

---

# Variance-Dependent Regret Bounds for Non-stationary Linear Bandits

---

Zhiyong Wang<sup>1</sup>      Jize Xie<sup>2</sup>      Yi Chen<sup>2</sup>      John C.S. Lui<sup>1</sup>      Dongruo Zhou<sup>3</sup>  
<sup>1</sup>The Chinese University of Hong Kong    <sup>2</sup>Hong Kong University of Science and Technology  
<sup>3</sup>Indiana University Bloomington

## Abstract

We investigate the non-stationary stochastic linear bandit problem where the reward distribution evolves each round. Existing algorithms characterize the non-stationarity by the *total variation budget*  $B_K$ , which is the summation of the change of the consecutive feature vectors of the linear bandits over  $K$  rounds. However, such a quantity only measures the non-stationarity with respect to the expectation of the reward distribution, which makes existing algorithms sub-optimal under the general non-stationary distribution setting. In this work, we propose algorithms that utilize the *variance* of the reward distribution as well as the  $B_K$ , and show that they can achieve tighter regret upper bounds. Specifically, we introduce two novel algorithms: Restarted WeightedOFUL<sup>+</sup> and Restarted SAVE<sup>+</sup>. These algorithms address cases where the variance information of the rewards is known and unknown, respectively. Notably, when the total variance  $V_K$  is much smaller than  $K$ , our algorithms outperform previous state-of-the-art results on non-stationary stochastic linear bandits under different settings. Experimental evaluations further validate the superior performance of our proposed algorithms over existing works.

## 1 Introduction

In this work, we study non-stationary stochastic bandits, which is a generalization of the classical stationary stochastic bandits, where the reward distribution is

non-stationary. The intuition about the non-stationary setting comes from real-world applications such as dynamic pricing and ads allocation, where the environment changes rapidly and deviates significantly from stationarity [Auer et al. \(2002\)](#); [Cheung et al. \(2018\)](#). Most of the existing works in stochastic bandits consider a stationary setting where the goal of the agent is to minimize the *static regret*, *i.e.*, the summation of sub-optimality gaps between the agent’s selected arm and the fixed, time-independent best arm that maximizes the expectation of the reward distribution. In contrast, for the non-stationary setting, the emphasis shifts to minimizing the *dynamic regret*, which represents the gap between the cumulative reward of selecting the time-dependent optimal arm at each time and that of the learner. As we can always treat a stationary bandit instance as a special case of the non-stationary bandit instance, designing algorithms that work well under the non-stationary setting is significantly more challenging.

There have been a series of works aiming to minimize the *dynamic regret* for non-stationary stochastic bandits, such as Multi-Armed Bandits (MAB) ([Auer et al., 2002](#); [Garivier and Moulines, 2011](#); [Besbes et al., 2014b](#); [Wei et al., 2016](#)), linear bandits ([Cheung et al., 2018, 2019](#); [Zhao et al., 2020b](#); [Wei and Luo, 2021](#); [Wang et al., 2023](#)), general function approximation ([Fauray et al., 2021](#); [Russac et al., 2020, 2021](#)), and the even more challenging reinforcement learning (RL) setting ([Mao et al., 2021](#); [Touati and Vincent, 2020](#); [Gajane et al., 2018](#); [Cheung et al., 2020](#); [Wei and Luo, 2021](#)). In this work, we mainly consider the linear bandit setting, where each arm is a contextual vector, and the expected reward of each arm is assumed to be the linear product of the arm with an unknown feature vector. Most existing *dynamic regret* results for non-stationary linear bandits depend on both the *non-stationarity measurement* and the number of interaction rounds. Specifically, assume  $K$  is the total number of rounds, and for each  $k \in [K]$ ,  $\mathbf{x}$  is one of the arms,  $\boldsymbol{\theta}_k$  and  $\boldsymbol{\theta}_{k+1}$  are the feature vectors at  $k$  and  $k+1$  rounds, satisfying  $\|\mathbf{x}\|_2 \leq 1$ . Then, the non-stationarity measurement is often defined as the summation of the changes in the

mean of the reward distribution, which is

$$B_K := \sum_{k=1}^K \max_{\mathbf{x} \in \mathbb{R}^d} |\langle \mathbf{x}, \boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1} \rangle| = \sum_{k=1}^K \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2. \quad (1)$$

Existing works for non-stationary linear bandits (Rusac et al., 2019; Kim and Tewari, 2020a; Zhao et al., 2020a; Touati and Vincent, 2020; Cheung et al., 2018; Zhao et al., 2020b) achieved a regret upper bound of  $\tilde{O}(d^{7/8} B_K^{1/4} K^{3/4})$ , where  $d$  is the problem dimension. A recent work by Wei and Luo (2021) proposed a black-box reduction method that can achieve a regret upper bound of  $\tilde{O}(dB_K^{1/3} K^{2/3})$  in the setting with a fixed arm set across all rounds. Such regret bounds clearly demonstrate that regret grows as long as the non-stationarity grows, which is aligned with intuition.

Although existing works clearly demonstrate the relationship between the  $B_K$  and the regret, we claim that it is not sufficient for us to fully characterize the non-stationary level of the reward distributions. Consider applications such as hyperparameter tuning in physical systems, the noise distribution may highly depend on the evaluation point since the measurement noise often largely varies with the chosen parameter settings Kirschner and Krause (2018). For linear bandits, such examples suggest that the non-stationarity not only consists of the change of the mean of the distribution, but also the variance of the distribution. However, none of the previous works on non-stationary linear bandits considered how to leverage the variance information to improve regret bounds in the above heteroscedastic noise setting. Therefore, an open question arises:

*Can we design even better algorithms for non-stationary linear bandits by considering its variance information?*

In this paper, we answer this question affirmatively. We assume that at the  $k$ -th round, the reward distribution of an arm  $\mathbf{x}$  satisfies  $r_k \sim \langle \boldsymbol{\theta}_k, \mathbf{x} \rangle + \epsilon_k$ , where  $\epsilon_k$  is a zero-mean noise variable with variance  $\sigma_k^2$ . Our contributions are:

- We establish the first variance-dependent regret lower bound for non-stationary linear bandits. This result captures the interplay between non-stationarity and variance, which is not addressed in existing literature for non-stationary linear bandits.
- For the case where the reward variance  $\sigma_k^2$  at round  $k$  can be observed and the *total variation budget*  $B_K$  is known, we propose the Restarted-WeightedOFUL<sup>+</sup> algorithm, which uses variance-based weighted linear regression to deal with heteroscedastic noises (Zhou

et al., 2021; Zhou and Gu, 2022) and a restarted scheme to forget some historical data to hedge against the non-stationarity. We prove that the regret upper bound of Restarted-WeightedOFUL<sup>+</sup> is  $\tilde{O}(d^{7/8} (B_K V_K)^{1/4} \sqrt{K} + d^{5/6} B_K^{1/3} K^{2/3})$ . Notably, our regret surpasses the best result for non-stationary linear bandits  $\tilde{O}(dB_K^{1/3} K^{2/3})$  (Wei and Luo, 2021) when the total variance  $V_K = \tilde{O}(1)$  is small, which indicates that additional variance information benefits non-stationary linear bandit algorithms.

- For the case where the reward variance  $\sigma_k^2$  is unknown but the total variance  $V_K$  and variation budget  $B_K$  are known, we propose the Restarted-SAVE<sup>+</sup> algorithm. It maintains a multi-layer weighted linear regression structure with carefully-designed weight within each layer to handle the unknown variances (Zhao et al., 2023). We prove that Restarted-SAVE<sup>+</sup> can achieve a regret upper bound of  $\tilde{O}(d^{4/5} V_K^{2/5} B_K^{1/5} K^{2/5} + d^{2/3} B_K^{1/3} K^{2/3})$ . Specifically, when  $V_K = \tilde{O}(1)$ , our regret is also better than the existing best result  $\tilde{O}(dB_K^{1/3} K^{2/3})$  (Wei and Luo, 2021), which again verifies the effect of the variance information.
- Lastly, we propose Restarted-SAVE<sup>+</sup>-BOB for the case where both the reward variance  $\sigma_k^2$  and  $B_K$  are unknown. Restarted-SAVE<sup>+</sup>-BOB equips a *bandit-over-bandit* (BOB) framework to handle the unknown  $B_K$  (Cheung et al., 2019), and also maintains a multi-layer structure as Restarted-SAVE<sup>+</sup>. We show that Restarted-SAVE<sup>+</sup>-BOB achieves a regret upper bound of  $\tilde{O}(d^{4/5} V_K^{2/5} B_K^{1/5} K^{2/5} + d^{2/3} B_K^{1/3} K^{2/3} + d^{1/5} K^{7/10})$ , and it behaves the same as Restarted-SAVE<sup>+</sup> when  $V_K = \tilde{O}(1)$  and  $B_K = \Omega(d^{-14} K^{1/10})$ .
- We also conduct experimental evaluations to validate the outperformance of our proposed algorithms over existing works.

**Notation** We use lower case letters to denote scalars, and use lower and upper case bold face letters to denote vectors and matrices respectively. We denote by  $[n]$  the set  $\{1, \dots, n\}$ . For a vector  $\mathbf{x} \in \mathbb{R}^d$  and a positive semi-definite matrix  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ , we denote by  $\|\mathbf{x}\|_2$  the vector's Euclidean norm and define  $\|\mathbf{x}\|_{\boldsymbol{\Sigma}} = \sqrt{\mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x}}$ . For two positive sequences  $\{a_n\}$  and  $\{b_n\}$  with  $n = 1, 2, \dots$ , we write  $a_n = O(b_n)$  if there exists an absolute constant  $C > 0$  such that  $a_n \leq C b_n$  holds for all  $n \geq 1$  and write  $a_n = \Omega(b_n)$  if there exists an absolute constant  $C > 0$  such that  $a_n \geq C b_n$  holds for all  $n \geq 1$ . We use  $\tilde{O}(\cdot)$  to further hide the polylogarithmic factors.

Model	Algorithm	Regret	Variance -Dependent	Varying Arm Set	Require $B_K$
Linear Bandit	SW-UCB (Cheung et al., 2018)	$\tilde{O}(d^{\frac{7}{8}} B_K^{\frac{1}{4}} K^{\frac{3}{4}})$	No	Yes	Yes
	BOB (Cheung et al., 2018)	$\tilde{O}(d^{\frac{7}{8}} B_K^{\frac{1}{4}} K^{\frac{3}{4}})$	No	Yes	No
	RestartUCB (Zhao et al., 2020b)	$\tilde{O}(d^{\frac{7}{8}} B_K^{\frac{1}{4}} K^{\frac{3}{4}})$	No	Yes	Yes
	RestartUCB-BOB (Zhao et al., 2020b)	$\tilde{O}(d^{\frac{7}{8}} B_K^{\frac{1}{4}} K^{\frac{3}{4}})$	No	Yes	No
	LB-WeightUCB (Wang et al., 2023)	$\tilde{O}(d^{\frac{3}{4}} B_K^{\frac{1}{4}} K^{\frac{3}{4}})$	No	Yes	Yes
	MASTER + OFUL (Wei and Luo, 2021)	$\tilde{O}(dB_K^{\frac{1}{3}} K^{\frac{2}{3}})$	No	No	No
	Restarted-WeightedOFUL <sup>+</sup> (Ours)	$\tilde{O}(d^{\frac{7}{8}} (B_K V_K)^{\frac{1}{4}} K^{\frac{1}{2}} + d^{\frac{5}{6}} B_K^{\frac{1}{3}} K^{\frac{2}{3}})$	Yes	Yes	Yes
	Restarted SAVE <sup>+</sup> (Ours)	$\tilde{O}(d^{\frac{4}{5}} V_K^{\frac{2}{5}} B_K^{\frac{1}{5}} K^{\frac{2}{5}} + d^{\frac{2}{3}} B_K^{\frac{1}{3}} K^{\frac{2}{3}})$	Yes	Yes	Yes
	Restarted SAVE <sup>+</sup> -BOB (Ours)	$\tilde{O}(d^{\frac{4}{5}} V_K^{\frac{2}{5}} B_K^{\frac{1}{5}} K^{\frac{2}{5}} + d^{\frac{2}{3}} B_K^{\frac{1}{3}} K^{\frac{2}{3}} + d^{\frac{1}{5}} K^{\frac{7}{10}})$	Yes	Yes	No
	Lower Bound (Ours)	$\tilde{\Omega}(d^{2/3} B_K^{1/3} V_K^{1/3} K^{1/3} \wedge V_K + \sqrt{B_K K})$	Yes	Yes	-
MAB	Rerun-UCB-V (Wei et al., 2016)	$\tilde{O}( \mathcal{A} ^{\frac{2}{3}} B_K^{\frac{1}{3}} V_K^{\frac{1}{3}} K^{\frac{1}{3}} +  \mathcal{A} ^{\frac{1}{2}} B_K^{\frac{1}{2}} K^{\frac{1}{2}})$	Yes	No	Yes
	Lower Bound (Wei et al., 2016)	$\tilde{\Omega}(B_K^{\frac{1}{3}} V_K^{\frac{1}{3}} K^{\frac{1}{3}} + B_K^{\frac{1}{2}} K^{\frac{1}{2}})$	Yes	No	-

Table 1: Comparison of non-stationary bandits in terms of regret guarantee.  $K$  is the total rounds,  $d$  is the problem dimension for linear bandits,  $B_K$  is the *total variation budget* defined in Section 3 (for the MAB setting,  $B_K = \sum_{k=1}^K \|\mu_k - \mu_{k+1}\|_\infty$ , where  $\mu_k$  is the mean of the reward distribution at round  $k$ ),  $V_K$  is the *total variance* defined in Section 3,  $|\mathcal{A}|$  is the number of arms for MAB.

## 2 Related Work

### 2.1 Non-stationary (Linear) Bandits

There have been a series of works about non-stationary bandits Auer et al. (2002); Garivier and Moulines (2011); Besbes et al. (2014b); Wei et al. (2016); Cheung et al. (2019); Russac et al. (2019); Auer et al. (2019); Chen et al. (2019); Russac et al. (2020); Zhao et al. (2020b); Kim and Tewari (2020b); Wei and Luo (2021); Russac et al. (2021); Chen et al. (2021); Deng et al. (2022); Suk and Kpotufe (2022); Liu et al. (2023); Abbasi-Yadkori et al. (2023); Clerici et al. (2023).

In non-stationary linear bandits, the unknown feature vector  $\theta_k$  can be dynamically and adversarially adjusted, with the total change upper bounded by

the *total variation budget*  $B_K$  over  $K$  rounds, *i.e.*,  $\sum_{k=1}^{K-1} \|\theta_{k+1} - \theta_k\|_2 \leq B_K$ . To tackle this problem, some works proposed forgetting strategies such as sliding window, restart, and weighted regression (Cheung et al., 2019; Russac et al., 2019; Zhao et al., 2020b). Kim and Tewari (2020b) also introduced the randomized exploration with weighting strategy. The regret upper bounds in these works are all of  $\tilde{O}(B_K^{\frac{1}{4}} K^{\frac{3}{4}})$ . A recent work by Wei and Luo (2021) proposed the MASTER-OFUL algorithm based on a black-box approach, which can achieve a regret bound of  $\tilde{O}(B_K^{\frac{1}{3}} K^{\frac{2}{3}})$  in the case where the arm set is fixed over  $K$  rounds. To the best of our knowledge, none of the existing works consider how to utilize the variance information to improve the regret bound in the case with time-dependent variances. The only exception of utilizing

the variance information in the non-stationary bandit setting is Wei et al. (2016), which proposed the Rerun-UCB-V algorithm for the non-stationary MAB setting with a regret dependent on the action set size  $|\mathcal{A}|$ . To compare with, the regret upper bounds of our algorithms are independent of the action set size, thus our algorithms are more efficient for the case where the number of actions is large.

## 2.2 Linear Bandits with Heteroscedastic Noises

Some recent works study the heteroscedastic linear bandit problem, where the noise distribution is assumed to vary over time. Kirschner and Krause (2018) first proposed the linear bandit model with heteroscedastic noise. In this model, the noise at round  $k \in [K]$  is assumed to be  $\sigma_k$ -sub-Gaussian. Some follow-up works relaxed the  $\sigma_k$ -sub-Gaussian assumption by assuming the noise at the  $k$ -th round to be of variance  $\sigma_k^2$  (Zhou et al., 2021; Zhang et al., 2021; Kim et al., 2022; Zhou and Gu, 2022; Dai et al., 2022; Zhao et al., 2023). Specifically, Zhou et al. (2021) and Zhou and Gu (2022) considered the case where  $\sigma_k$  is observed by the learner after the  $k$ -th round. Zhang et al. (2021) and Kim et al. (2022) proposed statistically efficient but computationally inefficient algorithms for the unknown-variance case. A recent work by Zhao et al. (2023) proposed an algorithm that achieves both statistical and computational efficiency in the unknown-variance setting. Dai et al. (2022) also considered a specific heteroscedastic linear bandit problem where the linear model is sparse.

## 3 Problem Setting

We consider a heteroscedastic variant of the classic non-stationary linear contextual bandit problem. Let  $K$  be the total number of rounds. At each round  $k \in [K]$ , the learner interacts with the environment as follows: (1) the environment generates an arbitrary arm set  $\mathcal{D}_k \subseteq \mathbb{R}^d$  where each element represents a feasible arm for the learner to choose, and also generates an *unknown* feature vector  $\theta_k$ ; (2) the learner observes  $\mathcal{D}_k$  and selects  $\mathbf{a}_k \in \mathcal{D}_k$ ; (3) the environment generates the stochastic noise  $\epsilon_k$  and reveals the stochastic reward  $r_k = \langle \theta_k, \mathbf{a}_k \rangle + \epsilon_k$  to the learner. We assume that for all  $k \geq 1$  and all  $\mathbf{a} \in \mathcal{D}_k$ ,  $\langle \mathbf{a}, \theta_k \rangle \in [-1, 1]$ ,  $\|\theta_k\|_2 \leq B$ ,  $\|\mathbf{a}\|_2 \leq A$ .

Following Zhou et al. (2021); Zhao et al. (2023), we assume the following condition on the random noise  $\epsilon_k$  at each round  $k$ :

$$\begin{aligned} \mathbb{P}(|\epsilon_k| \leq R) &= 1, \quad \mathbb{E}[\epsilon_k | \mathbf{a}_{1:k}, \epsilon_{1:k-1}] = 0, \\ \mathbb{E}[\epsilon_k^2 | \mathbf{a}_{1:k}, \epsilon_{1:k-1}] &\leq \sigma_k^2. \end{aligned} \quad (2)$$

Following Cheung et al. (2018, 2019); Russac et al. (2019); Zhao et al. (2020b), we assume the summation of  $\ell_2$  differences of consecutive  $\theta_k$ 's is upper bounded by the *total variation budget*  $B_K$ , i.e.,  $\sum_{k=1}^{K-1} \|\theta_{k+1} - \theta_k\|_2 \leq B_K$ , where the  $\theta_k$ 's can be adversarially chosen by an oblivious adversary. We also assume that the *total variance* is upper bounded by  $V_K$ , which is  $\sum_{k=1}^K \sigma_k^2 \leq V_K$ . The goal of the agent is to minimize the *dynamic regret* defined as follows:  $\text{Regret}(K) = \sum_{k \in [K]} (\langle \mathbf{a}_k^*, \theta_k \rangle - \langle \mathbf{a}_k, \theta_k \rangle)$ , where  $\mathbf{a}_k^* = \arg\max_{\mathbf{a} \in \mathcal{D}_k} \langle \mathbf{a}, \theta_k \rangle$  is the optimal arm at round  $k$  with the highest expected reward.

## 4 Lower Bound

In this section, we establish a novel variance-dependent regret lower bound for non-stationary linear bandits, which reveals new insights into the problem structure.

**Theorem 4.1.** *Given  $K > 0$ . For any bandit algorithm there exists  $\theta_1, \dots, \theta_K$  satisfying the problem setting denoted in Section 3, such that*

$$\begin{aligned} \text{Regret}(K) &\geq \Omega(\min\{d^{2/3} B_K^{1/3} V_K^{1/3} K^{1/3}, V_K\} + \sqrt{B_K K}). \end{aligned}$$

*Proof.* See Appendix C.  $\square$

*Remark 4.2.* Note that Cheung et al. (2019) proposed a lower bound of  $\Omega(d^{2/3} B_K^{1/3} K^{2/3})$  for general non-stationary linear bandits. However, their result applies only to cases without the variance restriction  $V_K$ , making it inapplicable to our setting.

Theorem 4.1 represents the first variance-dependent regret lower bound specifically tailored for non-stationary linear bandits. The bound highlights the inherent complexity of balancing variance and non-stationarity, offering a foundation for future work aimed at designing algorithms with matching upper bounds. Notably, our result improves the existing variance-dependent lower bound  $\Omega(B_K^{1/3} V_K^{1/3} K^{1/3} + B_K^{1/2} K^{1/2})$  (Wei et al., 2016) by a factor of  $d^{2/3}$  for the linear bandits setting.

## 5 Non-stationary Linear Contextual Bandit with Known Variance

In this section, we introduce our Algorithm 1 under the setting where the variance  $\sigma_k^2$  at  $k$ -th iteration is known to the agent in prior. We start from WeightedOFUL<sup>+</sup> (Zhou and Gu, 2022), an *weighted ridge regression*-based algorithm for heteroscedastic linear bandits under the stationary reward assumption. For our non-stationary linear bandit setting where  $\theta_k$  is changing over the round  $k$ , WeightedOFUL<sup>+</sup> aims to build an

**Algorithm 1** Restarted-WeightedOFUL<sup>+</sup>

**Require:** Regularization parameter  $\lambda > 0$ ;  $B$ , an upper bound on the  $\ell_2$ -norm of  $\theta_k$  for all  $k \in [K]$ ; confidence radius  $\hat{\beta}_k$ , variance parameters  $\alpha, \gamma$ ; restart window size  $w$ .

- 1:  $\hat{\Sigma}_1 \leftarrow \lambda \mathbf{I}$ ,  $\hat{\mathbf{b}}_1 \leftarrow \mathbf{0}$ ,  $\hat{\theta}_1 \leftarrow \mathbf{0}$ ,  $\hat{\beta}_1 = \sqrt{\lambda}B$
- 2: **for**  $k = 1, \dots, K$  **do**
- 3:   **if**  $k \% w == 0$  **then**
- 4:      $\hat{\Sigma}_k \leftarrow \lambda \mathbf{I}$ ,  $\hat{\mathbf{b}}_k \leftarrow \mathbf{0}$ ,  $\hat{\theta}_k \leftarrow \mathbf{0}$ ,  $\hat{\beta}_k = \sqrt{\lambda}B$
- 5:   **end if**
- 6:   Observe  $\mathcal{D}_k$  and choose  $\mathbf{a}_k \leftarrow \arg\max_{\mathbf{a} \in \mathcal{D}_k} \langle \mathbf{a}, \hat{\theta}_k \rangle + \hat{\beta}_k \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}}$
- 7:   Observe  $(r_k, \sigma_k)$ , set  $\bar{\sigma}_k$  as
 
$$\bar{\sigma}_k \leftarrow \max\{\sigma_k, \alpha, \gamma \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}}^{1/2}\} \quad (3)$$
- 8:    $\hat{\Sigma}_{k+1} \leftarrow \hat{\Sigma}_k + \mathbf{a}_k \mathbf{a}_k^\top / \bar{\sigma}_k^2$ ,  $\hat{\mathbf{b}}_{k+1} \leftarrow \hat{\mathbf{b}}_k + r_k \mathbf{a}_k / \bar{\sigma}_k^2$ ,  
 $\hat{\theta}_{k+1} \leftarrow \hat{\Sigma}_{k+1}^{-1} \hat{\mathbf{b}}_{k+1}$
- 9: **end for**

$\hat{\theta}_k$  which estimates the feature vector  $\theta_k$  by using the solution to the following regression problem:

$$\hat{\theta}_k \leftarrow \arg \min_{\theta} \sum_{t=1}^{k-1} \bar{\sigma}_t^{-2} (\langle \theta, \mathbf{a}_t \rangle - r_t)^2 + \lambda \|\theta\|_2^2, \quad (4)$$

where the weight is defined as in (3). After obtaining  $\hat{\theta}_k$ , WeightedOFUL<sup>+</sup> chooses arm  $\mathbf{a}_k$  by maximizing the upper confidence bound (UCB) of  $\langle \mathbf{a}, \hat{\theta} \rangle$ , with an exploration bonus  $\hat{\beta}_k \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}}$ , where  $\hat{\Sigma}_k$  is the covariance matrix over  $\mathbf{a}_k$ . The weight  $\bar{\sigma}_k^2$  is introduced to balance the different past examples based on their reward variance  $\sigma_k^2$ , and such a strategy has been proved as a state-of-the-art algorithm for the stationary heteroscedastic linear bandits (Zhou and Gu, 2022). However, the non-stationary nature of our setting prevents us from directly using  $\hat{\theta}_k$  defined in (4) as an estimate to  $\theta$ . Therefore, inspired by the *restarting* strategy which has been adopted by previous algorithms for non-stationary linear bandits (Zhao et al., 2020b), we propose Restarted-WeightedOFUL<sup>+</sup>, which periodically restarts itself and runs WeightedOFUL<sup>+</sup> as its submodule. The restart window size is set as  $w$ , which is used to balance the nonstationarity and the total regret and will be fine-tuned in the next steps. Combined with the restart window size  $w$ , we set  $\{\hat{\beta}_k\}_{k \geq 1}$  to

$$\begin{aligned} \hat{\beta}_k = & 12 \sqrt{d \log(1 + \frac{(k \% w) A^2}{\alpha^2 d \lambda}) \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{(k \% w)^2}{\delta}))} \\ & + 30 \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{(k \% w)^2}{\delta})) \frac{R}{\gamma^2} + \sqrt{\lambda} B. \end{aligned} \quad (5)$$

We now propose the theoretical guarantee for Algorithm 1. The following key lemma shows how nonstationarity affects our estimation of the reward of each arm.

**Lemma 5.1.** *Let  $0 < \delta < 1$ . Then with probability at least  $1 - \delta$ , for any action  $\mathbf{a} \in \mathbb{R}^d$ , we have*

$$\begin{aligned} |\mathbf{a}^\top (\hat{\theta}_k - \theta_k)| \leq & \underbrace{\frac{A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{t=w \cdot \lfloor k/w \rfloor + 1}^{k-1} \|\theta_t - \theta_{t+1}\|_2}_{\text{Drifting term}} \\ & + \underbrace{\hat{\beta}_k \|\mathbf{a}\|_{\hat{\Sigma}_k^{-1}}}_{\text{Stochastic term}}. \end{aligned}$$

*Proof.* See Appendix D for the full proof.  $\square$

Here we provide a proof sketch of Lemma 5.1 to show the technical challenge we need to overcome. Without loss of generality, we prove the lemma for  $k \in [1, w]$ . We have

$$\begin{aligned} |\mathbf{a}^\top (\hat{\theta}_k - \theta_k)| \leq & \left| \mathbf{a}^\top \hat{\Sigma}_k^{-1} \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} (\theta_t - \theta_k) \right| \\ & + \|\mathbf{a}\|_{\hat{\Sigma}_k^{-1}} \left\| \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \epsilon_t}{\bar{\sigma}_t^2} \right\|_{\hat{\Sigma}_k^{-1}} + \sqrt{\lambda} B \|\mathbf{a}\|_{\hat{\Sigma}_k^{-1}}, \end{aligned} \quad (6)$$

For the first term, it gets involved by the nonstationarity of  $\theta_k$ . By rearranging the summation orders and several calculation steps, we have

$$\begin{aligned} \left| \mathbf{a}^\top \hat{\Sigma}_k^{-1} \sum_{t=1}^k \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} (\theta_t - \theta_k) \right| & \leq \sum_{t=1}^{k-1} |\mathbf{a}^\top \hat{\Sigma}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t}| \cdot \left\| \frac{\mathbf{a}_t}{\bar{\sigma}_t} \right\|_2 \\ & \cdot \left\| \sum_{s=t}^{k-1} (\theta_s - \theta_{s+1}) \right\|_2 \leq \frac{A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{s=1}^{k-1} \|\theta_s - \theta_{s+1}\|_2, \end{aligned}$$

We would like to highlight the subtleties in both our algorithm design and analysis to get the desired improvement. First, from here, we can see the necessity of introducing  $\alpha$  in the design of  $\bar{\sigma}_k$  in Eq.(3), which makes it possible to upper bound  $\bar{\sigma}_k^{-1}$  and get a tunable  $\alpha$  in the drifting term, which can subsequently be used to optimize the regret bound. Second, we show that it is essential to split the term  $\bar{\sigma}_t^{-2}$  as how we did. Only by doing that can we bound the  $\sum_{t=1}^s \frac{\mathbf{a}_t}{\bar{\sigma}_t} \hat{\Sigma}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t}$  term by  $d$  with the elliptical potential lemma. Otherwise, we can get a  $1/\alpha^2$  term rather than the  $A/\alpha$  term, which will hurt the final regret bound. For the second term in Eq.(6), a vanilla way to control it is adopting a self-normalized concentration inequality from (Abbasi-Yadkori et al., 2011). However, it can not utilize variance information, but just the magnitude of the noise, which fails to get a tight bound with the variance information. Inspired by Zhou and Gu (2022); Zhou et al. (2021); Zhao et al. (2023), we adapt a variance-adaptive concentration inequality in Theorem H.1 to get a tighter bound. Similar arguments also hold for the proof of Theorem 6.1 for the unknown variance case. We refer to Appendix D for the full proof.



Lemma 5.1 suggests that under the non-stationary setting, the difference between the true expected reward and our estimated reward will be upper bounded by two separate terms. The first drifting term characterizes the error caused by the non-stationary environment, and the second stochastic term characterizes the error caused by the estimation of the stochastic environment. Note that similar bound has also been discovered in Touati and Vincent (2020). We want to emphasize that our bound differs from existing ones in 1) an additional variance parameter  $\alpha$  in the drifting term, and 2) a weighted covariance matrix  $\hat{\Sigma}$  rather than a vanilla covariance matrix.

Next we present our main theorem.

**Theorem 5.2.** *Let  $0 < \delta < 1$ . By treating  $A, \lambda, B, R$  as constants and setting  $\gamma^2 = R/\sqrt{d}$ , with probability at least  $1 - \delta$ , the regret of Restarted-WeightedOFUL<sup>+</sup> is bounded by*

$$\text{Regret}(K) = \tilde{O}(B_K w^{3/2} d^{1/2} \alpha^{-1} + dK\alpha/\sqrt{w} + d\sqrt{KV_K/w} + dK/w). \quad (7)$$

*Proof.* See Appendix E.  $\square$

**Remark 5.3.** For the stationary linear bandit case where  $B_K = 0$ , we can set the restart window size  $w = K$  and the variance parameter  $\alpha = 1/\sqrt{K}$ , then we obtain an  $\tilde{O}(d\sqrt{V_K} + d)$  regret for Algorithm 1, which is identical to the one in Zhou and Gu (2022).

Next, we aim to select parameters  $\alpha$  and  $w$  in order to optimize (7).

**Corollary 5.4.** *Assume that  $B_K, V_K \in [\Omega(1), O(K)]$ . Then by selecting*

$$\begin{aligned} w &= d^{1/4} \sqrt{V_K/B_K}, & dV_K^6 &\geq K^4 B_K^2, \\ w &= d^{1/6} (K/B_K)^{1/3} & \text{otherwise.} \end{aligned}$$

and  $\alpha = d^{-1/4} B_K^{1/2} w K^{-1/2}$ , the regret is in the order

$$\text{Regret}(K) = \tilde{O}(d^{7/8} (B_K V_K)^{1/4} \sqrt{K} + d^{5/6} B_K^{1/3} K^{2/3}). \quad (8)$$

**Remark 5.5.** We compare the regret of Algo.1 in Corollary 5.4 with previous results in the special cases below.

- In the worst case where  $V_K = O(K)$ , our result becomes  $\tilde{O}(d^{7/8} B_K^{1/4} K^{3/4})$ , matching the state-of-the-art results for restarting and sliding window strategies Cheung et al. (2018); Zhao et al. (2020b).
- In the case where the *total variance* is small, i.e.,  $V_K = \tilde{O}(1)$ , assuming that  $K^4 > d$ , our result becomes  $\tilde{O}(d^{5/6} B_K^{1/3} K^{2/3})$ , better than all the previous results Cheung et al. (2018); Zhao et al. (2020b); Wang et al. (2023); Wei and Luo (2021).

**Remark 5.6.** Wei et al. (2016) has studied non-stationary MAB with dynamic variance. With the knowledge of  $V_K$  and  $B_K$ , Wei et al. (2016) proposed a restart-based Rerun-UCB-V algorithm with a  $\tilde{O}(|\mathcal{A}|^{2/3} B_K^{1/3} V_K^{1/3} K^{1/3} + |\mathcal{A}|^{1/2} B_K^{1/2} K^{1/2})$  regret, where  $\mathcal{A}$  is the action set. Reduced to the MAB setting, our Restarted-WeightedOFUL<sup>+</sup> achieves an  $\tilde{O}(|\mathcal{A}|^{7/8} (B_K V_K)^{1/4} \sqrt{K} + |\mathcal{A}|^{5/6} B_K^{1/3} K^{2/3})$  regret, which is worse than Wei et al. (2016). We claim that this is due to the generality of the linear bandits, which brings us a looser bound to the drifting term in Lemma 5.1. When restricting to the MAB setting, our drifting term enjoys a tighter bound, which could further tighten our final regret. To develop an algorithm achieving the same regret as Wei et al. (2016) is beyond the scope of this work.

**Remark 5.7.** Wei et al. (2016) has established a lower bound  $\tilde{\Omega}(B_K^{1/3} V_K^{1/3} K^{1/3} + B_K^{1/2} K^{1/2})$  for MAB with total variance  $V_K$  and total variation budget  $B_K$ . There still exist gaps between our regret and their lower bound regarding the dependence of  $K, V_K, B_K$ , and we leave to fix the gaps as future work.

## 6 Non-stationary Linear Contextual Bandit with Unknown Variance and Total Variation Budget

By Theorem 5.2, we know that Algorithm 1 is able to utilize the total variance  $V_K$  and obtain a better regret result compared with existing algorithms which do not utilize  $V_K$ . However, the success of Algorithm 1 depends on the knowledge of the per-round variance  $\sigma_k$ , and it also depends on a good selection of restart window size  $w$ , whose optimal selection depends on both  $V_K$  and  $B_K$ . In this section, we aim to relax these two requirements with still better regret results.

### 6.1 Unknown Per-round Variance, Known $V_K$ and $B_K$

We first aim to relax the requirement that each  $\sigma_k^2$  is known to the agent at the beginning of  $k$ -th round. We follow the SAVE algorithm (Zhao et al., 2023) which introduces a multi-layer structure (Chu et al., 2011; He et al., 2021) to deal with unknown  $\sigma_k^2$ . In detail, SAVE maintains multiple estimates to the current feature vector  $\theta_k$ , which we denote them as  $\hat{\theta}_{k,1}, \dots, \hat{\theta}_{k,L}$  in line 2. Each  $\hat{\theta}_{k,\ell}$  is calculated based on a subset  $\hat{\Psi}_{k,\ell} \subseteq [k-1]$  of samples  $\{(\mathbf{a}_t, r_t)\}$ . The rule that whether to add the current  $k$  to some  $\hat{\Psi}_{k,\ell}$  is based on the uncertainty of  $\mathbf{a}_k$  with the sample set  $\{(\mathbf{a}_t, r_t)\}_{t \in \hat{\Psi}_{k,\ell}}$ . As long as  $\mathbf{a}_k$  is too uncertain w.r.t. some level  $\ell_k$  (line 2), we add  $k$  to  $\hat{\Psi}_{k,\ell}$  and update the estimate  $\hat{\theta}_{k,\ell_k}$  accordingly (line 2). Each  $\hat{\theta}_{k,\ell_k}$  is calculated as the

---

**Algorithm 2** Restarted SAVE<sup>+</sup>


---

**Require:**  $\alpha > 0$ ; the upper bound on the  $\ell_2$ -norm of  $\mathbf{a}$  in  $\mathcal{D}_k (k \geq 1)$ , i.e.,  $A$ ; the upper bound on the  $\ell_2$ -norm of  $\boldsymbol{\theta}_k (k \geq 1)$ , i.e.,  $B$ ; restart window size  $w$ .

- 1: Initialize  $L \leftarrow \lceil \log_2(1/\alpha) \rceil$ .
- 2: Initialize the estimators for all layers:  $\hat{\Sigma}_{1,\ell} \leftarrow 2^{-2\ell} \cdot \mathbf{I}$ ,  $\hat{\mathbf{b}}_{1,\ell} \leftarrow \mathbf{0}$ ,  $\hat{\boldsymbol{\theta}}_{1,\ell} \leftarrow \mathbf{0}$ ,  $\hat{\beta}_{1,\ell} \leftarrow 2^{-\ell+1}$ ,  $\hat{\Psi}_{1,\ell} \leftarrow \emptyset$  for all  $\ell \in [L]$ .
- 3: **for**  $k = 1, \dots, K$  **do**
- 4:   **if**  $k \% w == 0$  **then**
- 5:     Set  $\hat{\Sigma}_{k,\ell} \leftarrow 2^{-2\ell} \cdot \mathbf{I}$ ,  $\hat{\mathbf{b}}_{k,\ell} \leftarrow \mathbf{0}$ ,  $\hat{\boldsymbol{\theta}}_{k,\ell} \leftarrow \mathbf{0}$ ,  $\hat{\beta}_{k,\ell} \leftarrow 2^{-\ell+1}$ ,  $\hat{\Psi}_{k,\ell} \leftarrow \emptyset$  for all  $\ell \in [L]$ .
- 6:   **end if**
- 7:   Observe  $\mathcal{D}_k$ , choose  $\mathbf{a}_k \leftarrow \operatorname{argmax}_{\mathbf{a} \in \mathcal{D}_k} \min_{\ell \in [L]} \langle \mathbf{a}, \hat{\boldsymbol{\theta}}_{k,\ell} \rangle + \hat{\beta}_{k,\ell} \|\mathbf{a}\|_{\hat{\Sigma}_{k,\ell}^{-1}}$  and observe  $r_k$ .
- 8:   Set  $\ell_k \leftarrow L + 1$
- 9:   Let  $\mathcal{L}_k \leftarrow \{\ell \in [L] : \|\mathbf{a}_k\|_{\hat{\Sigma}_{k,\ell}^{-1}} \geq 2^{-\ell}\}$ , set  $\ell_k \leftarrow \min(\mathcal{L}_k)$  if  $\mathcal{L}_k \neq \emptyset$
- 10:    $\hat{\Psi}_{k,\ell_k} \leftarrow \hat{\Psi}_{k,\ell_k} \cup \{k\}$
- 11:   **if**  $\mathcal{L}_k \neq \emptyset$  **then**
- 12:     Set  $w_k \leftarrow \frac{2^{-\ell_k}}{\|\mathbf{a}_k\|_{\hat{\Sigma}_{k,\ell_k}^{-1}}}$  and update

$$\hat{\Sigma}_{k+1,\ell_k} \leftarrow \hat{\Sigma}_{k,\ell_k} + w_k^2 \mathbf{a}_k \mathbf{a}_k^\top, \hat{\mathbf{b}}_{k+1,\ell} \leftarrow \hat{\mathbf{b}}_{k,\ell} + w_k^2 \cdot r_k \mathbf{a}_k, \hat{\boldsymbol{\theta}}_{k+1,\ell_k} \leftarrow \hat{\Sigma}_{k+1,\ell_k}^{-1} \hat{\mathbf{b}}_{k+1,\ell_k}.$$

- 13:   Compute the adaptive confidence radius  $\hat{\beta}_{k+1,\ell}$  for the next round according to (9).
  - 14:   **end if**
  - 15:   For  $\ell \neq \ell_k$  let  $\hat{\Sigma}_{k+1,\ell} \leftarrow \hat{\Sigma}_{k,\ell}$ ,  $\hat{\mathbf{b}}_{k+1,\ell} \leftarrow \hat{\mathbf{b}}_{k,\ell}$ ,  $\hat{\boldsymbol{\theta}}_{k+1,\ell} \leftarrow \hat{\boldsymbol{\theta}}_{k,\ell}$ ,  $\hat{\beta}_{k+1,\ell} \leftarrow \hat{\beta}_{k,\ell}$ .
  - 16: **end for**
- 

solution of a weighted regression problem, where the weight  $w_k$  is selected as the inverse of the uncertainty of the arm  $\mathbf{a}_k$  w.r.t. the samples in the  $\ell$ -th layer. Maintaining  $L$  different  $\hat{\boldsymbol{\theta}}_{k,\ell}, \ell \in [L]$ , Algorithm 2 then calculates  $L$  number of UCB for each arm  $\mathbf{a}$  w.r.t.  $L$  different  $\hat{\boldsymbol{\theta}}_{k,\ell}$ , and selects the arm which maximizes the minimization of  $L$  UCBs (line 2). It has been shown in Zhao et al. (2023) that such a multilayer structure is able to utilize the  $V_K$  information without knowing the per-round variance  $\sigma_k^2$ . Similar to Algorithm 1, in order to deal with the nonstationarity issue, we introduce a restarting scheme that Algorithm 2 restarts itself by a restart window size  $w$  (line 2).

Next we show the theoretical guarantee of Algorithm 2. We call the restart time rounds *grids* and denote them by  $g_1, g_2, \dots, g_{\lceil \frac{K}{w} \rceil - 1}$ , where  $g_i \% w = 0$  for all  $i \in [\lceil \frac{K}{w} \rceil - 1]$ . Let  $i_k$  be the grid index of time round  $k$ , i.e.,  $g_{i_k} \leq k < g_{i_k+1}$ . We denote  $\hat{\Psi}_{k,\ell} := \{t : t \in [g_{i_k}, k-1], \ell_t = \ell\}$ . We define the confidence radius  $\hat{\beta}_{k,\ell}$  at round  $k$  and layer  $\ell$  as

$$\hat{\beta}_{k,\ell} := 16 \cdot 2^{-\ell} \sqrt{\left(8\hat{\text{Var}}_{k,\ell} + 6R^2 \log\left(\frac{4(w+1)^2 L}{\delta}\right) + 2^{-2\ell+4}\right)} \times \sqrt{\log\left(\frac{4w^2 L}{\delta}\right) + 6 \cdot 2^{-\ell} R \log\left(\frac{4w^2 L}{\delta}\right) + 2^{-\ell} B}, \quad (9)$$

where we set  $\hat{\text{Var}}_{k,\ell}$  as  $\sum_{i \in \hat{\Psi}_{k,\ell}} w_i^2 (r_i - \langle \hat{\boldsymbol{\theta}}_{k,\ell}, \mathbf{a}_i \rangle)^2$ , if

$2^\ell \geq 64 \sqrt{\log\left(\frac{4(w+1)^2 L}{\delta}\right)}$ , or  $R^2 |\hat{\Psi}_{k,\ell}|$  for the remaining cases.

Note that our selection of the confidence radius  $\hat{\beta}_{k,\ell}$  only depends on  $\hat{\text{Var}}_{k,\ell}$ , which serves as an estimate of the total variance of samples at  $\ell$ -th layer without knowing  $\sigma_k^2$ .

We build the theoretical guarantee of Algorithm 2 as follows.

**Theorem 6.1.** *Let  $0 < \delta < 1$ . Define  $\{\beta_{k,\ell}\}_{k \geq 1, \ell \in [L]}$  as in (9), regarding  $A, R$  as constants, we have*

$$\text{Regret}(K) = \tilde{O}(\sqrt{dw^{1.5}} B_K / \alpha + \alpha^2 (K^2 + \sqrt{wKV_K}) + d\sqrt{KV_K/w} + dK/w).$$

*Proof.* See Appendix F for the full proof.  $\square$

**Remark 6.2.** Like Remark 5.3, we consider the case where  $B_K = 0$ . We set  $w = K$  and  $\alpha^2 = 1/K\sqrt{V_K}$ , then we obtain a regret  $\tilde{O}(d\sqrt{V_K} + d)$ , which matches the regret of the SAVE algorithm in Zhao et al. (2023).

**Corollary 6.3.** *Assume that  $B_K, V_K \in [\Omega(1), O(K)]$ , then by selecting*

$$\begin{aligned} w &= d^{1/3} (K/B_K)^{1/3}, & K^2 &\geq V_K^3 d / B_K, \\ w &= d^{2/5} (KV_K)^{1/5} / B_K^{2/5}, & &\text{otherwise.} \end{aligned}$$

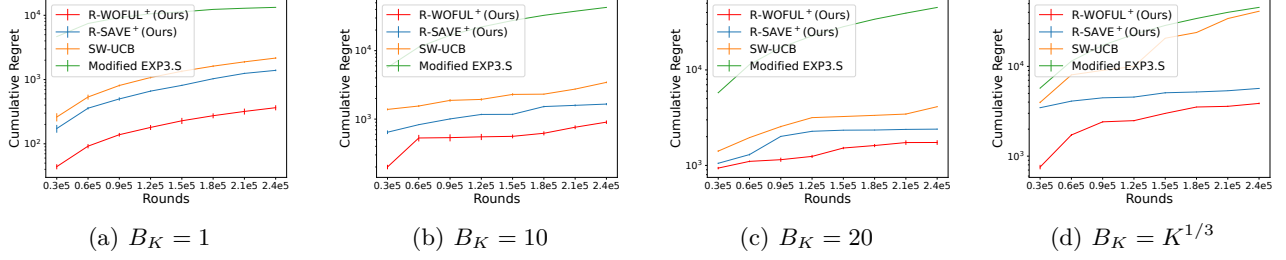


Figure 1: The regret of Restarted-WeightedOFUL<sup>+</sup>, Restarted SAVE<sup>+</sup>, SW-UCB and Modified EXP3.S under different total rounds.

and  $\alpha = d^{1/6} \sqrt{w} B_K^{1/3} / (K^{1/3} + (V_K K w)^{1/6})$ , we have

$$\text{Regret}(K) = \tilde{O}(d^{4/5} V_K^{2/5} B_K^{1/5} K^{2/5} + d^{2/3} B_K^{1/3} K^{2/3}).$$

*Remark 6.4.* We discuss the regret of Algo.2 in Corollary 6.3 in the following special cases. In the case where the *total variance* is small, *i.e.*,  $V_K = \tilde{O}(1)$ , assuming that  $K^2 > d$ , our result becomes  $\tilde{O}(d^{2/3} B_K^{1/3} K^{2/3})$ , better than all the previous results [Cheung et al. \(2018\)](#); [Zhao et al. \(2020b\)](#); [Wang et al. \(2023\)](#); [Wei and Luo \(2021\)](#). In the worst case where  $V_K = O(K)$ , our result becomes  $\tilde{O}(d^{4/5} B_K^{1/5} K^{4/5})$ .

**Unknown Per-round Variance, Unknown  $V_K$  and  $B_K$**  In Corollary 6.3, we need to know the *total variance*  $V_K$  and *total variation budget*  $B_K$  to select the optimal  $w$  and  $\alpha$ . To deal with the more general case where  $V_K$  and  $B_K$  are unknown, we can employ the *Bandits-over-Bandits* (BOB) mechanism ([Cheung et al. \(2019\)](#); [Wang et al. \(2023\)](#); [Zhao et al. \(2020b\)](#)). We name the Restarted SAVE<sup>+</sup> algorithm with BOB mechanism as “Restarted SAVE<sup>+</sup>-BOB”. Due to the space limit, we put the algorithm design, descriptions, and theoretical analysis of Restarted SAVE<sup>+</sup>-BOB (Algo.3) in Appendix A.

## 7 Experiments

To validate the effectiveness of our methods, we conduct a series of experiments on the synthetic data.

**Problem Setting and Baselines** Following the experimental set up in [Cheung et al. \(2019\)](#), we consider the 2-armed bandits setting, where the action set  $\mathcal{D}_k = \{(1, 0), (0, 1)\}$ , and

$$\theta_k = \begin{pmatrix} 0.5 + \frac{3}{10} \sin(5B_K \pi k / K) \\ 0.5 + \frac{3}{10} \sin(\pi + 5B_K \pi k / K) \end{pmatrix}.$$

It is easy to see that the total variation budget can be bounded as  $B_K$ . At each round  $k$ , the  $\epsilon_k$  satisfies the following distribution:

$$\epsilon_k \sim \text{Bernoulli}(0.5/k) - 0.5/k.$$

We can verify that under such a distribution for  $\epsilon_k$ , the variance of the reward distribution at  $k$ -th round is  $(1 - 0.5/k) \cdot 0.5/k$ , and the total variance  $V_K \sim \log K$ .

We compare the proposed Restarted-WeightedOFUL<sup>+</sup> and Restarted SAVE<sup>+</sup> with SW-UCB [Cheung et al. \(2019\)](#) and Modified EXP3.S [Besbes et al. \(2014a\)](#). We leave the detailed setup for the baselines in Appendix B.

**Result** We plot the results in Figure.1, where all the empirical results are averaged over ten independent trials and the error bar is the standard error divided by  $\sqrt{10}$ . The results are consistent with our theoretical findings. It is evident that our algorithms significantly outperform both SW-UCB and Modified EXP3.S. Among our proposed algorithms, Restarted-WeightedOFUL<sup>+</sup> achieves the best performance. This can be attributed to the fact that it knows the variance and can make more informed decisions. Although Restarted SAVE<sup>+</sup> performed slightly worse than Restarted-WeightedOFUL<sup>+</sup>, it still outperforms the baseline algorithms, particularly when  $B_K = K^{1/3}$ . These results highlight the superiority of our methods.

## 8 Conclusion and Future Work

We study non-stationary stochastic linear bandits in this work. We establish the first variance-dependent regret lower bound for non-stationary linear bandits, which captures the interplay between variance, non-stationarity, and dimensionality in the linear bandit setting, offering new insights into the complexity of this problem. We propose Restarted-WeightedOFUL<sup>+</sup> and Restarted SAVE<sup>+</sup>, two algorithms that utilize the dynamic variance information of the dynamic reward distribution. We show that both of our algorithms are able to achieve better dynamic regret compared with best existing results ([Wei and Luo, 2021](#)) under several parameter regimes, *e.g.*, when the total variance  $V_K$  is small. Experiment results backup our theoretical claim. It is worth noting there still exist gaps between our current obtained regret and the lower bound, and



to fix such a gap leaves as our future work.

## 9 Acknowledgement

The work of John C.S. Lui was supported in part by the RGC GRF-14202923.

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- Yasin Abbasi-Yadkori, András György, and Nevena Lazić. A new look at dynamic regret for non-stationary stochastic bandits. *Journal of Machine Learning Research*, 24(288):1–37, 2023.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1): 48–77, 2002.
- Peter Auer, Pratik Gajane, and Ronald Ortner. Adaptively tracking the best bandit arm with an unknown number of distribution changes. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 138–158. PMLR, 25–28 Jun 2019. URL <https://proceedings.mlr.press/v99/auer19a.html>.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Optimal exploration-exploitation in a multi-armed-bandit problem with non-stationary rewards. *SSRN Electronic Journal*, 2014a.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Stochastic multi-armed-bandit problem with non-stationary rewards. *Advances in neural information processing systems*, 27, 2014b.
- Wei Chen, Liwei Wang, Haoyu Zhao, and Kai Zheng. Combinatorial semi-bandit in the non-stationary environment. In *Uncertainty in Artificial Intelligence*, pages 865–875. PMLR, 2021.
- Yifang Chen, Chung-Wei Lee, Haipeng Luo, and Chen-Yu Wei. A new algorithm for non-stationary contextual bandits: Efficient, optimal and parameter-free. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 696–726. PMLR, 25–28 Jun 2019. URL <https://proceedings.mlr.press/v99/chen19b.html>.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Hedging the drift: Learning to optimize under non-stationarity. *Available at SSRN 3261050*, 2018.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Learning to optimize under non-stationarity. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1079–1087. PMLR, 2019.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Reinforcement learning for non-stationary markov

- decision processes: The blessing of (more) optimism. In *International Conference on Machine Learning*, pages 1843–1854. PMLR, 2020.
- Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214. JMLR Workshop and Conference Proceedings, 2011.
- Giulia Clerici, Pierre Laforgue, and Nicolò Cesa-Bianchi. Linear bandits with memory: from rotting to rising, 2023.
- Yan Dai, Ruosong Wang, and Simon S Du. Variance-aware sparse linear bandits. *arXiv preprint arXiv:2205.13450*, 2022.
- Yuntian Deng, Xingyu Zhou, Baekjin Kim, Ambuj Tewari, Abhishek Gupta, and Ness Shroff. Weighted gaussian process bandits for non-stationary environments. In *International Conference on Artificial Intelligence and Statistics*, pages 6909–6932. PMLR, 2022.
- Louis Faury, Yoan Russac, Marc Abeille, and Clément Calauzènes. Regret bounds for generalized linear bandits under parameter drift. *arXiv preprint arXiv:2103.05750*, 2021.
- David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pages 100–118, 1975.
- Pratik Gajane, Ronald Ortner, and Peter Auer. A sliding-window algorithm for markov decision processes with arbitrarily changing rewards and transitions. *arXiv preprint arXiv:1805.10066*, 2018.
- Aurélien Garivier and Eric Moulines. On upper-confidence bound policies for switching bandit problems. In Jyrki Kivinen, Csaba Szepesvári, Esko Ukkonen, and Thomas Zeugmann, editors, *Algorithmic Learning Theory*, pages 174–188, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg. ISBN 978-3-642-24412-4.
- Jiafan He, Dongruo Zhou, and Quanquan Gu. Uniform-pac bounds for reinforcement learning with linear function approximation. *Advances in Neural Information Processing Systems*, 34:14188–14199, 2021.
- Baekjin Kim and Ambuj Tewari. Randomized exploration for non-stationary stochastic linear bandits. In *Uncertainty in Artificial Intelligence*, 2020a.
- Baekjin Kim and Ambuj Tewari. Randomized exploration for non-stationary stochastic linear bandits. In *Conference on Uncertainty in Artificial Intelligence*, pages 71–80. PMLR, 2020b.
- Yeoneung Kim, Insoon Yang, and Kwang-Sung Jun. Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture mdps. *Advances in Neural Information Processing Systems*, 35:1060–1072, 2022.
- Johannes Kirschner and Andreas Krause. Information directed sampling and bandits with heteroscedastic noise. In *Conference On Learning Theory*, pages 358–384. PMLR, 2018.
- Yueyang Liu, Benjamin Van Roy, and Kuang Xu. A definition of non-stationary bandits. *arXiv preprint arXiv:2302.12202*, 2023.
- Weichao Mao, Kaiqing Zhang, Ruihao Zhu, David Simchi-Levi, and Tamer Basar. Near-optimal model-free reinforcement learning in non-stationary episodic mdps. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning Research*, pages 7447–7458. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/mao21b.html>.
- Yoan Russac, Claire Vernade, and Olivier Cappé. Weighted linear bandits for non-stationary environments. *Advances in Neural Information Processing Systems*, 2019.
- Yoan Russac, Olivier Cappé, and Aurélien Garivier. Algorithms for non-stationary generalized linear bandits. *arXiv preprint arXiv:2003.10113*, 2020.
- Yoan Russac, Louis Faury, Olivier Cappé, and Aurélien Garivier. Self-concordant analysis of generalized linear bandits with forgetting. In *International Conference on Artificial Intelligence and Statistics*, pages 658–666. PMLR, 2021.
- Joe Suk and Samory Kpotufe. Tracking most significant arm switches in bandits. In *Conference on Learning Theory*, pages 2160–2182. PMLR, 2022.
- Ahmed Touati and Pascal Vincent. Efficient learning in non-stationary linear markov decision processes. *arXiv preprint arXiv:2010.12870*, 2020.
- Jing Wang, Peng Zhao, and Zhi-Hua Zhou. Revisiting weighted strategy for non-stationary parametric bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 7913–7942. PMLR, 2023.
- Chen-Yu Wei and Haipeng Luo. Non-stationary reinforcement learning without prior knowledge: An optimal black-box approach. In *Conference on learning theory*, pages 4300–4354. PMLR, 2021.
- Chen-Yu Wei, Yi-Te Hong, and Chi-Jen Lu. Tracking the best expert in non-stationary stochastic environments. *Advances in neural information processing systems*, 29, 2016.
- Zihan Zhang, Jiaqi Yang, Xiangyang Ji, and Simon S Du. Improved variance-aware confidence sets for linear bandits and linear mixture mdp. *Advances*

in *Neural Information Processing Systems*, 34:4342–4355, 2021.

Heyang Zhao, Jiafan He, Dongruo Zhou, Tong Zhang, and Quanquan Gu. Variance-dependent regret bounds for linear bandits and reinforcement learning: Adaptivity and computational efficiency. *arXiv preprint arXiv:2302.10371*, 2023.

Peng Zhao, Lijun Zhang, Yuan Jiang, and Zhi-Hua Zhou. A simple approach for non-stationary linear bandits. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 746–755. PMLR, 26–28 Aug 2020a.

Peng Zhao, Lijun Zhang, Yuan Jiang, and Zhi-Hua Zhou. A simple approach for non-stationary linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 746–755. PMLR, 2020b.

Dongruo Zhou and Quanquan Gu. Computationally efficient horizon-free reinforcement learning for linear mixture mdps. *Advances in neural information processing systems*, 35:36337–36349, 2022.

Dongruo Zhou, Quanquan Gu, and Csaba Szepesvari. Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*, pages 4532–4576. PMLR, 2021.

## Checklist

1. For all models and algorithms presented, check if you include:
  - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. Yes
  - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. Yes
  - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes/No/Not Applicable]
2. For any theoretical claim, check if you include:
  - (a) Statements of the full set of assumptions of all theoretical results. Yes
  - (b) Complete proofs of all theoretical results. Yes
  - (c) Clear explanations of any assumptions. Yes
3. For all figures and tables that present empirical results, check if you include:
  - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). Yes
  - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). Yes
  - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). Yes
  - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). Yes
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
  - (a) Citations of the creator If your work uses existing assets. Not Applicable
  - (b) The license information of the assets, if applicable. [Yes/No/Not Applicable]
  - (c) New assets either in the supplemental material or as a URL, if applicable. Not Applicable
  - (d) Information about consent from data providers/curators. Not Applicable
  - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. Not Applicable
5. If you used crowdsourcing or conducted research with human subjects, check if you include:

- (a) The full text of instructions given to participants and screenshots. Not Applicable
- (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. Not Applicable
- (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. Not Applicable

## A Restarted SAVE<sup>+</sup>-BOB

In this section, we provide the details of our proposed Restarted SAVE<sup>+</sup>-BOB algorithm. The Restarted SAVE<sup>+</sup>-BOB algorithm is summarized in Algo.3. We divide the  $K$  rounds into  $\lceil \frac{K}{H} \rceil$  blocks, with each block having  $H$  rounds (except the last one may have less than  $H$ ). Within each block  $i$ , we use a fixed  $(\alpha_i, w_i)$  pair to run the Restarted SAVE<sup>+</sup> algorithm. To adaptively learn the optimal  $(\alpha, w)$  pair without the knowledge of  $V_K$  and  $B_K$ , we employ an adversarial bandit algorithm (Exp3 in Auer et al. (2002)) as the meta-learner to select  $\alpha_i, w_i$  over time for  $i \in \lceil \frac{K}{H} \rceil$  blocks. Specifically, in each block, the meta learner selects a  $(\alpha, w)$  pair from the candidate pool to feed to Restarted SAVE<sup>+</sup>, and the cumulative reward received by Restarted SAVE<sup>+</sup> within the block is fed to the meta-learner as the reward feedback to select a better pair for the next block.

We set  $H$  to be  $\lceil d^{\frac{2}{5}} K^{\frac{2}{5}} \rceil$ , and set the candidate pool of  $(\alpha, w)$  pairs for the Exp3 algorithm as:

$$\mathcal{P} = \{(w, \alpha) : w \in \mathcal{W}, \alpha \in \mathcal{J}\}, \quad (10)$$

where

$$\mathcal{W} = \{w_i = d^{\frac{1}{3}} 2^{i-1} | i \in \lceil \frac{1}{3} \log_2 K \rceil + 1\} \cup \{w_i = d^{\frac{2}{5}} 2^{i-1} | i \in \lceil \frac{2}{5} \log_2 K \rceil + 1\}, \quad (11)$$

and

$$\mathcal{J} = \{\alpha_i = d^{\frac{1}{3}} 2^{-i+1} | i \in \lceil \frac{1}{3} \log_2 K \rceil + 1\} \cup \{\alpha_i = d^{\frac{11}{30}} 2^{-i+1} | i \in \lceil \frac{11}{30} \log_2 K \rceil + 1\}. \quad (12)$$

The algorithm also labels all the  $|\mathcal{P}| = (\lceil \frac{1}{3} \log_2 K \rceil + \lceil \frac{2}{5} \log_2 K \rceil + 2) \cdot (\lceil \frac{1}{3} \log_2 K \rceil + \lceil \frac{11}{30} \log_2 K \rceil + 2)$  candidate pairs of parameters in  $\mathcal{P}$ , i.e.,  $\mathcal{P} = \{(w_i, \alpha_i)\}_{i=1}^{|\mathcal{P}|}$ . The algorithm initializes  $\{s_{j,1}\}_{j=1}^{|\mathcal{P}|}$  to be  $s_{j,1} = 1, \forall j = 0, 1, \dots, |\mathcal{P}|$ , which means that at the beginning, the algorithm selects a pair from  $\mathcal{P}$  uniformly at random. At the beginning of each block  $i \in \lceil [K/H] \rceil$ , the meta-learner (Exp3) calculates the distribution  $(p_{j,i})_{j=1}^{|\mathcal{P}|}$  over the candidate set  $\mathcal{P}$  by

$$p_{j,i} = (1 - \gamma) \frac{s_{j,i}}{\sum_{u=1}^{|\mathcal{P}|} s_{u,i}} + \frac{\gamma}{|\mathcal{P}| + 1}, \quad \forall j = 1, \dots, |\mathcal{P}|, \quad (13)$$

where  $\gamma$  is defined as

$$\gamma = \min \left\{ 1, \sqrt{\frac{(|\mathcal{P}| + 1) \ln(|\mathcal{P}| + 1)}{(e - 1) \lceil K/H \rceil}} \right\}. \quad (14)$$

Then, the meta-learner draws a  $j_i$  from the distribution  $(p_{j,i})_{j=1}^{|\mathcal{P}|}$ , and sets the pair of parameters in block  $i$  to be  $(w_{j_i}, \alpha_{j_i})$ , and runs the base algorithm Algo.2 from scratch in this block with  $(w_{j_i}, \alpha_{j_i})$ , then feeds the cumulative reward in the block  $\sum_{k=(i-1)H+1}^{\min\{i \cdot H, K\}} r_k$  to the meta-learner. The meta-learner rescales  $\sum_{k=(i-1)H+1}^{\min\{i \cdot H, K\}} r_k$  to

$\frac{\sum_{k=(i-1)H+1}^{\min\{i \cdot H, K\}} r_k}{H + R \sqrt{\frac{H}{2} \log \left( K \left( \frac{K}{H} + 1 \right) \right)} + \frac{2}{3} \cdot R \log \left( K \left( \frac{K}{H} + 1 \right) \right)}$  to make it in the range  $[0, 1]$  with high probability (supported by Lemma

H.7). The meta-learner updates the parameter  $s_{j_i, i+1}$  to be

$$s_{j_i, i+1} = s_{j_i, i} \cdot \exp \left( \frac{\gamma}{(|\mathcal{P}| + 1) p_{j_i, i}} \left( \frac{1}{2} + \frac{\sum_{k=(i-1)H+1}^{\min\{i \cdot H, K\}} r_k}{H + R \sqrt{\frac{H}{2} \log \left( K \left( \frac{K}{H} + 1 \right) \right)} + \frac{2}{3} \cdot R \log \left( K \left( \frac{K}{H} + 1 \right) \right)} \right) \right), \quad (15)$$

and keep others unchanged, i.e.,  $s_{u, i+1} = s_{u, i}, \forall u \neq j_i$ . After that, the algorithm will go to the next block, and repeat the same process in block  $i + 1$ .

We have the following theorem to bound the regret of Restarted SAVE<sup>+</sup>-BOB.

**Theorem A.1.** *By using the BOB framework with Exp3 as the meta-algorithm and Restarted SAVE<sup>+</sup> as the base algorithm, with the candidate pool  $\mathcal{P}$  for Exp3 specified as in Eq.(10), Eq.(11), Eq.(12), and  $H = \lceil d^{\frac{2}{5}} K^{\frac{2}{5}} \rceil$ , then the regret of Restarted SAVE<sup>+</sup>-BOB (Algo.3) satisfies*

$$\text{Regret}(K) = \tilde{O}(d^{4/5} V_K^{2/5} B_K^{1/5} K^{2/5} + d^{2/3} B_K^{1/3} K^{2/3} + d^{2/5} K^{7/10}). \quad (16)$$



**Algorithm 3** Restarted SAVE<sup>+</sup>-BOB

**Require:** total time rounds  $K$ ; problem dimension  $d$ ; noise upper bound  $R$ ;  $\alpha > 0$ ; the upper bound on the  $\ell_2$ -norm of  $\mathbf{a}$  in  $\mathcal{D}_k (k \geq 1)$ , i.e.,  $A$ ; the upper bound on the  $\ell_2$ -norm of  $\boldsymbol{\theta}_k (k \geq 1)$ , i.e.,  $B$ .

- 1: Initialize  $H = \lceil d^{\frac{2}{5}} K^{\frac{2}{5}} \rceil$ ;  $\mathcal{P}$  as defined in Eq.(10), and index the  $|\mathcal{P}| = (\lceil \frac{1}{3} \log_2 K \rceil + \lceil \frac{2}{5} \log_2 K \rceil + 2) \cdot (\lceil \frac{1}{3} \log_2 K \rceil + \lceil \frac{11}{30} \log_2 K \rceil + 2)$  items in  $\mathcal{P}$ , i.e.,  $\mathcal{P} = \{(w_i, \alpha_i)\}_{i=1}^{|\mathcal{P}|}$ ;  $\gamma = \min \left\{ 1, \sqrt{\frac{(|\mathcal{P}|+1) \ln(|\mathcal{P}|+1)}{(e-1) \lceil K/H \rceil}} \right\}$ ;  $\{s_{j,1}\}_{j=1}^{|\mathcal{P}|}$  is set to  $s_{j,1} = 1, \forall j = 0, 1, \dots, |\mathcal{P}|$ .
- 2: **for**  $i = 1, 2, \dots, \lceil K/H \rceil$  **do**
- 3: Calculate the distribution  $(p_{j,i})_{j=1}^{|\mathcal{P}|}$  by  $p_{j,i} = (1 - \gamma) \frac{s_{j,i}}{\sum_{u=1}^{|\mathcal{P}|} s_{u,i}} + \frac{\gamma}{|\mathcal{P}|+1}, \forall j = 1, \dots, |\mathcal{P}|$ .
- 4: Set  $j_i \leftarrow j$  with probability  $p_{j,i}$ , and  $(w_i, \alpha_i) \leftarrow (w_{j_i}, \alpha_{j_i})$ .
- 5: Run Algo.2 from scratch in block  $i$  (i.e., in rounds  $k = (i-1)H + 1, \dots, \min\{i \cdot H, K\}$ ) with  $(w, \alpha) = (w_i, \alpha_i)$ .
- 6: Update  $s_{j_i,i+1} = s_{j_i,i} \cdot \exp \left( \frac{\gamma}{(|\mathcal{P}|+1)p_{j_i,i}} \left( \frac{1}{2} + \frac{\sum_{k=(i-1)H+1}^{\min\{i \cdot H, K\}} r_k}{H + R \sqrt{\frac{H}{2} \log \left( K \left( \frac{K}{H} + 1 \right) \right) + \frac{2}{3} \cdot R \log \left( K \left( \frac{K}{H} + 1 \right) \right)}} \right) \right)$ , and keep all the others unchanged, i.e.,  $s_{u,i+1} = s_{u,i}, \forall u \neq j_i$ .
- 7: **end for**

*Proof.* See Appendix G for the full proof.  $\square$

*Remark A.2.* We discuss the regret of Algo.3 in Corollary 6.3 in the following special cases. In the case where the total variance is small, i.e.,  $V_K = \tilde{O}(1)$ , assuming  $K^2 > d$ , our result becomes  $\tilde{O}(d^{2/3} B_K^{1/3} K^{2/3} + d^{1/5} K^{7/10})$ , when  $d^{14} B_K^{10} > K$ , it becomes  $\tilde{O}(d^{2/3} B_K^{1/3} K^{2/3})$ , better than all the previous results Cheung et al. (2018); Zhao et al. (2020b); Wang et al. (2023); Wei and Luo (2021). In the worst case where  $V_K = O(K)$ , our result becomes  $\tilde{O}(d^{4/5} B_K^{1/5} K^{4/5})$ .

## B Additional Experiment Setup

For Restarted-WeightedOFUL<sup>+</sup>, we set  $\lambda = 1$ ,  $\hat{\beta}_k = 10$ ,  $w = 1000$ , and we grid search the variance parameters  $\alpha$  and  $\gamma$ , both among values  $[1, 1.5, 2, 2.5, 3]$ . Finally we set  $\alpha = 1$ , and  $\gamma = 2$ . For Restarted SAVE<sup>+</sup> we set  $w = 1000$ ,  $\hat{\beta}_{k,\ell} = 2^{-\ell+1}$ , and grid search  $L$  from 1 to 10 with stepsize of 1 and finally choose  $L = 6$ . For SW-UCB, we set  $\lambda = 1$ ,  $w = 1000$ ,  $\beta_k = 10$ . The Modified EXP3.S requires two parameters  $\bar{\alpha}$  and  $\bar{\gamma}$ , and we set  $\bar{\gamma} = 0.01$  and  $\bar{\alpha} = \frac{1}{K}$ .

To test the algorithms' performance under different total time horizons, we let  $K$  vary from  $3 \times 10^4$  to  $2.4 \times 10^5$ , with a stepsize of  $3 \times 10^4$ , and plot the cumulative regret  $\text{Regret}(K)$  for these different total time step  $K$ . We set  $B_K = 1, 10, 20$ , and  $K^{1/3}$  to observe their performance in different levels of  $B_K$ .

## C Proof of Theorem 4.1

We prove the lower bound in Theorem 4.1 here. We need the following lemma from Zhou et al. (2021).

**Lemma C.1** (Modification from Lemma 25, Zhou et al. 2021). *Fix a positive real  $0 < \delta \leq 1/3$ , and positive integers  $T, d$  and assume that  $T \geq d^2/(2\delta)$ . Let  $\Delta = \sqrt{d\delta/T}/(4\sqrt{2})$  and consider the linear bandit problems  $\mathcal{L}_\mu$  parameterized with a parameter vector  $\mu \in \{-\Delta, \Delta\}^d$  and action set  $\mathcal{A} = \{-1/\sqrt{d}, 1/\sqrt{d}\}^d$  so that the reward distribution for taking action  $\mathbf{a} \in \mathcal{A}$  is a Bernoulli distribution  $B(\delta + \langle \mu^*, \mathbf{a} \rangle)$ . Then for any bandit algorithm  $\mathcal{B}$  such that*

$$\mathbb{E}_{\mu \sim \text{Unif}\{-\Delta, \Delta\}^d} [\text{Regret}(T, \mathcal{L}_\mu)] \geq \frac{d\sqrt{T}\delta}{8\sqrt{2}}. \quad (17)$$

Here  $\text{Regret}(T, \mathcal{L}_\mu)$  represents the regret under algorithm  $\mathcal{B}$  on the instance  $\mathcal{L}_\mu$ .

Next we prove Theorem 4.1.

*Proof of Theorem 4.1.* Let  $T < K$  be some constant to be defined. Let  $\delta$  be a constant satisfying  $2\delta \leq d^2/T$ . We create  $w = K/T$  number of linear bandit instances with the linear parameter  $\mu_1, \dots, \mu_w$ , where  $\mu_i \sim$

$\{-\Delta, \Delta\}^d$ ,  $\Delta = \sqrt{d\delta/T}/4\sqrt{2}$ . Our nonstationary instance  $\mathcal{L}_{\mu_1, \dots, \mu_w}$  consists of  $\mathcal{L}_{\mu_1}, \dots, \mathcal{L}_{\mu_w}$ , where at the step  $i \cdot T + 1, \dots, i \cdot T + T$ ,  $\mathcal{L}_{\mu_1, \dots, \mu_w}$  follows  $\mathcal{L}_{\mu_i}$ . Then by the independence of  $\mu_i$ , we have

$$\mathbb{E}_{\mu_1, \dots, \mu_w \sim \text{Unif}\{-\Delta, \Delta\}^d} \text{Regret}(T, \mathcal{L}_{\mu_1, \dots, \mu_w}) = \sum_{i=1}^w \mathbb{E}_{\mu_i \sim \text{Unif}\{-\Delta, \Delta\}^d} [\text{Regret}(T, \mathcal{L}_{\mu_i})] \geq \frac{d\sqrt{T\delta}}{8\sqrt{2}} \cdot \frac{K}{T}. \quad (18)$$

Next we calculate the total variation and total variance for instance  $\mathcal{L}_{\mu_1, \dots, \mu_w}$ . For each step, the reward distribution is a Bernoulli distribution  $B(\delta + \langle \mu_i, \mathbf{a} \rangle)$ , whose variance is

$$(\delta + \langle \mu_i, \mathbf{a} \rangle)(1 - \delta - \langle \mu_i, \mathbf{a} \rangle) \leq (\delta + \langle \mu_i, \mathbf{a} \rangle) \leq 2\delta, \quad (19)$$

where we use the fact  $\sqrt{d}\Delta \leq \delta$ . Therefore, the total variance over  $K$  steps is bounded by

$$V \leq 2K\delta. \quad (20)$$

Next, for the total variation, we have for any  $k, k+1$  belong to the same  $\mu_i$ , the variation of  $\mu$  is 0. Note that for any two different  $\mu_i, \mu_j$ , their difference is at most  $\|\mu_i - \mu_j\| \leq 2\sqrt{d} \cdot \Delta^2$ , then the total variation is bounded by

$$B \leq \frac{K}{T} \cdot 2\Delta\sqrt{d} = \sqrt{d\delta/T}/(4\sqrt{2}) \cdot \frac{K}{T} \cdot 2\sqrt{d} = \frac{dK\sqrt{\delta}}{2\sqrt{2}T^3}. \quad (21)$$

Then we select  $\delta$  and  $T$  as

$$\delta = \frac{V_K}{2K}, \quad T = \max\left\{\left(\frac{KV_K d^2}{16B_K^2}\right)^{1/3}, d^2 K/V_K\right\}, \text{ satisfying } 2K\delta \leq V_K, \frac{dK\sqrt{\delta}}{2\sqrt{2}T^3} \leq B_K, \quad T \geq \frac{d^2}{2\delta}. \quad (22)$$

We have the lower bound as

$$\mathbb{E}_{\mu_1, \dots, \mu_w \sim \text{Unif}\{-\Delta, \Delta\}^d} \text{Regret}(T, \mathcal{L}_{\mu_1, \dots, \mu_w}) \geq \Omega(d^{2/3} B_K^{1/3} V_K^{1/3} K^{1/3} \wedge V_K). \quad (23)$$

Therefore, there must exists  $\mu_1^*, \dots, \mu_w^*$ , satisfying

$$\text{Regret}(T, \mathcal{L}_{\mu_1^*, \dots, \mu_w^*}) \geq \Omega(d^{2/3} B_K^{1/3} V_K^{1/3} K^{1/3} \wedge V_K). \quad (24)$$

Finally, combining (24) with the lower bound result in Wei et al. (2016) concludes our proof.  $\square$

## D Proof of Lemma 5.1

For simplicity, we denote

$$\hat{\beta} := 12\sqrt{d \log(1 + \frac{wA^2}{\alpha^2 d \lambda}) \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{w^2}{\delta}) + 30 \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{w^2}{\delta}) \frac{R}{\gamma^2} + \sqrt{\lambda} B)}. \quad (25)$$

It is obvious that  $\hat{\beta} \geq \hat{\beta}_k$  for all  $k \in [K]$ . We call the restart time rounds *grids* and denote them by  $g_1, g_2, \dots, g_{\lceil \frac{K}{w} \rceil - 1}$ , where  $g_i \% w = 0$  for all  $i \in [\lceil \frac{K}{w} \rceil - 1]$ . Let  $i_k$  be the grid index of time round  $k$ , i.e.,  $g_{i_k} \leq k < g_{i_k+1}$ .

For ease of exposition and without loss of generality, we prove the lemma for  $k \in [1, w]$ . We calculate the estimation difference  $|\mathbf{a}^\top(\hat{\theta}_k - \theta_k)|$  for any  $\mathbf{a} \in \mathbb{R}^d$ ,  $\|\mathbf{a}\|_2 \leq A$ ,  $k \in [1, w]$ . By definition:

$$\hat{\theta}_k = \hat{\Sigma}_k^{-1} \mathbf{b}_k = \hat{\Sigma}_k^{-1} \left( \sum_{t=1}^{k-1} \frac{r_t \mathbf{a}_t}{\bar{\sigma}_t^2} \right) = \hat{\Sigma}_k^{-1} \left( \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top \theta_t}{\bar{\sigma}_t^2} + \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \epsilon_t}{\bar{\sigma}_t^2} \right), \quad (26)$$

where  $\hat{\Sigma}_k = \lambda \mathbf{I} + \sum_{t=g_{i_k}}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2}$ .

Then we have

$$\hat{\theta}_k - \theta_k = \hat{\Sigma}_k^{-1} \left( \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} (\theta_t - \theta_k) + \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \epsilon_t}{\bar{\sigma}_t^2} \right) - \lambda \hat{\Sigma}_k^{-1} \theta_k. \quad (27)$$

Therefore

$$|\mathbf{a}^\top (\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k)| \leq \left| \mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} (\boldsymbol{\theta}_t - \boldsymbol{\theta}_k) \right| + \|\mathbf{a}\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} \left\| \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \epsilon_t}{\bar{\sigma}_t^2} \right\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} + \lambda \|\mathbf{a}\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} \|\hat{\boldsymbol{\Sigma}}_k^{-\frac{1}{2}} \boldsymbol{\theta}_k\|_2, \quad (28)$$

where we use the Cauchy-Schwarz inequality.

For the first term, we have that for any  $k \in [1, w]$

$$\begin{aligned} \left| \mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \sum_{t=1}^k \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} (\boldsymbol{\theta}_t - \boldsymbol{\theta}_k) \right| &\leq \sum_{t=1}^{k-1} |\mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t}| \cdot \left| \frac{\mathbf{a}_t^\top}{\bar{\sigma}_t} \left( \sum_{s=t}^{k-1} (\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}) \right) \right| && \text{(triangle inequality)} \\ &\leq \sum_{t=1}^{k-1} |\mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t}| \cdot \left\| \frac{\mathbf{a}_t}{\bar{\sigma}_t} \right\|_2 \cdot \left\| \sum_{s=t}^{k-1} (\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}) \right\|_2 && \text{(Cauchy-Schwarz)} \\ &\leq \frac{A}{\alpha} \sum_{t=1}^{k-1} |\mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t}| \cdot \left\| \sum_{s=t}^{k-1} (\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}) \right\|_2 && (\|\mathbf{a}_t\| \leq A, \bar{\sigma}_t \geq \alpha) \\ &\leq \frac{A}{\alpha} \sum_{s=1}^{k-1} \sum_{t=1}^s |\mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t}| \cdot \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 && (\sum_{t=1}^{k-1} \sum_{s=t}^{k-1} = \sum_{s=1}^{k-1} \sum_{t=1}^s) \\ &\leq \frac{A}{\alpha} \sum_{s=1}^{k-1} \sqrt{\left[ \sum_{t=1}^s \mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \mathbf{a} \right] \cdot \left[ \sum_{t=1}^s \frac{\mathbf{a}_t^\top}{\bar{\sigma}_t} \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t} \right]} \cdot \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 && \text{(Cauchy-Schwarz)} \\ &\leq \frac{A}{\alpha} \sum_{s=1}^{k-1} \sqrt{\left[ \sum_{t=1}^s \mathbf{a}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \mathbf{a} \right] \cdot d} \cdot \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 && ((\star)) \\ &\leq \frac{A \|\mathbf{a}\|_2}{\alpha} \sqrt{d} \sum_{s=1}^{k-1} \sqrt{\frac{\sum_{t=1}^{k-1} 1}{\lambda}} \cdot \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 && (\lambda_{\max}(\hat{\boldsymbol{\Sigma}}_k^{-1}) \leq \frac{1}{\lambda}) \\ &\leq \frac{A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2, && (29) \end{aligned}$$

where the inequality  $(\star)$  follows from the fact that  $\sum_{t=1}^s \frac{\mathbf{a}_t^\top}{\bar{\sigma}_t} \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t} \leq d$  that can be proved as follows. We have  $\sum_{t=1}^{k-1} \frac{\mathbf{a}_t^\top}{\bar{\sigma}_t} \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t} = \sum_{t=1}^{k-1} \text{tr} \left( \frac{\mathbf{a}_t^\top}{\bar{\sigma}_t} \hat{\boldsymbol{\Sigma}}_k^{-1} \frac{\mathbf{a}_t}{\bar{\sigma}_t} \right) = \text{tr} \left( \hat{\boldsymbol{\Sigma}}_k^{-1} \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} \right)$ . Given the eigenvalue decomposition  $\sum_{t=1}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} = \text{diag}(\lambda_1, \dots, \lambda_d)^\top$ , we have  $\hat{\boldsymbol{\Sigma}}_k = \text{diag}(\lambda_1 + \lambda, \dots, \lambda_d + \lambda)^\top$ , and  $\text{tr} \left( \hat{\boldsymbol{\Sigma}}_k^{-1} \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \mathbf{a}_t^\top}{\bar{\sigma}_t^2} \right) = \sum_{i=1}^d \frac{\lambda_i}{\lambda_j + \lambda} \leq d$ .

For the second term, by the assumption on  $\epsilon_k$ , we know that

$$\begin{aligned} |\epsilon_k / \bar{\sigma}_k| &\leq R / \alpha, \\ |\epsilon_k / \bar{\sigma}_k| \cdot \min\{1, \|\mathbf{a}_k / \bar{\sigma}_k\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}}\} &\leq R \|\mathbf{a}_k\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} / \bar{\sigma}_k^2 \leq R / \gamma^2, \\ \mathbb{E}[\epsilon_k | \mathbf{a}_{1:k}, \epsilon_{1:k-1}] &= 0, \quad \mathbb{E}[(\epsilon_k / \bar{\sigma}_k)^2 | \mathbf{a}_{1:k}, \epsilon_{1:k-1}] \leq 1, \quad \|\mathbf{a}_k / \bar{\sigma}_k\|_2 \leq A / \alpha, \end{aligned}$$

Therefore, setting  $\mathcal{G}_k = \sigma(\mathbf{a}_{1:k}, \epsilon_{1:k-1})$ , and using that  $\sigma_k$  is  $\mathcal{G}_k$ -measurable, applying Theorem H.1 to  $(\mathbf{x}_k, \eta_k) = (\mathbf{a}_k / \bar{\sigma}_k, \epsilon_k / \bar{\sigma}_k)$  with  $\epsilon = R / \gamma^2$ , we get that with probability at least  $1 - \delta$ , for all  $k \in [1, w]$ ,

$$\left\| \sum_{t=1}^{k-1} \frac{\mathbf{a}_t \epsilon_t}{\bar{\sigma}_t^2} \right\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} \leq 12 \sqrt{d \log(1 + \frac{(k\%w)A^2}{\alpha^2 d \lambda}) \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{(k\%w)^2}{\delta})) + 30 \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{(k\%w)^2}{\delta})) \frac{R}{\gamma^2}}. \quad (30)$$

For the last term

$$\lambda \|\mathbf{a}\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} \|\hat{\boldsymbol{\Sigma}}_k^{-\frac{1}{2}} \boldsymbol{\theta}_k\|_2 \leq \lambda \|\mathbf{a}\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} \|\hat{\boldsymbol{\Sigma}}_k^{-\frac{1}{2}}\|_2 \|\boldsymbol{\theta}_k\|_2 \leq \lambda \|\mathbf{a}\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}} \frac{1}{\sqrt{\lambda_{\min}(\hat{\boldsymbol{\Sigma}}_k)}} \|\boldsymbol{\theta}_k\|_2 \leq \sqrt{\lambda} B \|\mathbf{a}\|_{\hat{\boldsymbol{\Sigma}}_k^{-1}}, \quad (31)$$

where we use the fact that  $\lambda_{\min}(\hat{\Sigma}_k) \geq \lambda$ .

Therefore, with probability at least  $1 - \delta$ , we have

$$\begin{aligned}
 & |\mathbf{a}^\top (\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k)| \\
 & \leq \frac{A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{t=1}^{k-1} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}\|_2 \\
 & \quad + \|\mathbf{a}\|_{\hat{\Sigma}_k^{-1}} \left( 12 \sqrt{d \log(1 + \frac{(k\%w)A^2}{\alpha^2 d \lambda}) \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{(k\%w)^2}{\delta}))} + 30 \log(32(\log(\frac{\gamma^2}{\alpha} + 1) \frac{(k\%w)^2}{\delta})) \frac{R}{\gamma^2} + \sqrt{\lambda} B \right) \\
 & = \frac{A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{t=1}^{k-1} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}\|_2 + \hat{\beta}_k \|\mathbf{a}\|_{\hat{\Sigma}_k^{-1}}, \tag{32}
 \end{aligned}$$

where  $\hat{\beta}_k$  is defined in Eq.(5).

## E Proof for Theorem 5.2

For simplicity of analysis, we only analyze the regret over the first grid, *i.e.*, we try to analyze  $\text{Regret}(\tilde{K})$  for  $\tilde{K} \in [1, w]$ . Denote  $\mathcal{E}_1$  as the event when Lemma 5.1 holds. Therefore, under event  $\mathcal{E}_1$ , for any  $\tilde{K} \in [1, w]$ , the regret can be bounded by

$$\begin{aligned}
 \text{Regret}(\tilde{K}) &= \sum_{k=1}^{\tilde{K}} [\langle \mathbf{a}_k^* - \mathbf{a}_k, \boldsymbol{\theta}_k \rangle] \\
 &= \sum_{k=1}^{\tilde{K}} [\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k - \hat{\boldsymbol{\theta}}_k \rangle + (\langle \mathbf{a}_k^*, \hat{\boldsymbol{\theta}}_k \rangle + \hat{\beta}_k \|\mathbf{a}_k^*\|_{\hat{\Sigma}_k^{-1}}) - (\langle \mathbf{a}_k, \hat{\boldsymbol{\theta}}_k \rangle + \hat{\beta}_k \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}}) \\
 & \quad + \langle \mathbf{a}_k, \hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k \rangle + \hat{\beta}_k \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}} - \hat{\beta}_k \|\mathbf{a}_k^*\|_{\hat{\Sigma}_k^{-1}}] \\
 &\leq \frac{2A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{k=1}^{\tilde{K}} \sum_{t=1}^{k-1} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}\|_2 + 2 \sum_{k=1}^{\tilde{K}} \min \left\{ 1, \hat{\beta}_k \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}} \right\}, \tag{33}
 \end{aligned}$$

where in the last inequality we use the definition of event  $\mathcal{E}_1$ , the arm selection rule in Line 7 of Algo.1, and  $0 \leq \langle \mathbf{a}_k^*, \boldsymbol{\theta}^* \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}^* \rangle \leq 2$ .

Then we will bound the two terms in Eq.(33).

For the first term, we have

$$\begin{aligned}
 & \frac{2A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{k=1}^{\tilde{K}} \sum_{t=1}^{k-1} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}\|_2 \\
 &= \frac{2A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} \sum_{t=1}^{\tilde{K}-1} \sum_{k=t}^{\tilde{K}} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}\|_2 \\
 &\leq \frac{2A^2}{\alpha} \sqrt{\frac{dw}{\lambda}} w \sum_{t=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}\|_2. \tag{34}
 \end{aligned}$$

To bound the second term in Eq.(33), we decompose the set  $[\tilde{K}]$  into a union of two disjoint subsets  $[K] = \mathcal{I}_1 \cup \mathcal{I}_2$ .

$$\mathcal{I}_1 = \left\{ k \in [\tilde{K}] : \left\| \frac{\mathbf{a}_k}{\sigma_k} \right\|_{\hat{\Sigma}_k^{-1}} \geq 1 \right\}, \quad \mathcal{I}_2 = \left\{ k \in [\tilde{K}] : \left\| \frac{\mathbf{a}_k}{\sigma_k} \right\|_{\hat{\Sigma}_k^{-1}} < 1 \right\}. \tag{35}$$

Then the following upper bound of  $|\mathcal{I}_1|$  holds:

$$\begin{aligned}
 |\mathcal{I}_1| &= \sum_{k \in \mathcal{I}_1} \min \left\{ 1, \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}}^2 \right\} \\
 &\leq \sum_{k=1}^{\tilde{K}} \min \left\{ 1, \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}}^2 \right\} \\
 &\leq 2d\iota,
 \end{aligned} \tag{36}$$

where  $\iota = \log(1 + \frac{wA^2}{d\lambda\alpha^2})$ , the first equality holds since  $\|\frac{\mathbf{x}_k}{\bar{\sigma}_k}\|_{\hat{\Sigma}_k^{-1}} \geq 1$  for  $k \in \mathcal{I}_1$ , the last inequality holds due to Lemma H.2 together with the fact  $\|\frac{\mathbf{a}_k}{\bar{\sigma}_k}\|_2 \leq \frac{A}{\alpha}$  since  $\bar{\sigma}_k \geq \alpha$  and  $\|\mathbf{a}_k\|_2 \leq A$ .

Then, we have

$$\begin{aligned}
 &\sum_{k=1}^{\tilde{K}} \min \left\{ 1, \hat{\beta}_k \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}} \right\} \\
 &= \sum_{k \in \mathcal{I}_1} \min \left\{ 1, \bar{\sigma}_k \hat{\beta}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} \right\} + \sum_{k \in \mathcal{I}_2} \min \left\{ 1, \bar{\sigma}_k \hat{\beta}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} \right\} \\
 &\leq \left[ \sum_{k \in \mathcal{I}_1} 1 \right] + \sum_{k \in \mathcal{I}_2} \bar{\sigma}_k \hat{\beta}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} \\
 &\leq 2d\iota + \hat{\beta} \sum_{k \in \mathcal{I}_2} \bar{\sigma}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}},
 \end{aligned} \tag{37}$$

where the first inequality holds since  $\min\{1, x\} \leq 1$  and also  $\min\{1, x\} \leq x$ , the second inequality holds by Eq.(36), and the fact the  $\hat{\beta} \geq \hat{\beta}_k$  for all  $k \in [K]$  ( $\hat{\beta}$  is defined in Eq.(25)). Next we further bound the second summation term in (37). We decompose  $\mathcal{I}_2 = \mathcal{J}_1 \cup \mathcal{J}_2$ , where

$$\mathcal{J}_1 = \left\{ k \in \mathcal{I}_2 : \bar{\sigma}_k = \sigma_k \cup \bar{\sigma}_k = \alpha \right\}, \quad \mathcal{J}_2 = \left\{ k \in \mathcal{I}_2 : \bar{\sigma}_k = \gamma \sqrt{\|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}}} \right\}.$$

Then  $\sum_{k \in \mathcal{I}_2} \bar{\sigma}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} = \sum_{k \in \mathcal{J}_1} \bar{\sigma}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} + \sum_{k \in \mathcal{J}_2} \bar{\sigma}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}}$ . First, for  $k \in \mathcal{J}_1$ , we have

$$\begin{aligned}
 \sum_{k \in \mathcal{J}_1} \bar{\sigma}_k \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} &\leq \sum_{k \in \mathcal{J}_1} (\sigma_k + \alpha) \min \left\{ 1, \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} \right\} \\
 &\leq \sqrt{\sum_{k=1}^{\tilde{K}} (\sigma_k + \alpha)^2} \sqrt{\sum_{k=1}^{\tilde{K}} \min \left\{ 1, \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}} \right\}^2} \\
 &\leq \sqrt{2 \sum_{k=1}^{\tilde{K}} (\sigma_k^2 + \alpha^2)} \sqrt{\sum_{k=1}^{\tilde{K}} \min \left\{ 1, \left\| \frac{\mathbf{a}_k}{\bar{\sigma}_k} \right\|_{\hat{\Sigma}_k^{-1}}^2 \right\}} \\
 &\leq 2 \sqrt{\sum_{k=1}^{\tilde{K}} \sigma_k^2 + \tilde{K} \alpha^2} \sqrt{d\iota},
 \end{aligned} \tag{38}$$

where the first inequality holds since  $\bar{\sigma}_k \leq \sigma_k + \alpha$  for  $k \in \mathcal{J}_1$  and  $\|\frac{\mathbf{a}_k}{\bar{\sigma}_k}\|_{\hat{\Sigma}_k^{-1}} \leq 1$  since  $k \in \mathcal{J}_1 \subseteq \mathcal{I}_2$ , the second inequality holds by Cauchy-Schwarz inequality, the third inequality holds due to  $(a+b)^2 \leq 2(a^2 + b^2)$ , and the last inequality holds due to Lemma H.2.



Finally we bound the summation for  $k \in \mathcal{J}_2$ . When  $k \in \mathcal{J}_2$ , we have  $\bar{\sigma}_k = \gamma^2 \|\frac{\mathbf{a}_k}{\bar{\sigma}_k}\|_{\hat{\Sigma}_k^{-1}}$ . Therefore we have

$$\begin{aligned} \sum_{k \in \mathcal{J}_2} \bar{\sigma}_k \|\frac{\mathbf{a}_k}{\bar{\sigma}_k}\|_{\hat{\Sigma}_k^{-1}} &= \sum_{k \in \mathcal{J}_2} \gamma^2 \|\frac{\mathbf{a}_k}{\bar{\sigma}_k}\|_{\hat{\Sigma}_k^{-1}}^2 \\ &\leq \sum_{k=1}^{\tilde{K}} \gamma^2 \min \left\{ 1, \|\frac{\mathbf{a}_k}{\bar{\sigma}_k}\|_{\hat{\Sigma}_k^{-1}}^2 \right\} \\ &\leq 2\gamma^2 d\iota, \end{aligned} \quad (39)$$

where in the first inequality we use the fact that  $\|\frac{\mathbf{a}_k}{\bar{\sigma}_k}\|_{\hat{\Sigma}_k^{-1}} \leq 1$  since  $k \in \mathcal{J}_2 \subseteq \mathcal{I}_2$ , and in the last inequality we use Lemma H.2.

Therefore, with Eq.(33), Eq.(34), Eq.(37), Eq.(38), Eq.(39), we can get the regret upper bound for  $\tilde{K} \in [1, w]$

$$\text{Regret}(\tilde{K}) \leq \frac{2A^2 w^{\frac{3}{2}}}{\alpha} \sqrt{\frac{d}{\lambda}} \sum_{k=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + 4\hat{\beta} \sqrt{d\iota} \sqrt{\sum_{k \in [\tilde{K}]} \sigma_k^2 + w\alpha^2 + 4d\iota\gamma^2\hat{\beta} + 4d\iota}. \quad (40)$$

Therefore, by the same deduction, we can get that

$$\text{Regret}([g_i, g_{i+1}]) \leq \frac{2A^2 w^{\frac{3}{2}}}{\alpha} \sqrt{\frac{d}{\lambda}} \sum_{k=g_i}^{g_{i+1}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + 4\hat{\beta} \sqrt{d\iota} \sqrt{\sum_{k=g_i}^{g_{i+1}} \sigma_k^2 + w\alpha^2 + 4d\iota\gamma^2\hat{\beta} + 4d\iota}, \quad (41)$$

where we use  $\text{Regret}([g_i, g_{i+1}])$  to denote the regret accumulated in the time period  $[g_i, g_{i+1}]$ .

Finally, without loss of generality, we assume  $K \% w = 0$ . Then we have

$$\begin{aligned} \text{Regret}(\tilde{K}) &= \sum_{i=0}^{\frac{\tilde{K}}{w}-1} \text{Regret}([g_i, g_{i+1}]) \\ &\leq \frac{2A^2 w^{\frac{3}{2}}}{\alpha} \sqrt{\frac{d}{\lambda}} \sum_{i=0}^{\frac{\tilde{K}}{w}-1} \sum_{k=g_i}^{g_{i+1}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + 4\hat{\beta} \sqrt{d\iota} \sum_{i=0}^{\frac{\tilde{K}}{w}-1} \sqrt{\sum_{k=g_i}^{g_{i+1}} \sigma_k^2 + w\alpha^2 + \frac{4d\iota\gamma^2\hat{\beta}K}{w} + \frac{4dK\iota}{w}} \\ &\leq \frac{2A^2 w^{\frac{3}{2}}}{\alpha} \sqrt{\frac{d}{\lambda}} \sum_{k=1}^{K-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + 4\hat{\beta} \sqrt{d\iota} \sqrt{\frac{K}{w} \sum_{i=0}^{\frac{\tilde{K}}{w}-1} (\sum_{k=g_i}^{g_{i+1}} \sigma_k^2 + w\alpha^2) + \frac{4d\iota\gamma^2\hat{\beta}K}{w} + \frac{4dK\iota}{w}} \\ &\leq \frac{2A^2 w^{\frac{3}{2}} B_K}{\alpha} \sqrt{\frac{d}{\lambda}} + 4\hat{\beta} \sqrt{\frac{Kd\iota}{w}} \sqrt{\sum_{k=1}^K \sigma_k^2 + K\alpha^2 + \frac{4d\iota\gamma^2\hat{\beta}K}{w} + \frac{4dK\iota}{w}}, \end{aligned}$$

where in the second inequality we use Cauchy-Schwarz inequality, and the last inequality holds due to  $\sum_{k \in [K-1]} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 \leq B_K$ .

## F Proof for Theorem 6.1

Recall that we call the restart time rounds *grids* and denote them by  $g_1, g_2, \dots, g_{\lceil \frac{\tilde{K}}{w} \rceil - 1}$ , where  $g_i \% w = 0$  for all  $i \in [\lceil \frac{\tilde{K}}{w} \rceil - 1]$ . Let  $i_k$  be the grid index of time round  $k$ , i.e.,  $g_{i_k} \leq k < g_{i_k+1}$ . We denote  $\hat{\Psi}_{k,\ell} := \{t : t \in [g_{i_k}, k-1], \ell_t = \ell\}$ .

For simplicity of analysis, we first try to bound the regret over the first grid, i.e., we try to analyze  $\text{Regret}(\tilde{K})$  for  $\tilde{K} \in [1, w]$ . Note that in this case, for any  $k \in [\tilde{K}]$  with  $\tilde{K} \in [1, w]$ , we have  $g_{i_k} = 1$ , so  $\hat{\Psi}_{k,\ell} := \{t : t \in [1, k-1], \ell_t = \ell\}$ .

First, we calculate the estimation difference  $|\mathbf{a}^\top (\hat{\boldsymbol{\theta}}_{k,\ell} - \boldsymbol{\theta}_k)|$  for any  $\mathbf{a} \in \mathbb{R}^d$ ,  $\|\mathbf{a}\|_2 \leq A$ . Recall that by definition,  $\hat{\Sigma}_{k,\ell} = 2^{-2\ell} \mathbf{I} + \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top$ ,  $\hat{\mathbf{b}}_{k,\ell} = \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 r_t \mathbf{a}_t$ , and

$$\hat{\boldsymbol{\theta}}_{k,\ell} = \hat{\Sigma}_{k,\ell}^{-1} \hat{\mathbf{b}}_{k,\ell} = \hat{\Sigma}_{k,\ell}^{-1} \left( \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 r_t \mathbf{a}_t \right) = \hat{\Sigma}_{k,\ell}^{-1} \left( \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top \boldsymbol{\theta}_t + \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \epsilon_t \right).$$

Then we have

$$\hat{\theta}_{k,\ell} - \theta_k = \hat{\Sigma}_{k,\ell}^{-1} \left( \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top (\theta_t - \theta_k) + \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \epsilon_t \right) - 2^{-2\ell} \hat{\Sigma}_{k,\ell}^{-1} \theta_k. \quad (42)$$

Therefore, we can get

$$|\mathbf{a}^\top (\hat{\theta}_{k,\ell} - \theta_k)| \leq \left| \mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top (\theta_t - \theta_k) \right| + \|\mathbf{a}\|_{\hat{\Sigma}_{k,\ell}^{-1}} \left\| \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \epsilon_t \right\|_{\hat{\Sigma}_{k,\ell}^{-1}} + 2^{-2\ell} \|\mathbf{a}\|_{\hat{\Sigma}_{k,\ell}^{-1}} \|\hat{\Sigma}_{k,\ell}^{-\frac{1}{2}} \theta_k\|_2, \quad (43)$$

where we use the Cauchy-Schwarz inequality.

For the first term, we have that for any  $k \in [1, w]$

$$\begin{aligned} \left| \mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top (\theta_t - \theta_k) \right| &\leq \sum_{t \in \hat{\Psi}_{k,\ell}} |\mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t| \cdot \left| w_t \mathbf{a}_t^\top \left( \sum_{s=t}^{k-1} (\theta_s - \theta_{s+1}) \right) \right| && \text{(triangle inequality)} \\ &\leq \sum_{t \in \hat{\Psi}_{k,\ell}} |\mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t| \cdot \|w_t \mathbf{a}_t\|_2 \cdot \left\| \sum_{s=t}^{k-1} (\theta_s - \theta_{s+1}) \right\|_2 && \text{(Cauchy-Schwarz)} \\ &\leq A \sum_{t \in \hat{\Psi}_{k,\ell}} |\mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t| \cdot \left\| \sum_{s=t}^{k-1} (\theta_s - \theta_{s+1}) \right\|_2 \quad (\|\mathbf{a}_t\| \leq A, w_t = \frac{2^{-\ell_t}}{\|\mathbf{a}_t\|_{\hat{\Sigma}_{t,\ell_t}^{-1}}} \leq 1) \\ &\leq A \sum_{s=1}^{k-1} \sum_{t \in \hat{\Psi}_{k,\ell}} |\mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t| \cdot \|\theta_s - \theta_{s+1}\|_2 \\ &\leq A \sum_{s=1}^{k-1} \sqrt{\left[ \sum_{t \in \hat{\Psi}_{k,\ell}} \mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} \mathbf{a} \right] \cdot \left[ \sum_{t \in \hat{\Psi}_{k,\ell}} w_t \mathbf{a}_t^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t \right]} \cdot \|\theta_s - \theta_{s+1}\|_2 && \text{(Cauchy-Schwarz)} \\ &\leq A \sum_{s=1}^{k-1} \sqrt{\left[ \sum_{t \in \hat{\Psi}_{k,\ell}} \mathbf{a}^\top \hat{\Sigma}_{k,\ell}^{-1} \mathbf{a} \right] \cdot d} \cdot \|\theta_s - \theta_{s+1}\|_2 && ((\star)) \\ &\leq A \|\mathbf{a}\|_2 \sqrt{d} \sum_{s=1}^{k-1} \sqrt{2^{2\ell} \sum_{t \in \hat{\Psi}_{k,\ell}} 1} \cdot \|\theta_s - \theta_{s+1}\|_2 && (\lambda_{\max}(\hat{\Sigma}_{k,\ell}^{-1}) \leq \frac{1}{2^{-2\ell}} = 2^{2\ell}) \\ &\leq A^2 2^\ell \sqrt{dw} \sum_{s=1}^{k-1} \|\theta_s - \theta_{s+1}\|_2, && (44) \end{aligned}$$

where the inequality  $(\star)$  follows from the fact that  $\sum_{t \in \hat{\Psi}_{k,\ell}} w_t \mathbf{a}_t^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t \leq d$  that can be proved as follows. We have  $\sum_{t \in \hat{\Psi}_{k,\ell}} w_t \mathbf{a}_t^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t = \sum_{t \in \hat{\Psi}_{k,\ell}} \text{tr} \left( w_t \mathbf{a}_t^\top \hat{\Sigma}_{k,\ell}^{-1} w_t \mathbf{a}_t \right) = \text{tr} \left( \hat{\Sigma}_{k,\ell}^{-1} \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top \right)$ . Given the eigenvalue decomposition  $\sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top = \text{diag}(\lambda_1, \dots, \lambda_d)^\top$ , we have  $\hat{\Sigma}_{k,\ell} = \text{diag}(\lambda_1 + \lambda, \dots, \lambda_d + \lambda)^\top$ , and  $\text{tr} \left( \hat{\Sigma}_{k,\ell}^{-1} \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \mathbf{a}_t^\top \right) = \sum_{j=1}^d \frac{\lambda_j}{\lambda_j + \lambda} \leq d$ .

For the second term in Eq.(43), we can apply Theorem H.3 for the layer  $\ell$ . In detail, for any  $k \in [K]$ , for each  $t \in \hat{\Psi}_{k,\ell}$ , we have

$$\|w_t \mathbf{a}_t\|_{\hat{\Sigma}_{t,\ell}^{-1}} = 2^{-\ell}, \quad \mathbb{E}[w_t^2 \epsilon_t^2 | \mathcal{F}_t] \leq w_t^2 \mathbb{E}[\epsilon_t^2 | \mathcal{F}_t] \leq w_t^2 \sigma_t^2, \quad |w_t \epsilon_t| \leq |\epsilon_t| \leq R,$$

where the last inequality holds due to the fact that  $w_t = \frac{2^{-\ell_t}}{\|\mathbf{a}_t\|_{\hat{\Sigma}_{t,\ell_t}^{-1}}} \leq 1$ . According to Theorem H.3, and taking a

union bound, we can deduce that with probability at least  $1 - \delta$ , for all  $\ell \in [L]$ , for all round  $k \in \Psi_{K+1,\ell}$ ,

$$\left\| \sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \mathbf{a}_t \epsilon_t \right\|_{\hat{\Sigma}_{k,\ell}^{-1}} \leq 16 \cdot 2^{-\ell} \sqrt{\sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \sigma_t^2 \log\left(\frac{4w^2 L}{\delta}\right)} + 6 \cdot 2^{-\ell} R \log\left(\frac{4w^2 L}{\delta}\right). \quad (45)$$

For simplicity, we denote  $\mathcal{E}_{\text{conf}}$  as the event such that Eq.(45) holds.

For the third term in Eq.(43), we have

$$2^{-2\ell} \|\mathbf{a}\|_{\hat{\Sigma}_{k,\ell}^{-1}} \|\hat{\Sigma}_{k,\ell}^{-\frac{1}{2}} \boldsymbol{\theta}_k\|_2 \leq 2^{-2\ell} \|\mathbf{a}\|_{\hat{\Sigma}_{k,\ell}^{-1}} \|\hat{\Sigma}_k^{-\frac{1}{2}}\|_2 \|\boldsymbol{\theta}_k\|_2 \leq 2^{-2\ell} \|\mathbf{a}\|_{\hat{\Sigma}_{k,\ell}^{-1}} \frac{1}{\sqrt{\lambda_{\min}(\hat{\Sigma}_{k,\ell})}} \|\boldsymbol{\theta}_k\|_2 \leq 2^{-\ell} B \|\mathbf{a}\|_{\hat{\Sigma}_k^{-1}}, \quad (46)$$

where we use the fact that  $\lambda_{\min}(\hat{\Sigma}_{k,\ell}) \geq 2^{-2\ell}$ .

For simplicity, we denote  $\ell^* = \lceil \frac{1}{2} \log_2 \log(4(w+1)^2 L / \delta) \rceil + 8$ . Then, under  $\mathcal{E}_{\text{conf}}$ , by the definition of  $\hat{\beta}_{k,\ell}$  in Eq.(9), Lemma H.4 and Lemma H.5, with probability at least  $1 - \delta$ , we have for all  $\ell^* + 1 \leq \ell \leq L$ ,

$$\hat{\beta}_{k,\ell} \geq 16 \cdot 2^{-\ell} \sqrt{\sum_{t \in \hat{\Psi}_{k,\ell}} w_t^2 \sigma_t^2 \log\left(\frac{4w^2 L}{\delta}\right)} + 6 \cdot 2^{-\ell} R \log\left(\frac{4w^2 L}{\delta}\right) + 2^{-\ell} B. \quad (47)$$

Therefore, with Eq.(43), Eq.(44), Eq.(45), Eq.(46), Eq.(47), with probability at least  $1 - 3\delta$ , for all  $\ell^* + 1 \leq \ell \leq L$  we have

$$|\mathbf{a}^\top (\hat{\boldsymbol{\theta}}_{k,\ell} - \boldsymbol{\theta}_k)| \leq A^2 2^\ell \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 + \hat{\beta}_{k,\ell} \|\mathbf{a}\|_{\hat{\Sigma}_{k,\ell}^{-1}}. \quad (48)$$

Then for all  $k \in [K]$  such that  $\ell^* + 1 \leq \ell_k \leq L$ , with probability at least  $1 - 3\delta$  we have

$$\begin{aligned} \langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle &\leq \min_{\ell \in [L]} \langle \mathbf{a}_k^*, \hat{\boldsymbol{\theta}}_{k,\ell} \rangle + A^2 2^\ell \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 + \hat{\beta}_{k,\ell} \|\mathbf{a}_k^*\|_{\hat{\Sigma}_{k,\ell}^{-1}} \\ &\leq A^2 2^L \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 + \min_{\ell \in [L]} \langle \mathbf{a}_k^*, \hat{\boldsymbol{\theta}}_{k,\ell} \rangle + \hat{\beta}_{k,\ell} \|\mathbf{a}_k^*\|_{\hat{\Sigma}_{k,\ell}^{-1}} \\ &\leq A^2 2^L \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 + \min_{\ell \in [L]} \langle \mathbf{a}_k^*, \hat{\boldsymbol{\theta}}_{k,\ell} \rangle + \hat{\beta}_{k,\ell} \|\mathbf{a}_k\|_{\hat{\Sigma}_{k,\ell}^{-1}} \\ &\leq A^2 2^L \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 + \langle \mathbf{a}_k^*, \hat{\boldsymbol{\theta}}_{k,\ell_{k-1}} \rangle + \hat{\beta}_{k,\ell_{k-1}} \|\mathbf{a}_k\|_{\hat{\Sigma}_{k,\ell_{k-1}}^{-1}}, \end{aligned} \quad (49)$$

where the first inequality holds because of Eq.(48), the third inequality holds because of the arm selection rule in Line 8 of Algo.2.

We decompose the regret for  $\tilde{K} \in [1, w]$  as follows

$$\begin{aligned} \text{Regret}(\tilde{K}) &= \sum_{k \in [\tilde{K}]} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle) \\ &= \sum_{\ell \in [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1,\ell}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle) + \sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1,\ell}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle) \\ &\quad + \sum_{k \in \hat{\Psi}_{\tilde{K}+1,L+1}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle). \end{aligned} \quad (50)$$

We will bound the three terms separately. For the first term, we have for layer  $\ell \in [\ell^*]$  and round  $k \in \hat{\Psi}_{\tilde{K}+1,\ell}$ , we

have

$$\begin{aligned}
 \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}^* \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}^* \rangle) &\leq 2 |\Psi_{K+1, \ell}| \\
 &= 2^{2\ell+1} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \|w_k \mathbf{a}_k\|_{\hat{\Sigma}_{k, \ell}^{-1}}^2 \\
 &\leq 2 \cdot 128^2 \log\left(\frac{4(w+1)^2 L}{\delta}\right) \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \|w_k \mathbf{a}_k\|_{\hat{\Sigma}_{k, \ell}^{-1}}^2 \\
 &\leq 2 \cdot 128^2 \log\left(\frac{4(w+1)^2 L}{\delta}\right) \cdot 2d \log\left(1 + \frac{2^{2\ell} w A^2}{d}\right) \\
 &= \tilde{O}(d),
 \end{aligned} \tag{51}$$

where the first inequality holds because the reward is in  $[-1, 1]$ , the equation follows from the fact that  $\|w_k \mathbf{a}_k\|_{\hat{\Sigma}_{k, \ell}^{-1}} = 2^{-\ell}$  holds for all  $k \in \Psi_{K+1, \ell}$ , the second inequality holds due to the fact that  $2^{\ell^*} \leq 128\sqrt{\log(4(w+1)^2 L/\delta)}$ , and the last inequality holds due to Lemma H.2.

Therefore

$$\sum_{\ell \in [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle) = \tilde{O}(d). \tag{52}$$

For the second part in Eq.(50), we have

$$\begin{aligned}
 &\sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle) \\
 &\leq \sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \left( \langle \mathbf{a}_k, \hat{\boldsymbol{\theta}}_{k, \ell-1} \rangle + \hat{\beta}_{k, \ell-1} \|\mathbf{a}_k\|_{\hat{\Sigma}_{k, \ell-1}^{-1}} \right. \\
 &\quad \left. + A^2 2^L \sqrt{dw} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle \right) \\
 &\leq 2 \sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \hat{\beta}_{k, \ell-1} \|\mathbf{a}_k\|_{\hat{\Sigma}_{k, \ell-1}^{-1}} + A^2 \sqrt{dw} \sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} 2^L \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2,
 \end{aligned} \tag{53}$$

where the inequality holds due to Eq.(49), the second inequality holds due to Eq.(48). We then try to bound the two terms.

For the first term in Eq.(53), we have

$$\begin{aligned}
 \sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \hat{\beta}_{k, \ell-1} \|\mathbf{a}_k\|_{\hat{\Sigma}_{k, \ell-1}^{-1}} &\leq \sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \hat{\beta}_{k, \ell-1} \cdot 2^{-\ell} \\
 &\leq \sum_{\ell \in [L] \setminus [\ell^*]} \hat{\beta}_{\tilde{K}, \ell-1} \cdot 2^{-\ell} |\hat{\Psi}_{\tilde{K}+1, \ell}| \\
 &= \sum_{\ell \in [L] \setminus [\ell^*]} \hat{\beta}_{\tilde{K}, \ell-1} \cdot 2^\ell \sum_{k \in \hat{\Psi}_{\tilde{K}+1, \ell}} \|w_k \mathbf{a}_k\|_{\hat{\Sigma}_{k, \ell}^{-1}}^2 \\
 &\leq \sum_{\ell \in [L] \setminus [\ell^*]} \hat{\beta}_{\tilde{K}, \ell-1} \cdot 2^\ell \cdot 2d \log\left(1 + \frac{2^{2\ell} \tilde{K} A^2}{d}\right) \\
 &= \tilde{O}(d \cdot 2^\ell \cdot \hat{\beta}_{\tilde{K}, \ell-1}) \\
 &= \tilde{O}\left(d \left(\sqrt{\sum_{k=1}^{\tilde{K}} \sigma_k^2} + R + 1\right)\right),
 \end{aligned} \tag{54}$$

where the first inequality holds because by the algorithm design, we have for all  $k \in \hat{\Psi}_{\tilde{K}+1,\ell}$ :  $\|\mathbf{a}_k\|_{\hat{\Sigma}_{k,\ell-1}^{-1}} \leq 2^{-\ell}$ ; the second inequality holds because for all  $k \in \hat{\Psi}_{\tilde{K}+1,\ell}$ ,  $\hat{\beta}_{k,\ell-1} \leq \hat{\beta}_{\tilde{K},\ell-1}$ ; the first equality holds because for all  $k \in \hat{\Psi}_{\tilde{K}+1,\ell}$ ,  $\|w_k \mathbf{a}_k\|_{\Sigma_{k,\ell}^{-1}}^2 = 2^{-2\ell}$ ; the third inequality holds by Lemma H.2; the last two equalities hold because by Lemma H.4 and Lemma H.5, we have  $\hat{\beta}_{\tilde{K},\ell-1} = \tilde{O}\left(2^{-\ell}(\sqrt{\sum_{k=1}^{\tilde{K}} \sigma_k^2} + R + 1)\right)$ .

For the second term in Eq.(53), we have

$$\begin{aligned} A^2 \sqrt{dw} \sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1,\ell}} 2^L \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 &\leq A^2 2^L \sqrt{dw} \sum_{k \in [\tilde{K}-1]} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 \\ &\leq \frac{A^2 \sqrt{dw}^{\frac{3}{2}}}{\alpha} \sum_{k=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 \end{aligned} \quad (55)$$

Therefore, with this, Eq.(53), and Eq.(54), we have

$$\sum_{\ell \in [L] \setminus [\ell^*]} \sum_{k \in \hat{\Psi}_{\tilde{K}+1,\ell}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle) \leq \frac{A^2 \sqrt{dw}^{\frac{3}{2}}}{\alpha} \sum_{k=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + \tilde{O}\left(d \left(\sqrt{\sum_{k=1}^{\tilde{K}} \sigma_k^2} + R + 1\right)\right). \quad (56)$$

Finally, for the last term in Eq.(50), we have

$$\begin{aligned} \sum_{k \in \hat{\Psi}_{\tilde{K}+1,L+1}} (\langle \mathbf{a}_k^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle) &\leq \sum_{k \in \hat{\Psi}_{\tilde{K}+1,L+1}} \left( \langle \mathbf{a}_k, \hat{\boldsymbol{\theta}}_{k,L} \rangle + \hat{\beta}_{k,L} \|\mathbf{a}_k\|_{\hat{\Sigma}_{k,L}^{-1}} + A^2 2^L \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 - \langle \mathbf{a}_k, \boldsymbol{\theta}_k \rangle \right) \\ &\leq \sum_{k \in \hat{\Psi}_{\tilde{K}+1,L+1}} \left( 2\hat{\beta}_{k,L} \|\mathbf{a}_k\|_{\hat{\Sigma}_{k,L}^{-1}} + A^2 2^{L+1} \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 \right) \\ &\leq \sum_{k \in \hat{\Psi}_{\tilde{K}+1,L+1}} \left( 2^{-L+1} \hat{\beta}_{k,L} + A^2 2^{L+1} \sqrt{dw} \sum_{s=1}^{k-1} \|\boldsymbol{\theta}_s - \boldsymbol{\theta}_{s+1}\|_2 \right) \\ &\leq \frac{2A^2 \sqrt{dw}^{\frac{3}{2}}}{\alpha} \sum_{k=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + \sum_{k \in \hat{\Psi}_{\tilde{K}+1,L+1}} 2^{-L+1} \hat{\beta}_{\tilde{K},L} \\ &\leq \frac{2A^2 \sqrt{dw}^{\frac{3}{2}}}{\alpha} \sum_{k=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + w \cdot 2\alpha \cdot \hat{\beta}_{\tilde{K},L} \\ &= \frac{2A^2 \sqrt{dw}^{\frac{3}{2}}}{\alpha} \sum_{k=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + \tilde{O}\left(w\alpha^2 \cdot \left(\sqrt{\sum_{k=1}^{\tilde{K}} \sigma_k^2} + R + 1\right)\right), \end{aligned} \quad (57)$$

where the first inequality holds due to Eq.(49), the second inequality holds due to Eq.(48), the third inequality holds because by the algorithm design, we have for all  $k \in \hat{\Psi}_{\tilde{K}+1,L+1}$ :  $\|\mathbf{a}_k\|_{\hat{\Sigma}_{k,L}^{-1}} \leq 2^{-L}$ , the fourth inequality holds due to the same reasons as before, and the fact that  $\hat{\beta}_{\tilde{K},L} \geq \hat{\beta}_{k,L}$  for all  $k \in \hat{\Psi}_{\tilde{K},L}$ ; the last inequality holds due to  $\hat{\beta}_{\tilde{K},\ell-1} = \tilde{O}\left(\alpha(\sqrt{\sum_{k=1}^{\tilde{K}} \sigma_k^2} + R + 1)\right)$ .

Plugging Eq.(56), Eq.(57), and Eq.(52) into Eq.(50), we can get that for  $\tilde{K} \in [1, w]$

$$\text{Regret}(\tilde{K}) = \tilde{O}\left(\frac{A^2 \sqrt{dw}^{\frac{3}{2}}}{\alpha} \sum_{k=1}^{\tilde{K}-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + (w\alpha^2 + d) \cdot \left(\sqrt{\sum_{k=1}^{\tilde{K}} \sigma_k^2} + R + 1\right)\right). \quad (58)$$



By the same deduction we can get

$$\text{Regret}([g_i, g_{i+1}]) = \tilde{O}\left(\frac{A^2\sqrt{d}w^{\frac{3}{2}}}{\alpha} \sum_{k=g_i}^{g_{i+1}} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + (w\alpha^2 + d) \cdot \left(\sqrt{\sum_{k=g_i}^{g_{i+1}} \sigma_k^2} + R + 1\right)\right). \quad (59)$$

Finally, without loss of generality, we assume  $K\%w = 0$ . Then we have

$$\begin{aligned} \text{Regret}(K) &= \sum_{i=0}^{\frac{K}{w}-1} \text{Regret}([g_i, g_{i+1}]) \\ &= \tilde{O}\left(\frac{A^2\sqrt{d}w^{\frac{3}{2}}}{\alpha} \sum_{i=0}^{\frac{K}{w}-1} \sum_{k=g_i}^{g_{i+1}} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + (w\alpha^2 + d) \cdot \sum_{i=0}^{\frac{K}{w}-1} \left(\sqrt{\sum_{k=g_i}^{g_{i+1}} \sigma_k^2} + R + 1\right)\right) \\ &\leq \tilde{O}\left(\frac{A^2\sqrt{d}w^{\frac{3}{2}}}{\alpha} \sum_{k=1}^{K-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 + (w\alpha^2 + d) \cdot \left(\sqrt{\frac{K}{w} \sum_{i=0}^{\frac{K}{w}-1} \sum_{k=g_i}^{g_{i+1}} \sigma_k^2} + \frac{KR}{w} + \frac{K}{w}\right)\right) \\ &\leq \tilde{O}\left(\frac{A^2\sqrt{d}w^{\frac{3}{2}}B_K}{\alpha} + (w\alpha^2 + d) \cdot \sqrt{\frac{K}{w} \sum_{k=1}^K \sigma_k^2} + (1+R) \cdot \left(K\alpha^2 + \frac{Kd}{w}\right)\right), \end{aligned}$$

where the first inequality holds due to the Cauchy-Schwarz inequality, the last inequality holds because  $\sum_{k=1}^{K-1} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}\|_2 \leq B_K$ .

## G Proof of Theorem A.1

With the candidate pool set  $\mathcal{P}$  designed as in Eq.(10), Eq.(11), Eq.(12), and  $H = \lceil d^{\frac{2}{5}} K^{\frac{2}{5}} \rceil$ , we have  $|\mathcal{P}| = O(\log K)$ , and for any  $w \in \mathcal{W}$ ,  $w \leq H$ .

We denote the optimal  $(w, \alpha)$  with the knowledge of  $V_K$  and  $B_K$  in Corollary 6.3 as  $(w^*, \alpha^*)$ . We denote the best approximation of  $(w^*, \alpha^*)$  in the candidate set  $\mathcal{P}$  as  $(w^+, \alpha^+)$ . Then we can decompose the regret as follows

$$\begin{aligned} \text{Regret}(K) &= \sum_{k=1}^K \langle \mathbf{a}_t^*, \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_t, \boldsymbol{\theta}_k \rangle = \underbrace{\sum_{k=1}^K \langle \mathbf{a}_t^*, \boldsymbol{\theta}_k \rangle - \sum_{i=1}^{\lceil \frac{K}{H} \rceil} \sum_{k=(i-1)H+1}^{iH} \langle \mathbf{a}_t(w^+, \alpha^+), \boldsymbol{\theta}_k \rangle}_{(1)} \\ &\quad + \underbrace{\sum_{i=1}^{\lceil \frac{K}{H} \rceil} \sum_{k=(i-1)H+1}^{iH} \langle \mathbf{a}_t(w^+, \alpha^+), \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_t(w_i, \alpha_i), \boldsymbol{\theta}_k \rangle}_{(2)}. \quad (60) \end{aligned}$$

The first term (1) is the dynamic regret of Restarted SAVE<sup>+</sup> with the best parameters in the candidate pool  $\mathcal{P}$ . The second term (2) is the regret overhead of meta-algorithm due to adaptive exploration of unknown optimal parameters.

By the design of the candidate pool set  $\mathcal{P}$  in Eq.(10), Eq.(11), Eq.(12), we have that there exists a pair  $(w^+, \alpha^+) \in \mathcal{P}$  such that  $w^+ < w^* < 2w^+$ , and  $\alpha^+ < \alpha^* < 2\alpha^+$ . Therefore, employing the regret bound in

Theorem 6.1, we can get

$$\begin{aligned}
 (1) &\leq \sum_{i=1}^{\lceil \frac{K}{H} \rceil} \tilde{O}(\sqrt{d}w^{+1.5}B_i/\alpha^+ + \alpha^{+2}(H + \sqrt{w^+HV_i}) + d\sqrt{HV_i/w^+} + dH/w^+) \\
 &\leq \tilde{O}(\sqrt{d}w^{+1.5}B_K/\alpha^+ + \alpha^{+2}(K + \sqrt{w^+H \sum_{i=1}^{\lceil \frac{K}{H} \rceil} V_i}) + d\sqrt{H \sum_{i=1}^{\lceil \frac{K}{H} \rceil} V_i/w^+} + dK/w^+) \\
 &= \tilde{O}(\sqrt{d}w^{+1.5}B_K/\alpha^+ + \alpha^{+2}(K + \sqrt{w^+KV_K}) + d\sqrt{KV_K/w^+} + dK/w^+) \\
 &= \tilde{O}(\sqrt{d}w^{*1.5}B_K/\alpha^* + \alpha^{*2}(K + \sqrt{w^*KV_K}) + d\sqrt{KV_K/w^*} + dK/w^*) \\
 &= \tilde{O}(d^{4/5}V_K^{2/5}B_K^{1/5}K^{2/5} + d^{2/3}B_K^{1/3}K^{2/3}), \tag{61}
 \end{aligned}$$

where we denote  $B_i$  as the total variation budget in block  $i$ ,  $V_i$  is the total variance in block  $i$ , the second inequality is by Cauchy-Schwarz inequality, the first equality holds due to  $\sum_{i=1}^{\lceil \frac{K}{H} \rceil} B_i = B_K$ ,  $\sum_{i=1}^{\lceil \frac{K}{H} \rceil} V_i = V_K$ , the second equality holds due to  $w^+ < w^* < 2w^+$  and  $\alpha^+ < \alpha^* < 2\alpha^+$ , the last equality holds by Corollary 6.3.

We then try to bound the second term (2). We denote by  $\mathcal{E}$  the event such that Lemma H.7 holds, and denote by  $R_i := \sum_{k=(i-1)H+1}^{iH} \langle \mathbf{a}_t(w^+, \alpha^+), \boldsymbol{\theta}_k \rangle - \langle \mathbf{a}_t(w_i, \alpha_i), \boldsymbol{\theta}_k \rangle$  the instantaneous regret of the meta learner in the block  $i$ . Then we have

$$\begin{aligned}
 (2) &= \mathbb{E} \left[ \sum_{i=1}^{\lceil \frac{K}{H} \rceil} R_i \right] \\
 &= \mathbb{E} \left[ \sum_{i=1}^{\lceil \frac{K}{H} \rceil} R_i | \mathcal{E} \right] P(\mathcal{E}) + \mathbb{E} \left[ \sum_{i=1}^{\lceil \frac{K}{H} \rceil} R_i | \bar{\mathcal{E}} \right] P(\bar{\mathcal{E}}) \\
 &\leq \tilde{O} \left( L_{\max} \sqrt{\frac{K}{H} |\mathcal{P}|} \right) \cdot \left(1 - \frac{2}{K}\right) + \tilde{O}(K) \cdot \frac{2}{K} \\
 &= \tilde{O}(\sqrt{H |\mathcal{P}| K}) \\
 &= \tilde{O}(d^{\frac{1}{5}} K^{\frac{7}{10}}), \tag{62}
 \end{aligned}$$

where  $L_{\max} := \max_{i \in [\lceil \frac{K}{H} \rceil]} L_i$ , the first inequality holds due to the standard regret upper bound result for Exp3 Auer et al. (2002), the third equality holds due to Lemma H.7, the last equality holds since  $H = \lceil d^{\frac{2}{5}} K^{\frac{2}{5}} \rceil$ , and  $|\mathcal{P}| = O(\log K)$ .

Finally, combining the above results for term (1) and term (2), we have

$$\text{Regret}(K) = \tilde{O}(d^{4/5}V_K^{2/5}B_K^{1/5}K^{2/5} + d^{2/3}B_K^{1/3}K^{2/3} + d^{\frac{1}{5}}K^{\frac{7}{10}}). \tag{63}$$

## H Technical Lemmas

**Theorem H.1** (Theorem 4.3, Zhou and Gu (2022)). *Let  $\{\mathcal{G}_k\}_{k=1}^\infty$  be a filtration, and  $\{\mathbf{x}_k, \eta_k\}_{k \geq 1}$  be a stochastic process such that  $\mathbf{x}_k \in \mathbb{R}^d$  is  $\mathcal{G}_k$ -measurable and  $\eta_k \in \mathbb{R}$  is  $\mathcal{G}_{k+1}$ -measurable. Let  $L, \sigma, \lambda, \epsilon > 0$ ,  $\boldsymbol{\mu}^* \in \mathbb{R}^d$ . For  $k \geq 1$ , let  $y_k = \langle \boldsymbol{\mu}^*, \mathbf{x}_k \rangle + \eta_k$  and suppose that  $\eta_k, \mathbf{x}_k$  also satisfy*

$$\mathbb{E}[\eta_k | \mathcal{G}_k] = 0, \quad \mathbb{E}[\eta_k^2 | \mathcal{G}_k] \leq \sigma^2, \quad |\eta_k| \leq R, \quad \|\mathbf{x}_k\|_2 \leq L. \tag{64}$$

For  $k \geq 1$ , let  $\mathbf{Z}_k = \lambda \mathbf{I} + \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top$ ,  $\mathbf{b}_k = \sum_{i=1}^k y_i \mathbf{x}_i$ ,  $\boldsymbol{\mu}_k = \mathbf{Z}_k^{-1} \mathbf{b}_k$ , and

$$\begin{aligned}
 \beta_k &= 12\sqrt{\sigma^2 d \log(1 + kL^2/(d\lambda)) \log(32(\log(R/\epsilon) + 1)k^2/\delta)} \\
 &\quad + 24\log(32(\log(R/\epsilon) + 1)k^2/\delta) \max_{1 \leq i \leq k} \{|\eta_i| \min\{1, \|\mathbf{x}_i\|_{\mathbf{Z}_{i-1}^{-1}}\}\} + 6\log(32(\log(R/\epsilon) + 1)k^2/\delta)\epsilon.
 \end{aligned}$$

Then, for any  $0 < \delta < 1$ , we have with probability at least  $1 - \delta$  that,

$$\forall k \geq 1, \quad \left\| \sum_{i=1}^k \mathbf{x}_i \eta_i \right\|_{\mathbf{Z}_k^{-1}} \leq \beta_k, \quad \|\boldsymbol{\mu}_k - \boldsymbol{\mu}^*\|_{\mathbf{Z}_k} \leq \beta_k + \sqrt{\lambda} \|\boldsymbol{\mu}^*\|_2.$$

**Lemma H.2** (Lemma 11, Abbasi-Yadkori et al. (2011)). For any  $\lambda > 0$  and sequence  $\{\mathbf{x}_k\}_{k=1}^K \subset \mathbb{R}^d$  for  $k \in [K]$ , define  $\mathbf{Z}_k = \lambda \mathbf{I} + \sum_{i=1}^{k-1} \mathbf{x}_i \mathbf{x}_i^\top$ . Then, provided that  $\|\mathbf{x}_k\|_2 \leq L$  holds for all  $k \in [K]$ , we have

$$\sum_{k=1}^K \min \{1, \|\mathbf{x}_k\|_{\mathbf{Z}_k^{-1}}^2\} \leq 2d \log(1 + KL^2/(d\lambda)).$$

**Theorem H.3** (Theorem 2.1, Zhao et al. (2023)). Let  $\{\mathcal{G}_k\}_{k=1}^\infty$  be a filtration, and  $\{\mathbf{x}_k, \eta_k\}_{k \geq 1}$  be a stochastic process such that  $\mathbf{x}_k \in \mathbb{R}^d$  is  $\mathcal{G}_k$ -measurable and  $\eta_k \in \mathbb{R}$  is  $\mathcal{G}_{k+1}$ -measurable. Let  $L, \sigma, \lambda, \epsilon > 0$ ,  $\boldsymbol{\mu}^* \in \mathbb{R}^d$ . For  $k \geq 1$ , let  $y_k = \langle \boldsymbol{\mu}^*, \mathbf{x}_k \rangle + \eta_k$ , where  $\eta_k, \mathbf{x}_k$  satisfy

$$\mathbb{E}[\eta_k | \mathcal{G}_k] = 0, \quad |\eta_k| \leq R, \quad \sum_{i=1}^k \mathbb{E}[\eta_i^2 | \mathcal{G}_i] \leq v_k, \quad \text{for } \forall k \geq 1$$

For  $k \geq 1$ , let  $\mathbf{Z}_k = \lambda \mathbf{I} + \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top$ ,  $\mathbf{b}_k = \sum_{i=1}^k y_i \mathbf{x}_i$ ,  $\boldsymbol{\mu}_k = \mathbf{Z}_k^{-1} \mathbf{b}_k$ , and

$$\beta_k = 16\rho \sqrt{v_k \log(4w^2/\delta)} + 6\rho R \log(4w^2/\delta),$$

where  $\rho \geq \sup_{k \geq 1} \|\mathbf{x}_k\|_{\mathbf{Z}_{k-1}^{-1}}$ . Then, for any  $0 < \delta < 1$ , we have with probability at least  $1 - \delta$  that,

$$\forall k \geq 1, \quad \left\| \sum_{i=1}^k \mathbf{x}_i \eta_i \right\|_{\mathbf{Z}_k^{-1}} \leq \beta_k, \quad \|\boldsymbol{\mu}_k - \boldsymbol{\mu}^*\|_{\mathbf{Z}_k} \leq \beta_k + \sqrt{\lambda} \|\boldsymbol{\mu}^*\|_2.$$

**Lemma H.4** (Adopted from Lemma B.4, Zhao et al. (2023)). Let weight  $w_i$  be defined in Algorithm 2. With probability at least  $1 - 2\delta$ , for all  $k \geq 1$ ,  $\ell \in [L]$ , the following two inequalities hold simultaneously:

$$\begin{aligned} \sum_{i \in \hat{\Psi}_{k+1, \ell}} w_i^2 \sigma_i^2 &\leq 2 \sum_{i \in \hat{\Psi}_{k+1, \ell}} w_i^2 \epsilon_i^2 + \frac{14}{3} R^2 \log(4w^2 L/\delta), \\ \sum_{i \in \hat{\Psi}_{k+1, \ell}} w_i^2 \epsilon_i^2 &\leq \frac{3}{2} \sum_{i \in \hat{\Psi}_{k+1, \ell}} w_i^2 \sigma_i^2 + \frac{7}{3} R^2 \log(4w^2 L/\delta). \end{aligned}$$

For simplicity, we denote  $\mathcal{E}_V$  as the event such that the two inequalities in Lemma H.4 holds.

**Lemma H.5** (Adopted from Lemma B.5, Zhao et al. (2023)). Suppose that  $\|\boldsymbol{\theta}^*\|_2 \leq B$ . Let weight  $w_i$  be defined in Algorithm 2. On the event  $\mathcal{E}_{\text{conf}}$  and  $\mathcal{E}_V$  (defined in Eq.(45), Lemma H.4), for all  $k \geq 1$ ,  $\ell \in [L]$  such that  $2^\ell \geq 64\sqrt{\log(4(w+1)^2 L/\delta)}$ , we have the following inequalities:

$$\begin{aligned} \sum_{i \in \Psi_{k+1, \ell}} w_i^2 \sigma_i^2 &\leq 8 \sum_{i \in \Psi_{k+1, \ell}} w_i^2 \left( r_i - \langle \hat{\boldsymbol{\theta}}_{k+1, \ell}, \mathbf{a}_i \rangle \right)^2 + 6R^2 \log(4(w+1)^2 L/\delta) + 2^{-2\ell+2} B^2, \\ \sum_{i \in \Psi_{k+1, \ell}} w_i^2 \left( r_i - \langle \hat{\boldsymbol{\theta}}_{k+1, \ell}, \mathbf{a}_i \rangle \right)^2 &\leq \frac{3}{2} \sum_{i \in \Psi_{k+1, \ell}} w_i^2 \sigma_i^2 + \frac{7}{3} R^2 \log(4w^2 L/\delta) + 2^{-2\ell} B^2. \end{aligned}$$

**Lemma H.6** (Freedman (1975)). Let  $M, v > 0$  be fixed constants. Let  $\{x_i\}_{i=1}^n$  be a stochastic process,  $\{\mathcal{G}_i\}_i$  be a filtration so that for all  $i \in [n]$ ,  $x_i$  is  $\mathcal{G}_i$ -measurable, while almost surely

$$\mathbb{E}[x_i | \mathcal{G}_{i-1}] = 0, \quad |x_i| \leq M, \quad \sum_{i=1}^n \mathbb{E}[x_i^2 | \mathcal{G}_{i-1}] \leq v.$$

Then for any  $\delta > 0$ , with probability at least  $1 - \delta$ , we have

$$\sum_{i=1}^n x_i \leq \sqrt{2v \log(1/\delta)} + 2/3 \cdot M \log(1/\delta).$$

**Lemma H.7.** Let  $N = \lceil \frac{K}{H} \rceil$ . Denote by  $L_i$  the absolute value of cumulative rewards for episode  $i$ , i.e.,  $L_i = \sum_{k=(i-1)H+1}^{iH} r_k$ , then

$$\mathbb{P} \left[ \forall i \in [N], L_i \leq H + R \sqrt{\frac{H}{2} \log \left( K \left( \frac{K}{H} + 1 \right) \right)} + \frac{2}{3} \cdot R \log \left( K \left( \frac{K}{H} + 1 \right) \right) \right] \geq 1 - \frac{1}{K}. \quad (65)$$

*Proof.* By Lemma H.6, we have that with probability at least  $1 - 1/K$

$$\begin{aligned}
 \sum_{k=(i-1) \cdot H+1}^{i \cdot H} \epsilon_i &\leq \sqrt{2 \sum_{k=(i-1) \cdot H+1}^{i \cdot H} \sigma_k^2 \log(NK) + 2/3 \cdot R \log(NK)} \\
 &\leq \sqrt{2H \frac{R^2}{4} \log(NK) + 2/3 \cdot R \log(NK)} \\
 &\leq R \sqrt{\frac{H}{2} \log \left( K \cdot \left( \frac{K}{H} + 1 \right) \right) + \frac{2}{3} \cdot R \log \left( K \cdot \left( \frac{K}{H} + 1 \right) \right)}, \tag{66}
 \end{aligned}$$

where we use union bound, and in the second inequality we use the fact that since  $|\epsilon_k| \leq R$ , we have  $\sigma_k^2 \leq \frac{R^2}{4}$ . Finally, together with the assumption that  $r_k \leq 1$  for all  $k \in [K]$ , we complete the proof.  $\square$