{H}(p')\|_F \leq 7\epsilon \quad (33)$$

$$\|\mathbf{H}(p_t) - \mathbf{H}(p')\|_F \leq \epsilon, \quad \|\hat{\mathbf{H}}(p_t) - \hat{\mathbf{H}}(p')\|_F \leq \epsilon \quad (34)$$

673 On the other hand, using [Lemma 12](#) and union bound, with probability  $1 - \delta$ , we have

$$\hat{H}(p') + \frac{4d \log(6d|\mathbf{P}'|/\delta)}{n} \mathbf{I} \succeq \frac{1}{2} H(p'), \quad (35)$$

$$\hat{\mathbf{H}}(p') + \frac{3d \log(d|\mathbf{P}'|/\delta)}{n} \mathbf{I} \succeq \frac{1}{2} \mathbf{H}(p'). \quad (36)$$

674 Combining [Eq. \(35\)](#) and [Eq. \(33\)](#), we get

$$\hat{H}(p_t) + 7\epsilon \mathbf{I} + \frac{4d \log(6d|\mathbf{P}'|/\delta)}{n} \mathbf{I} \succeq \hat{H}(p') + \frac{4d \log(6d|\mathbf{P}'|/\delta)}{n} \mathbf{I} \succeq \frac{1}{2} H(p') \succeq \frac{1}{2} H(p_t) - \frac{7}{2} \epsilon \mathbf{I}$$

675 which implies the first inequality in the lemma by plugging in the choice of  $\epsilon = \frac{1}{T^3}$  and the upper  
676 bound of  $\log |\mathbf{P}'|$  in [Proposition 2](#). The second inequality in the lemma can be obtained similarly by  
677 combining [Eq. \(34\)](#) and [Eq. \(36\)](#).

678 □

679 **Lemma 16.** *With probability of at least  $1 - \delta$ , for all  $t = 1, \dots, T$ ,*

$$\|x_t - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 \leq \mathcal{O}\left(\frac{d^3 \log(dT/\delta)}{t}\right)$$

680 *Proof.* Notice that  $x_t, \hat{x}_t, \hat{\Sigma}_t$  corresponds to  $x(p_t), \hat{x}(p_t), \hat{\Sigma}(p_t)$  defined in [Definition 11](#) with  $n =$   
681  $t - 1$ . To show the lemma, our strategy is to argue the following two facts: 1) the two desired  
682 inequalities hold for all policies in the cover  $\mathbf{P}'$  with high probability. This is simply by applying  
683 [Lemma 13](#) with an union bound over policies in  $\mathbf{P}'$ . 2)  $p_t$  is sufficiently close to the nearest element  
684 in  $\mathbf{P}'$  so the desired inequalities still approximately hold.

685 By [Proposition 1](#), we can find  $p' \in \mathbf{P}'$  such that for all  $\mathcal{A}$ ,

$$\left\| \widehat{\text{Cov}}(p_t^{\mathcal{A}}) - \widehat{\text{Cov}}(p'^{\mathcal{A}}) \right\|_F \leq \epsilon.$$

686 By [Proposition 2](#), we have

$$\|x(p') - x(p_t)\| \leq \epsilon, \quad \|\hat{x}(p') - \hat{x}(p_t)\| \leq \epsilon, \quad \|\hat{\Sigma}(p') - \hat{\Sigma}(p_t)\|_F \leq 7\epsilon \quad (37)$$

687 Thus,

$$\begin{aligned}
& \|x(p_t) - \hat{x}(p_t)\|_{\hat{\Sigma}(p_t)^{-1}}^2 \\
&= \left( \|x(p_t) - \hat{x}(p_t)\|_{\hat{\Sigma}(p_t)^{-1}}^2 - \|x(p') - \hat{x}(p')\|_{\hat{\Sigma}(p')^{-1}}^2 \right) + \|x(p') - \hat{x}(p')\|_{\hat{\Sigma}(p')^{-1}}^2 \\
&\leq \left( \|x(p_t) - \hat{x}(p_t)\|_{\hat{\Sigma}(p_t)^{-1}}^2 - \|x(p') - \hat{x}(p')\|_{\hat{\Sigma}(p')^{-1}}^2 \right) + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right) \\
&\quad \text{(by Lemma 13 with an union bound over } \mathbf{P}'\text{)} \\
&= \theta_t^\top \hat{\Sigma}(p_t)^{-1} \theta_t - \theta'^\top \hat{\Sigma}(p')^{-1} \theta' + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right) \\
&\quad \text{(define } \theta_t = x(p_t) - \hat{x}(p_t) \text{ and } \theta' = x(p') - \hat{x}(p')\text{)} \\
&= (\theta_t - \theta')^\top \hat{\Sigma}(p_t)^{-1} \theta_t + \theta'^\top (\hat{\Sigma}(p_t)^{-1} - \hat{\Sigma}(p')^{-1}) \theta_t + \theta'^\top \hat{\Sigma}(p')^{-1} (\theta_t - \theta') + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right) \\
&\leq (\theta_t - \theta')^\top (\hat{\Sigma}(p_t)^{-1} \theta_t + \hat{\Sigma}(p')^{-1} \theta') + \theta'^\top \hat{\Sigma}(p')^{-1} (\hat{\Sigma}(p') - \hat{\Sigma}(p_t)) \hat{\Sigma}(p_t)^{-1} \theta_t + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right)
\end{aligned}$$

688 The first two terms above can be bounded by the order of  $\mathcal{O}(\epsilon t^2)$  by Eq. (37). Using the choice  
689  $\epsilon = \frac{1}{T^3}$  and recalling that  $\log |\mathbf{P}'| = \mathcal{O}(d^2 \log(d/\epsilon))$  finishes the proof.

690 □

691 **Lemma 17.** *With probability of at least  $1 - \delta$ , for all  $t = 1, 2, \dots, T$ ,*

$$\|(\hat{\Sigma}_t - H_t)y_t\|_{\hat{\Sigma}_t^{-1}}^2 \leq \mathcal{O}\left(\frac{d^3 \log(dT/\delta)}{t}\right)$$

692 *Proof.* Notice that  $x_t, \hat{x}_t, \hat{\Sigma}_t$  corresponds to  $x(p_t), \hat{x}(p_t), \hat{\Sigma}(p_t)$  defined in Definition 11 with  $n =$   
693  $t - 1$ . To show the lemma, our strategy is to argue the following two facts: 1) the two desired  
694 inequalities hold for all policies in the cover  $\mathbf{P}'$  with high probability. This is simply by applying  
695 Lemma 13 with an union bound over policies in  $\mathbf{P}'$ . 2)  $p_t$  is sufficiently close to the nearest element  
696 in  $\mathbf{P}'$  so the desired inequalities still approximately hold.

697 By Proposition 1, we can find  $p' \in \mathbf{P}'$  such that for all  $\mathcal{A}$ ,

$$\left\| \widehat{\text{Cov}}(p_t^{\mathcal{A}}) - \widehat{\text{Cov}}(p'^{\mathcal{A}}) \right\|_F \leq \epsilon.$$

698 By Proposition 2, we have

$$\|x(p') - x(p_t)\| \leq \epsilon, \quad \|\hat{x}(p') - \hat{x}(p_t)\| \leq \epsilon, \quad \|\hat{\Sigma}(p') - \hat{\Sigma}(p_t)\|_F \leq 7\epsilon \quad (38)$$

699 Thus, for any  $\|y_t\|_2 \leq 1$ ,

$$\begin{aligned}
& \|(\hat{\Sigma}(p_t) - H(p_t))y_t\|_{\hat{\Sigma}(p_t)^{-1}}^2 \\
&= \left( \|(\hat{\Sigma}(p_t) - H(p_t))y_t\|_{\hat{\Sigma}(p_t)^{-1}}^2 - \|(\hat{\Sigma}(p') - H(p'))y_t\|_{\hat{\Sigma}(p')^{-1}}^2 \right) + \|(\hat{\Sigma}(p') - H(p'))y_t\|_{\hat{\Sigma}(p')^{-1}}^2 \\
&\leq \left( \|(\hat{\Sigma}(p_t) - H(p_t))y_t\|_{\hat{\Sigma}(p_t)^{-1}}^2 - \|(\hat{\Sigma}(p') - H(p'))y_t\|_{\hat{\Sigma}(p')^{-1}}^2 \right) + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right) \\
&\quad \text{(by Lemma 14 with an union bound over } \mathbf{P}'\text{)} \\
&= \theta_t^\top \hat{\Sigma}(p_t)^{-1} \theta_t - \theta'^\top \hat{\Sigma}(p')^{-1} \theta' + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right) \\
&\quad \text{(define } \theta_t = (\hat{\Sigma}(p_t) - H(p_t))y_t \text{ and } \theta' = (\hat{\Sigma}(p') - H(p'))y_t\text{)} \\
&= (\theta_t - \theta')^\top \hat{\Sigma}(p_t)^{-1} \theta_t + \theta'^\top (\hat{\Sigma}(p_t)^{-1} - \hat{\Sigma}(p')^{-1}) \theta_t + \theta'^\top \hat{\Sigma}(p')^{-1} (\theta_t - \theta') + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right) \\
&\leq (\theta_t - \theta')^\top (\hat{\Sigma}(p_t)^{-1} \theta_t + \hat{\Sigma}(p')^{-1} \theta') + \theta'^\top \hat{\Sigma}(p')^{-1} (\hat{\Sigma}(p') - \hat{\Sigma}(p_t)) \hat{\Sigma}(p_t)^{-1} \theta_t + \mathcal{O}\left(\frac{d \log(d|\mathbf{P}'|/\delta)}{t-1}\right)
\end{aligned}$$

700 The first two terms above can be bounded by the order of  $\mathcal{O}(\epsilon t^2)$  by Eq. (38). Plugging in the choice  
701 of  $\epsilon = \frac{1}{T^3}$  and recalling that  $\log |\mathbf{P}'| = \mathcal{O}(d^2 \log(d/\epsilon))$  finishes the proof.

702 □

703 **D Regret Analysis**

704 Consider the regret decomposition in [Section 3.5](#).

$$\begin{aligned} \text{Reg}(u) &= \mathbb{E} \left[ \sum_{t=1}^T \langle a_t - u^{\mathcal{A}_t}, y_t \rangle \right] = \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_t} - \mathbf{U}^{\mathcal{A}_t}, \gamma_t \rangle \right] = \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \gamma_t \rangle \right] \\ &\leq \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \gamma_t - \hat{\gamma}_t \rangle \right]}_{\text{Bias}} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \rangle \right]}_{\text{Bonus}} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right]}_{\text{FTRL-Reg}} \end{aligned}$$

705 where  $\mathcal{A}_0$  is drawn from  $D$  and is independent from the interaction between the learning and the  
706 environment. Recall that our algorithm is FTRL:

$$\mathbf{H}_t^{\mathcal{A}_0} = \underset{\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}}{\text{argmin}} \left\{ \sum_{s=1}^{t-1} \langle \mathbf{H}, \hat{\gamma}_s - \alpha_s \hat{\Sigma}_s^{-1} \rangle + \frac{F(\mathbf{H})}{\eta_t} \right\}.$$

707 The **FTRL-Reg** term can be handled by the standard FTRL analysis ([Lemma 5](#)). In order to deal  
708 with the issue that  $F$  can be unbounded on the boundary of  $\mathcal{H}^{\mathcal{A}_0}$ , we apply [Lemma 5](#) with the regret  
709 comparator  $\bar{\mathbf{U}}^{\mathcal{A}_0}$  defined as

$$\bar{\mathbf{U}}^{\mathcal{A}_0} = \left(1 - \frac{1}{T^2}\right) \mathbf{U}^{\mathcal{A}_0} + \frac{1}{T^2} \mathbf{H}_*^{\mathcal{A}_0}$$

710 where  $\mathbf{H}_*^{\mathcal{A}_0} \triangleq \underset{\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}}{\text{argmin}} F(\mathbf{H})$ . Thus,

**FTRL-Reg**

$$\begin{aligned} &\leq \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \bar{\mathbf{U}}^{\mathcal{A}_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right] + \mathbb{E} \left[ \sum_{t=1}^T \langle \bar{\mathbf{U}}^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right] \\ &\leq \underbrace{\mathbb{E} \left[ \frac{F(\bar{\mathbf{U}}^{\mathcal{A}_0}) - \min_{\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}} F(\mathbf{H})}{\eta_T} \right]}_{\text{Penalty}} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \max_{\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}} \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \right]}_{\text{Stability-1}} \\ &\quad + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \max_{\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}} \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, -\alpha_t \hat{\Sigma}_t^{-1} \rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \right]}_{\text{Stability-2}} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \bar{\mathbf{U}}^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right]}_{\text{Error}} \end{aligned} \tag{39}$$

711 In the rest of this section, we bound the following terms individually: **Bias**, **Bonus**, **Penalty**,  
712 **Stability-1**, **Stability-2**, **Error**.

713 For any  $t = 2, \dots, T$ , let  $\mathcal{E}_{t-1}$  be the event that the high-probability event in [Lemma 15](#), [Lemma 16](#),  
714 and [Lemma 17](#) happens for all  $1, \dots, t-1$  and  $\bar{\mathcal{E}}_{t-1}$  be the opposite event of  $\mathcal{E}_{t-1}$  (i.e. any of these  
715 three lemmas fails for any  $1, \dots, t-1$ ). We have  $\mathcal{P}[\mathcal{E}_{t-1}] = 1 - \mathcal{O}(\delta)$  and  $\mathcal{P}[\bar{\mathcal{E}}_{t-1}] = \mathcal{O}(\delta)$ . Let  
716  $\mathbb{E}[\cdot \mid \mathcal{E}_{t-1}]$  be the conditional expectation that event  $\mathcal{E}_{t-1}$  happens and let  $\mathbb{E}_t^{\mathcal{E}} = \mathbb{E}[\cdot \mid \mathcal{F}_{t-1}, \mathcal{E}_{t-1}]$

717 **D.1 Bounding the Bias term**

**Lemma 18.**

$$\text{Bias} = \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \gamma_t - \hat{\gamma}_t \rangle \right] \leq \frac{1}{4} \sum_{t=1}^T \alpha_t \|x_t - u\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \delta T^2 + \sum_{t=1}^T \frac{d^3 \log(T/\delta)}{\alpha_t t} \right)$$

718 *Proof.* For any  $t$ , we have

$$\begin{aligned}
& \mathbb{E}_t^{\mathcal{E}} \left[ \left\langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \right\rangle \right] \\
&= \mathbb{E}_t^{\mathcal{E}} \left[ \left\langle \mathbf{H}_t - \mathbf{U}, \gamma_t - \hat{\gamma}_t \right\rangle \right] && \text{(taking expectation over } \mathcal{A}_0) \\
&= \mathbb{E}_t^{\mathcal{E}} \left[ \left\langle x_t - u, y_t - \hat{y}_t \right\rangle \right] && \text{(by the definition of lifting)} \\
&= \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \left( y_t - \hat{\Sigma}_t^{-1} (a_t - \hat{x}_t) a_t^\top y_t \right) \right] && \text{(by the definition of } \hat{y}_t) \\
&= \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \left( y_t - \hat{\Sigma}_t^{-1} (a_t - \hat{x}_t) (a_t - \hat{x}_t)^\top y_t \right) \right] - \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \hat{\Sigma}_t^{-1} (a_t - \hat{x}_t) \hat{x}_t^\top y_t \right] \\
&= \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \left( I - \hat{\Sigma}_t^{-1} \mathbb{E}_{\mathcal{A} \sim \mathcal{D}} \mathbb{E}_{a_t \sim p_t^A} \left[ (a_t - \hat{x}_t) (a_t - \hat{x}_t)^\top \right] \right) y_t \right] \\
&\quad - \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \hat{\Sigma}_t^{-1} \left( \mathbb{E}_{\mathcal{A} \sim \mathcal{D}} \mathbb{E}_{a_t \sim p_t^A} [a_t] - \hat{x}_t \right) \hat{x}_t^\top y_t \right] && \text{(taking expectation over } \mathcal{A}_t \text{ and } a_t) \\
&= \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \hat{\Sigma}_t^{-1} \left( \hat{\Sigma}_t - H_t \right) y_t \right] - \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \hat{\Sigma}_t^{-1} (x_t - \hat{x}_t) \hat{x}_t^\top y_t \right] \\
&\hspace{15em} \text{(by the definition of } H_t \text{ and } x_t) \\
&\leq \mathbb{E}_t^{\mathcal{E}} \left[ (x_t - u)^\top \hat{\Sigma}_t^{-1} \left( \hat{\Sigma}_t - H_t \right) y_t \right] + \mathbb{E}_t^{\mathcal{E}} \left[ \left| (x_t - u)^\top \hat{\Sigma}_t^{-1} (x_t - \hat{x}_t) \right| \right] && (|\hat{x}_t^\top y_t| \leq 1) \\
&\leq \mathbb{E}_t^{\mathcal{E}} \left[ \|x_t - u\|_{\hat{\Sigma}_t^{-1}} \left( \|(\hat{\Sigma}_t - H_t)y_t\|_{\hat{\Sigma}_t^{-1}} + \|x_t - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}} \right) \right] && \text{(Cauchy-Schwarz)} \\
&\leq \mathcal{O} \left( \sqrt{\frac{d^3 \log(T/\delta)}{t}} \|x_t - u\|_{\hat{\Sigma}_t^{-1}} \right) && \text{(Lemma 17 and Lemma 16 given } \mathcal{E}_{t-1}) \\
&\leq \frac{\alpha_t}{4} \|x_t - u\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \frac{d^3 \log(T/\delta)}{\alpha_t t} \right) && \text{(AM-GM inequality)}
\end{aligned}$$

719 On the other hand, since  $\hat{\Sigma}_t \succeq \frac{1}{t} I \succeq \frac{1}{T} I$ , for any  $t = 1, \dots, T$ ,

$$\|\hat{y}_t\|_2 = \|\Sigma_t^{-1} (a_t - \hat{x}_t) a_t^\top y_t\|_2 \leq \|\Sigma_t^{-1} (a_t - \hat{x}_t)\|_2 \leq \mathcal{O}(T)$$

720 Thus, we have trivial bound

$$\mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \right\rangle \mid \overline{\mathcal{E}_{t-1}} \right] = \mathbb{E}_t \left[ \left\langle \mathbf{H}_t - \mathbf{U}, \gamma_t - \hat{\gamma}_t \right\rangle \mid \overline{\mathcal{E}_{t-1}} \right] = \mathbb{E}_t \left[ \left\langle x_t - u, y_t - \hat{y}_t \right\rangle \mid \overline{\mathcal{E}_{t-1}} \right] \leq \mathcal{O}(T)$$

721 Therefore, we have

$$\begin{aligned}
\text{Bias} &= \mathbb{E} \left[ \sum_{t=1}^T \left\langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \right\rangle \right] \\
&= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \right\rangle \right] \right] \\
&= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \right\rangle \mid \mathcal{E}_{t-1} \right] \mathbb{I}\{\mathcal{E}_{t-1}\} \right] + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \right\rangle \mid \overline{\mathcal{E}_{t-1}} \right] \mathbb{I}\{\overline{\mathcal{E}_{t-1}}\} \right] \\
&\leq \frac{1}{4} \sum_{t=1}^T \alpha_t \|x_t - u\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \log(T/\delta)}{\alpha_t t} + \delta T^2 \right)
\end{aligned}$$

722

□

## 723 D.2 Bounding the Bonus term

724 We first prove the following useful technique lemma to bound the inner product of lifted matrices.

725 **Lemma 19.** Let  $\mathbf{G} = \begin{bmatrix} G + gg^\top & g \\ g^\top & 1 \end{bmatrix}$ ,  $\mathbf{H} = \begin{bmatrix} H + hh^\top & h \\ h^\top & 1 \end{bmatrix}$  where  $G$  and  $H$  are positive semi-

726 definite, and  $\mathbf{H}' = \mathbf{H} + vv^\top$  where  $v = \begin{bmatrix} 0 \\ \sqrt{\beta} \end{bmatrix} \in \mathbb{R}^{d+1}$ . Then we have

727 1.  $\text{Tr}(\mathbf{H}^{-1}\mathbf{G}) = \text{Tr}(H^{-1}G) + \|g - h\|_{H^{-1}}^2 + 1$

728 2.  $\text{Tr}((\mathbf{H}')^{-1}\mathbf{G}) \geq \frac{1}{2(1+\frac{\beta}{1+\beta}\|h\|_{H^{-1}}^2)}\|g - h\|_{H^{-1}}^2 - \frac{\beta^2}{(1+\beta)^2}\|h\|_{H^{-1}}^2$

729 *Proof.* From Theorem 2.1 of [LS02], for any block matrix  $R = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$  if  $A$  is invertible and its

730 Schur complement  $S_A = D - CA^{-1}B$  is invertible, then

$$R^{-1} = \begin{bmatrix} A^{-1} + A^{-1}BS_A^{-1}CA^{-1} & -A^{-1}BS_A^{-1} \\ -S_A^{-1}CA & S_A^{-1} \end{bmatrix}$$

731 Using above equation, for the first equation, Since  $(H + hh^\top)^{-1} = H^{-1} - \frac{H^{-1}hh^\top H^{-1}}{1+h^\top H^{-1}h}$ . The inverse

732 Schur complement of  $H + hh^\top$  is  $1 + h^\top H^{-1}h$ . Thus

$$\mathbf{H}^{-1} = \begin{bmatrix} (I + H^{-1}hh^\top)(H + hh^\top)^{-1} & -H^{-1}h \\ -h^\top H^{-1} & 1 + h^\top H^{-1}h \end{bmatrix} = \begin{bmatrix} H^{-1} & -H^{-1}h \\ -h^\top H^{-1} & 1 + h^\top H^{-1}h \end{bmatrix}$$

733 and

$$\begin{aligned} \text{Tr}(\mathbf{H}^{-1}\mathbf{G}) &= \text{Tr}(H^{-1}G + H^{-1}gg^\top - H^{-1}hg^\top) - h^\top H^{-1}g + 1 + h^\top H^{-1}h \\ &= \text{Tr}(H^{-1}G) + g^\top H^{-1}g - 2g^\top H^{-1}h + h^\top H^{-1}h + 1 \\ &= \text{Tr}(H^{-1}G) + \|g - h\|_{H^{-1}}^2 + 1. \end{aligned}$$

734 For the second equation, observe that

$$\mathbf{H}' = \begin{bmatrix} H + hh^\top & h \\ h^\top & 1 + \beta \end{bmatrix} = (1 + \beta) \begin{bmatrix} \frac{1}{1+\beta}(H + hh^\top) & \frac{1}{1+\beta}h \\ \frac{1}{1+\beta}h^\top & 1 \end{bmatrix} = (1 + \beta) \begin{bmatrix} H' + h'h'^\top & h' \\ h'^\top & 1 \end{bmatrix}$$

735 where  $h' = \frac{1}{1+\beta}h$  and  $H' = \frac{1}{1+\beta}H + (\frac{1}{1+\beta} - \frac{1}{(1+\beta)^2})hh^\top = \frac{1}{1+\beta}H + \frac{\beta}{(1+\beta)^2}hh^\top \succeq 0$ .

736 Applying the first equality, we have

$$\text{Tr}((\mathbf{H}')^{-1}\mathbf{G}) = \frac{1}{1 + \beta} (\text{Tr}((H')^{-1}G) + \|g - h'\|_{H'^{-1}}^2 + 1) \geq \frac{1}{1 + \beta} \|g - h'\|_{H'^{-1}}^2.$$

737 Below, we continue to lower bound this term. By the same formula above, we have

$$H'^{-1} = \left( \frac{1}{1 + \beta}H + \frac{\beta}{(1 + \beta)^2}hh^\top \right)^{-1} = (1 + \beta)H^{-1} - \frac{\beta H^{-1}hh^\top H^{-1}}{1 + \frac{\beta}{1 + \beta}h^\top H^{-1}h}.$$

738 Thus

$$\begin{aligned} &\frac{1}{1 + \beta} \|g - h'\|_{H'^{-1}}^2 \\ &\geq \frac{1}{2(1 + \beta)} \|g - h\|_{H^{-1}}^2 - \frac{1}{1 + \beta} \|h - h'\|_{H'^{-1}}^2 \quad (\text{using } \|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2) \\ &= \frac{1}{2}(g - h)^\top \left( H^{-1} - \frac{\frac{\beta}{1 + \beta}H^{-1}hh^\top H^{-1}}{1 + \frac{\beta}{1 + \beta}h^\top H^{-1}h} \right) (g - h) - (h - h')^\top \left( H^{-1} - \frac{\frac{\beta}{1 + \beta}H^{-1}hh^\top H^{-1}}{1 + \frac{\beta}{1 + \beta}h^\top H^{-1}h} \right) (h - h') \\ &\geq \frac{1}{2} \|g - h\|_{H^{-1}}^2 - \frac{\frac{\beta}{1 + \beta}((g - h)^\top H^{-1}h)^2}{2 \left( 1 + \frac{\beta}{1 + \beta} \|h\|_{H^{-1}}^2 \right)} - \frac{\beta^2}{(1 + \beta)^2} \|h\|_{H^{-1}}^2 \quad (\text{using } h - h' = \frac{\beta}{1 + \beta}h) \\ &\geq \frac{1}{2} \|g - h\|_{H^{-1}}^2 - \frac{\frac{\beta}{1 + \beta} \|h\|_{H^{-1}}^2}{2 \left( 1 + \frac{\beta}{1 + \beta} \|h\|_{H^{-1}}^2 \right)} \|g - h\|_{H^{-1}}^2 - \frac{\beta^2}{(1 + \beta)^2} \|h\|_{H^{-1}}^2 \quad (\text{Cauchy-Schwarz}) \\ &= \frac{1}{2 \left( 1 + \frac{\beta}{1 + \beta} \|h\|_{H^{-1}}^2 \right)} \|g - h\|_{H^{-1}}^2 - \frac{\beta^2}{(1 + \beta)^2} \|h\|_{H^{-1}}^2. \end{aligned}$$

739

□

740 Using Lemma 19, we are able to show Corollary 20 which bound part of the second term.

741 **Corollary 20.**  $\text{Tr}(\mathbf{U}\hat{\Sigma}_t^{-1}) \geq \frac{1}{4}\|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 - \frac{1}{4}$ .

742 *Proof.* From Lemma 19, we have

$$\text{Tr}(\mathbf{U}\hat{\Sigma}_t^{-1}) \geq \frac{1}{2\left(1 + \frac{\beta_t}{1+\beta_t}\|\hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2\right)}\|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 - \frac{\beta_t^2}{(1+\beta_t)^2}\|\hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2.$$

743 Since  $\hat{\Sigma}_t \succeq \beta_t \mathbf{I}$ ,  $\hat{\Sigma}_t^{-1} \preceq \frac{1}{\beta_t} \mathbf{I}$ . Since  $\|\hat{x}_t\|_2 \leq 1$ , we have  $\|\hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 \leq \frac{1}{\beta_t}$ . Then

$$\begin{aligned} \text{Tr}(\mathbf{U}\hat{\Sigma}_t^{-1}) &\geq \frac{1}{2\left(1 + \frac{1}{1+\beta_t}\right)}\|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 - \frac{\beta_t}{(1+\beta_t)^2} \\ &\geq \frac{1}{4}\|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 - \frac{\beta_t}{(2\sqrt{\beta_t})^2} && (\beta_t \geq 0) \\ &= \frac{1}{4}\|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 - \frac{1}{4}. \end{aligned}$$

744

□

**Lemma 21.**

$$\begin{aligned} \mathbf{Bonus} &= \mathbb{E} \left[ \sum_{t=1}^T \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \right\rangle \right] \\ &\leq 2(d+2) \sum_{t=1}^T \alpha_t - \frac{1}{4} \sum_{t=1}^T \alpha_t \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \alpha_t \log(T/\delta)}{t} + \delta T \sum_{t=1}^T \alpha_t \right). \end{aligned}$$

745 *Proof.* For any  $t$ , we have

$$\begin{aligned} &\mathbb{E}_t^{\mathcal{E}} \left[ \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \right\rangle \right] \\ &= \mathbb{E}_t^{\mathcal{E}} \left[ \text{Tr} \left( \alpha_t (\mathbf{H}_t - \mathbf{U}) \hat{\Sigma}_t^{-1} \right) \right] && \text{(taking expectation over } \mathcal{A}_0) \\ &= \mathbb{E}_t^{\mathcal{E}} \left[ \alpha_t \text{Tr} \left( \mathbf{H}_t \hat{\Sigma}_t^{-1} \right) - \alpha_t \text{Tr} \left( \mathbf{U} \hat{\Sigma}_t^{-1} \right) \right] \\ &\leq \alpha_t \text{Tr} \left( \mathbb{E}_t^{\mathcal{E}} [\mathbf{H}_t] \hat{\Sigma}_t^{-1} \right) - \mathbb{E}_t^{\mathcal{E}} \left[ \frac{\alpha_t}{4} \|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 \right] + \frac{1}{4} \alpha_t && \text{(Corollary 20)} \\ &\leq 2\alpha_t(d+2) - \mathbb{E}_t^{\mathcal{E}} \left[ \frac{\alpha_t}{4} \|u - \hat{x}_t\|_{\hat{\Sigma}_t^{-1}}^2 \right] \\ &\leq 2\alpha_t(d+2) - \mathbb{E}_t^{\mathcal{E}} \left[ \frac{\alpha_t}{4} \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2 - \frac{\alpha_t}{4} \|\hat{x}_t - x_t\|_{\hat{\Sigma}_t^{-1}}^2 \right] \\ &\leq 2\alpha_t(d+2) - \frac{\alpha_t}{4} \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \frac{d^3 \alpha_t \log(T/\delta)}{t} \right) && \text{(Lemma 16)} \end{aligned}$$

746 On the other hand, since  $\hat{\Sigma}_t \succeq \frac{1}{t} \mathbf{I} \succeq \frac{1}{T} \mathbf{I}$ , we have trivial bound

$$\mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \right\rangle \mid \overline{\mathcal{E}_{t-1}} \right] \leq \mathcal{O}(\alpha_t T)$$

747 Therefore, we have

$$\begin{aligned}
\text{Bonus} &= \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \rangle \right] \\
&= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \rangle \right] \right] \\
&= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \rangle \mid \mathcal{E}_{t-1} \right] \mathbb{I}\{\mathcal{E}_{t-1}\} \right] + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{U}^{\mathcal{A}_0}, \alpha_t \hat{\Sigma}_t^{-1} \rangle \mid \overline{\mathcal{E}_{t-1}} \right] \mathbb{I}\{\overline{\mathcal{E}_{t-1}}\} \right] \\
&\leq 2(d+2) \sum_{t=1}^T \alpha_t - \frac{(1-\delta)}{4} \sum_{t=1}^T \alpha_t \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \alpha_t \log(T/\delta)}{t} + \delta T \sum_{t=1}^T \alpha_t \right) \\
&\leq 2(d+2) \sum_{t=1}^T \alpha_t - \frac{1}{4} \sum_{t=1}^T \alpha_t \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \alpha_t \log(T/\delta)}{t} + \delta T \sum_{t=1}^T \alpha_t \right)
\end{aligned}$$

748

□

### 749 D.3 Bounding the Penalty term

750 **Lemma 22.**  $\bar{\mathbf{U}}^{\mathcal{A}_0}$ , we have

$$\frac{F(\bar{\mathbf{U}}^{\mathcal{A}_0}) - \min_{\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}} F(\mathbf{H})}{\eta_T} \leq \frac{2d \log(T)}{\eta_T}$$

751 *Proof.* Since  $\bar{\mathbf{U}}^{\mathcal{A}_0} = (1 - \frac{1}{T^2}) \mathbf{U}^{\mathcal{A}_0} + \frac{1}{T^2} \mathbf{H}_*^{\mathcal{A}_0}$ , we have  $\bar{\mathbf{U}}^{\mathcal{A}_0} \succeq \frac{1}{T^2} \mathbf{H}_*^{\mathcal{A}_0}$ . Then

$$\frac{F(\bar{\mathbf{U}}^{\mathcal{A}_0}) - \min_{\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}} F(\mathbf{H})}{\eta_T} = \frac{1}{\eta_T} \log \frac{\det(\mathbf{H}_*^{\mathcal{A}_0})}{\det(\bar{\mathbf{U}}^{\mathcal{A}_0})} \leq \frac{2d \log(T)}{\eta_T}.$$

752

□

### 753 D.4 Bounding the Stability-1 term

754 [ZL22] gave a useful identity to bound the Bregman divergence. We restate it in Lemma 23 for  
755 completeness.

756 **Lemma 23.** Let  $\mathbf{G} = \begin{bmatrix} G + gg^\top & g \\ g^\top & 1 \end{bmatrix}$  and  $\mathbf{H} = \begin{bmatrix} H + hh^\top & h \\ h^\top & 1 \end{bmatrix}$ , we have

$$D(\mathbf{G}, \mathbf{H}) = D(G, H) + \|g - h\|_{H^{-1}}^2 \geq \|g - h\|_{H^{-1}}^2$$

*Proof.*

$$\begin{aligned}
D(\mathbf{G}, \mathbf{H}) &= F(\mathbf{G}) - F(\mathbf{H}) - \langle \nabla F(\mathbf{H}), \mathbf{G} - \mathbf{H} \rangle \\
&= \log \left( \frac{\det(\mathbf{H})}{\det(\mathbf{G})} \right) + \text{Tr}(\mathbf{H}^{-1}(\mathbf{G} - \mathbf{H})) \\
&= \log \left( \frac{\det(\mathbf{H})}{\det(\mathbf{G})} \right) + \text{Tr}(\mathbf{H}^{-1}\mathbf{G}) - d - 1 \\
&= \log \left( \frac{\det(\mathbf{H})}{\det(\mathbf{G})} \right) + \text{Tr}(\mathbf{H}^{-1}\mathbf{G}) - d - 1 \\
&= \log \left( \frac{\det(H)}{\det(G)} \right) + \text{Tr}(H^{-1}G) + \|g - h\|_{H^{-1}}^2 - d \quad (\text{Lemma 19}) \\
&= D(G, H) + \|g - h\|_{H^{-1}}^2 \\
&\geq \|g - h\|_{H^{-1}}^2
\end{aligned}$$

758 **Lemma 24.** For any  $\mathbf{H} \in \mathcal{H}^{\mathcal{A}_0}$ , we have

$$\text{Stability-1} = \mathbb{E} \left[ \sum_{t=1}^T \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \right] \leq 2d \sum_{t=1}^T \eta_t + \mathcal{O}(\delta T^2)$$

759 *Proof.* Recall that  $\mathbf{H}_t^{\mathcal{A}_0} = \widehat{\text{Cov}}(p_t^{\mathcal{A}_0})$  and  $\widehat{\text{Cov}}(p) = \begin{bmatrix} \text{Cov}(p) + \mu(p)\mu(p)^\top & \mu(p) \\ \mu(p)^\top & 1 \end{bmatrix}$ , we have

$$\begin{aligned} \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} &\leq \left\langle x_t^{\mathcal{A}_0} - \mu(p), \hat{y}_t \right\rangle - \frac{\|\mu(p) - x_t^{\mathcal{A}_0}\|_{\text{Cov}(p_t^{\mathcal{A}_0})^{-1}}^2}{2\eta_t} \quad (\text{Lemma 23}) \\ &\leq \|x_t^{\mathcal{A}_0} - \mu(p)\|_{\text{Cov}(p_t^{\mathcal{A}_0})^{-1}} \|\hat{y}_t\|_{\text{Cov}(p_t^{\mathcal{A}_0})} - \frac{\|\mu(p) - x_t^{\mathcal{A}_0}\|_{\text{Cov}(p_t^{\mathcal{A}_0})^{-1}}^2}{2\eta_t} \\ &\leq \frac{\eta_t}{2} \|\hat{y}_t\|_{\text{Cov}(p_t^{\mathcal{A}_0})}^2 \quad (\text{AM-GM inequality}) \\ &= \frac{\eta_t}{2} \|\hat{\Sigma}_t^{-1}(a_t - \hat{x}_t)\ell_t\|_{\text{Cov}(p_t^{\mathcal{A}_0})}^2 \\ &\leq \frac{\eta_t}{2} (a_t - \hat{x}_t)^\top \hat{\Sigma}_t^{-1} \text{Cov}(p_t^{\mathcal{A}_0}) \hat{\Sigma}_t^{-1} (a_t - \hat{x}_t) \quad (|\ell_t| \leq 1) \\ &= \frac{\eta_t}{2} \text{Tr} \left( (a_t - \hat{x}_t)(a_t - \hat{x}_t)^\top \hat{\Sigma}_t^{-1} \text{Cov}(p_t^{\mathcal{A}_0}) \hat{\Sigma}_t^{-1} \right) \end{aligned}$$

760 Since  $\mathbb{E}_{\mathcal{A} \sim \mathcal{D}} \mathbb{E}_{a \sim p^{\mathcal{A}}} [(a - \hat{x}_t)(a - \hat{x}_t)^\top] = H_t$ , taking expectations over  $\mathcal{A}_t$ ,  $a_t$  and  $\mathcal{A}_0$  conditioned  
761 on  $\mathcal{E}_{t-1}$ , we have

$$\begin{aligned} \mathbb{E}_t^\mathcal{E} \left[ \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \right] &\leq \mathbb{E}_t^\mathcal{E} \left[ \frac{\eta_t}{2} \text{Tr} \left( (a_t - \hat{x}_t)(a_t - \hat{x}_t)^\top \hat{\Sigma}_t^{-1} \text{Cov}(p_t^{\mathcal{A}_0}) \hat{\Sigma}_t^{-1} \right) \right] \\ &= \mathbb{E}_t^\mathcal{E} \left[ \frac{\eta_t}{2} \text{Tr} \left( H_t \hat{\Sigma}_t^{-1} \mathbb{E}_{\mathcal{A}_0 \sim D} [\text{Cov}(p_t^{\mathcal{A}_0})] \hat{\Sigma}_t^{-1} \right) \right]. \end{aligned}$$

762 Notice that given  $\mathcal{E}_{t-1}$ ,

$$\hat{\Sigma}_t \succeq \frac{1}{2} H_t = \frac{1}{2} \mathbb{E}_{\mathcal{A} \sim D} [\text{Cov}(p_t^{\mathcal{A}})] + \frac{1}{2} (\hat{x}_t - x_t)(\hat{x}_t - x_t)^\top \succeq \frac{1}{2} \mathbb{E}_{\mathcal{A} \sim D} [\text{Cov}(p_t^{\mathcal{A}})]$$

763 Hence we continue to upper bound the last expression by

$$\mathbb{E}_t^\mathcal{E} \left[ \eta_t \text{Tr} \left( H_t \hat{\Sigma}_t^{-1} \hat{\Sigma}_t \hat{\Sigma}_t^{-1} \right) \right] \leq \mathbb{E}_t^\mathcal{E} \left[ \eta_t \text{Tr} \left( H_t \hat{\Sigma}_t^{-1} \right) \right] \leq 2\eta_t d.$$

764 On the other hand, since  $\hat{\Sigma}_t \succeq \frac{1}{t} I \succeq \frac{1}{T} I$ , we have trivial bound

$$\mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \mid \overline{\mathcal{E}_{t-1}} \right] \leq \mathcal{O}(T)$$

765 Combining everything, we get

$$\begin{aligned} \text{Stability-1} &= \mathbb{E} \left[ \sum_{t=1}^T \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \right] \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \mid \mathcal{E}_{t-1} \right] \mathbb{I}\{\mathcal{E}_{t-1}\} \right] \\ &\quad + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \left\langle \mathbf{H}_t^{\mathcal{A}_0} - \mathbf{H}, \hat{\gamma}_t \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathcal{A}_0})}{2\eta_t} \mid \overline{\mathcal{E}_{t-1}} \right] \mathbb{I}\{\overline{\mathcal{E}_{t-1}}\} \right] \\ &\leq 2d \sum_{t=1}^T \eta_t + \mathcal{O}(\delta T^2). \end{aligned}$$

## 767 D.5 Bounding the Stability-2 term

768 Note that Lemma 8 does not require matrix  $A, B$  to be positive semi-definite. We will use it to prove  
769 the following lemma based on Lemma 34 in [DWZ23b].

770 **Lemma 25.** *If  $\eta_t \alpha_t \leq \frac{1}{64t}$ , then*

$$\text{Stability-2} = \mathbb{E} \left[ \sum_{t=1}^T \max_{\mathbf{H} \in \mathcal{H}^{\mathbf{A}_0}} \left\langle \mathbf{H}_t^{\mathbf{A}_0} - \mathbf{H}, -\alpha_t \hat{\Sigma}_t^{-1} \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathbf{A}_0})}{2\eta_t} \right] \leq d \sum_{t=1}^T \alpha_t + \mathcal{O}(\delta T^2)$$

771 *Proof.* We first show that  $\max_{\mathbf{H} \in \mathcal{H}^{\mathbf{A}_0}} \left\langle \mathbf{H}_t^{\mathbf{A}_0} - \mathbf{H}, -\alpha_t \hat{\Sigma}_t^{-1} \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathbf{A}_0})}{2\eta_t} \leq \frac{\alpha_t}{2} \|\hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})}$ .

772 Define

$$G(\mathbf{H}) = \left\langle \mathbf{H}_t^{\mathbf{A}_0} - \mathbf{H}, -\alpha_t \hat{\Sigma}_t^{-1} \right\rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{\mathbf{A}_0})}{2\eta_t}$$

773 and  $\lambda = \|\alpha_t \hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})}$ . Since  $\hat{\Sigma}_t \succeq \frac{1}{t}I$ ,  $\mathbf{H}_t^{\mathbf{A}_0} \preceq 2I$ ,  $\eta_t \alpha_t \leq \frac{1}{64t}$ , we have

$$\eta_t \lambda = \eta_t \|\alpha_t \hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})} = \eta_t \alpha_t \sqrt{\text{Tr}(\mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1} \mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1})} \leq 2\eta_t \alpha_t t \leq \frac{1}{32}.$$

774 Let  $\mathbf{H}'$  be the maximizer of  $G$ . Since  $G(\mathbf{H}_t^{\mathbf{A}_0}) = 0$ , we have  $G(\mathbf{H}') \geq 0$ . It suffices to show  
775  $\|\mathbf{H}' - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} \leq 16\eta_t \lambda$  because from Lemma 8, it leads to

$$G(\mathbf{H}') \leq \|\mathbf{H}_t^{\mathbf{A}_0} - \mathbf{H}'\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} \|\alpha_t \hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})} \leq 16\eta_t \lambda \alpha_t \|\hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})} = \frac{\alpha_t}{2} \|\hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})}$$

776 To show  $\|\mathbf{H}' - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} \leq 16\eta_t \lambda$ , it suffices to show that for all  $\mathbf{U}$  such that  $\|\mathbf{U} -$   
777  $\mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} = 16\eta_t \lambda$ ,  $G(\mathbf{U}) \leq 0$ . This is because given this condition, if  $\|\mathbf{H}' -$   
778  $\mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} > 16\eta_t \lambda$ , then there is a  $\mathbf{U}$  in the line segment between  $\mathbf{H}_t^{\mathbf{A}_0}$  and  $\mathbf{H}'$  such that  
779  $\|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} = 16\eta_t \lambda$ . From the condition,  $G(\mathbf{U}) \leq 0 \leq \min\{G(\mathbf{H}_t^{\mathbf{A}_0}), G(\mathbf{H}')\}$  which  
780 contradicts to the strictly concave of  $G$ .

781 Now consider any  $\mathbf{U}$  such that  $\|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} = 16\eta_t \lambda$ . By Taylor expansion, there exists  $\mathbf{U}'$   
782 in the line segment between  $\mathbf{U}$  and  $\mathbf{H}_t^{\mathbf{A}_0}$  such that

$$G(\mathbf{U}) \leq \|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} \|\alpha_t \hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})} - \frac{1}{4\eta_t} \|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{U}')}^2$$

783 We have  $\|\mathbf{U}' - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} \leq \|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} = 16\eta_t \lambda \leq \frac{1}{2}$ . From the Equation 2.2 in  
784 page 23 of [Nem04] (also appear in Eq.(5) of [AHR09]) and log det is a self-concordant function,  
785 we have  $\|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{U}')}^2 \geq \frac{1}{4} \|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})}^2$ . Thus, we have

$$G(\mathbf{U}) \leq \|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{\nabla^2F(\mathbf{H}_t^{\mathbf{A}_0})} \|\alpha_t \hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})} - \frac{1}{16\eta_t} \|\mathbf{U} - \mathbf{H}_t^{\mathbf{A}_0}\|_{(\mathbf{H}_t^{\mathbf{A}_0})^{-1}}^2 = 16\eta_t \lambda^2 - \frac{(16\eta_t \lambda)^2}{16\eta_t} = 0$$

786 We have  $\|\hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{\mathbf{A}_0})} = \sqrt{\text{Tr}(\mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1} \mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1})} = \sqrt{\text{Tr}((\mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1})^2)}$ . Observe the fol-  
787 lowing two facts: 1) all eigenvalues of  $\mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1}$  are non-negative since  $\mathbf{H}_t^{\mathbf{A}_0}$  and  $\hat{\Sigma}_t^{-1}$  are both  
788 positive semi-definite, 2) for a square matrix  $A$  with all non-negative eigenvalues,  $\text{Tr}(A^2) \leq \text{Tr}(A)^2$   
789 because  $\text{Tr}(A^2) = \sum_i \lambda_i(A)^2 = \sum_i \lambda_i(A) \lambda_i(A) \leq (\sum_i \lambda_i(A))^2$ . We have

$$\sqrt{\text{Tr}((\mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1})^2)} \leq \text{Tr}(\mathbf{H}_t^{\mathbf{A}_0} \hat{\Sigma}_t^{-1}).$$

790 This allows us to conclude

$$\mathbb{E}_t^\mathcal{E} \left[ \frac{\alpha_t}{2} \|\hat{\Sigma}_t^{-1}\|_{\nabla^{-2}F(\mathbf{H}_t^{A_0})} \right] \leq \frac{\alpha_t}{2} \mathbb{E}_t^\mathcal{E} \left[ \text{Tr}(\mathbf{H}_t^{A_0} \hat{\Sigma}_t^{-1}) \right] \leq \alpha_t d$$

791 where we use that  $\hat{\Sigma}_t \succeq \frac{1}{2} \mathbb{E}_{\mathcal{A}_0 \sim D} [\mathbf{H}_t^{A_0}]$  given  $\mathcal{E}_{t-1}$ .

792 On the other hand, since  $\hat{\Sigma}_t \succeq \frac{1}{t} \mathbf{I} \succeq \frac{1}{T} \mathbf{I}$ , for any  $t = 1, \dots, T$ , we have trivial bound

$$\mathbb{E}_t \left[ \max_{\mathbf{H} \in \mathcal{H}^{A_0}} \langle \mathbf{H}_t^{A_0} - \mathbf{H}, -\alpha_t \hat{\Sigma}_t^{-1} \rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{A_0})}{2\eta_t} \middle| \overline{\mathcal{E}_{t-1}} \right] \leq \mathcal{O}(T)$$

793 Overall,

$$\begin{aligned} \text{Stability-2} &= \mathbb{E} \left[ \sum_{t=1}^T \max_{\mathbf{H} \in \mathcal{H}^{A_0}} \langle \mathbf{H}_t^{A_0} - \mathbf{H}, -\alpha_t \hat{\Sigma}_t^{-1} \rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{A_0})}{2\eta_t} \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \max_{\mathbf{H} \in \mathcal{H}^{A_0}} \langle \mathbf{H}_t^{A_0} - \mathbf{H}, -\alpha_t \hat{\Sigma}_t^{-1} \rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{A_0})}{2\eta_t} \right] \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \max_{\mathbf{H} \in \mathcal{H}^{A_0}} \langle \mathbf{H}_t^{A_0} - \mathbf{H}, \hat{\gamma}_t \rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{A_0})}{2\eta_t} \middle| \mathcal{E}_{t-1} \right] \mathbb{I}\{\mathcal{E}_{t-1}\} \right] \\ &\quad + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \max_{\mathbf{H} \in \mathcal{H}^{A_0}} \langle \mathbf{H}_t^{A_0} - \mathbf{H}, \hat{\gamma}_t \rangle - \frac{D(\mathbf{H}, \mathbf{H}_t^{A_0})}{2\eta_t} \middle| \overline{\mathcal{E}_{t-1}} \right] \mathbb{I}\{\overline{\mathcal{E}_{t-1}}\} \right] \\ &\leq d \sum_{t=1}^T \alpha_t + \mathcal{O}(\delta T^2). \end{aligned}$$

794

□

## 795 D.6 Bounding the Error term

**Lemma 26.**

$$\mathbf{Error} = \mathbb{E} \left[ \sum_{t=1}^T \langle \bar{\mathbf{U}}^{A_0} - \mathbf{U}^{A_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right] \leq \mathcal{O}(1).$$

796 *Proof.* Since  $\bar{\mathbf{U}}^{A_0} = (1 - \frac{1}{T^2}) \mathbf{U}^{A_0} + \frac{1}{T^2} \mathbf{H}_*^{A_0}$ , and  $\hat{\Sigma}_t \succeq \frac{1}{T} \mathbf{I}$ ,  $\hat{\Sigma}_t \succeq \frac{1}{T} \mathbf{I}$  we have

$$\begin{aligned} \mathbf{Error} &= \mathbb{E} \left[ \sum_{t=1}^T \langle \bar{\mathbf{U}}^{A_0} - \mathbf{U}^{A_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right] \\ &= \mathbb{E} \left[ \frac{1}{T^2} \sum_{t=1}^T \langle \bar{\mathbf{U}}^{A_0} - \mathbf{H}_*^{A_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right] \\ &\leq \mathcal{O}(1). \end{aligned}$$

797

□

## 798 D.7 Finishing up

799 Recall the regret decomposition at the beginning of [Appendix D](#). From [Lemma 22](#), [Lemma 24](#),  
800 [Lemma 25](#), and [Lemma 26](#), we have

$$\begin{aligned} \mathbf{FTRL-Reg} &= \mathbf{Penalty} + \mathbf{Stability-1} + \mathbf{Stability-2} + \mathbf{Error} \\ &\leq \mathcal{O} \left( \frac{d \log(T)}{\eta_T} + d \sum_{t=1}^T \eta_t + d \sum_{t=1}^T \alpha_t + \delta T^2 \right) \end{aligned}$$

801 From [Lemma 18](#) and [Lemma 21](#), we can cancel out the additional regret induced by bias through the  
 802 well-designed bonus term. Namely,

$$\begin{aligned}
 \mathbf{Bias} + \mathbf{Bonus} &= \frac{1}{4} \sum_{t=1}^T \alpha_t \|x_t - u\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \log(T/\delta)}{\alpha_t t} + \delta T^2 \right) \\
 &\quad + 2(d+2) \sum_{t=1}^T \alpha_t - \frac{1}{4} \sum_{t=1}^T \alpha_t \|u - x_t\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \alpha_t \log \frac{T}{\delta}}{t} + \delta \sum_{t=1}^T \alpha_t T \right) \\
 &= \mathcal{O} \left( d \sum_{t=1}^T \alpha_t + \sum_{t=1}^T \frac{d^3 \log(T/\delta)}{\alpha_t t} + \sum_{t=1}^T \frac{d^3 \alpha_t \log(T/\delta)}{t} + \delta T^2 \right)
 \end{aligned}$$

803 Thus, we have

$$\begin{aligned}
 \mathbf{Reg} &= \mathbf{Bias} + \mathbf{Bonus} + \mathbf{FTRL-Reg} \\
 &= \mathcal{O} \left( \frac{d \log(T)}{\eta_T} + d \sum_{t=1}^T \eta_t + d \sum_{t=1}^T \alpha_t + \sum_{t=1}^T \frac{d^3 \log(T/\delta)}{\alpha_t t} + \sum_{t=1}^T \frac{d^3 \alpha_t \log(T/\delta)}{t} + \delta T^2 \right)
 \end{aligned}$$

804 Recall that we have an additional condition in [Lemma 25](#) such that for any  $t$ ,  $\eta_t \alpha_t \leq \frac{1}{64t}$ . Picking  
 805  $\alpha_t = \frac{d}{\sqrt{t}}$ ,  $\eta_t = \frac{1}{64d\sqrt{t}}$  and  $\delta = \frac{1}{T^2}$ , we get

$$\mathbf{Reg} = \mathcal{O} \left( d^2 \sqrt{T} \log(T) + d^4 \log(T) \right) = \mathcal{O}(d^2 \sqrt{T} \log(T))$$

806 where we assume  $d^2 \leq \sqrt{T}$  without loss of generality (otherwise the bound is vacuous).

## 807 E Handling Misspecification

808 In this section, we discuss how to handle misspecification as defined in [Section 3.6](#). In [Appendix E.1](#),  
 809 we study the case where the amount of misspecification  $\varepsilon$  is known by the learner. In [Appendix E.2](#),  
 810 we use a blackbox approach to turn it into an algorithm that achieves almost the same regret bound  
 811 (up to  $\log T$  factors) without knowing  $\varepsilon$ .

### 812 E.1 Known misspecification

813 As discussed in [Section 3.6](#), when the amount of misspecification  $\varepsilon$  is known, we still use [Algorithm 1](#),  
 814 but with different  $\alpha_t$  and  $\eta_t$ . Throughout this subsection, we let  $\alpha_t = \frac{d}{\sqrt{t} + \frac{\varepsilon}{\sqrt{d}}}$  and  $\eta_t = \frac{1}{64 \left( d\sqrt{t} + \frac{\varepsilon}{\sqrt{d}} t \right)}$ ,  
 815 and point out the modifications of the analysis from [Appendix D](#).

816 We start with the regret decomposition similar to that in [Appendix D](#), but here we define

$$\begin{aligned}
 y_t &= \operatorname{argmin}_{y \in \mathbb{B}_2^d} \max_{\mathcal{A} \in \operatorname{supp}(D)} \max_{a \in \mathcal{A}} |f_t(a) - \langle a, y \rangle|, \\
 \varepsilon_t &= \max_{\mathcal{A} \in \operatorname{supp}(D)} \max_{a \in \mathcal{A}} |f_t(a) - \langle a, y_t \rangle|, \\
 c_t(a) &= f_t(a) - \langle a, y_t \rangle.
 \end{aligned}$$

817 The regret decomposition goes as follows:

$$\begin{aligned}
\text{Reg}(u) &= \mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - f_t(u^{A_t})) \right] \\
&\leq \mathbb{E} \left[ \sum_{t=1}^T \langle a_t - u^{A_t}, y_t \rangle \right] + \sum_{t=1}^T \varepsilon_t \\
&\leq \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{A_t} - \mathbf{U}^{A_t}, \gamma_t \rangle \right] + \varepsilon T = \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t \rangle \right] + \varepsilon T \\
&\leq \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \rangle \right]}_{\text{Bias}} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \alpha_t \hat{\Sigma}_t^{-1} \rangle \right]}_{\text{Bonus}} \\
&\quad + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \hat{\gamma}_t - \alpha_t \hat{\Sigma}_t^{-1} \rangle \right]}_{\text{FTRL-Reg}} + \varepsilon T.
\end{aligned}$$

818 Now  $\hat{y}_t = \hat{\Sigma}_t^{-1}(a_t - \hat{x}_t)\ell_t$  with  $\mathbb{E}[\ell_t] = a_t^\top y_t + c_t(a_t)$ .

819 For the **Bias** term, the proof is almost the same as [Lemma 18](#). The only difference is that from the  
820 fourth line, we have

$$\mathbb{E}_t \left[ (x_t - u)^\top \left( y_t - \hat{\Sigma}_t^{-1}(a_t - \hat{x}_t) (a_t^\top y_t + c_t(a_t)) \right) \right]$$

821 for some  $c_t(a_t)$  such that  $|c_t(a_t)| \leq \varepsilon_t$ . This leads to an additional term of

$$\begin{aligned}
&\mathbb{E}_t^\mathcal{E} \left[ -(x_t - u)^\top \hat{\Sigma}_t^{-1}(a_t - \hat{x}_t)c_t(a_t) \right] \\
&\leq \mathbb{E}_t^\mathcal{E} \left[ \sqrt{(x_t - u)^\top \hat{\Sigma}_t^{-1} c_t(a_t)^2 (a_t - \hat{x}_t)(a_t - \hat{x}_t)^\top \hat{\Sigma}_t^{-1} (x_t - u)} \right] \\
&\leq \mathbb{E}_t^\mathcal{E} \left[ \sqrt{(x_t - u)^\top \hat{\Sigma}_t^{-1} \mathbb{E}_{\mathcal{A}_t, a_t} [c_t(a_t)^2 (a_t - \hat{x}_t)(a_t - \hat{x}_t)^\top] \hat{\Sigma}_t^{-1} (x_t - u)} \right] \\
&\leq \mathbb{E}_t^\mathcal{E} \left[ \varepsilon_t \sqrt{(x_t - u)^\top \hat{\Sigma}_t^{-1} (\mathbb{E}_{\mathcal{A}_t, a_t} [(a_t - \hat{x}_t)(a_t - \hat{x}_t)^\top]) \hat{\Sigma}_t^{-1} (x_t - u)} \right] \\
&\leq \mathbb{E}_t^\mathcal{E} \left[ \varepsilon_t \sqrt{(x_t - u)^\top \hat{\Sigma}_t^{-1} H_t \hat{\Sigma}_t^{-1} (x_t - u)} \right] \\
&\leq \varepsilon_t \|x_t - u\|_{\hat{\Sigma}_t^{-1}}
\end{aligned}$$

822 Plugging it into the proof of [Lemma 18](#), we have

$$\begin{aligned}
\mathbb{E}_t^\mathcal{E} \left[ \langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \rangle \right] &\leq \mathcal{O} \left( \sqrt{\frac{d^3 \log(T/\delta)}{t}} + \varepsilon_t \right) \|x_t - u\|_{\hat{\Sigma}_t^{-1}} \\
&\leq \frac{\alpha_t}{4} \|x_t - u\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \frac{d^3 \log(T/\delta)}{\alpha_t t} + \frac{\varepsilon_t^2}{\alpha_t} \right)
\end{aligned}$$

823 Other parts of the proof follow those in [Lemma 18](#). Finally, we get

$$\begin{aligned}
\text{Bias} &= \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{H}_t^{A_0} - \mathbf{U}^{A_0}, \gamma_t - \hat{\gamma}_t \rangle \right] \\
&\leq \frac{1}{4} \sum_{t=1}^T \alpha_t \|x_t - u\|_{\hat{\Sigma}_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \log(T/\delta)}{\alpha_t t} + \sum_{t=1}^T \frac{\varepsilon_t^2}{\alpha_t} + \delta T^2 \right)
\end{aligned}$$

824 The **Bonus** term will not be affected, according to [Lemma 21](#), we have

$$\mathbf{Bonus} \leq 2(d+2) \sum_{t=1}^T \alpha_t - \frac{1}{4} \sum_{t=1}^T \alpha_t \|u - x_t\|_{\Sigma_t^{-1}}^2 + \mathcal{O} \left( \sum_{t=1}^T \frac{d^3 \alpha_t \log(T/\delta)}{t} + \delta T^2 \right)$$

825 The **Penalty** term will not be affected, according to [Lemma 22](#), we have

$$\frac{F(\bar{U}^{A_0}) - \min_{\mathbf{H} \in \mathcal{H}^{A_0}} F(\mathbf{H})}{\eta_T} \leq \frac{2d \log(T)}{\eta_T}$$

826 **Stability-1** term is also unchanged, as we assume that  $\ell_t$  still lies in  $[-1, 1]$  even under misspecifica-  
827 tion. We still have

$$\mathbf{Stability-1} \leq \mathcal{O} \left( d \sum_{t=1}^T \eta_t + \delta T^2 \right)$$

828 The **Stability-2** term will not be affected as long as  $\eta_t \alpha_t \leq \frac{1}{64t}$ . According to [Lemma 25](#), we have

$$\mathbf{Stability-2} \leq \mathcal{O} \left( d \sum_{t=1}^T \alpha_t + \delta T^2 \right)$$

829 The **Error** term is also unaffected. We still have  $\mathbf{Error} = \mathcal{O}(1)$ .

830 Adding these terms together, the regret caused by bias and the negative term induced by bonus cancel  
831 out. We have

$$\mathbf{Reg} = \mathcal{O} \left( \frac{d \log(T)}{\eta_T} + d \sum_{t=1}^T (\eta_t + \alpha_t) + \sum_{t=1}^T \frac{d^3 \log(T/\delta)}{\alpha_t t} + \sum_{t=1}^T \frac{d^3 \alpha_t \log(T/\delta)}{t} + \sum_{t=1}^T \frac{\varepsilon_t^2}{\alpha_t} + \delta T^2 \right)$$

832 Recall that we pick  $\alpha_t = \frac{d}{\sqrt{t}} + \frac{\varepsilon}{\sqrt{d}} \cdot \eta_t = \frac{1}{64d\sqrt{t} + 64\frac{\varepsilon}{\sqrt{d}}t}$  and  $\delta = \frac{1}{T^2}$ . This gives

$$\mathbf{Reg} = \mathcal{O}(d^2 \sqrt{T} \log(T) + d^4 \log(T) + \sqrt{d} \varepsilon T) = \mathcal{O}(d^2 \sqrt{T} \log(T) + \sqrt{d} \varepsilon T)$$

833 where we assume  $d^2 \leq \sqrt{T}$  without loss of generality.

## 834 E.2 Unknown misspecification

835 In this subsection, we use a model selection technique to convert the algorithm in [Appendix E.1](#) which  
836 requires knowledge on  $\varepsilon$  into an algorithm that achieves a similar regret bound without knowing  $\varepsilon$ .  
837 Such a procedure to handle unknown misspecification/corruption has appeared in several previous  
838 works [[FGMZ20](#), [WDZ22](#)], though we adopt the technique in an unpublished concurrent work  
839 [[Ano23](#)] to handle the adversarial case.<sup>3</sup>

840 The idea here is a black-box reduction which turns an algorithm that only deals with known  $\varepsilon$  to one  
841 that handles unknown  $\varepsilon$ . This is similar to [[WDZ22](#)] but additionally handles adversarial losses using  
842 a different approach.

843 More specifically, the reduction has two layers. The bottom layer takes as input an arbitrary  
844 misspecification-robust algorithm that operates under known  $\varepsilon$  (e.g., [Algorithm 1](#)), and outputs  
845 a *stable* misspecification-robust algorithm (formally defined later) that still operates under known  
846  $\varepsilon$ . The top layer follows the standard Corral idea and takes as input a stable algorithm that operates  
847 under known  $\varepsilon$ , and outputs an algorithm that operates under unknown  $\varepsilon$ . Below, we explain these  
848 two layers of reduction in details.

<sup>3</sup>Since [[Ano23](#)] has not been published, for completeness, we restate all their results in [Appendix E.2](#). The goal is to use their reduction idea to handle the unknown misspecification case. We do not claim our contribution in the reduction idea.

---

**Algorithm 3** STable Algorithm By Independent Learners and Instance SElection (STABILISE)

---

**Input:**  $\varepsilon$  and a base algorithm satisfying [Definition 27](#).

**Initialize:**  $\lceil \log_2 T \rceil$  instances of the base algorithm  $\text{ALG}_1, \dots, \text{ALG}_{\lceil \log_2 T \rceil}$ , where  $\text{ALG}_j$  is configured with the parameter

$$\theta = \theta_j \triangleq 2^{-j}\varepsilon T + 4\sqrt{2^{-j}T \log T} + 8 \log(T).$$

**for**  $t = 1, 2, \dots$  **do**

    Receive  $w_t$ .

**if**  $w_t \leq \frac{1}{T}$  **then**

        play an arbitrary policy  $\pi_t$

**continue** (without updating any instances)

    Let  $j_t$  be such that  $w_t \in (2^{-j_t-1}, 2^{-j_t}]$ .

    Let  $\pi_t$  be the policy suggested by  $\text{ALG}_{j_t}$ .

    Output  $\pi_t$ .

**If** feedback is received, send it to  $\text{ALG}_{j_t}$  with probability  $\frac{2^{-j_t-1}}{w_t}$ , and discard it otherwise.

---

849 **Bottom Layer (from an Arbitrary Algorithm to a Stable Algorithm)** The input of the bottom  
850 layer is an arbitrary misspecification-robust algorithm, formally defined as:

851 **Definition 27.** An algorithm is misspecification-robust if it takes  $\theta$  as input, and achieves the following  
852 regret for any random stopping time  $t' \leq T$  and any policy  $u$ :

$$\mathbb{E} \left[ \sum_{t=1}^{t'} (f_t(a_t) - f_t(u^{A_t})) \right] \leq \mathbb{E} [c_1 \sqrt{t'} + c_2 \theta] + \Pr [\varepsilon_{1:t'} > \theta] T$$

853 for problem-dependent and  $\log(T)$  factors  $c_1, c_2 \geq 1$  and  $\varepsilon_{1:t'} \triangleq \sqrt{t' \sum_{\tau=1}^{t'} \varepsilon_\tau^2}$ .

854 In our case,  $c_1 = \Theta(d^2 \log T)$  and  $c_2 = \Theta(\sqrt{d})$ . While the regret bound in [Definition 27](#) might look  
855 cumbersome, it is in fact fairly reasonable: if the guess  $\theta$  is not smaller than the true amount of  $\varepsilon_{1:t'}$ ,  
856 the regret should be of order  $d^2 \sqrt{t'} + \sqrt{d}\theta$ ; otherwise, the regret bound is vacuous since  $T$  is its  
857 largest possible value. The only extra requirement is that the algorithm needs to be *anytime* (i.e., the  
858 regret bound holds for any stopping time  $t'$ ), but even this is known to be easily achievable by using a  
859 doubling trick over a fixed-time algorithm. It is then clear that [Algorithm 1](#) (together with a doubling  
860 trick) indeed satisfies [Definition 27](#).

861 As mentioned, the output of the bottom layer is a stable robust algorithm. To characterize stability,  
862 we follow [\[ALNS17\]](#) and define a new learning protocol that abstracts the interaction between the  
863 output algorithm of the bottom layer and the master algorithm from the top layer:

864 **Protocol 1.** In every round  $t$ , before the learner makes a decision, a probability  $w_t \in [0, 1]$  is revealed  
865 to the learner. After making a decision, the learner sees the desired feedback from the environment  
866 with probability  $w_t$ , and sees nothing with probability  $1 - w_t$ .

867 One can convert any misspecification-robust algorithm (defined in [Definition 27](#)) into a stable  
868 misspecification-robust algorithm (characterized in [Theorem 28](#)).

869 This conversion is achieved by a procedure that called STABILISE (see [Algorithm 3](#) for details). The  
870 high-level idea of STABILISE is as follows. Noticing that the challenge when learning in [Protocol 1](#)  
871 is that  $w_t$  varies over time, we discretize the value of  $w_t$  and instantiate one instance of the input  
872 algorithm to deal with one possible discretized value, so that it is learning in [Protocol 1](#) but with a  
873 fixed  $w_t$ , making it straightforward to bound its regret based on what it promises in [Definition 27](#).

874 More concretely, STABILISE instantiates  $\mathcal{O}(\log_2 T)$  instances  $\{\text{ALG}_j\}_{j=0}^{\lceil \log_2 T \rceil}$  of the input algorithm  
875 that satisfies [Definition 27](#), each with a different parameter  $\theta_j$ . Upon receiving  $w_t$  from the environ-  
876 ment, it dispatches round  $t$  to the  $j$ -th instance where  $j$  is such that  $w_t \in (2^{-j-1}, 2^{-j}]$ , and uses the  
877 policy generated by  $\text{ALG}_j$  to interact with the environment (if  $w_t \leq \frac{1}{T}$ , simply ignore this round).  
878 Based on [Protocol 1](#), the feedback for this round is received with probability  $w_t$ . To *equalize* the  
879 probability of  $\text{ALG}_j$  receiving feedback as mentioned in the high-level idea, when the feedback is

880 actually obtained, STABILISE sends it to  $\text{ALG}_j$  only with probability  $\frac{2^{-j-1}}{w_t}$  (and discards it other-  
881 wise). This way, every time  $\text{ALG}_j$  is assigned to a round, it always receives the desired feedback with  
882 probability  $w_t \cdot \frac{2^{-j-1}}{w_t} = 2^{-j-1}$ . This equalization step allows us to use the original guarantee of the  
883 base algorithm ([Definition 27](#)) and run it as it is, without requiring it to perform extra importance  
884 weighting steps as in [[ALNS17](#)].

885 The choice of  $\theta_j$  is crucial in making sure that STABILISE only has  $\varepsilon T$  regret overhead instead of  
886  $\frac{\varepsilon T}{\min_{t \in [T]} w_t}$ . Since  $\text{ALG}_j$  only receives feedback with probability  $2^{-j-1}$ , the expected total misspeci-  
887 fication it experiences is on the order of  $2^{-j-1} \varepsilon T$ . Therefore, its input parameter  $\theta_j$  only needs to be  
888 of this order instead of the total amount of misspecification  $\varepsilon T$ .

889 The formal guarantee of the conversion is stated in the following [Theorem 28](#).

890 **Theorem 28.** *If an algorithm is misspecification robust according to [Definition 27](#) for some constants*  
891 *( $c_1, c_2$ ), then [Algorithm 3](#) ensures*

$$\text{Reg} \leq \mathcal{O} \left( \mathbb{E} \left[ c'_1 \sqrt{T \rho_T} \right] + c'_2 \varepsilon T \right)$$

892 *under [Protocol 1](#), where  $\rho_T = \frac{1}{\min_{t \in [T]} w_t}$ , with  $c'_1 = \Theta((c_1 + c_2) \sqrt{\log T})$ .*

893 *Proof of [Theorem 28](#).* Define indicators

$$g_{t,j} = \mathbb{I}\{w_t \in (2^{-j-1}, 2^{-j}]\}$$

$$h_{t,j} = \mathbb{I}\{\text{ALG}_j \text{ receives the feedback for episode } t\}.$$

894 Now we consider the regret of  $\text{ALG}_j$ . Notice that  $\text{ALG}_j$  makes an update only when  $g_{t,j} h_{t,j} = 1$ . By  
895 the guarantee of the base algorithm ([Definition 27](#)), we have

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - f_t(u^{A_t})) g_{t,j} h_{t,j} \right] \\ & \leq \mathbb{E} \left[ c_1 \sqrt{\sum_{t=1}^T g_{t,j} h_{t,j}} + c_2 \theta_j \max_{t \leq T} g_{t,j} \right] + \Pr \left[ \sqrt{\left( \sum_{t=1}^T g_{t,j} h_{t,j} \right) \left( \sum_{t=1}^T \varepsilon_t^2 g_{t,j} h_{t,j} \right)} > \theta_j \right] T. \end{aligned} \quad (40)$$

896 We first bound the last term: Notice that  $\mathbb{E}[h_{t,j} | g_{t,j}] = 2^{-j-1} g_{t,j}$  by [Algorithm 3](#). Therefore,

$$\sum_{t=1}^T \varepsilon_t^2 g_{t,j} \mathbb{E}[h_{t,j} | g_{t,j}] = 2^{-j-1} \sum_{t=1}^T \varepsilon_t^2 g_{t,j} \leq 2^{-j-1} \varepsilon^2 T \quad (41)$$

$$\sum_{t=1}^T g_{t,j} \mathbb{E}[h_{t,j} | g_{t,j}] = 2^{-j-1} \sum_{t=1}^T g_{t,j} \leq 2^{-j-1} T \quad (42)$$

897 By Freedman's inequality, with probability at least  $1 - \frac{1}{T^2}$ ,

$$\begin{aligned} & \sum_{t=1}^T \varepsilon_t^2 g_{t,j} h_{t,j} - \sum_{t=1}^T \varepsilon_t^2 g_{t,j} \mathbb{E}[h_{t,j} | g_{t,j}] \\ & \leq 2 \sqrt{\sum_{t=1}^T (\varepsilon_t)^4 g_{t,j} \mathbb{E}[h_{t,j} | g_{t,j}] \log(T) + 4 \log(T)} \\ & \leq 4 \sqrt{\sum_{t=1}^T \varepsilon_t^2 g_{t,j} \mathbb{E}[h_{t,j} | g_{t,j}] \log(T) + 4 \log(T)} \\ & \leq \sum_{t=1}^T \varepsilon_t^2 g_{t,j} \mathbb{E}[h_{t,j} | g_{t,j}] + 8 \log(T) \quad (\text{AM-GM inequality}) \end{aligned}$$

898 which gives

$$\sum_{t=1}^T \varepsilon_t^2 g_{t,j} h_{t,j} \leq 2 \sum_{t=1}^T \varepsilon_t^2 g_{t,j} \mathbb{E}[h_{t,j}|g_{t,j}] + 8 \log(T) \leq 2^{-j} \varepsilon^2 T + 8 \log(T)$$

899 with probability at least  $1 - \frac{1}{T^2}$  using Eq. (41). Similarly,

$$\sum_{t=1}^T g_{t,j} h_{t,j} \leq 2 \sum_{t=1}^T g_{t,j} \mathbb{E}[h_{t,j}|g_{t,j}] + 8 \log(T) \leq 2^{-j} T + 8 \log(T)$$

900 with probability at least  $1 - \frac{1}{T^2}$ . Therefore, with probability at least  $1 - \frac{2}{T^2}$ ,

$$\begin{aligned} \sqrt{\left( \sum_{t=1}^T g_{t,j} h_{t,j} \right) \left( \sum_{t=1}^T \varepsilon_t^2 g_{t,j} h_{t,j} \right)} &\leq \sqrt{2^{-2j} \varepsilon^2 T^2 + 16 \cdot 2^{-j} T \log T + 64 \log^2 T} \\ &\leq 2^{-j} \varepsilon T + 4 \sqrt{2^{-j} T \log T} + 8 \log(T) \\ &\leq \theta_j \end{aligned}$$

901 Therefore, the last term in Eq. (40) is bounded by  $\frac{2}{T^2} T \leq \frac{2}{T}$ .

902 Next, we deal with other terms in Eq. (40). Again, by  $\mathbb{E}[h_{t,j}|g_{t,j}] = 2^{-j-1} g_{t,j}$ , Eq. (40) implies

$$2^{-j-1} \mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - f_t(u^{A_t})) g_{t,j} \right] \leq \mathbb{E} \left[ c_1 \sqrt{2^{-j-1} \sum_{t=1}^T g_{t,j}} + c_2 \theta_j \max_{t \leq T} g_{t,j} \right] + \frac{2}{T}.$$

903 which implies after rearranging:

$$\begin{aligned} &\mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - f_t(u^{A_t})) g_{t,j} \right] \\ &\leq \mathbb{E} \left[ c_1 \sqrt{\frac{1}{2^{-j-1}} \sum_{t=1}^T g_{t,j}} + \left( \frac{c_2 \theta_j}{2^{-j-1}} \right) \max_{t \leq T} g_{t,j} \right] + \frac{2}{T 2^{-j-1}} \\ &\leq \mathbb{E} \left[ c_1 \sqrt{\sum_{t=1}^T \frac{2g_{t,j}}{w_t}} + 4c_2 \left( \varepsilon T + \sqrt{\frac{T \log T}{2^{-j}}} + \log T \right) \max_{t \leq T} g_{t,j} \right] + \frac{2}{T 2^{-j-1}}. \end{aligned}$$

(using that when  $g_{t,j} = 1$ ,  $\frac{1}{2^{-j-1}} \leq \frac{2}{w_t}$ , and the definition of  $\theta_j$ )

904 Now, summing this inequality over all  $j \in \{0, 1, \dots, \lceil \log_2 T \rceil\}$ , we get

$$\begin{aligned} &\mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - f_t(u^{A_t})) \mathbb{I} \left\{ w_t > \frac{1}{T} \right\} \right] \\ &\leq \mathcal{O} \left( \mathbb{E} \left[ c_1 \sqrt{N \sum_{t=1}^T \frac{1}{w_t}} + N c_2 \varepsilon T + c_2 \sqrt{\frac{T \log T}{\min_{t \leq T} w_t}} + c_2 N \log T \right] + 1 \right) \\ &\leq \mathcal{O} \left( \mathbb{E} \left[ (c_1 + c_2) \sqrt{T \log(T) \rho_T} \right] + c_2 \varepsilon T \log T \right) \end{aligned}$$

905 where  $N \leq \mathcal{O}(\log T)$  is the number of  $\text{ALG}_j$ 's that has been executed at least once.

906 On the other hand,

$$\mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - f_t(u^{A_t})) \mathbb{I} \left\{ w_t \leq \frac{1}{T} \right\} \right] < T \mathbb{E} [\mathbb{I} \{\rho_T \geq T\}] \leq \mathbb{E} [\rho_T].$$

907 Combining the two parts and using the assumption  $c_2 \geq 1$  finishes the proof.  $\square$

---

**Algorithm 4** (A Variant of) Corral

---

**Initialize:** a log-barrier algorithm with each arm being an instance of an algorithm satisfying the guarantee in [Theorem 28](#). The hypothesis on  $\varepsilon T$  is set to  $2^i$  for arm  $i$  ( $i = 1, 2, \dots, M \triangleq \lceil \log_2 T \rceil$ ).  
**Initialize:**  $\rho_{0,i} = M, \forall i$ .

**for**  $t = 1, 2, \dots, T$  **do**

Let

$$w_t = \operatorname{argmin}_{w \in \Delta(M), w_i \geq \frac{1}{T}, \forall i} \left\{ \left\langle w, \sum_{\tau=1}^{t-1} (\hat{z}_\tau - r_\tau) \right\rangle + \frac{1}{\eta} \sum_{i=1}^M \log \frac{1}{w_i} \right\}$$

where  $\eta = \frac{1}{4c'_1 \sqrt{T}}$ .

For all  $i$ , send  $w_{t,i}$  to instance  $i$ .

Draw  $i_t \sim w_t$ .

Execute the  $a_t$  output by instance  $i_t$

Receive the loss  $z_{t,i_t}$  for action  $a_t$  (whose expectation is  $f_t(a_t)$ ) and send it to instance  $i_t$ .

Define for all  $i$ :

$$\hat{z}_{t,i} = \frac{z_{t,i} \mathbb{I}[i_t = i]}{w_{t,i}},$$

$$\rho_{t,i} = \min_{\tau \leq t} \frac{1}{w_{\tau,i}},$$

$$r_{t,i} = c'_1 \left( \sqrt{\rho_{t,i} T} - \sqrt{\rho_{t-1,i} T} \right).$$

908 **Top Layer (from Known  $\varepsilon$  to Unknown  $\varepsilon$ )** In this subsection, we use the algorithm that we  
909 construct in [Theorem 28](#) as a base algorithm, and further construct an algorithm with  $\sqrt{T} + \varepsilon$  regret  
910 under unknown  $\varepsilon$ . The idea is to run multiple base algorithms, each with a different hypothesis  
911 on  $\varepsilon$ ; on top of them, run another multi-armed bandit algorithm to adaptively choose among them.  
912 The goal is to let the top-level bandit algorithm perform almost as well as the best base algorithm.  
913 This is the Corral idea outlined in [[ALNS17](#), [FGMZ20](#), [LZZZ22](#)], and the algorithm is presented in  
914 [Algorithm 4](#).

915 **Theorem 29.** *Using an algorithm constructed in [Theorem 28](#) as a base algorithm, [Algorithm 4](#)*  
916 *ensures  $\operatorname{Reg} = \mathcal{O} \left( c'_1 \sqrt{T \log^3 T} + c'_2 \varepsilon T \right)$  without knowing  $\varepsilon$ .*

917 The top-level bandit algorithm is an FTRL with log-barrier regularizer. We first state the standard  
918 regret bound of FTRL under log-barrier regularizer, whose proof can be found in, e.g., [Theorem 7](#) of  
919 [[WL18](#)].

920 **Lemma 30.** *The FTRL algorithm over a convex subset  $\Omega$  of the  $(M-1)$ -dimensional simplex  $\Delta(M)$ :*

$$w_t = \operatorname{argmin}_{w \in \Omega} \left\{ \left\langle w, \sum_{\tau=1}^{t-1} \ell_\tau \right\rangle + \frac{1}{\eta} \sum_{i=1}^M \log \frac{1}{w_i} \right\}$$

921 *ensures for all  $u \in \Omega$ ,*

$$\sum_{t=1}^T \langle w - u, \ell_t \rangle \leq \frac{M \log T}{\eta} + \eta \sum_{t=1}^T \sum_{i=1}^M w_{t,i}^2 \ell_{t,i}^2$$

922 *as long as  $\eta w_{t,i} |\ell_{t,i}| \leq \frac{1}{2}$  for all  $t, i$ .*

923 *Proof of [Theorem 29](#).* The Corral algorithm is essentially an FTRL with log-barrier regularizer. To  
924 apply [Lemma 30](#), we first verify the condition  $\eta w_{t,i} |\ell_{t,i}| \leq \frac{1}{2}$  where  $\ell_{t,i} = \hat{z}_{t,i} - r_{t,i}$ . By our choice

925 of  $\eta$ ,

$$\begin{aligned} \eta w_{t,i} |\hat{z}_{t,i}| &\leq \eta z_{t,i} \leq \frac{1}{4}, && \text{(because } c'_1 \geq 1) \\ \eta w_{t,i} r_{t,i} &= \eta c'_1 \sqrt{T} w_{t,i} (\sqrt{\rho_{t,i}} - \sqrt{\rho_{t-1,i}}). \end{aligned}$$

926 The right-hand side of the last equality is non-zero only when  $\rho_{t,i} > \rho_{t-1,i}$ , implying that  $\rho_{t,i} = \frac{1}{w_{t,i}}$ .  
 927 Therefore, we further bound it by

$$\begin{aligned} \eta w_{t,i} r_{t,i} &\leq \eta c'_1 \sqrt{T} \frac{1}{\rho_{t,i}} (\sqrt{\rho_{t,i}} - \sqrt{\rho_{t-1,i}}) \\ &= \eta c'_1 \sqrt{T} \left( \frac{1}{\sqrt{\rho_{t,i}}} - \frac{\sqrt{\rho_{t-1,i}}}{\rho_{t,i}} \right) \\ &\leq \eta c'_1 \sqrt{T} \left( \frac{1}{\sqrt{\rho_{t-1,i}}} - \frac{1}{\sqrt{\rho_{t,i}}} \right) && \left( \frac{1}{\sqrt{a}} - \frac{\sqrt{b}}{a} \leq \frac{1}{\sqrt{b}} - \frac{1}{\sqrt{a}} \text{ for } a, b > 0 \right) \\ &\leq \eta c'_1 \sqrt{T} && (\rho_{t,i} \geq 1) \\ &= \frac{1}{4} && \text{(definition of } \eta) \end{aligned} \tag{43}$$

928 which can be combined to get the desired property  $\eta w_{t,i} |\hat{z}_{t,i} - r_{t,i}| \leq \frac{1}{2}$ .

929 Hence, by the regret guarantee of log-barrier FTRL (Lemma 30), we have

$$\begin{aligned} &\mathbb{E} \left[ \sum_{t=1}^T (z_{t,i_t} - z_{t,i^*}) \right] \\ &\leq \mathcal{O} \left( \frac{M \log T}{\eta} + \eta \mathbb{E} \left[ \underbrace{\sum_{t=1}^T \sum_{i=1}^M w_{t,i}^2 (\hat{z}_{t,i} - r_{t,i})^2}_{\text{term}_1} \right] \right) + \mathbb{E} \left[ \underbrace{\sum_{t=1}^T \left( \sum_{i=1}^M w_{t,i} r_{t,i} - r_{t,i^*} \right)}_{\text{term}_2} \right] \end{aligned}$$

930 where  $i^*$  is the smallest  $i$  such that  $2^i$  upper bounds the true total misspecification amount  $\varepsilon T$ .

931 **Bounding term<sub>1</sub>:**

$$\text{term}_1 \leq 2\eta \sum_{t=1}^T \sum_{i=1}^M w_{t,i}^2 (\hat{z}_{t,i}^2 + r_{t,i}^2)$$

932 where

$$2\eta \sum_{t=1}^T \sum_{i=1}^M w_{t,i}^2 \hat{z}_{t,i}^2 = 2\eta \sum_{t=1}^T \sum_{i=1}^M z_{t,i}^2 \mathbb{I}\{i_t = i\} \leq \mathcal{O}(\eta T)$$

933 and

$$\begin{aligned} 2\eta \sum_{t=1}^T \sum_{i=1}^M w_{t,i}^2 r_{t,i}^2 &\leq 4\eta \sum_{t=1}^T \sum_{i=1}^M (c'_1 \sqrt{T})^2 \left( \frac{1}{\sqrt{\rho_{t-1,i}}} - \frac{1}{\sqrt{\rho_{t,i}}} \right)^2 && \text{(continue from Eq. (43))} \\ &\leq 4\eta c_1'^2 T \times \sum_{t=1}^T \sum_{i=1}^M \left( \frac{1}{\sqrt{\rho_{t-1,i}}} - \frac{1}{\sqrt{\rho_{t,i}}} \right) \\ &\hspace{15em} \left( \frac{1}{\sqrt{\rho_{t-1,i}}} - \frac{1}{\sqrt{\rho_{t,i}}} \leq 1 \text{ and } 1 - a \leq -\ln a \right) \\ &\leq 4\eta c_1'^2 T M^{\frac{3}{2}}. && \text{(telescoping and using } \rho_{0,i} = M \text{ and } \rho_{T,i} \leq T) \end{aligned}$$

934 **Bounding term<sub>2</sub>:**

$$\begin{aligned}
\text{term}_2 &= \sum_{t=1}^T \sum_{i=1}^M w_{t,i} r_{t,i} - \sum_{t=1}^T r_{t,i^*} \\
&\leq c'_1 \sqrt{T} \sum_{t=1}^T \sum_{i=1}^M \left( \frac{1}{\sqrt{\rho_{t-1,i}}} - \frac{1}{\sqrt{\rho_{t,i}}} \right) - \left( c'_1 \sqrt{\rho_{T,i^*} T} - c'_1 \sqrt{\rho_{0,i^*} T} \right) \\
&\quad \text{(continue from Eq. (43) and using } 1 - a \leq -\ln a) \\
&\leq \mathcal{O} \left( c'_1 \sqrt{T} M^{\frac{3}{2}} \right) - c'_1 \sqrt{\rho_{T,i^*} T}.
\end{aligned}$$

935 Combining the two terms and using  $\eta = \Theta \left( \frac{1}{c'_1 \sqrt{T} + c'_2} \right)$ ,  $M = \Theta(\log T)$ , we get

$$\begin{aligned}
\mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - z_{t,i^*}) \right] &= \mathbb{E} \left[ \sum_{t=1}^T (z_{t,i_t} - z_{t,i^*}) \right] \\
&= \mathcal{O} \left( c'_1 \sqrt{T \log^3 T} \right) - \mathbb{E} \left[ c'_1 \sqrt{\rho_{T,i^*} T} \right] \tag{44}
\end{aligned}$$

936 On the other hand, by the guarantee of the base algorithm (Theorem 28) and that  $\varepsilon T \in [2^{i^* - 1}, 2^{i^*}]$ ,  
937 we have

$$\mathbb{E} \left[ \sum_{t=1}^T (z_{t,i^*} - f_t(u^{A_t})) \right] \leq \mathbb{E} \left[ c'_1 \sqrt{\rho_{T,i^*} T} \right] + c'_2 \varepsilon T. \tag{45}$$

938 Combining Eq. (44) and Eq. (45), we get

$$\mathbb{E} \left[ \sum_{t=1}^T (f_t(a_t) - f_t(u^{A_t})) \right] \leq \mathcal{O} \left( c'_1 \sqrt{T \log^3 T} \right) + c'_2 \varepsilon T,$$

939 which finishes the proof.  $\square$

940 *Proof of Theorem 3.* As shown in Appendix E.1, our Algorithm 1 can be adapted to satisfy Defini-  
941 tion 27 with  $c_1 = \Theta(d^2 \log T)$  and  $c_2 = \Theta(\sqrt{d})$ . By a concatenation of Theorem 28 and Theorem 29,  
942 we conclude that there is an algorithm that achieves

$$\mathcal{O} \left( (c_1 + c_2) \sqrt{T} \log^2 T + c_2 \varepsilon T \log T \right) = \mathcal{O} \left( d^2 \sqrt{T} \log^2 T + \sqrt{d} \varepsilon T \log T \right).$$

943 regret under unknown  $\varepsilon$ .  $\square$

## 944 **F Analysis for Linear EXP4**

945 *Proof of Theorem 4.* We first show that

$$\forall \pi \in \Pi : \text{Reg}(\pi) \triangleq \mathbb{E} \left[ \sum_{t=1}^T a_t^\top y_t - \sum_{t=1}^T \pi(\mathcal{A}_t)^\top y_t \right] \leq \mathcal{O} \left( \gamma T + \frac{\ln |\Pi|}{\eta} + \eta d T \right). \tag{46}$$

946 The magnitude of the loss is bounded by

$$\begin{aligned}
|\hat{\ell}_{t,\pi}| &= \left| \left\langle \pi(\mathcal{A}_t), \tilde{H}_t^{-1} a_t \ell_t \right\rangle \right| \\
&\leq \|\pi(\mathcal{A}_t)\|_{\tilde{H}_t^{-1}} \|a_t\|_{\tilde{H}_t} \\
&\leq \frac{1}{\gamma} \|\pi(\mathcal{A}_t)\|_{G_t^{-1}} \|a_t\|_{G_t} \leq \frac{d}{\gamma}.
\end{aligned}$$

947 If  $\gamma \geq 2d\eta$ , then we have  $|\hat{\ell}_{t,\pi}| \leq \frac{1}{2}$  and we can use the standard regret bound of exponential weights:

$$\forall \pi \in \Pi : \quad \text{Reg}(\pi) \leq \gamma T + \frac{\ln |\Pi|}{\eta} + \eta \sum_{t=1}^T \mathbb{E} \left[ \mathbb{E}_{a_t \sim p_t} \left[ \sum_{\pi \in \Pi} P_{t,\pi} \hat{\ell}_{t,\pi}^2 \right] \right].$$

948 Let  $H_t = \mathbb{E}_{a \sim p_t} [aa^\top]$ . Then we have  $\tilde{H}_t^{-1} \preceq \frac{1}{1-\gamma} H_t^{-1}$ , and thus

$$\begin{aligned} \mathbb{E}_{a_t \sim p_t} \left[ \sum_{\pi \in \Pi} P_{t,\pi} \hat{\ell}_{t,\pi}^2 \right] &\leq \mathbb{E}_{a_t \sim p_t} \left[ \sum_{\pi \in \Pi} P_{t,\pi} \cdot \langle \pi(\mathcal{A}_t), \tilde{H}_t^{-1} a_t \rangle^2 \right] \\ &= \mathbb{E}_{a_t \sim p_t} \mathbb{E}_{a \sim p_t} \left[ \langle a, \tilde{H}_t^{-1} a_t \rangle^2 \right] \quad (\text{by the definition of } p_{t,a}) \\ &\leq \frac{1}{(1-\gamma)^2} \text{Tr} (H_t H_t^{-1} H_t H_t^{-1}) = \mathcal{O}(d). \end{aligned}$$

949 Combining all proves [Eq. \(46\)](#).

950 Next, we show that there exists  $\theta \in \Theta$  such that

$$\mathbb{E}_{\mathcal{A} \sim D} \left[ \sum_{t=1}^T (\pi_\theta(\mathcal{A}) - \pi^*(\mathcal{A}))^\top y_t \right] \leq \mathcal{O}(1). \quad (47)$$

951 Let  $\hat{\theta}$  be the closest element in  $\Theta$  to  $\sum_{t=1}^T y_t$ . By the definition of  $\Theta$  and the assumption that  $\|y_t\| \leq 1$ ,  
952 we have  $\|\hat{\theta} - \sum_{t=1}^T y_t\| \leq \epsilon$ . Thus, for any  $\mathcal{A}$ ,

$$\sum_{t=1}^T (\pi_{\hat{\theta}}(\mathcal{A}) - \pi^*(\mathcal{A}))^\top y_t \leq \sum_{a \in \mathcal{A}} (\pi_{\hat{\theta}}(\mathcal{A}) - \pi^*(\mathcal{A}))^\top \hat{\theta} + \epsilon \leq \epsilon$$

953 where the last inequality is by the fact that  $\pi_{\hat{\theta}}(\mathcal{A}) = \operatorname{argmin}_{a \in \mathcal{A}} a^\top \hat{\theta}$ . Taking expectation over  $\mathcal{A}$   
954 gives [Eq. \(47\)](#).

955 Finally, combining [Eq. \(46\)](#) and [Eq. \(47\)](#), choosing  $\epsilon = 1$  and  $\gamma = 2d\eta = 2d\sqrt{\frac{\log T}{T}}$ , we get

$$\begin{aligned} \text{Reg} &= \mathbb{E} \left[ \sum_{t=1}^T a_t^\top y_t - \sum_{t=1}^T \pi^*(\mathcal{A}_t)^\top y_t \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T a_t^\top y_t - \sum_{t=1}^T \pi_{\hat{\theta}}(\mathcal{A}_t)^\top y_t \right] + \mathbb{E}_{\mathcal{A} \sim D} \left[ \sum_{t=1}^T (\pi_{\hat{\theta}}(\mathcal{A}) - \pi^*(\mathcal{A}))^\top y_t \right] \\ &= \mathcal{O} \left( \gamma T + \frac{\ln((2T)^d)}{\eta} + \eta d T + 1 \right) \\ &= \mathcal{O} \left( d \sqrt{T \log T} \right), \end{aligned}$$

956 finishing the proof. □

## 957 G Comparison with [\[DLWZ23, SKM23\]](#)

958 We state the exponential weight algorithm adopted by [\[LWL21, DLWZ23, SKM23\]](#) in [Algorithm 5](#),  
959 which is an algorithm that we know to achieve the prior-art regret bound in our setting (though they  
960 studied a more general MDP setting).

961 Their algorithm proceeds in *epochs* (indexed by  $k$ ), where every epoch consists of  $W$  rounds. The  
962 policy on action set  $\mathcal{A}$  in the  $k$ -th epoch is defined as

$$p_k^{\mathcal{A}}(a) \propto \exp \left( -\eta \sum_{s=1}^{k-1} (a^\top \hat{y}_s - b_s(a)) \right)$$

963 where  $\hat{y}_k$  is the loss estimator for epoch  $k$ , and  $b_k(a)$  is a (non-linear) bonus. In all  $W$  rounds in  
964 epoch  $k$ , the same policy is executed. The samples obtained in these  $W$  rounds are randomly divided  
965 into two halves. One half is used to estimate the covariance matrix  $\hat{\Sigma}_k$ , and the other half is used to  
966 construct the loss estimator  $\hat{y}_k$  (see [Line 5](#) of [Algorithm 5](#)).

---

**Algorithm 5** Exponential weights with magnitude-reduced loss estimators
 

---

1 **for**  $k = 1, 2, \dots, \frac{T}{W}$  **do**  
 2   For all  $\mathcal{A}$ , define

$$p_k^{\mathcal{A}}(a) = \frac{\exp\left(-\eta \sum_{s=1}^{k-1} (a^\top \hat{y}_s - b_s(a))\right)}{\sum_{a' \in \mathcal{A}} \exp\left(-\eta \sum_{s=1}^{k-1} (a'^\top \hat{y}_s - b_s(a'))\right)} \quad \text{for all } a \in \mathcal{A}.$$

3   Randomly partition  $\{(k-1)W + 1, \dots, kW\}$  into two equal parts  $\mathcal{T}_k, \mathcal{T}'_k$ .  
 4   **for**  $t = (k-1)W + 1, \dots, kW$  **do**  
   | receive  $\mathcal{A}_t$ , sample  $a_t \sim p_k^{\mathcal{A}_t}$ , and receive  $\ell_t$ .  
 5   Define

$$\hat{\Sigma}_k = \beta I + \frac{1}{|\mathcal{T}_k|} \sum_{t \in \mathcal{T}_k} a_t a_t^\top$$

$$\hat{y}_k = \hat{\Sigma}_k^{-1} \left( \frac{1}{|\mathcal{T}'_k|} \sum_{t \in \mathcal{T}'_k} a_t \ell_t \right)$$

$$b_k(a) = \alpha \|a\|_{\hat{\Sigma}_k^{-1}}.$$


---

**G.1 Regret Analysis Sketch**

968 The regret analysis starts with a standard decomposition that is similar to ours. We abuse the notation  
 969 by defining  $y_k = \frac{1}{W} \sum_{t=(k-1)W}^{kW} y_t$ . Then

$$\begin{aligned} \text{Reg} &= W \mathbb{E} \left[ \sum_{k=1}^{T/W} p_k^{\mathcal{A}_0}(a) \langle a - u^{\mathcal{A}_0}, y_k \rangle \right] \\ &= W \underbrace{\mathbb{E} \left[ \sum_{k=1}^{T/W} p_k^{\mathcal{A}_0}(a) \left( \langle a, \hat{y}_k \rangle - b_k(a) \right) - \left( u^{\mathcal{A}_0} - b_k(u^{\mathcal{A}_0}) \right) \right]}_{\text{EW-Reg}} + W \underbrace{\mathbb{E} \left[ \sum_{k=1}^{T/W} p_k^{\mathcal{A}_0}(a) b_k(a) - b_k(u^{\mathcal{A}_0}) \right]}_{\text{Bonus}} \\ &\quad + W \underbrace{\mathbb{E} \left[ \sum_{k=1}^{T/W} p_k^{\mathcal{A}_0}(a) \langle a - u^{\mathcal{A}_0}, y_k - \hat{y}_k \rangle \right]}_{\text{Bias}}. \end{aligned}$$

970 Bounding the regret term follows the standard analysis of exponential weight:

$$\begin{aligned} \text{EW-Reg} &\leq W \mathbb{E} \left[ \frac{\ln |\mathcal{A}_0|}{\eta} + \eta \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) \langle a, \hat{y}_k \rangle^2 + \eta \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) b_k(a)^2 \right] \\ &\leq W \mathbb{E} \left[ \frac{\ln |\mathcal{A}_0|}{\eta} + \eta \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) a^\top \hat{\Sigma}_k^{-1} H_k \hat{\Sigma}_k^{-1} a + \eta \sum_{k=1}^{T/W} \frac{\alpha^2}{\beta} \right] \end{aligned}$$

971 where  $H_k = \mathbb{E}_{\mathcal{A} \sim D} \mathbb{E}_{a \sim p_k^{\mathcal{A}}} [a a^\top]$ . Then they use the following fact to bound the stability term: as  
 972 long as  $W \geq \frac{d}{\beta^2}$ , it holds with high probability that  $\hat{\Sigma}_k^{-1} H_k \hat{\Sigma}_k^{-1} \preceq 2\hat{\Sigma}_k^{-1}$ . Thus **EW-Reg** can be

973 further bounded by

$$\begin{aligned} \mathbf{EW-Reg} &\lesssim W \left( \frac{\ln |\mathcal{A}_0|}{\eta} + \eta \mathbb{E} \left[ \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) \|a\|_{\hat{\Sigma}_k^{-1}}^2 \right] + \eta \frac{T}{W} \frac{\alpha^2}{\beta} \right) \\ &\leq \frac{W \ln |\mathcal{A}_0|}{\eta} + \eta dT + \eta T \frac{\alpha^2}{\beta}. \end{aligned}$$

974 By the definition of the bonus function  $b_t$ , it holds that

$$\mathbf{Bonus} = W \mathbb{E} \left[ \alpha \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) \|a\|_{\hat{\Sigma}_k^{-1}} \right] - W \mathbb{E} \left[ \alpha \sum_{k=1}^{T/W} \|u^{\mathcal{A}_0}\|_{\hat{\Sigma}_k^{-1}} \right].$$

975 Finally, the bias term can be bounded as follows:

$$\begin{aligned} \mathbf{Bias} &= W \mathbb{E} \left[ \sum_{k=1}^{T/W} p_k^{\mathcal{A}_0}(a) (a - u^{\mathcal{A}_0})^\top (y_k - \hat{\Sigma}_k^{-1} H_k y_k) \right] \\ &= W \mathbb{E} \left[ \sum_{k=1}^{T/W} p_k^{\mathcal{A}_0}(a) (a - u^{\mathcal{A}_0})^\top \hat{\Sigma}_k^{-1} (\hat{\Sigma}_k - H_k) y_k \right] \\ &\leq W \mathbb{E} \left[ \sum_{k=1}^{T/W} p_k^{\mathcal{A}_0}(a) \|a - u^{\mathcal{A}_0}\|_{\hat{\Sigma}_k^{-1}} \|(\hat{\Sigma}_k - H_k) y_k\|_{\hat{\Sigma}_k^{-1}} \right]. \end{aligned}$$

976 The bias here has a similar form as in our case. They use the following fact to bound the bias: as  
977 long as  $W \geq \frac{d}{\beta^2}$ , it holds that  $\|(\hat{\Sigma}_k - H_k) y_k\|_{\hat{\Sigma}_k^{-1}} \leq \sqrt{\beta d}$ . Therefore, the bias can further be upper  
978 bounded by

$$\mathbf{Bias} \leq W \mathbb{E} \left[ \sqrt{\beta d} \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) \|a\|_{\hat{\Sigma}_k^{-1}} + \sqrt{\beta d} \sum_{k=1}^{T/W} \|u^{\mathcal{A}_0}\|_{\hat{\Sigma}_k^{-1}} \right].$$

979 Combining the three parts, we get that the overall regret is of order

$$\mathbb{E} \left[ \frac{W \ln |\mathcal{A}_0|}{\eta} + \eta dT + \eta T \frac{\alpha^2}{\beta} + W(\alpha + \sqrt{\beta d}) \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) \|a\|_{\hat{\Sigma}_k^{-1}} + W(\sqrt{\beta d} - \alpha) \sum_{k=1}^{T/W} \|u^{\mathcal{A}_0}\|_{\hat{\Sigma}_k^{-1}} \right].$$

980 Choosing  $\alpha \approx \sqrt{\beta d}$ , we further bound it by

$$\begin{aligned} &\mathbb{E} \left[ \frac{W \ln |\mathcal{A}_0|}{\eta} + \eta dT + W \sqrt{\beta d} \sum_{k=1}^{T/W} \sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) \|a\|_{\hat{\Sigma}_k^{-1}} \right] \\ &\leq \mathbb{E} \left[ \frac{W \ln |\mathcal{A}_0|}{\eta} + \eta dT + W \sqrt{\beta d} \sum_{k=1}^{T/W} \sqrt{\sum_{a \in \mathcal{A}_0} p_k^{\mathcal{A}_0}(a) \|a\|_{\hat{\Sigma}_k^{-1}}^2} \right] \\ &\leq \frac{W \ln |\mathcal{A}_0|}{\eta} + \eta dT + \sqrt{\beta d} T. \end{aligned}$$

981 Recall the constraint  $W \geq \frac{d}{\beta^2}$ . Choosing  $W = \frac{d}{\beta^2}$  gives

$$\frac{d \ln |\mathcal{A}_0|}{\eta \beta^2} + \eta dT + \sqrt{\beta d} T \quad (48)$$

982 which gives  $d(\ln |\mathcal{A}_0|)^{\frac{1}{6}} T^{\frac{5}{6}}$  with the optimally chosen  $\eta$  and  $\beta$ .

983 **Remark** Due to the restrictions on the magnitude of the loss estimator required by the exponential  
984 weight algorithm, there is actually another constraint  $\frac{\eta}{\beta} \leq 1$ , which makes Eq. (48) be  $d(\ln |\mathcal{A}_0|)^{\frac{1}{7}} T^{\frac{6}{7}}$   
985 at best. This is exactly the bound obtained by [SKM23]. A more sophisticated way to construct  $\hat{y}_k$   
986 developed by [DLWZ23] removes this additional requirement and allows a bound of  $d(\ln |\mathcal{A}_0|)^{\frac{1}{6}} T^{\frac{5}{6}}$ .  
987 The sub-optimal bound  $T^{\frac{8}{5}}$  reported in [DLWZ23] is due to issues related to MDPs, which is not  
988 presented in the contextual bandit case here.