

MBZUAI

Digital.Commons@MBZUAI

Machine Learning Faculty Publications

Scholarly Works

6-6-2022

Learning to Control under Time-Varying Environment

Yuzhen Han

Ruben Solozabal

Jing Dong

Xingyu Zhou

Martin Takac

See next page for additional authors

Follow this and additional works at: <https://dclibrary.mbzuai.ac.ae/mlfp>



Part of the [Artificial Intelligence and Robotics Commons](#)

Open access version thanks to arXiv

License: CC BY 4.0

Uploaded July 05, 2022

Authors

Yuzhen Han, Ruben Solozabal, Jing Dong, Xingyu Zhou, Martin Takac, and Bin Gu

Learning to Control under Time-Varying Environment

Yuzhen Han

Department of Mechanical & Industrial Engineering
University of Toronto
Toronto, ON CANADA

Ruben Solozabal

MBZUAI
Masdar City, Abu Dhabi, UAE
ruben.solozabal@mbzuai.ac.ae

Jing Dong

The Chinese University of Hong Kong
Shenzhen, China
jingdong@link.cuhk.edu.cn

Xingyu Zhou

Wayne State University
Detroit, USA
xingyu.zhou@wayne.edu

Martin Takáč

MBZUAI
Masdar City, Abu Dhabi, UAE
Takac.MT@gmail.com

Bin Gu

MBZUAI
Masdar City, Abu Dhabi, UAE
bin.gu@mbzuai.ac.ae

Abstract

This paper investigates the problem of regret minimization in linear time-varying (LTV) dynamical systems. Due to the simultaneous presence of uncertainty and non-stationarity, designing online control algorithms for unknown LTV systems remains a challenging task. At a cost of NP-hard offline planning, prior works have introduced online convex optimization algorithms, although they suffer from nonparametric rate of regret. In this paper, we propose the first computationally tractable online algorithm with regret guarantees that avoids offline planning over the state linear feedback policies. Our algorithm is based on the optimism in the face of uncertainty (OFU) principle in which we optimistically select the best model in a high confidence region. Our algorithm is then more explorative when compared to previous approaches. To overcome non-stationarity, we propose either a restarting strategy (R-OFU) or a sliding window (SW-OFU) strategy. With proper configuration, our algorithm attains sublinear regret $O(T^{2/3})$. These algorithms utilize data from the current phase for tracking variations on the system dynamics. We corroborate our theoretical findings with numerical experiments, which highlight the effectiveness of our methods. To the best of our knowledge, our study establishes the first model-based online algorithm with regret guarantees under LTV dynamical systems.

1 Introduction

Regret minimization in online control has been extensively investigated in the context of either unknown time-invariant or known time-varying dynamics systems. Yet real applications such as dynamic pricing and ad allocations call for the need for an unknown time-varying system. Under such a setting, the problems become significantly more challenging due to the coexistence of non-stationarity and uncertainty. Despite previous attempts on unknown LTV on stable controllers [1] or system identification [2], it remains open whether an algorithm can achieve meaningful regret guarantees in this scenario. This paper thus addresses the problem of minimizing the cumulative

regret in unknown LTV systems

$$\mathcal{R}_T = \sum_{t=1}^T (\min_a f(a) - f(a_t)), \quad (1)$$

where f is the (cost) function being optimized, and a_t is the action chosen at time t . To the best of our knowledge, this is the first work that achieves regret guarantees with a computationally tractable algorithm.

When the system is time-invariant, the regret minimization (1) problem has been well studied. With offline simulations, numerous existing results achieve sublinear regret (i.e. $O(\sqrt{T})$) [3–8]. By further encouraging exploration with intrinsic noise from the system dynamics, [9, 10] achieve a logarithmic regret of $O(\log^2 T)$. Recent work [11] presents a finite-time sublinear regret from a single chain of black-box interactions without access to offline simulations. With online convex optimization (OCO) and structured memory [12] achieves a constant, dimension-free competitive ratio of regret.

The regret minimization problem is complicated when the system is time-varying. This encapsulates a wide range of possible scenarios, the dynamics can be slowly changing or abruptly changing. With the knowledge of the system dynamic, several approaches that have investigated (1) in LTV systems, summarized in Table 1. Specifically, [13] studies regret for predictive control of LTV systems. The iGPC [14] utilizes a nested-OCO formulation to design an iterative algorithm for minimizing planning regret in the presence of a time-varying system. Similarly, [15] adopts OCO with memory to minimize the adaptive regret, which is the supremum of the local regret (with respect to the local optimal comparator) over all contiguous intervals in time.

While promising results are presented in LTV with known system dynamics, such requirement is often too stringent, if not impossible to fulfill in real applications. Yet the uncertainty of the system dynamic poses new challenges for algorithm designs and regret guarantees. To the best of our knowledge, there is only one work that addresses (1) on unknown LTV environments [16]. The work achieves a sublinear regret bound for convex parametrization policies. However, for the class of linear state policies ($u = Kx$), the regret is proportional to $\exp(\Omega(nm))$, where n, m are the dimension of state and action spaces. This reveals the impractical nature of the algorithm as it can be intractable for a wide range of problems. Furthermore, the algorithms rely on an offline planning procedure over the entire state linear feedback policies. While this is possible in a linear time-invariant system, which admits efficient convex relaxations, this is NP-hard in LTV with unknown dynamics. The following question remains open.

Does there exist a computationally tractable regret minimization algorithm for LTV with unknown system dynamic?

In this paper, we study the regret minimization problem on LTV with unknown system dynamics and answer the above question affirmatively. We propose Restarting based OFU (R-OFU) and Sliding Window based OFU (SW-OFU) algorithms to find a class of linear feedback policies for minimizing the long-term cumulative dynamic regret [13] across episodes. We note that our objective is thus different from adaptive regret, which focuses on (worst case) regret over intervals [15, 16], and remains different from planning regret [14]. Both of our algorithms are based on the optimism in the face of uncertainty (OFU) principle [17, 18]. This encourages our algorithms to explore for optimal solutions given the current estimation of the system dynamics. We further verify that R-OFU and SW-OFU are computationally tractable. This is because only a mini-batch of historical data in the current epoch (or sliding window) is utilized for online planning. With proper configuration, our algorithm attains sublinear regret $O(T^{2/3})$.

We further demonstrate the versatility and practicality of our algorithm with extensive experiments on switching and time-variant systems. Our empirical results corroborate our theoretical findings with respect to the regret and cost.

Paper Structure. We first present the necessary definitions and problem formulation in Section 2. Then, we present the detailed algorithms in Section 3. In Section 4, we provide the main results with respect to the R-OFU and SW-OFU algorithms. The analysis as well as proof sketches of the proposed algorithms are presented in Section 5. Lastly, in Section 6, the results and details of the experiments are provided. Proofs and other details are deferred to the Appendix.

Table 1: Representative (online) control algorithms for regret minimization.

Ref.	Dynamic	Environment	Type of Regret	Knowledge	Regret Bound (in terms of T)
[9]	Linear	Time-Inv.	Cumulative Regret	Partial	$O(\log^2 T)$
[6]	Linear	Time-Inv.	Cumulative Regret	No	$O(\sqrt{T})$
[11]	Linear	Time-Inv.	Cumulative Regret	Yes	$\tilde{O}(T^{2/3})$
[19]	Nonlinear	Time-Inv.	Cumulative Regret	Partial	$O(\sqrt{T})$
[20]	Nonlinear	Time-Inv.	Mistake	Partial	-----
[13]	Linear	Time-Var.	Dynamic Regret	Yes	$O(\lambda^k T)$
[15]	Linear	Time-Var.	Adaptive Regret	Yes	-----
[14]	Nonlinear	Time-Var.	Planning Regret	Partial	-----
[16]	Linear	Time-Var.	Adaptive Regret	No	$O(e^{\Omega(dx du)} T^{1 - \frac{1}{2(dx du + 3)}})$
Ours	Linear	Time-Var.	Dynamic Regret	No	$\tilde{O}(T)$ (epoch < episode length)
Ours	Linear	Time-Var.	Dynamic Regret	No	$\tilde{O}(T^{2/3})$ (epoch \geq episode length)

2 Problem Setting

Notation We use $\|A\|_F = \sqrt{\langle A, A \rangle_F} = \sqrt{\text{Trace}\langle A * A \rangle}$ to denote the Frobenius norm of matrix A . For two matrices X and Y , we also define $\|X\|_Y^2 = \text{Trace}(X^\top Y X)$. $\mathbb{E}[X]$ denotes the expectation of a random variable X and $x \vee y$ denotes the maximum between $x, y \in \mathbb{R}$.

2.1 Problem Formulation

We consider the episodic time-varying linear quadratic regulator (LQR) setting with K episodes and H steps. We let $x \in \mathcal{X} \in \mathbb{R}^n$ denotes the vector of system state, $u \in \mathcal{U} \in \mathbb{R}^m$ denotes the vector of control input, and w_t denotes the system noise, which is zero-mean. In each episode k , the agent starts from a random initial state sampled from the initial distribution $x_{k,h=1} \sim \rho$, and executes $H - 1$ control steps to finish the episode. Then the agent starts over from $h = 1$ for the $(k + 1)$ -th episode with a new initial state $x_{k+1,h=1}$ sampled from ρ . This process repeats for the specified number of episodes K . The dynamic of the k -th episode on the time-varying LQR system is described as

$$x_{k,h+1} = A_{k,h}x_{k,h} + B_{k,h}u_{k,h} + w_{k,h}, \quad (2)$$

with a quadratic cost function $c_{k,h} = x_{k,h}^\top Q x_{k,h} + u_{k,h}^\top R u_{k,h}$. This dynamic is governed by unknown time-varying matrices $A_{k,h}$ and $B_{k,h}$, while Q, R are known positive definite matrices. The key difference between the non-stationary equation (2) and existing stationary LQR learning systems (i.e. $\hat{x}_{k,h+1} = A\hat{x}_{k,h} + Bu_{k,h} + w_{k,h}$) [21, 22] is that the transition matrix $A_{k,h}$ and $B_{k,h}$ evolve with the time step h and the episode k . Remark that the dynamics vary between different episodes, which makes information non-transferable between them.

The goal is to design a control policy $\pi : [H] \times \mathcal{X} \rightarrow \mathcal{U}$ that minimizes the accumulated expected cost within each episode $k \in [K]$

$$J_{k,h}^\pi(x) = \mathbb{E}_\pi \left[\sum_{h'=h}^H c_{k,h'} | x_{k,h} = x \right], \quad (3)$$

where \mathbb{E}_π denotes expectation over the random trajectories generated by π starting from x at (k, h) .

Let $\Theta_{k,h}^* = [A_{k,h}, B_{k,h}]^\top \in \mathbb{R}^{(n+m) \times n}$. The optimal policy π^* can then be expressed as

$$\pi_{k,h}^* = K_{k,h}(\Theta_{k,h}^*)x_{k,h}, \quad (4)$$

where $K_{k,h}(\Theta_{k,h}^*)$ is the gain of the control policy

$$K_{k,h}(\Theta_{k,h}^*) = -(R + B_{k,h}^\top P_{k,h}(\Theta_{k,h}^*) B_{k,h})^{-1} B_{k,h}^\top P_{k,h}(\Theta_{k,h}^*) A_{k,h}, \quad (5)$$

and $P_{k,h}(\Theta_{k,h}^*)$ is the solution to the Riccati equation [23].

The optimal cost is thus given by

$$J_{k,h}^*(\Theta_{k,h}^*, x) = x^\top P_{k,h}(\Theta_{k,h}^*) x + \sum_{h'=h}^H \mathbb{E} \left[w_{k,h'}^\top P_{k,h'+1}(\Theta_{k,h}^*) w_{k,h'} \right].$$

¹We make Θ^* and Θ_* equivalent and interchangeable through the whole paper.

Intuitively, controlling such system is intractable, the natural choice is playing zero control input $u_{k,h} = 0$. However, we assume that the system dynamic evolves slowly according to the following Assumption 2.3. Therefore, at each time step, an optimal controller is optimistically computed based on the current estimation, similar to a LTI system.

The agent's performance over K episodes is measured by the cumulative (pseudo) dynamic regret \mathcal{R} with respect to the true system dynamics of the model, which is also time dependent. Formally, this is referred to as the dynamic regret.

Definition 2.1 (Dynamic regret). Over K episodes, the dynamic regret of an agent is

$$\mathcal{R}(K) = \sum_{k=1}^K J_1^{\pi_k}(\Theta_k^*, x_{k,1}) - J_1^*(\Theta_k^*, x_{k,1}), \quad (6)$$

where $J_1^{\pi_k}(\Theta_k^*, x_{k,1})$ is the expected cost under chosen π_k at the episode k , $\Theta_k^* = [\Theta_{k,1}^*, \Theta_{k,2}^*, \dots, \Theta_{k,H}^*]$, and $J_1^*(\Theta_k^*, x_{k,1})$ is the expected cost under optimal control policy for the episode k .

We make the following assumptions on controllability and boundedness to make this problem tractable. Note that similar assumptions are also used in the literature [21, 22, 8].

Assumption 2.2. The true system Θ^* is controllable and open-loop stable (i.e., $\text{Rank} \left(\begin{bmatrix} B_{k,h} & A_{k,h} B_{k,h} & A_{k,h}^2 B_{k,h} & \dots & A_{k,h}^{n-1} B_{k,h} \end{bmatrix} \right) = n$) and bounded $\|\Theta^*\|_F \leq 1$. There also exists constants v, v_A, v_B , and v_w such that $\|A_{k,h}\| \leq v_A < 1$, $\|B_{k,h}\| \leq v_B < 1$, $\|w_{k,h}\|_2 \leq v_w < 1$, and $\|R\|, \|Q\| \leq v$. For $k \geq 1$, the states $\|x_{k,1}\| \leq 1$. Further, $v_w + \Upsilon v_B + v_A \leq 1$ with Υ being a constant.

Assumption 2.3. We assume that the total system variability on every episode is bounded,

$$\sum_{h=1}^{H-1} \|\Theta_{k,h+1} - \Theta_{k,h}\|_F \leq \mathcal{B}_H, \quad \forall k \in K.$$

Assumption 2.4. Let $\{\mathcal{F}_{k,h}\}_{h=0}^\infty$ be a filtration generated by the random variables $\{x_{k,h}, u_{k,h}\}_{h=1}^\infty$. We assume that $\{w_{k,h}\}_{h \geq 1}$ is a vector valued martingale process adapted to filtration $\{\mathcal{F}_{k,h}\}_{h \geq 0}$. Further, let η_t be a sub-Gaussian random vector with a fixed constant $R > 0$, and for any $\chi \in \mathbb{R}^n$,

$$\mathbb{E} [\exp(\chi^\top w_{k,h}) | \mathcal{F}_{k,h-1}] \leq \exp\left(\frac{R^2 \|\chi\|^2}{2}\right), \quad \forall h \geq 1.$$

3 Algorithms

In this section, we propose R-OFU and SW-OFU algorithms to minimize dynamic regret \mathcal{R} under LTV systems. Both algorithms conduct planning and policy execution in a fully online fashion. In the online planning step, the algorithm estimates the current Θ_h^* based on historical data from the current phase with restarting (R) or sliding Window (SW) strategies. In the policy execution step, we apply greedy policy search with optimism in the face of uncertainty (OFU). Specifically, a better model estimation (in terms of cost) is searched under a confidence region and the model with the best estimated dynamics is selected for solving the Riccati equation.

3.1 Online Planning

The key ingredients of our online planning phase are the restarting and sliding window strategy, which allows us to only use the data from the current epoch for estimating Θ_h^* . This thus greatly reduces the computation overhead and allows for a tractable algorithm.

To shorthand the notation, we write the system parameters $z_{k,h} = [x_{k,h}^\top, u_{k,h}^\top]^\top$, $Z_{k,h} = [z_{k,h}^\top]$, $X_{k,h}^{next} = [x_{k,h+1}^\top]$, and $W_{k,h} = [w_{k,h}^\top]$ for step $h \in [H]$ in the episode $k \in [K]$. Also, in the following paper we abbreviate the nomenclature when referring to any episode k as $x_h = x_{k,h}$, similarly we define $z_h, X_h, X_h^{next}, Z_h$ and W_h .

Restarting (R): Within each episode, the restarting least-square ridge regression estimator is implemented using the historical data in the current epoch,

$$\Theta_h = \arg \min_{\Theta} \|\Theta\|_{\lambda I}^2 + \sum_{s=h_0}^{h-1} \|X_s^{next} - Z_s \Theta\|_F^2, \quad (7)$$

where h_0 is the starting point of the current epoch. Then, Θ_h admits a closed-form solution

$$\Theta_h = \mathcal{V}_h^{-1} \mathcal{U}_h = \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top X_s^{\text{next}} \right).$$

Sliding Window (SW): Consider a sliding window of length \mathcal{W} , $(1 \vee (h - \mathcal{W})) : (h - 1)$, with observation history $\{(Z_s, X_s^{\text{next}})\}_{s=1 \vee (h - \mathcal{W})}^{h-1}$, the sliding window least-square ridge regression estimator is defined as

$$\Theta_h = \arg \min_{\Theta} \|\Theta\|_{\lambda I}^2 + \sum_{s=1 \vee (h - \mathcal{W})}^{h-1} \|X_s^{\text{next}} - Z_s \Theta\|_F^2. \quad (8)$$

Similar to the closed form solution of (8), the solution of the SW estimator is $\tilde{\mathcal{V}}_h^{-1} \tilde{\mathcal{U}}_h$ where

$$\tilde{\mathcal{V}}_h = \sum_{s=1 \vee (h - \mathcal{W})}^{h-1} Z_s^\top Z_s + \lambda I, \quad \tilde{\mathcal{U}}_h = \sum_{s=1 \vee (h - \mathcal{W})}^{h-1} Z_s^\top X_s^{\text{next}}$$

Comparing the restarting and sliding window strategy. The restarting and sliding window strategies are two common strategies used in non-stationary online estimation literature [24, 25, 16]. Both strategies are depicted in Figure 1. Specifically, the restarting R strategy within each epoch, it discards data and re-identifies the model.

In contrast, SW draws and throws out data continuously using a sliding window. Therefore, the R adapts better in abruptly changing systems, especially with a given detecting mechanism [25]. Though the sliding window strategy achieves better performance in slowly changing dynamics. This phenomenon will be further discussed throughout the experiments in Section 6.

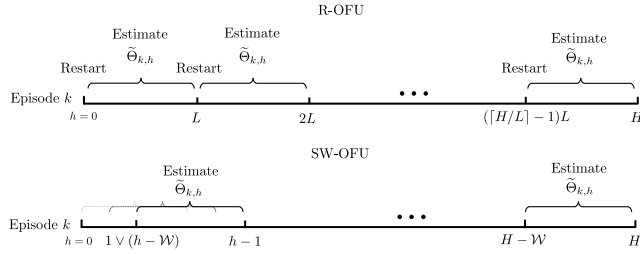


Figure 1: Comparison between R-OFU and SW-OFU.

3.2 Policy Execution

We integrate into our algorithms the OFU principle during the policy execution step. Therefore, after estimating the dynamics of the system, we optimistically select the best model within a confidence interval around the initial estimation. This allows our algorithm to explore in the uncertain environment [26]. We then greedily select our action with respect to our chosen model using (4). The overall description of the methods is summarized in Algorithm 1 and Algorithm 2.

High Probability Confidence Set. Based on the estimation of Θ_h obtained in the planning step (section 3.1), we construct a high probability confidence set for the system model Θ^* . Inspired by the analysis of [22], we design the confidence set as follows.

Lemma 3.1. *For any $h \in [H]$ of an episode and $\delta \in (0, 1)$, with probability at least $1 - \delta$, the estimation error is upper bounded as*

$$\|\Theta_h^* - \Theta_h\|_{\mathcal{V}_h} \leq \zeta_h(\delta), \quad \|\Theta_h^* - \Theta_h\|_{\tilde{\mathcal{V}}_h} \leq \tilde{\zeta}_h(\delta). \quad (9)$$

where

$$\begin{aligned} \zeta_h(\delta) &= \sqrt{\lambda} + v_w \sqrt{2 \ln\left(\frac{2}{\delta}\right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}} + \frac{\sqrt{L(m+n)}}{\sqrt{\lambda}} \mathcal{B}_H, \\ \tilde{\zeta}_h(\delta) &= \sqrt{\lambda} + v_w \sqrt{2 \ln\left(\frac{2H}{\delta}\right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}} + \frac{\sqrt{\mathcal{W}(m+n)}}{\sqrt{\lambda}} \mathcal{B}_H. \end{aligned} \quad (10)$$

are the radius of confidence region for R with length L , SW with length \mathcal{W} respectively and v_w is from Assumption 2.2.

With Lemma 3.1 and the estimator Θ_h from the online planning step, our algorithm maintains confidence radius,

$$\mathcal{C}_h(\delta) = \{\Theta : \|\Theta - \Theta_h\|_{\mathcal{V}_h} \leq \zeta_h(\delta)\}, \quad \tilde{\mathcal{C}}_h(\delta) = \{\Theta : \|\Theta - \Theta_h\|_{\tilde{\mathcal{V}}_h} \leq \tilde{\zeta}_h(\delta)\}. \quad (11)$$

OFU-Based Action Search.

Algorithm 1 R-OFU based online control algorithm

Require: Number of episodes K , time horizon H , epoch size L , regularization strength λ

```
1: for Episode  $k = 1, 2, \dots, K$ ; do
2:   Set epoch counter  $j = 1$ 
3:   while  $j \leq \lceil H/L \rceil$  do
4:     Set  $\kappa = (j - 1)L$  and initialize  $\mathcal{V}_\kappa = \lambda I$ 
5:     for  $h = \kappa + 1, \dots, \kappa + L - 1$  do
6:       Compute  $\Theta_h = \mathcal{V}_h^{-1} \mathcal{U}_h$  with  $\zeta_h(\delta)$  computed from (10)
7:       Construct high confidence set  $\mathcal{C}_h(\delta)$  and select  $\Theta_h \in \arg \min_{\Theta \in \mathcal{C}_h(\delta)} J_1^*(\Theta, x_{k,1})$ 
8:       Implement control  $u_{k,h} = K_h(\Theta_h)x_{k,h}$  and observe cost  $c_{k,h}$ ,  $Z_{k,h}$  and  $X_{k,h}^{next}$ 
9:     end for
10:    Set  $j = j + 1$ 
11:  end while
12: end for
```

Algorithm 2 SW-OFU based online control algorithm

Require: Number of episodes K , time horizon H , sliding window size \mathcal{W} , regularization strength λ

```
1: for Episodes  $k = 1, 2, \dots, K$ ; do
2:   Initialize  $\tilde{\mathcal{V}}_{k,0} = \lambda I$ 
3:   for  $h = 1, \dots, H$  do
4:     Compute  $\Theta_h = \tilde{\mathcal{V}}_h^{-1} \tilde{\mathcal{U}}_h$  with set  $\tilde{\zeta}_h$  computed from (10)
5:     Construct high confidence set  $\mathcal{C}_h(\delta)$  and select  $\Theta_h \in \arg \min_{\Theta \in \mathcal{C}_h(\delta)} J_1^*(\Theta, x_{k,1})$ 
6:     Implement control  $u_{k,h} = K_h(\Theta_h)x_{k,h}$  and observe cost  $c_{k,h}$ ,  $Z_{k,h}$  and  $X_{k,h}^{next}$ 
7:   end for
8: end for
```

Within the confidence set $\mathcal{C}_h(\delta)$ or $\tilde{\mathcal{C}}_h(\delta)$, we adopt the OFU principle to compute an optimistic estimate of $\tilde{\Theta}_h$,

$$\tilde{\Theta}_h \in \arg \min_{\Theta \in \mathcal{C}_h(\delta)} J_1^*(\Theta, x_{k,1}) \quad (12)$$

where $J_1^*(\Theta, x_{k,1})$ is the optimal cost when the true dynamic is Θ . Then, the agent computes the control following the policy

$$u_h = \pi_h(x_{k,h}) = K_h(\tilde{\Theta}_h)x_{k,h}, \quad (13)$$

where the gain $K_h(\tilde{\Theta}_h)$ can be calculated through (5).

To ensure that equation (13) is well-defined and satisfies the stability condition, we establish the following propositions. The detailed proofs and discussion is deferred to the Appendix A.

Proposition 3.2. *The region encompassed by the high probability confidence set (11) is closed and bounded.*

Proposition 3.3. *Given any $\tilde{\Theta}_h$ in (12), the gain of the controller $K_h(\tilde{\Theta}_h)$ is well defined.*

4 Main Results and Analysis

With the R-OFU and SW-OFU algorithms, we obtain the following dynamic guarantees for the unknown time-varying LTV system.

Theorem 4.1 (Dynamic regret with R-OFU). *Algorithm (1) achieves a high probability dynamic regret bound*

$$\begin{aligned} \mathcal{R}(K) = & O\left(H^{\frac{3}{2}}\sqrt{K}\right) + O\left(HK\vartheta(\delta)\sqrt{(n+m)\ln\left(1 + \frac{HK}{(n+m)\lambda}\right)}\right) \\ & + O\left(HK\left(\ln\frac{1}{\delta} + n(n+m)\ln\left(1 + \frac{HK}{(n+m)\lambda}\right)\right)\right), \end{aligned}$$

where $\delta \in (0, 1)$ is the probability parameter, $\vartheta(\delta) = \sqrt{\lambda} + \sqrt{\frac{L(m+n)}{\lambda}} \mathcal{B}_H$, and L is the epoch size.

Theorem 4.2 (Dynamic regret with SW-OFU). *Algorithm (2) achieves a high probability dynamic regret bound*

$$\begin{aligned} \mathcal{R}(K) = & O\left(H^{\frac{3}{2}}\sqrt{K}\right) + O\left(HK\tilde{\vartheta}(\delta)\sqrt{(n+m)\ln\left(1+\frac{HK}{(n+m)\lambda}\right)}\right) \\ & + O\left(HK\left(\ln\frac{H}{\delta} + n(n+m)\ln\left(1+\frac{HK}{(n+m)\lambda}\right)\right)\right), \end{aligned}$$

where $\tilde{\vartheta}(\delta) = \sqrt{\lambda} + \sqrt{\frac{\mathcal{W}(m+n)}{\lambda}} \mathcal{B}_H$, $\delta \in (0, 1)$ is the probability parameter, and \mathcal{W} is the sliding window size.

Corollary 4.3. *From Theorems 4.1 and 4.2, the dynamic regret is sublinear in K when the sliding window size or the restarting epoch length is set to be larger than H .*

Theorem 4.4 (Dynamic regret with with a larger size of epoch). *Our R-OFU algorithm achieves a high probability dynamic regret bound with a larger size of $L \geq H$*

$$\mathcal{R}(K) = \tilde{O}\left(L\mathcal{B}_{HK} + HK\sqrt{\frac{1}{L}} + \sqrt{HK}\sqrt{\mathcal{B}_{HK}L^{1/4}}\right),$$

where \mathcal{B}_{HK} is the total variation budget along the whole steps. By setting $L = L^* = (HK)^{2/3}\mathcal{B}_{HK}^{-2/3}$, we achieve a minmax near-optimal dynamic regret $\tilde{O}\left((HK)^{2/3}\mathcal{B}_{HK}^{1/3}\right)$.

Remark 4.5. From Theorems 4.1 and 4.2, it is noted that R-OFU and SW-OFU achieve the same order of regret in terms of $T = HK$, while R-OFU is slightly better with a factor of $\ln H$. The additional $\ln H$ factor comes from the information loss due to the SW.

Remark 4.6. When compared with prior results, our regret bound is much more practical. In previous work [16], the regret is $\Omega\left(\exp(nm)T^{1-\frac{1}{2(nm+3)}}\right)$, which scales exponentially with the dimension of state and action space. In practice, achieving the regret bound can be computationally intractable. In contrast, the Theorems 4.1 and 4.2 achieves regret independent of the state and action space size. Considering that attaining polynomial regret (e.g., $T^{1-\alpha}$, $\alpha > 0$) may not be even possible for unknown LTV, our results achieve a reasonable order of $T^{1.5}$ when $L < H$. Additionally, this regret bound can be further improved to $T^{2/3}$ under proper configurations.

5 Analysis

In this section, we present the analysis of Lemma 3.1 and Theorems 4.1 and 4.2.

5.1 Analysis of High Confidence Set

As shown in the closed-form expression of R and SW regressors, the key difference in the solutions is the term $h_0 = \max(1, h - \mathcal{W})$ in SW. We present details on the construction of a high confidence set for the restarting strategy, since it can also be applied to sliding window case with minor changes.

Proposition 5.1. *From the closed-form solution of (7), the estimate error can be decomposed as,*

$$\|\Theta_h^* - \Theta_h\|_{\mathcal{V}_h} \leq \underbrace{\left\|(\lambda I)^{\frac{1}{2}}\right\|_2}_{\ell_1} + \underbrace{\left\|\sum_{s=h_0}^{h-1} Z_s^\top W_s\right\|_{\mathcal{V}_h^{-1}}}_{\ell_2} + \underbrace{\left\|\sum_{s=h_0}^{h-1} Z_s^\top Z_s(\Theta_s^* - \Theta_h^*)\right\|_{\mathcal{V}_h^{-1}}}_{\ell_3}. \quad (14)$$

The detailed result is described in Appendix B.1.

Remark 5.2. The term ℓ_1 and ℓ_2 are the estimate errors caused by the regularizer and random noise; while the last term ℓ_3 is due to the time-varying property. Both R and SW have these three sources of estimate errors. The first and third terms are from the same bound for both R and SW. However, the bound for the second term is different among them.

The terms ℓ_1 , ℓ_2 and ℓ_3 can be bounded separately, as we summarized in the following lemmas 5.3 and 5.4. The proof of the lemmas can be found in Appendix B.

Lemma 5.3. *For any $h \in H$ in an episode and $\delta \in (0, 1)$ in R , with probability at least $1 - \delta$, the following holds*

$$\ell_1 = \sqrt{\lambda}, \quad \ell_2 \leq v_w \sqrt{2 \ln \left(\frac{1}{\delta} \right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}}, \quad \ell_3 \leq \frac{\sqrt{L(m+n)}}{\sqrt{\lambda}} \mathcal{B}_H.$$

Lemma 5.4. *For any $h \in H$ in an episode and $\delta \in (0, 1)$ in SW , with probability at least $1 - \delta$, the following holds*

$$\ell_1 = \sqrt{\lambda}, \quad \ell_2 \leq v_w \sqrt{2 \ln \left(\frac{H}{\delta} \right) + n \ln \frac{\det(\tilde{\mathcal{V}}_h)}{\det(\lambda I)}}, \quad \ell_3 \leq \frac{\sqrt{\mathcal{W}(m+n)}}{\sqrt{\lambda}} \mathcal{B}_H.$$

5.2 Analysis of Dynamic Regret

Armed with our analysis of high confidence set, we are now able to give an upper bound of the dynamic regret with R-OFU and SW-OFU algorithms. We first start with a careful decomposition of the dynamic regret under the good event where $\varepsilon_K(\delta) = \{\Theta_* \in \mathcal{C}_h(\delta), \forall h \in [H]\}$.

Lemma 5.5. *Let $\tilde{P}_{k,h} = P_h(\tilde{\Theta}_{k,h})$ and $\mathcal{F}_{k,h}$ denotes all randomness before the step (k, h) . Under a 'good' event $\varepsilon_K(\delta) = \{\Theta_* \in \mathcal{C}_h(\delta), \forall h \in [H]\}$, the dynamic regret $\mathcal{R}(K)$ in (6) is decomposed as*

$$\mathcal{R}(K) \leq \sum_{k=1}^K \sum_{h=1}^H \varsigma_{h,k},$$

where

$$\begin{aligned} \varsigma_{k,h} &= \mathbb{E}[J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) | \mathcal{F}_{k,h}] - J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) + \|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} - \|\tilde{\Theta}_{k,h}^T z_{k,h}\|_{\tilde{P}_{k,h+1}} \\ &\quad - \mathbb{E} \left[\|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} | \mathcal{F}_{k,h} \right] + \|\Theta_*^T z_{k,h}\|_{\tilde{P}_{k,h+1}}. \end{aligned}$$

The detailed proof of Lemma 5.5 is presented in Appendix C. Based on this decomposition, one can bound (6) separately using the following lemma for R strategy.

Lemma 5.6. *Under Assumption 1 and event $\varepsilon_K(\delta)$, we have the following dynamic regret bound with at least probability $1 - 2\delta$,*

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}[J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) | \mathcal{F}_{k,h}] - J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) &\leq O \left(\sqrt{KH^3 \ln \frac{2}{\delta}} \right), \\ \sum_{k=1}^K \sum_{h=1}^H \|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} - \mathbb{E} \left[\|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} | \mathcal{F}_{k,h} \right] &\leq O \left(\sqrt{KH \ln \frac{2}{\delta}} \right), \\ \sum_{k=1}^K \sum_{h=1}^H \|\Theta_*^T z_{k,h}\|_{\tilde{P}_{k,h+1}} - \|\tilde{\Theta}_{k,h}^T z_{k,h}\|_{\tilde{P}_{k,h+1}} &\leq O \left(HK \zeta_h(\delta) \sqrt{\ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}} \right). \end{aligned}$$

We present the proof sketch for Theorem 4.1 with algorithm R-OFU. In the case of the SW-OFU algorithm, the Proof of Theorem 4.2 is similar but with the difference in the radius term $\tilde{\zeta}_h(\delta)$.

Proof. By the boundness results in Appendix A, we have

$$\ln \det(\mathcal{V}_h) \leq (n+m) \ln \left(\lambda + \frac{HK(1+\gamma)^2}{n+m} \right).$$

Therefore, $\zeta_h(\delta)$ can be rewritten as

$$\zeta_h(\delta) = v_w \sqrt{2 \ln \left(\frac{2}{\delta} \right) + (n+nm) \ln \left(1 + \frac{HK(1+\gamma)^2}{\lambda(n+m)} \right)} + \sqrt{\lambda} + \frac{\sqrt{L(m+n)}}{\sqrt{\lambda}} \mathcal{B}_H$$

Replacing $\zeta_h(\delta)$ into the third inequality in Lemma 5.6 and putting everything together yields the final result. \square

Followed by the Proof of Theorem 4.1 and 4.2, we present the Proof of Theorem 4.4 in the Appendix E.

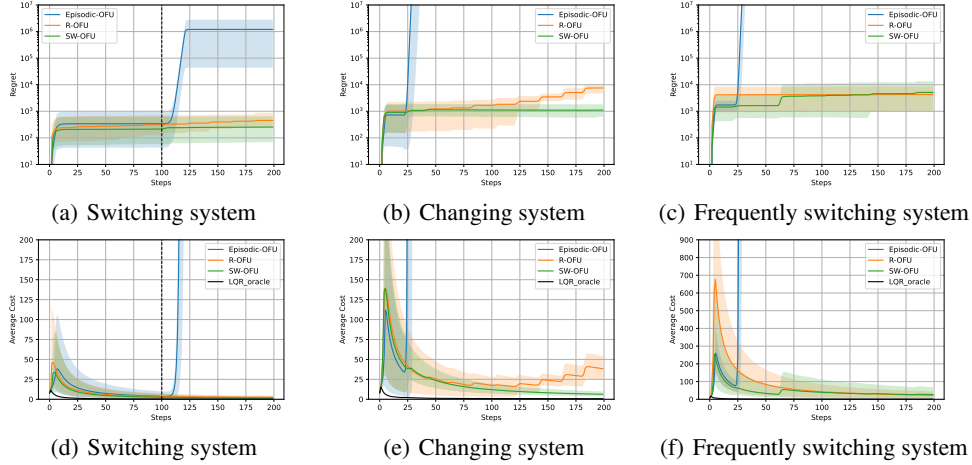


Figure 2: Performance comparison in the time-variant linear systems. On the top, the regret for each of the controllers (from left to right) on the switching system, slowly changing and frequently switching environments. All experiments are simulated under Gaussian perturbations, i.e. $w_t = \mathcal{N}(0, 0.1^2)$. The averaged results over 5 runs are plotted and the confidence intervals are shadowed in the picture.

6 Experiments

In this section, we provide empirical analysis of our algorithms under the following time-variant linear systems with an oracle LQR controller:

Switching system. In the first scenario we consider a linear system whose dynamics are defined by $A_1 = \begin{bmatrix} 1 & 0.5 \\ 0 & 1 \end{bmatrix}$ and $B_1 = \begin{bmatrix} 0 \\ 1.2 \end{bmatrix}$ for the first $H/2$ timesteps in the episode. Then, the system switches to $A_2 = \begin{bmatrix} 1 & 1.5 \\ 0 & 1 \end{bmatrix}$ and $B_2 = \begin{bmatrix} 0 \\ 0.9 \end{bmatrix}$ for the last $H/2$ timesteps in the episode.

Slowly changing system. For the second experiment we consider the slowly changing system defined by $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ and $B_h = \begin{bmatrix} 0 \\ h/20 \end{bmatrix}$. In this case, B_h constantly evolves with h .

Frequently switching system. On the frequently switching model, the dynamics changes every 20 steps. Specifically, the dynamics is randomly selected between a set of configurations whose controllability has been previously tested. The system configurations used are:

$$A_1 = A_3 = \begin{bmatrix} 1 & 0.5 \\ 0 & 1 \end{bmatrix}, B_1 = -B_3 = \begin{bmatrix} 0 \\ 1.2 \end{bmatrix}, A_2 = A_4 = \begin{bmatrix} 1 & 1.5 \\ 0 & 1 \end{bmatrix}, B_2 = -B_4 = \begin{bmatrix} 0 \\ 0.9 \end{bmatrix}.$$

We consider that all systems are perturbed under i.i.d. Gaussian noise i.e. $w_t = \mathcal{N}(0, 0.1^2)$. The performance of the algorithms is measured with quadratic cost function $c_h = \|x_h\|^2 + \|u_h\|^2$.

In order to control the proposed unknown LTV systems, we use the R-OFU and SW-OFU algorithms. In the case of R-OFU we set the length of the epoch to $L = 20$ whereas the size of the SW-OFU window is $\mathcal{W} = 20$. According to the OFU principle, we select the best model from a set of $m = 50$ candidates generated using random noise $\mathcal{U}_{[-0.5, 0.5]}$ along each of the search directions.

The results of the experiments are summarized in Fig. 2. Regarding the regret (6) and average cost the R-OFU performs better in scenarios in abruptly changing systems (switching and frequently switching), whereas the SW-OFU is better under slowly changing dynamics. This is due to the fact that R-OFU can adapt to changes more rapidly, as it discards the previous history at the start of a new epoch. While the SW-OFU takes advantage of the recent history to derive a control policy that performs better on slowly changing scenarios as it does not experience the aggregated cost of restarting the estimation on every epoch (observe Fig. 2b and 2e). Lastly, and as expected, we observe that the oracle LQR is not capable of adapting to the time-varying scenarios, getting non-stable.

7 Conclusion and Future Work

This paper studies the regret minimization problem for LTV control, where the system dynamics are unknown and change over time along with the episodes. We propose two practical algorithms based on the R and SW strategies, that notwithstanding their simplicity achieve sublinear regret under the proper configuration. To the best of our knowledge, this is the first work to obtain a tight theoretical regret bound on this setting while being computationally tractable. An interesting alternative direction in the future is to incorporate detection mechanisms to adaptively sense variations of the environment. This may provide an indication of how much the size of R and SW changes, potentially paving the way towards improved algorithms with tighter dynamic regret bounds. On the algorithm side, the OFU strategy has is shown to be computational tractable but inefficient. A promising direction is using nonconvex optimization [27], which has shown to be efficient in LTI systems.

References

- [1] Richard H Middleton and Graham C Goodwin. Adaptive control of time-varying linear systems. *IEEE Transactions on Automatic Control*, 33(2):150–155, 1988.
- [2] Tuhin Sarkar, Alexander Rakhlin, and Munther Dahleh. Nonparametric system identification of stochastic switched linear systems. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 3623–3628. IEEE, 2019.
- [3] Asaf Cassel and Tomer Koren. Online policy gradient for model free learning of linear quadratic regulators with \sqrt{T} regret. In *International Conference on Machine Learning*. PMLR, 2021.
- [4] Naman Agarwal, Brian Bullins, Elad Hazan, Sham Kakade, and Karan Singh. Online control with adversarial disturbances. In *International Conference on Machine Learning*, pages 111–119. PMLR, 2019.
- [5] Alon Cohen, Tomer Koren, and Yishay Mansour. Learning linear-quadratic regulators efficiently with only \sqrt{t} regret. In *International Conference on Machine Learning*, pages 1300–1309. PMLR, 2019.
- [6] Max Simchowitz and Dylan Foster. Naive exploration is optimal for online lqr. In *International Conference on Machine Learning*, pages 8937–8948. PMLR, 2020.
- [7] Horia Mania, Stephen Tu, and Benjamin Recht. Certainty equivalence is efficient for linear quadratic control. In *Advances in Neural Information Processing Systems*, pages 10154–10164, 2019.
- [8] Marc Abeille and Alessandro Lazaric. Efficient optimistic exploration in linear-quadratic regulators via lagrangian relaxation. In *International Conference on Machine Learning*, pages 23–31. PMLR, 2020.
- [9] Asaf Cassel, Alon Cohen, and Tomer Koren. Logarithmic regret for learning linear quadratic regulators efficiently. In *International Conference on Machine Learning*, pages 1328–1337. PMLR, 2020.
- [10] Dylan Foster and Max Simchowitz. Logarithmic regret for adversarial online control. In *International Conference on Machine Learning*, pages 3211–3221. PMLR, 2020.
- [11] Xinyi Chen and Elad Hazan. Black-box control for linear dynamical systems. In *Conference on Learning Theory*, pages 1114–1143. PMLR, 2021.
- [12] Guanya Shi, Yiheng Lin, Soon-Jo Chung, Yisong Yue, and Adam Wierman. Online optimization with memory and competitive control. In *Advances in Neural Information Processing Systems*, 2020.
- [13] Yiheng Lin, Yang Hu, Haoyuan Sun, Guanya Shi, Guannan Qu, and Adam Wierman. Perturbation-based regret analysis of predictive control in linear time varying systems. In *Advances in Neural Information Processing Systems*, 2021.

- [14] Naman Agarwal, Elad Hazan, Anirudha Majumdar, and Karan Singh. A regret minimization approach to iterative learning control. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 100–109. PMLR, 2021.
- [15] Paula Gradu, Elad Hazan, and Edgar Minasyan. Adaptive regret for control of time-varying dynamics. *arXiv preprint arXiv:2007.04393*, 2020.
- [16] Edgar Minasyan, Paula Gradu, Max Simchowitz, and Elad Hazan. Online control of unknown time-varying dynamical systems. volume 34, 2021.
- [17] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24:2312–2320, 2011.
- [18] Sayak Ray Chowdhury, Xingyu Zhou, and Ness Shroff. Adaptive control of differentially private linear quadratic systems. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pages 485–490. IEEE, 2021.
- [19] Sham M. Kakade, Akshay Krishnamurthy, Kendall Lowrey, Motoya Ohnishi, and Wen Sun. Information theoretic regret bounds for online nonlinear control. *Advances in neural information processing systems*, abs/2006.12466, 2020.
- [20] Dimitar Ho, Hoang Le, John Doyle, and Yisong Yue. Online robust control of nonlinear systems with large uncertainty. In *International Conference on Artificial Intelligence and Statistics*, pages 3475–3483. PMLR, 2021.
- [21] Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 1–26. JMLR Workshop and Conference Proceedings, 2011.
- [22] Tianyu Wang and Lin F Yang. Episodic linear quadratic regulators with low-rank transitions. *arXiv preprint arXiv:2011.01568*, 2020.
- [23] Dimitri Bertsekas. *Reinforcement learning and optimal control*. Athena Scientific, 2019.
- [24] Jin Jiang and Youmin Zhang. A novel variable-length sliding window blockwise least-squares algorithm for on-line estimation of time-varying parameters. *International journal of adaptive control and signal processing*, 18(6):505–521, 2004.
- [25] 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.
- [26] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [27] Asaf Cassel, Alon Cohen, and Tomer Koren. Efficient online linear control with stochastic convex costs and unknown dynamics. In <https://arxiv.org/abs/2203.01170>.
- [28] Yoan Russac, Claire Vernade, and Olivier Cappé. Weighted linear bandits for non-stationary environments. In *Advances in Neural Information Processing Systems*, pages 12017–12026, 2019.

A Proof for Well-Definedness in Section 3.2

Proposition A.1. *The region encompassed by high the probability confidence sets (11) are closed and bounded.*

Proof. By definition both \mathcal{C}_h and $\tilde{\mathcal{C}}_h$ are closed and bounded in the region with a constant radius $\zeta_h(\delta)$ and $\tilde{\zeta}_h(\delta)$. \square

Proposition A.2. *Given any $\tilde{\Theta}_h$ in (12), the gain of the controller $K_h(\tilde{\Theta}_h)$ is well defined.*

Proof. Please refer to Lemma 1 in [22]. \square

Based on the propositions A.1 and A.2, we also present several boundeness results through the following corollaries.

Corollary A.3. *Under Assumption 2.2, the following holds,*

$$\|x_{k,h}\|_2 \leq 1, \quad \|u_{k,h}\|_2 \leq \gamma, \quad \|z_{k,h}\|_2 \leq 1 + \gamma,$$

for all $k \geq$ and $h \in [H]$.

Corollary A.4. *Under Assumption 2.2, there exists a constant D such that,*

$$\|P_{k,h}(\Theta)\|_2 \leq D,$$

for all $k \geq$ and $h \in [H]$.

Corollary A.5. *The Spectrum of matrices \mathcal{V}_h and $\tilde{\mathcal{V}}_h$ is bounded, i.e.,*

$$\rho(\mathcal{V}_h), \quad \rho(\tilde{\mathcal{V}}_h) \leq M.$$

Corollaries A.3-A.5 are consistent to the boundness results in [22]. Thus, we can adopt the same reasoning to prove it.

B Proof for Proposition 5.1 and Lemmas 5.3, 5.4 in Section 5.1

B.1 Proof of Proposition 5.1

From the closed-form solution of (7), one can verify that the estimate error can be decomposed as,

$$\begin{aligned} & \Theta_h^* - \Theta_h \\ &= \Theta_h^* - \mathcal{V}_h^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top X_s^{next} \right) \\ &= \Theta_h^* - \mathcal{V}_h^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s \Theta_s^* + \sum_{s=h_0}^{h-1} Z_s^\top W_s \right) \\ &= \Theta_h^* - \mathcal{V}_h^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) - \sum_{s=h_0}^{h-1} Z_s^\top W_s - \sum_{s=h_0}^{h-1} Z_s^\top Z_s \Theta_h^* \right) \\ &= \mathcal{V}_h^{-1} \left(\lambda \Theta_h^* - \sum_{s=h_0}^{h-1} Z_s^\top W_s - \sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right). \end{aligned} \quad (15)$$

Therefore, based on (15), the following holds,

$$\begin{aligned} \|\Theta_h^* - \Theta_h\|_{\mathcal{V}_h} &\stackrel{(a)}{\leq} \|\lambda I \Theta_s^*\|_{\mathcal{V}_h^{-1}} + \left\| \sum_{s=h_0}^{h-1} Z_s^\top W_s \right\|_{\mathcal{V}_h^{-1}} + \left\| \sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right\|_{\mathcal{V}_h^{-1}} \\ &\leq \left\| (\lambda I)^{\frac{1}{2}} \Theta_s^* \right\|_F + \left\| \sum_{s=h_0}^{h-1} Z_s^\top W_s \right\|_{\mathcal{V}_h^{-1}} + \left\| \sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right\|_{\mathcal{V}_h^{-1}} \\ &\stackrel{(b)}{\leq} \underbrace{\left\| (\lambda I)^{\frac{1}{2}} \right\|_2}_{\ell_1} + \underbrace{\left\| \sum_{s=h_0}^{h-1} Z_s^\top W_s \right\|_{\mathcal{V}_h^{-1}}}_{\ell_2} + \underbrace{\left\| \sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right\|_{\mathcal{V}_h^{-1}}}_{\ell_3}. \end{aligned} \quad (16)$$

where (a) follows from triangle inequality and (b) holds from fact that $\|\mathcal{A}\mathcal{B}\|_F \leq \|\mathcal{A}\|_2 \|\mathcal{B}\|_F$ for any two matrices \mathcal{A} and \mathcal{B} .

B.1.1 Proof for bound of ℓ_2 in Lemma 5.3

Let $\{\mathcal{F}_t\}_{t=0}^\infty$ be a filtration generated by the random variables $\{s_{t+1}, a_{t+1}\}_{t=0}^\infty$. Let $\{\eta_t\}_{t \geq 1}$ be a vector-valued martingale difference process adapted to the filtration $\{\mathcal{F}_t\}_{t=0}^\infty$ be a sub-Gaussian random vector, i.e., it satisfies for some $R \geq 0$, the following holds

$$\mathbb{E}[\exp(\alpha^\top \eta_t) | \mathcal{F}_{t-1}] \leq \exp\left(\frac{R^2 \|\alpha\|^2}{2}\right), \quad \forall t \geq 1, \forall \alpha \in \mathbb{R}^n. \quad (17)$$

Let $V_t = \sum_{i=1}^t Z_i^\top Z_i$, $\mathcal{V}_t = V_t + \lambda I_m$, $S_t = \sum_{s=1}^t Z_s^\top W_s$ and $d = m + n$. For any $0 < \delta \leq 1$, with probability at least $1 - \delta$, uniformly over all $t \geq 1$, it holds that

$$\left\| \mathcal{V}_t^{-1/2} S_t \right\|_F \leq R \sqrt{2 \ln\left(\frac{1}{\delta}\right) + n \ln \frac{\det(\mathcal{V}_t)}{\det(\lambda I_d)}}. \quad (18)$$

Proof. For any $\gamma \in \mathbb{R}^{d \times n}$ and $t \geq 0$, let us define

$$M_t^\gamma = \exp\left(\frac{1}{R} \text{tr}(\gamma^\top S_t) - \frac{1}{2} \text{tr}(\gamma^\top S_t \gamma)\right). \quad (19)$$

From where we derive that $M_t^\gamma = \prod_{i=1}^t D_i^\gamma$, where

$$\begin{aligned} D_i^\gamma &= \exp\left(\frac{1}{R} \text{tr}(\gamma^\top Z_i^\top W_i) - \frac{1}{2} \text{tr}(\gamma^\top Z_i^\top Z_i \gamma)\right) \\ &= \exp\left(\frac{1}{R} \text{tr}(W_i \gamma^\top Z_i^\top) - \frac{1}{2} \text{tr}(Z_i \gamma \gamma^\top Z_i^\top)\right) \\ &= \exp\left(\frac{W_i \gamma^\top Z_i^\top}{R} - \frac{1}{2} \|\gamma^\top Z_i^\top\|^2\right). \end{aligned} \quad (20)$$

Note that $M_t^\gamma \geq 0$ and D_i^γ is \mathcal{F}_t measurable, as is $M_t^\gamma \geq 0$. Further, due to the conditional sub-Gaussian property, it holds that $\mathbb{E}[D_i^\gamma | \mathcal{F}_{t-1}] \leq 1$, and thus that $\mathbb{E}[M_t^\gamma | \mathcal{F}_{t-1}] \leq M_{t-1}^\gamma$. Therefore, $\{M_t^\gamma\}_{t=0}^\infty$ is a super-martingale w.r.t the filtration $\{\mathcal{F}_t\}_{t=0}^\infty$ satisfying $\mathbb{E}[M_t^\gamma] \leq 1$.

Let τ be a stopping time w.r.t filtration $\{\mathcal{F}_t^\gamma\}_{t=0}^\infty$. By the convergence theorem for non-negative super-martingales, $M_\infty^\gamma = \lim_{t \rightarrow \infty} M_t^\gamma$ is almost surely well-defined. Thus M_τ^γ is well-defined as well, irrespective of whether τ is finite or not. Let $Q_t^\gamma = M_{\min\{\tau, t\}}^\gamma$ be a stop version of $\{M_t^\gamma\}_t$. By Fatou's lemma,

$$\mathbb{E}[M_t^\gamma] = \mathbb{E}\left[\liminf_{t \rightarrow \infty} Q_t^\gamma\right] \leq \liminf_{t \rightarrow \infty} \mathbb{E}[Q_t^\gamma] = \liminf_{t \rightarrow \infty} \mathbb{E}\left[M_{\min\{\tau, t\}}^\gamma\right] \leq 1. \quad (21)$$

since the stopped super-martingale $\left\{M_{\min\{\tau, t\}}^\gamma\right\}_t$ is also super-martingale.

Let \mathcal{F}_∞ be σ -algebra generated by $\{\mathcal{F}_t^\gamma\}_{t=0}^\infty$, and $\Gamma \in \mathbb{R}^{d \times n}$ be a random matrix with its entries being i.i.d. according to $\mathcal{N}(0, \lambda^{-1})$ independent of \mathcal{F}_∞ . Define a mixture of super-martingale $M_t = \mathbb{E}[M_t^\Gamma | \mathcal{F}_\infty]$, and it is immediate to see that $\{M_t\}_t$ is also a non-negative super-martingale w.r.t. the filtration $\{\mathcal{F}_t\}_t$. Hence, by similar argument, M_τ is well-defined and following holds

$$\mathbb{E}[M_t] = \mathbb{E}[M_t^\Gamma] = \mathbb{E}[\mathbb{E}[M_t^\Gamma | \mathcal{F}_\infty]] \leq \mathbb{E}[1] = 1 \quad (22)$$

Now let us start to compute M_t . To this end, we define $\mathcal{S}_t = \frac{S_t}{R}$ and $V = \lambda I_d$. The joint probability density function of Γ is given by

$$f(\gamma) = \left(\sqrt{\lambda/2\pi}\right)^{dn} \exp\left(-\frac{\lambda}{2} \sum_{i=1}^d \sum_{j=1}^n \gamma_{i,j}^2\right) = \left(\frac{\det(V)^{\frac{1}{2}}}{(2\pi)^{\frac{d}{2}}}\right)^n \exp\left(-\frac{1}{2} \text{tr}(\gamma^\top V \gamma)\right). \quad (23)$$

For any p.d. matrix M , define $c(M) = \left(\frac{\det(M)^{\frac{1}{2}}}{(2\pi)^{\frac{m}{2}}}\right)^n$. Then, we have ²

$$\begin{aligned} M_t &= \int_{\mathbb{R}^{m \times n}} \exp\left(\text{tr}(\gamma^\top \mathcal{S}_t) - \frac{1}{2} \text{tr}(\gamma^\top V_t \gamma)\right) f(\gamma) d\gamma \\ &= \int_{\mathbb{R}^{m \times n}} \exp\left(-\frac{1}{2} \text{tr}\left((\gamma - V_t^{-1} \mathcal{S}_t)^\top V_t (\gamma - V_t^{-1} \mathcal{S}_t)\right) + \frac{1}{2} \text{tr}(\mathcal{S}_t^\top V_t^{-1} \mathcal{S}_t)\right) f(\gamma) d\gamma \\ &= c(V) \exp\left(\frac{1}{2} \text{tr}(\mathcal{S}_t^\top V_t^{-1} \mathcal{S}_t)\right) \int_{\mathbb{R}^{m \times n}} \exp\left(-\frac{1}{2} \left\{ \text{tr}\left((\gamma - V_t^{-1} \mathcal{S}_t)^\top V_t (\gamma - V_t^{-1} \mathcal{S}_t)\right) + \frac{1}{2} \text{tr}(\gamma^\top V \gamma) \right\}\right) d\gamma \\ &= c(V) \exp\left(\frac{1}{2} \text{tr}(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t)\right) \int_{\mathbb{R}^{m \times n}} \exp\left(-\frac{1}{2} \text{tr}\left((\gamma - \mathcal{V}_t^{-1} \mathcal{S}_t)^\top \mathcal{V}_t (\gamma - \mathcal{V}_t^{-1} \mathcal{S}_t)\right)\right) d\gamma. \end{aligned} \quad (24)$$

where in the last step we use the fact that $\mathcal{V}_t = V_t + V$ and

$$\begin{aligned} &\text{tr}\left((\gamma - V_t^{-1} \mathcal{S}_t)^\top V_t (\gamma - V_t^{-1} \mathcal{S}_t)\right) + \text{tr}(\gamma^\top V \gamma) \\ &= \text{tr}\left((\gamma - \mathcal{V}_t^{-1} \mathcal{S}_t)^\top \mathcal{V}_t (\gamma - \mathcal{V}_t^{-1} \mathcal{S}_t)\right) + \text{tr}(\mathcal{S}_t^\top V_t^{-1} \mathcal{S}_t) - \text{tr}(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t). \end{aligned} \quad (25)$$

Let $P(1), \dots, P(n)$ denote the columns a m -by- n matrix P , and $A_t = \mathcal{V}_t^{-1} \mathcal{S}_t$, then

$$\begin{aligned} \text{tr}\left((\gamma - \mathcal{V}_t^{-1} \mathcal{S}_t)^\top \mathcal{V}_t (\gamma - \mathcal{V}_t^{-1} \mathcal{S}_t)\right) &= \sum_{i=1}^n (\gamma(i) - A_t(i))^\top \mathcal{V}_t (\gamma(i) - A_t(i)) \\ &= \sum_{i=1}^n \|\gamma(i) - A_t(i)\|_{\mathcal{V}_t}^2, \end{aligned} \quad (26)$$

which yields

$$\begin{aligned} M_t &= c(V) \exp\left(\frac{1}{2} \text{tr}(\mathcal{S}_t^\top V_t^{-1} \mathcal{S}_t)\right) \int_{\mathbb{R}^{d \times n}} \prod_{i=1}^n \exp\left(-\frac{1}{2} \|\gamma(i) - A_t(i)\|_{\mathcal{V}_t}^2\right) d\gamma \\ &= c(V) \exp\left(\frac{1}{2} \text{tr}(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t)\right) \prod_{i=1}^n \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2} \|\gamma(i) - A_t(i)\|_{\mathcal{V}_t}^2\right) d\gamma(i) \\ &= c(V) \exp\left(\frac{1}{2} \text{tr}(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t)\right) \prod_{i=1}^n \frac{(2\pi)^{\frac{m}{2}}}{\det(\mathcal{V}_t)^{\frac{1}{2}}} \\ &= \left(\frac{\det(\lambda I_d)}{\det(\lambda I_d + V_t)}\right)^{\frac{n}{2}} \exp\left(\frac{1}{2} \text{tr}(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t)\right). \end{aligned} \quad (27)$$

²We make tr and Trace interchangeable.

In the above steps, we rely in the fact that $\int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\|\gamma(i) - A_t(i)\|_{V_t}^2\right) = \sqrt{\frac{(2\pi)^d}{\det(V_t)}}$ since $V_t \in \mathbb{R}^{d \times d}$ is a positive definite matrix.

Now for any $0 < \delta \leq 1$, the following holds by Markov's inequality

$$\begin{aligned} \mathbb{P}\left[tr(\mathcal{S}_\tau^\top \mathcal{V}_\tau^{-1} \mathcal{S}_\tau) > 2 \log\left(\frac{\det(V_\tau + \lambda I_d)^{\frac{n}{2}}}{\delta \det(\lambda I_d)^{\frac{n}{2}}}\right)\right] &= \mathbb{P}\left[\frac{\exp\left(\frac{1}{2}tr(\mathcal{S}_\tau^\top \mathcal{V}_\tau^{-1} \mathcal{S}_\tau)\right)}{\frac{1}{\delta}\left(\frac{\det(V_\tau + \lambda I_d)^{\frac{n}{2}}}{\delta \det(\lambda I_d)^{\frac{n}{2}}}\right)^{\frac{n}{2}}} > 1\right] \\ &\leq \delta \mathbb{E}[M_\tau] \leq \delta. \end{aligned} \quad (28)$$

To complete the proof, we define τ as

$\tau = \min\left\{t \geq 0 : tr(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t) > 2 \log\left(\frac{\det(V_t + \lambda I_d)^{\frac{n}{2}}}{\delta \det(\lambda I_d)^{\frac{n}{2}}}\right)\right\}$, where $\min\{\emptyset\} = \infty$ by convention. Clearly, τ is a random stopping time and

$$\begin{aligned} &\mathbb{P}\left[\forall t \geq 0, tr(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t) > 2 \log\left(\frac{\det(V_t + \lambda I_d)^{\frac{n}{2}}}{\delta \det(\lambda I_d)^{\frac{n}{2}}}\right)\right] \\ &= \mathbb{P}[\tau < \infty] \leq \mathbb{P}\left[tr(\mathcal{S}_\tau^\top \mathcal{V}_\tau^{-1} \mathcal{S}_\tau) > 2 \log\left(\frac{\det(V_\tau + \lambda I_d)^{\frac{n}{2}}}{\delta \det(\lambda I_d)^{\frac{n}{2}}}\right)\right] \leq \delta \end{aligned} \quad (29)$$

The proof concludes due to the fact that $tr(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t) = \frac{1}{R^2} tr(\mathcal{S}_t^\top \mathcal{V}_t^{-1} \mathcal{S}_t) = \left(\frac{1}{R} \|\mathcal{V}_t^{-1/2} \mathcal{S}_t\|_F\right)^2$. \square

B.1.2 Proof for bound of ℓ_2 in Lemma 5.4

Proof. The key difference when Lemma 5.3 is compared to ℓ_2 is a the bound on $\left\|\sum_{s=h_0}^{h-1} Z_s^\top W_s\right\|_{\hat{V}_h^{-1}}$, where $h_0 = \max(1, h - \mathcal{W})$. Following the trick used in [28] to handle the information loss during the sliding window, we use the union bound instead of the 'stopping time' trick to get $\left\|\mathcal{V}_t^{-1/2} \mathcal{S}_t\right\|_F \leq R\sqrt{2 \ln\left(\frac{H}{\delta}\right)} + n \ln \frac{\det(V_t)}{\det(\lambda I_d)}$. \square

B.2 Proof for bound of ℓ_3

B.2.1 Proof for bound of ℓ_3 in Lemma 5.3

Now we begin to prove the bound of ℓ_3 . Firstly, we bound ℓ_3 in Lemma 5.3 (using R strategy). We first propose the following lemma.

Lemma B.1. *For any $h \in [H]$, we have*

$$\left\|\left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I\right)^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*)\right)\right\|_F \leq (1 + \gamma) \sqrt{\frac{L(m+n)}{\lambda}} \mathcal{B}_H,$$

where the L is the size of the restart epoch.

Proof. See Appendix F.1. \square

Now we are ready to prove the bound of ℓ_3 in Lemma 5.3 .

Proof. For any $h \in [H]$, one has,

$$\begin{aligned}
& \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right) \right\|_{\mathcal{V}_h^{-1}}^2 \\
&= \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right) \right\|_{\left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)}^2 \\
&\leq \lambda_{\max} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right) \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right) \right\|_F^2. \quad (30)
\end{aligned}$$

Then putting Corollary A.5 and Lemma B.1 into (30), we finish the proof for bound of ℓ_3 in Lemma 5.3. \square

B.2.2 Proof for bound of ℓ_3 in Lemma 5.4

From the mathematical conduction of Section B.2.1, one can easily verify that the results in Section B.2.1 can be directly applied to Lemma 5.4 by replacing L with sliding window size \mathcal{W} .

C Proof for Lemma 5.5 in Section 5.2

Proof. We begin with the following decomposition of the dynamic regret.

$$\begin{aligned}
\mathcal{R}(K) &= \sum_{k=1}^K J_1^{\pi_k}(\Theta_*, x_{k,1}) - J_1^*(\Theta_*, x_{k,1}) \\
&\leq \sum_{k=1}^K J_1^{\pi_k}(\Theta_*, x_{k,1}) - J_1^*(\tilde{\Theta}_k, x_{k,1}) \\
&= \sum_{k=1}^K \Gamma_{k,1}, \quad (31)
\end{aligned}$$

where the inequality holds due to optimistic algorithm (i.e., (12)) under the event $\varepsilon_K(\delta)$. To bound this, we first investigate $\Gamma_{k,h}$. Note that the action $u_{k,h}$ under π_k is the same as that under an optimal policy when the true dynamic is $\tilde{\Theta}$, hence

$$\begin{aligned}
\Gamma_{k,h} &= \|x_{k,h}\|_Q + \|u_{k,h}\|_R + \mathbb{E}[J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) | \mathcal{F}_{k,h}] \\
&\quad - \|x_{k,h}\|_Q - \|u_{k,h}\|_R - \sum_{h'=h+1}^H \mathbb{E}[\|w_{h'}\|_{P_{h'+1}}(\tilde{\Theta}_k)] - \mathbb{E}[\|x_{k,h+1}\|_{P_{h'+1}} | \mathcal{F}_{k,h}]. \quad (32)
\end{aligned}$$

Denote $\Gamma_{k,h} = \Delta_{k,h} + J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) - \psi_{k,h+1} - \mathbb{E}[\|\tilde{\Theta}_k^\top z_{k,h} + w_{k,h}\|_{P_{h'+1}} | \mathcal{F}_{k,h}]$, where $\psi_{k,h+1} = \sum_{h'=h+1}^H \mathbb{E}[\|w_{h'}\|_{P_{h'+1}}(\tilde{\Theta}_k)]$ and $\Delta_{k,h} = \mathbb{E}[J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) | \mathcal{F}_{k,h}] - J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1})$. Now, we rewrite $\Gamma_{k,h}$ as follows,

$$\begin{aligned}
\Gamma_{k,h} &\stackrel{(a)}{=} \Delta_{k,h} + J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) - \psi_{k,h+1} - \|\tilde{\Theta}_{k,h}^\top z_{k,h}\|_{\tilde{P}_{k,h+1}} \\
&\quad - \mathbb{E}[\|w_{k,h}\|_{\tilde{P}_{k,h+1}} | \mathcal{F}_{k,h}]
\end{aligned}$$

$$\begin{aligned}
&= \Delta_{k,h} + J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) - \psi_{k,h+1} - \|\tilde{\Theta}_{k,h}^\top z_{k,h}\|_{\tilde{P}_{k,h+1}} \\
&\quad - \mathbb{E} \left[\|x_{k,h+1} - \Theta_*^\top z_{k,h}\|_{\tilde{P}_{k,h+1}} | \mathcal{F}_{k,h} \right] \\
&\stackrel{(b)}{=} \Delta_{k,h} + J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) - \psi_{k,h+1} - \|\tilde{\Theta}_{k,h}^\top z_{k,h}\|_{\tilde{P}_{k,h+1}} \\
&\quad - \mathbb{E} \left[\|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} | F_{k,h} \right] + \|\Theta_*^\top z_{k,h}\|_{\tilde{P}_{k,h+1}}, \tag{33}
\end{aligned}$$

where in (a) and (b), we use the independence and mean zero properties of $w_{k,h}$. Notice that $\psi_{k,h+1} = J_h^*(\tilde{\Theta}_{k,h+1}, x_{k,h+1}) - \|x_{k,h+1}\|_{\tilde{P}_{k,h+1}}$, then, (32) can be expressed as

Then,

$$\begin{aligned}
\Gamma_{k,h} &= \Delta_{k,h} + J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) - J_h^*(\tilde{\Theta}_{k,h+1}, x_{k,h+1}) + \|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} \\
&\quad - \|\tilde{\Theta}_{k,h}^\top z_{k,h}\|_{\tilde{P}_{k,h+1}} - \mathbb{E} \left[\|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} | F_{k,h} \right] + \|\Theta_*^\top z_{k,h}\|_{\tilde{P}_{k,h+1}}. \tag{34}
\end{aligned}$$

Due to the fact that the cost of $H + 1$ step and beyond are 0, summarize (34) yields

$$\mathcal{R}(K) \leq \sum_{k=1}^K \sum_{h=1}^{H-1} \varsigma_{h,k}. \tag{35}$$

Thus, we finish the proof of Lemma 5.5. \square

D Proof of Lemma 5.6 in Section 5.2

Proof. We begin with the following decomposition of the dynamic regret.

For the sake of convenient, denote $I_{k,h} = \|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} - \mathbb{E} \left[\|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} | \mathcal{F}_{k,h} \right]$, from where one can verify that the sequence of term $\Delta_{k,h}$ and I_1 form a martingale difference sequence. Meanwhile

$$\mathbb{E} [I_1 | \mathcal{F}_{k,h}] = 0, \quad \mathbb{E} [\Delta_{k,h} | \mathcal{F}_{k,h}] = 0,$$

holds since $\mathcal{F}_{k,h}$ is all randomness before the step (k, h) . In order to bound the term I_1 and $\Delta_{k,h}$, we resort to Corollary A.3 as well as Assumption 2.2. From where the following inequalities hold

$$\begin{aligned}
|I_2| &= \left| \|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} - \mathbb{E} \left[\|x_{k,h+1}\|_{\tilde{P}_{k,h+1}} | \mathcal{F}_{k,h} \right] \right| \\
&\leq 2v, \tag{36}
\end{aligned}$$

and to bound $\Delta_{k,h}$, we bound it backwards by using Assumption 2.2, first notice that

$$|J_H^{\pi_k}(\Theta_*, x_{k,H})| = \|x_{k,H}\|_Q + \|u_{k,H}\|_R = (1 + \gamma^2)v. \tag{37}$$

Thus, for $h \in [H]$, we have

$$|J_h^{\pi_k}(\Theta_*, x_{k,h})| = \|x_{k,h}\|_Q + \|u_{k,h}\|_R + |\mathbb{E} [J_{h+1}^{\pi_k}(\Theta_*, x_{k,h+1}) | \mathcal{F}_{k,h}]| \leq H(1 + \gamma^2)v. \tag{38}$$

Finally, we apply the Azuma-Hoeffding inequality (Lemma F.1) to (36) and (38)

$$\begin{aligned}
\sum_{k=1}^K \sum_{h=1}^H I_1 &\leq 4C \sqrt{KH \ln \frac{1}{\delta}} \leq O \left(\sqrt{KH \ln \frac{1}{\delta}} \right), \\
\sum_{k=1}^K \sum_{h=1}^H \Delta_{k,h} &\leq 2H(1 + \gamma^2)C \sqrt{KH \ln \frac{1}{\delta}} \leq O \left(\sqrt{KH^3 \ln \frac{1}{\delta}} \right).
\end{aligned}$$

Now we investigate the third term in Lemma 5.6. Note that

$$\begin{aligned}
& \sum_{k=1}^K \sum_{h=1}^H \|\Theta_*^\top z_{k,h}\|_{\tilde{P}_{k,h+1}}^2 - \|\tilde{\Theta}^\top z_{k,h}\|_{\tilde{P}_{k,h+1}}^2 \\
& \leq \sum_{k=1}^K \sum_{h=1}^H \left| \|\Theta_*^\top z_{k,h}\|_{\tilde{P}_{k,h+1}} - \|\tilde{\Theta}^\top z_{k,h}\|_{\tilde{P}_{k,h+1}} \right| \\
& = \sum_{k=1}^K \sum_{h=1}^H \left| \left(\|\tilde{P}_{k,h+1}^{1/2} \Theta_*^\top z_{k,h}\|_2 - \|\tilde{P}_{k,h+1}^{1/2} \tilde{\Theta}^\top z_{k,h}\|_2 \right) \left(\|\tilde{P}_{k,h+1}^{1/2} \Theta_*^\top z_{k,h}\|_2 + \|\tilde{P}_{k,h+1}^{1/2} \tilde{\Theta}^\top z_{k,h}\|_2 \right) \right| \\
& \leq \underbrace{\left[\sum_{k=1}^K \sum_{h=1}^H \left(\|\tilde{P}_{k,h+1}^{1/2} \Theta_*^\top z_{k,h}\|_2 + \|\tilde{P}_{k,h+1}^{1/2} \tilde{\Theta}^\top z_{k,h}\|_2 \right)^2 \right]^{1/2}}_{\ell_{\mathcal{A}}} \\
& \quad \times \underbrace{\left[\sum_{k=1}^K \sum_{h=1}^H \left(\|\tilde{P}_{k,h+1}^{1/2} \Theta_*^\top z_{k,h}\|_2 - \|\tilde{P}_{k,h+1}^{1/2} \tilde{\Theta}^\top z_{k,h}\|_2 \right)^2 \right]^{1/2}}_{\ell_{\mathcal{B}}}. \tag{39}
\end{aligned}$$

We bound the $\ell_{\mathcal{A}}$ and $\ell_{\mathcal{B}}$ respectively. We first bound the $\ell_{\mathcal{A}}$ as follows

$$\begin{aligned}
& \left[\sum_{k=1}^K \sum_{h=1}^H \left(\|\tilde{P}_{k,h+1}^{1/2} \Theta_*^\top z_{k,h}\|_2 + \|\tilde{P}_{k,h+1}^{1/2} \tilde{\Theta}^\top z_{k,h}\|_2 \right)^2 \right]^{1/2} \\
& \leq \left[\sum_{k=1}^K \sum_{h=1}^H \left(\|\tilde{P}_{k,h+1}^{1/2} \Theta_*^\top z_{k,h}\|_2 + \|\tilde{P}_{k,h+1}^{1/2} \tilde{\Theta}^\top z_{k,h}\|_2 \right)^2 \right]^{1/2} \\
& = \left[\sum_{k=1}^K \sum_{h=1}^H \left(\|\tilde{P}_{k,h+1}^{1/2} z_{k,h+1}\|_2 + \|\tilde{P}_{k,h+1}^{1/2} \tilde{z}_{k,h+1}\|_2 \right)^2 \right]^{1/2} \\
& \leq \left[\sum_{k=1}^K \sum_{h=1}^H \left(\|\tilde{P}_{k,h+1}^{1/2}\|_2 \|z_{k,h+1}\|_2 + \|\tilde{P}_{k,h+1}^{1/2}\|_2 \|\tilde{z}_{k,h+1}\|_2 \right)^2 \right]^{1/2} \\
& \leq \left[\sum_{k=1}^K \sum_{h=1}^H (2D(1+\gamma)) \right]^{1/2} \\
& = 2D\sqrt{KH}(1+\gamma). \tag{40}
\end{aligned}$$

Then, under the Assumption 2.2 and event $\varepsilon_K(\delta)$, we have

$$\begin{aligned}
\left(\|\tilde{P}_{k,h+1}^{1/2} \Theta_*^\top z_{k,h}\|_2 - \|\tilde{P}_{k,h+1}^{1/2} \tilde{\Theta}^\top z_{k,h}\|_2 \right)^2 & \leq \left\| \tilde{P}_{k,h+1}^{1/2} (\Theta_* - \tilde{\Theta}_{k,h})^\top z_{k,h} \right\|_2^2 \\
& \leq D \left\| (\Theta_* - \tilde{\Theta}_{k,h})^\top z_{k,h} \right\|_2^2 \\
& \leq D \left\| (\Theta_* - \tilde{\Theta}_k)^\top \mathcal{V}_{h,k}^{1/2} \right\|_2^2 \left\| \mathcal{V}_{h,k}^{-1/2} z_{k,h} \right\|_2^2 \\
& \leq D \left\| \Theta_* - \tilde{\Theta}_k \right\|_{\mathcal{V}_{h,k}}^2 \left\| \mathcal{V}_{h,k}^{-1/2} z_{k,h} \right\|_2^2 \\
& \leq D\zeta_h^2(\delta) \left\| \mathcal{V}_{h,k}^{-1/2} z_{k,h} \right\|_2^2. \tag{41}
\end{aligned}$$

Then, using (41) and fact that $\left(\left\|\tilde{P}_{k,h+1}^{1/2}\Theta_*^\top z_{k,h}\right\|_2 - \left\|\tilde{P}_{k,h+1}^{1/2}\tilde{\Theta}_{k,h}^\top z_{k,h}\right\|_2\right)^2 \leq 2D(1+\gamma)^2$ we derive that

$$\begin{aligned} \left(\left\|\tilde{P}_{k,h+1}^{1/2}\Theta_*^\top z_{k,h}\right\|_2 - \left\|\tilde{P}_{k,h+1}^{1/2}\tilde{\Theta}_{k,h}^\top z_{k,h}\right\|_2\right)^2 &\leq 2D(1+\gamma)^2\zeta_h^2(\delta)\min\left\{\left\|\mathcal{V}_h^{-1/2}z_{k,h}\right\|_2^2, 1\right\} \\ &= 2D(1+\gamma)^2\zeta_k^2(\delta)\min\left\{\|z_{k,h}\|_{\mathcal{V}_h^{-1}}^2, 1\right\} \\ &\leq 2D(1+\gamma)^2\zeta_h^2(\delta)\ln\left(\|z_{k,h}\|_{\mathcal{V}_h^{-1}}^2 + 1\right). \end{aligned} \quad (42)$$

Therefore, we have

$$\sum_{k=1}^K \sum_{h=1}^H \left(\left\|\tilde{P}_{k,h+1}^{1/2}\Theta_*^\top z_{k,h}\right\|_2 - \left\|\tilde{P}_{k,h+1}^{1/2}\tilde{\Theta}_{k,h}^\top z_{k,h}\right\|_2\right)^2 \leq 4KHD(1+\gamma)^2\zeta_k^2(\delta)\ln\left(\frac{\det(\mathcal{V}_h)}{\det(\lambda I)}\right), \quad (43)$$

where the last step in (43) is derived applying Lemma F.2 in Appendix F. Finally, taking the square root of (43) and multiplying it by (40) yields to the final bound for the third term in Lemma 5.6. \square

E Proof of Theorem 4.4 in Section 4

Proof. The regret $\mathcal{R}(K)$ under the case $L \geq H$ can be decomposed as shown in (35). Recall that

$$\begin{aligned} \mathcal{R}(K) &\stackrel{(I)}{=} \sum_{k=1}^K \sum_{h=1}^H I_1 + \sum_{k=1}^K \sum_{h=1}^H \Delta_{k,h} + \sum_{k=1}^K \sum_{h=1}^H \left\|\Theta_*^\top z_{k,h}\right\|_{\tilde{P}_{k,h+1}}^2 - \left\|\tilde{\Theta}^\top z_{k,h}\right\|_{\tilde{P}_{k,h+1}}^2, \\ &\leq O\left(\sqrt{KH \ln \frac{1}{\delta}}\right) + O\left(\sqrt{KH^3 \ln \frac{1}{\delta}}\right) + \underbrace{\ell_A \times \ell_B}_{\mathcal{R}'''(K)}. \end{aligned}$$

Where the first two terms in (I) enjoy the same regret bound under the Corollaries and Propositions in Section A. And the third term in (I) could be decomposed as $\mathcal{R}'''(K) = \ell_A \times \ell_B$ according to (39). The only difference is the regret bound w.r.t ℓ_B . To simplify the notations for the later proof, we denote ℓ_B as

$$\ell_B = \left[\sum_{k=1}^K \sum_{h=1}^H \left(\left\|\tilde{P}_{k,h+1}^{1/2}\Theta_*^\top z_{k,h}\right\|_2 - \left\|\tilde{P}_{k,h+1}^{1/2}\tilde{\Theta}^\top z_{k,h}\right\|_2\right)^2 \right]^{1/2} = [\mathcal{R}_{\ell_B}(K)]^{1/2}$$

Recall (41), the following holds

$$\left(\left\|\tilde{P}_{k,h+1}^{1/2}\Theta_*^\top z_{k,h}\right\|_2 - \left\|\tilde{P}_{k,h+1}^{1/2}\tilde{\Theta}_{k,h}^\top z_{k,h}\right\|_2\right)^2 \leq D\zeta_h^2(\delta)\left\|\mathcal{V}_{h,k}^{-1/2}z_{k,h}\right\|_2^2.$$

Hence, dynamic regret $\mathcal{R}(\mathcal{L})$ within the epoch \mathcal{L} is bounded by

$$\begin{aligned} \mathcal{R}(\mathcal{L}) &= \sum_{h \in \mathcal{L}} 2D(1+\gamma)^2\zeta_h^2(\delta)\left\|\mathcal{V}_{h,k}^{-1/2}z_{k,h}\right\| \\ &\leq \sum_{h \in \mathcal{L}} 2D(1+\gamma)^2(\ell_1 + \ell_2 + \ell_3)^2\left\|\mathcal{V}_{h,k}^{-1/2}z_{k,h}\right\| \\ &\leq \sum_{h \in \mathcal{L}} \Upsilon_1 L\mathcal{B}_L^2 + \Upsilon_2\sqrt{L}\mathcal{B}_L + \Upsilon_3 \\ &\leq L(\Upsilon_1 L\mathcal{B}_L^2 + \Upsilon_2\sqrt{L}\mathcal{B}_L + \Upsilon_3) \\ &= \underbrace{L^2\mathcal{B}_L^2\Upsilon_1}_{\ell'} + \underbrace{L^{3/2}\mathcal{B}_L\Upsilon_2}_{\ell''} + \underbrace{L\Upsilon_3}_{\ell'''} \end{aligned} \quad (44)$$

where $B_L = \sum_{p=h_0}^{h-1} \|\Theta_p - \Theta_{p+1}\|_F$ is the total variability in one epoch \mathcal{L} ,
 $\Upsilon_1 = \frac{(m+n)}{\lambda} \Upsilon_4$, $\Upsilon_2 = \left(\sqrt{\lambda} + v_w \sqrt{2 \ln \left(\frac{1}{\delta} \right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}} \right) \Upsilon_4$, $\Upsilon_3 =$
 $\Upsilon_4 \left(\lambda + v_w^2 \left(2 \ln \left(\frac{1}{\delta} \right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)} \right) + \sqrt{\lambda} v_w \sqrt{2 \ln \left(\frac{1}{\delta} \right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}} \right)$ and $\Upsilon_4 =$
 $2D(1+\gamma)^2 \left\| \mathcal{V}_{h,k}^{-1/2} z_{k,h} \right\|_2^2$.

Then, we can bound ℓ_B term ℓ' , ℓ'' and ℓ''' using Lemma F.2 as

$$\ell' \leq \underbrace{\frac{4(m+n)^2(1+\gamma)^2 D}{\lambda} \log \left(1 + \frac{(1+\gamma)^2 L}{\lambda(m+n)} \right)}_{\chi'} \mathcal{B}_L^2 L$$

$$\ell'' \leq \underbrace{4(m+n)(1+\gamma)^2 D \left(\sqrt{\lambda} + v_w \sqrt{2 \ln \left(\frac{1}{\delta} \right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}} \right) \log \left(1 + \frac{(1+\gamma)^2 L}{\lambda(m+n)} \right)}_{\chi''} \mathcal{B}_L L^{1/2}$$

$$\ell''' \leq \underbrace{4(m+n)(1+\gamma)^2 D \left(\lambda + v_w^2 \left(2 \ln \left(\frac{1}{\delta} \right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)} \right) + \sqrt{\lambda} v_w \sqrt{2 \ln \left(\frac{1}{\delta} \right) + n \ln \frac{\det(\mathcal{V}_h)}{\det(\lambda I)}} \right) \log \left(1 + \frac{(1+\gamma)^2 L}{\lambda(m+n)} \right)}_{\chi'''} \mathcal{B}_L L$$

By taking the union bound over the dynamic regret of all $\lceil HK/L \rceil$ epochs, we know that the following holds with probability at least $1 - 2/HK$

$$\begin{aligned} \mathcal{R}_{\ell_B}(K) &= \sum_{s=1}^{\lceil HK/L \rceil} \mathcal{R}(\mathcal{L}_s) \\ &\leq \frac{HK}{L} (\ell' + \ell'' + \ell''') \\ &= \frac{L^2}{HK} \frac{(HK)^2}{L^2} \mathcal{B}_L^2 \chi' + L^{1/2} \frac{HK}{L} \mathcal{B}_L \chi'' + \frac{HK}{L} \chi''' \\ &= \frac{L^2}{HK} \mathcal{B}_{HK}^2 \chi' + L^{1/2} \mathcal{B}_{HK} \chi'' + \frac{HK}{L} \chi''' \end{aligned} \quad (45)$$

Putting $\mathcal{R}_{\ell_B}(K)$ to ℓ_B , we obtain the bound for $\mathcal{R}'''(K)$ as

$$\begin{aligned} \mathcal{R}'''(K) &= \ell_A \times \ell_B = 2D\sqrt{KH}(1+\gamma) \sqrt{\frac{L^2}{HK} \mathcal{B}_{HK}^2 \chi' + L^{1/2} \mathcal{B}_{HK} \chi'' + \frac{HK}{L} \chi'''} \\ &\leq 2D\sqrt{KH}(1+\gamma) \left(\sqrt{\frac{L^2}{HK} \mathcal{B}_{HK}^2 \chi'} + \sqrt{L^{1/2} \mathcal{B}_{HK} \chi''} + \sqrt{\frac{HK}{L} \chi'''} \right). \end{aligned} \quad (46)$$

Ignoring logarithmic factors, we finally obtain that

$$\mathcal{R}(K) \leq \tilde{O} \left(L \mathcal{B}_{HK} + HK \sqrt{\frac{1}{L}} + \sqrt{HK} \sqrt{\mathcal{B}_{HK} L^{1/4}} \right) + \tilde{O}(\sqrt{KH}) + \tilde{O}(\sqrt{KH^3})$$

□

F Auxiliary Proof and Lemma

In this section, we provide several technical lemmas frequently used in the proofs.

Lemma F.1. (Azuma-Hoeffding inequality) Let $\{X_k\}_{k=0}^\infty$ be a discrete-parameter real-valued martingale sequence such that for every $k \in \mathbb{N}$, the condition $|X_k - X_{k-1}| \leq \mu$ holds for some non-negative constant μ . Then with probability at least $1 - \delta$, we have

$$|X_n - X_0| \leq 2\mu\sqrt{n \log \frac{1}{\delta}}.$$

Lemma F.2. [17] For any $\{x_t\}_{t=1}^T \in \mathbb{R}^d$ satisfying that $\|x_t\|_2 \leq L$, let $A_0 = \lambda I$ and $A_t = A_0 + \sum_{i=1}^{t-1} x_i x_i^\top$, then the following inequality holds

$$\sum_{t=1}^T \min\{1, \|x_t\|_{A_{t-1}^{-1}}\}^2 \leq 2d \log \frac{d\lambda + TL^2}{d\lambda}.$$

Proof. Notice that the following holds

$$A_t = A_{t-1} + x_t x_t^\top = A_{t-1}^{1/2} \left(I + A_{t-1}^{-1/2} x_t x_t^\top A_{t-1}^{-1/2} \right) A_{t-1}^{1/2}.$$

and taking the determinant yields

$$\det(A_t) = \det(A_{t-1}) \det\left(I + A_{t-1}^{-1/2} x_t x_t^\top A_{t-1}^{-1/2}\right).$$

Note the fact $\det(I + x x^\top) = 1 + \|x\|_2^2$, we have

$$\det(A_t) = \det(A_{t-1}) \left(1 + \|A_{t-1}^{-1/2} x_t\|_2^2\right) \geq \det(A_{t-1}) \exp\left(\frac{\|A_{t-1}^{-1/2} x_t\|_2^2}{2}\right).$$

where the inequality holds based on fact $1 + x \geq \exp(x/2)$ holds for $x \in [0, 1]$. Finally, by utilizing telescope structure, we get

$$\sum_{t=1}^T \|A_{t-1}^{-1/2} x_t\|_2^2 \leq 2 \log \frac{\det(A_T)}{\det(A_0)} \leq 2d \log \left(1 + \frac{L^2 T}{\lambda d}\right).$$

□

F.1 Proof of Lemma B.1

In this section, we provide the proof of Lemma B.1.

Proof. For any $h_0 \in [H]$, one has

$$\begin{aligned} & \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s (\Theta_s^* - \Theta_h^*) \right) \right\|_F \\ &= \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s \left(\sum_{p=s}^{h-1} (\Theta_p - \Theta_{p+1}) \right) \right) \right\|_F \\ &= \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{p=h_0}^{h-1} \left(\sum_{s=h_0}^p Z_s^\top Z_s (\Theta_p - \Theta_{p+1}) \right) \right) \right\|_F \\ &\leq \sum_{p=h_0}^{h-1} \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^p Z_s^\top Z_s \right) (\Theta_p - \Theta_{p+1}) \right\|_F \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{p=h_0}^{h-1} \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^p Z_s^\top Z_s \right) (\Theta_p - \Theta_{p+1}) \right\|_F \\
&\leq \sum_{p=h_0}^{h-1} \left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \sum_{s=h_0}^p Z_s^\top Z_s \right\|_2 \|\Theta_p - \Theta_{p+1}\|_F, \tag{47}
\end{aligned}$$

where the last inequality holds due to fact that $\|AB\|_F \leq \|A\|_2 \|B\|_F$ for any matrix A and B . Since

$$\|v\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \leq \frac{\|v\|_2}{\sqrt{\lambda}} \text{ as } \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right) \succeq \lambda I \text{ holds, thus we have}$$

$$\begin{aligned}
\left\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \sum_{s=h_0}^p Z_s^\top Z_s \right\|_2 &= \sup_{v \in \mathcal{B}(1)} \left| v^\top \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \sum_{s=h_0}^p Z_s^\top Z_s v \right| \\
&\stackrel{(a)}{\leq} \left| v_*^\top \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \sum_{s=h_0}^p Z_s^\top Z_s v_* \right| \\
&\leq \|v_*\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left\| \sum_{s=h_0}^p Z_s^\top Z_s v_* \right\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \\
&\leq \|v_*\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left\| \sum_{s=h_0}^p Z_s^\top \|Z_s\|_2 \|v_*\|_2 \right\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \\
&\leq \frac{1+\gamma}{\sqrt{\lambda}} \left\| \sum_{s=h_0}^p Z_s \right\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \\
&\leq \frac{1+\gamma}{\sqrt{\lambda}} \sum_{s=h_0}^p \|Z_s\| \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \\
&\stackrel{(b)}{\leq} \frac{1+\gamma}{\sqrt{\lambda}} \sqrt{L} \sqrt{\sum_{s=h_0}^p \|Z_s\|^2} \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \\
&\stackrel{(c)}{\leq} (1+\gamma) \sqrt{\frac{L(m+n)}{\lambda}}, \tag{48}
\end{aligned}$$

where v_* in (a) denotes the optimizer; (b) holds by Cauchy-Schwarz inequality. The inequality (c) makes use of the following algebra formulation: for $p \in \{h_0, \dots, h-1\}$

$$\begin{aligned}
&\sum_{s=h_0}^p \|Z_s\|^2 \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \\
&= \text{tr} \left[Z_s \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} Z_s^\top \right] \\
&= \text{tr} \left[\left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^p Z_s^\top Z_s \right) \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \text{tr} \left[\left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^p Z_s^\top Z_s \right) \right] \\
&\quad + \sum_{s=p+1}^{h-1} Z_s^\top \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} Z_s + \lambda \sum_{i=1}^d e_i^\top \left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} e_i \\
&= \text{tr} \left[\left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=h_0}^p Z_s^\top Z_s \right) \right] + \text{tr} \left[\left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \left(\sum_{s=p+1}^{h-1} Z_s^\top Z_s \right) \right] \\
&\quad + \text{tr} \left[\left(\sum_{s=h_0}^{h-1} Z_s^\top Z_s + \lambda I \right)^{-1} \lambda \sum_{i=1}^d e_i^\top e_i \right] \\
&= \text{tr}(I_{n+m}) = n + m. \tag{49}
\end{aligned}$$

Finally, putting Assumption 2.3 and (48) into (47) finishes our proof. \square