

---

# Dynamical Linear Bandits for Long-Lasting Vanishing Rewards

---

Marco Mussi<sup>1</sup> Alberto Maria Metelli<sup>1</sup> Marcello Restelli<sup>1</sup>

## Abstract

In many real-world sequential decision-making problems, an action might not immediately reflect on the feedback and spread its effects over a long time horizon. For instance, in online advertising, investing in a platform produces an increase in *awareness*, but the actual reward (e.g., a *conversion*) might occur far in the future. Furthermore, whether a conversion takes place depends on several factors: how fast the awareness grows; possible awareness vanishing effects; synergy or interference with other advertising platforms. Previous work has investigated the Multi-Armed Bandit framework with the possibility of delayed and aggregated feedback, without a particular structure on how an action propagates into the future, disregarding possible hidden dynamical effects. In this paper, we introduce a novel setting, Dynamical Linear Bandits (DLB), an extension of linear bandits characterized by a hidden state. When an action is performed, the learner observes a noisy reward whose mean is a linear function of the hidden state and the action, and the hidden state evolves according to a linear dynamics. In this way, the effects of each action are delayed by the system evolution, persist over time, and the interplay between the action components is taken into account. After introducing the setting and discussing the notion of optimal policy, we provide an any-time optimistic regret minimization algorithm Dynamical Linear Upper Confidence Bound (DynLin-UCB) that suffers regret of order  $\tilde{O}(cd\sqrt{T})$ , where  $c$  is a constant dependent on the properties of the linear dynamical evolution. Finally, we conduct a numerical validation on a synthetic environment to show the effectiveness of DynLin-UCB in comparison with bandit baselines.

## 1. Introduction

In a large variety of sequential decision-making problems, a learner is required choosing an action that, when executed, determines an evolution of the underlying system state, that is hidden to the learner. In these partially-observable problems, the learner might observe a reward (a.k.a. feedback) that is the combined effect of multiple actions played in the past and its realization might span a large time horizon. For instance, in online advertising campaigns, the process that leads to a *conversion*, also known as the marketing funnel (Court et al., 2009), is characterized by complex dynamics and comprises several phases. When multiple campaigns and heterogeneous platforms are involved, a profitable budget investment policy has to account for relationships between campaigns/platforms and how their effects are combined. In this scenario, the conversion attribution problem (Berman, 2018) consists in assigning the credit of a conversion (e.g., a user’s purchase of a promoted product) not only to the latest ad the user was exposed to, but also to the previous ones which contributed to generate such a conversion. The *joint* consideration of each funnel phase is a fundamental step towards an optimal investment solution, while considering the advertising campaigns/platforms *independently* leads to sub-optimal solutions. Consider, for instance, a simplified version of the funnel with two types of campaigns: *awareness* (i.e., impression) ads and *conversion* ads. The first kind of ad aims at improving the brand awareness, while the latter aims at creating the actual conversion. If we evaluate the performances in terms of conversions only, we will observe that impression ads are not instantaneously effective in creating conversions, so we will be tempted to reduce the amount of budget invested in such a campaign. However, this approach is clearly sub-optimal because, as demonstrated in several works (e.g., Hoban & Bucklin, 2015), impression ads increase the chance to convert when a conversion ad is shown after the impression. In addition, the effect of some ads, especially impression ads delivered via television, may be delayed. It has been demonstrated (Chapelle, 2014) that users remember advertising over time in a vanishing way, leading to delayed consequences that simple models cannot capture. It is worth noting that this kind of interplay comprises more general scenarios than simple delay (which are anyway present), including cases where the interaction is governed by dynamics

---

<sup>1</sup>DEIB, Politecnico di Milano, Milan, Italy. Correspondence to: Marco Mussi <marco.mussi@polimi.it>.

*hidden* to the observer.

While this problem can be indubitably modeled as a partially-observable Markov Decision Process (POMDP, Åström, 1965), the complexity of the framework and its generality are often not required to capture the main features of the problem. Indeed, for specific classes of problems, the Multi-Armed Bandit (MAB, Lattimore & Szepesvári, 2020) literature has explored the possibility of experiencing delayed reward (a.k.a. feedback) either assuming that the actual reward will be observed, individually, in the future (e.g., Joulani et al., 2013) or with the more realistic assumption that an aggregated feedback is available (e.g., Pike-Burke et al., 2018), with also specific applications for online advertising (Vernade et al., 2017). Although effective in dealing with delay effects and the possibility of a reward spread in the future (Cesa-Bianchi et al., 2018), they do not account for additional, more complex, dynamical effects, which can be regarded as the evolution of a hidden state.

In this work, we take a different perspective. We propose to model the non-observable dynamical effects, underlying the phenomena as a Linear Time-Invariant (LTI) system (Hespanha, 2018). In particular, the system is characterized by an hidden internal state  $\mathbf{x}_t$  (e.g., awareness) which evolves via a linear dynamics fed by the action  $\mathbf{u}_t$  (e.g., the amount invested on each platform) and is affected by noise. At each round, the learner experiences a reward  $y_t$  (e.g., conversions) which is a noisy observation that linearly combines the state  $\mathbf{x}_t$  and the action  $\mathbf{u}_t$ . Our goal consists in learning an optimal policy so as to maximize the expected cumulative reward. We call this setting *Dynamical Linear Bandits* (DLBs) that, as we shall see, reduce to linear bandits when no dynamics is involved. Coming back to the application scenario, the state allows encoding the awareness that accumulates and/or vanishes over time. Furthermore, thanks to the dynamics, the model allows for representing interference and synergy phenomena between platforms. Due to the dynamical nature of the system, the effect of each action persists over time indefinitely but, under stability conditions, it vanishes asymptotically.

**Contributions** The contributions of this paper are theoretical, algorithmic, and experimental and can be summarized as follows. In Section 2, we introduce the Dynamical Linear Bandit (DLB) setting to represent sequential decision-making problems characterized by a hidden state that evolves linearly according to an *unknown* dynamics and the learner observes a noisy reward obtained as a linear combination of the current state and action played. We formally define the hidden-state linear dynamics, the computation of the reward, the notion of policy, and the regret definition. In particular, under stability conditions, we show that the optimal policy corresponds to playing the *constant action* that leads the system to the most profitable steady-

state. In Section 3, we propose a novel anytime optimistic regret minimization algorithm, *Dynamical Linear Upper Confidence Bound* (DynLin-UCB) for the DLB setting. DynLin-UCB takes inspiration from Lin-UCB and subdivides the optimization horizon  $T$  into increasing-length epochs. In each epoch, an action is selected and kept constant (persisted) so that the system approximately reaches the steady-state. Now, a reliable observation of the reward is collected to update the estimates and optimistically select the next action to play. We provide a regret analysis for DynLin-UCB showing that, under certain assumptions, it enjoys  $\tilde{O}(cd\sqrt{T})$  expected regret, where  $c$  is a constant that accounts for the “speed” at which the system reaches the steady-state and  $d$  is the dimensionality of the action  $\mathbf{u}$ . In Section 5, we provide a numerical validation to highlight the properties of our approach in comparison with bandit baselines. The proof of all the results are reported in Appendix A.

**Notation** Let  $a, b \in \mathbb{N}$  with  $a \leq b$ , we denote with  $J_{a,b} := \{a, \dots, b\}$ , with  $Jb := J1, b$ , and with  $J_{a, \infty} := \{a, a+1, \dots\}$ . Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  be real-valued vectors, we denote with  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y} = \sum_{j=1}^n x_j y_j$  the inner product. For a positive semidefinite matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , we denote with  $\|\mathbf{x}\|_{\mathbf{A}}^2 = \mathbf{x}^\top \mathbf{A} \mathbf{x}$  the weighted 2-norm. The *spectral radius*  $\rho(\mathbf{A})$  is the largest absolute value of the eigenvalues of  $\mathbf{A}$ , the *spectral norm*  $\|\mathbf{A}\|_2$  is the square root of the maximum eigenvalue of  $\mathbf{A}^\top \mathbf{A}$ , the *Frobenius norm*  $\|\mathbf{A}\|_F$  is the trace of  $\mathbf{A}^\top \mathbf{A}$ . We introduce the maximum spectral norm to spectral radius ratio of the powers of  $\mathbf{A}$  defined as  $\kappa(\mathbf{A}) = \sup_{\neq 0} \|\mathbf{A}\|_2 / \rho(\mathbf{A})$  (Oymak & Ozay, 2019). We denote with  $\mathbf{I}_n$  the identity matrix of order  $n$  and with  $\mathbf{0}_n$  the vector of all zeros of dimension  $n$ . A random vector  $\mathbf{x} \in \mathbb{R}^n$  is  $\sigma^2$ -subgaussian, in the sense of (Hsu et al., 2012), if for every vector  $\mathbf{v} \in \mathbb{R}^n$  it holds that  $\mathbb{E}[\exp(\langle \mathbf{v}, \mathbf{x} \rangle)] \leq \exp(\|\mathbf{v}\|_2^2 \sigma^2 / 2)$ .

## 2. Problem Formulation

In this section, we introduce the *Dynamical Linear Bandits* (DLBs), formulate the learning problem, focusing on the learner-environment interaction, assumptions, policies, and definition of regret (Section 2.1). Then, we derive a closed-form expression for the optimal policy for DLBs (Section 2.2).

### 2.1. Setting

We consider the sequential interaction between a learner and an environment. In a Dynamical Linear Bandit (DLB), the environment is characterized by a *hidden* state, i.e., a  $n$ -dimensional real vector, initialized to  $\mathbf{x}_1 \in \mathcal{X}$ , where  $\mathcal{X} \subseteq \mathbb{R}^n$  is the state space. At each round  $t \in \mathbb{N}$ , the environment is in the hidden state  $\mathbf{x}_t \in \mathcal{X}$ , the learner chooses

an action, i.e., a  $d$ -dimensional real vector  $\mathbf{u}_t \in \mathcal{U}$ , where  $\mathcal{U} \subseteq \mathbb{R}^d$  is the action space. Then, the learner receives a noisy reward  $y_t = \langle \mathbf{I}, \mathbf{x}_t \rangle + \langle \mathbf{I}, \mathbf{u}_t \rangle + \eta_t \in \mathcal{Y}$ , where  $\mathcal{Y} \subseteq \mathbb{R}$  is the reward space,  $\mathbf{I} \in \mathbb{R}^n$ ,  $\mathbf{I} \in \mathbb{R}^d$  are unknown parameters, and  $\eta_t$  is a zero-mean  $\sigma^2$ -subgaussian random noise, conditioned to the past. Then, the environment evolves to the new state according to the unknown linear dynamics  $\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \mathbf{t}$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is the dynamic matrix,  $\mathbf{B} \in \mathbb{R}^{n \times d}$  is the action-state matrix, and  $\mathbf{t}$  is a zero-mean  $\sigma^2$ -subgaussian random noise, conditioned to the past, independent of  $\eta_t$ .<sup>1</sup>

**Remark 2.1.** *The setting proposed above is a particular case of a POMDP (Littman et al., 1995), in which the state  $\mathbf{x}_t$  is non-observable, while the learner has access to the noisy observation  $y_t$  that, in our setting, corresponds to the noisy reward too. Furthermore, the setting can be viewed as a MISO (Multiple Input Single Output) discrete-time LTI system (Kalman, 1963). Finally, the DLB reduces to (non-contextual) linear bandit (Lattimore & Szepesvári, 2020) when the hidden state does not affect the reward, i.e., when  $\mathbf{I} = \mathbf{0}$ .*

**Markov Parameters** We revise a useful representation, widely employed in the LTI literature, that allows expressing  $y_t$  in terms of the sequence of the most recent  $H + 1 \in \mathbb{N}$  actions  $(\mathbf{u}_s)_{s \in \mathbb{P}_t^H}$ , reward noise  $\eta_t$ ,  $H$  state noises  $(\mathbf{t}_s)_{s \in \mathbb{P}_t^H}$ , and starting state  $\mathbf{x}_{t-H}$  (Ho & Kálmán, 1966; Oymak & Ozay, 2019; Tsiamis & Pappas, 2019; Sarkar et al., 2021):

$$y_t = \sum_{s=0}^H \langle \mathbf{h}^{t, s}, \mathbf{u}_t \rangle + \underbrace{\mathbf{I}^\top \mathbf{A}^H \mathbf{x}_{t-H}}_{\text{starting state}} + \underbrace{\eta_t + \sum_{s=1}^H \mathbf{I}^\top \mathbf{A}^{s-1} \mathbf{t}_s}_{\text{noise}}, \quad (1)$$

where the sequence of vectors  $\mathbf{h}^{t, s} \in \mathbb{R}^d$  for every  $s \in \mathbb{N} \cup \{0\}$  are called *Markov parameters* and are defined as:  $\mathbf{h}^{t, 0} = \mathbf{I}$  and  $\mathbf{h}^{t, s} = \mathbf{B}^\top (\mathbf{A}^{s-1})^\top \mathbf{I}$  if  $s \geq 1$ . Furthermore, we introduce the *cumulative Markov parameters*, defined for every  $s, s^1 \in \mathbb{N} \cup \{0\}$  with  $s \leq s^1$  as  $\mathbf{h}^{J_s; s^1} = \sum_{l=s}^{s^1} \mathbf{h}^{t, l}$  and the corresponding limit as  $s^1 \rightarrow +\infty$ , i.e.,  $\mathbf{h}^{J_s; \infty} = \sum_{l=s}^{\infty} \mathbf{h}^{t, l}$ . Finally, we use the abbreviation  $\mathbf{h} = \mathbf{h}^{J_0; \infty} = \mathbf{B}^\top (\mathbf{I} - \mathbf{A})^{-1} \mathbf{I}$ . We will make use of the following standard assumptions.

**Assumption 2.1 (Boundedness).** *The following inequalities hold:  $\|\mathbf{I}\|_2 \leq \beta$ ,  $\|\mathbf{I}\|_2 \leq \beta$ ,  $\|\mathbf{B}\|_2 \leq B$ ,  $\|\mathbf{u}\|_2 \leq U$  with  $\mathbf{u} \in \mathcal{U}$ , and  $\|\mathbf{x}\|_2 \leq X$  with  $\mathbf{x} \in \mathcal{X}$ , and  $\sup_{\mathbf{u}, \mathbf{u}^1 \in \mathcal{U}} \langle \mathbf{u}, \mathbf{u}^1 \rangle \leq 1$ .<sup>2</sup>*

<sup>1</sup> $n$  is the order of the LTI system (Kalman, 1963). We make no assumption on the value of  $n$  and on its knowledge.

<sup>2</sup>The assumption of the bounded state norm  $\|\mathbf{x}\|_2 \leq X$  can

**Assumption 2.2 (Stability).** *The spectral radius of  $\mathbf{A}$  is strictly smaller than 1, i.e.,  $\rho(\mathbf{A}) < 1$ , and the maximum spectral norm to spectral radius ratio of the powers of  $\mathbf{A}$  is bounded, i.e.,  $\|\mathbf{A}^k\|_2 < +\infty$ .<sup>3</sup>*

The former assumption requires the boundedness of the norms of the relevant vectors, matrices, as well as states and actions; whereas the latter is related to the *stability* of the dynamic matrix  $\mathbf{A}$ , which is widely employed in discrete-time LTI literature (Oymak & Ozay, 2019; Lale et al., 2020b;a).

**Policies and Performance** The learner's behavior is modeled by a deterministic *policy*  $\pi = (\pi_t)_{t \in \mathbb{N}}$  defined, for every round  $t \in \mathbb{N}$ , as  $\pi_t : \mathcal{H}_{t-1} \rightarrow \mathcal{U}$ , mapping the history of observations  $H_{t-1} = (\mathbf{u}_1, y_1, \dots, \mathbf{u}_{t-1}, y_{t-1}) \in \mathcal{H}_{t-1}$  to an action  $\mathbf{u}_t = \pi_t(H_{t-1}) \in \mathcal{U}$ , where  $\mathcal{H}_{t-1} = (\mathcal{U} \times \mathcal{Y})^{t-1}$  is the set of histories of length  $t-1$ . The performance of a policy  $\pi$  is evaluated in terms of the *infinite-horizon expected average reward*:<sup>4</sup>

$$J(\pi) := \liminf_{H \in \mathbb{N}} \mathbb{E} \left[ \frac{1}{H} \sum_{t=1}^H y_t \right] \quad (2)$$

$$\text{where } \begin{cases} \mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \mathbf{t} \\ y_t = \langle \mathbf{I}, \mathbf{x}_t \rangle + \langle \mathbf{I}, \mathbf{u}_t \rangle + \eta_t \\ \mathbf{u}_t = \pi_t(H_{t-1}) \end{cases}, \quad t \in \mathbb{N},$$

where the expectation is taken w.r.t. the randomness of the state noise  $\mathbf{t}$  and reward noise  $\eta_t$ . If a policy  $\pi$  is *constant*, i.e.,  $\pi_t(H_{t-1}) = \mathbf{u}$  for every  $t \in \mathbb{N}$ , we abbreviate  $J(\mathbf{u}) = J(\pi)$ . A policy  $\pi$  is an *optimal infinite-horizon policy* if it maximizes the infinite-horizon expected average reward, i.e.,  $\pi \in \arg \max_{\pi} J(\pi)$ , whose performance is denoted as  $J^* := J(\pi^*)$ .

**Regret** The goal of the learner is to minimize the *online expected (policy) regret* by playing a policy  $\pi$ , competing against the optimal infinite-horizon policy  $\pi^*$  over a *learning horizon*  $T \in \mathbb{N}$ :  $\mathbb{E} R(\pi, T) := \mathbb{E} [\sum_{t=1}^T J^* - y_t]$ , where  $y_t$  is the sequence of rewards collected by playing  $\pi$  as in Equation (2). Furthermore, we introduce a different notion of regret, that will turn useful for analysis purposes, that we name *offline expected (policy) regret* that compares  $J^*$  with the infinite-horizon performance of the

be replaced with the assumption of bounded state noise. As shown in (Agarwal et al., 2019), this assumption can be relaxed by conditioning to the event that none of the noise vectors are ever large at the cost of an additional  $\log T$  factor in the regret.

<sup>3</sup>The latter is a mild assumption: if  $\mathbf{A}$  is diagonalizable as  $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$ , then  $\|\mathbf{A}^k\|_2 \leq \|\mathbf{Q}\|_2 \|\mathbf{\Lambda}^k\|_2 \|\mathbf{Q}^{-1}\|_2$  and it is finite. In particular, if  $\mathbf{A}$  is symmetric then  $\|\mathbf{A}^k\|_2 \leq \rho(\mathbf{A})^k$ .

<sup>4</sup>In Appendix B, we show that, under Assumption 2.2, the infinite-horizon setting is not dissimilar from the finite-horizon setting, provided that the horizon is sufficiently large.

action  $u_t \in \mathcal{U}$  played at each round  $t \in \{1, \dots, T\}$  by the agent. Let  $R_t = \sum_{q \in \mathcal{Q}} R_{t,q} \mathbf{1}_{\{u_t = q\}}$ . Clearly, the two regret definitions coincide when the system does not have a dynamics, i.e.,  $A = 0$ .

## 2.2. Optimal Policy

In this section, we derive a closed-form expression for the optimal policy  $u^*$  for the finite-horizon objective function, as introduced in Equation (2).

**Theorem 2.1 (Optimal Policy)** Under Assumptions 2.1 and 2.2, an optimal policy  $u^*$  maximizing the finite-horizon expected average reward  $\bar{r}_T$  as in Equation (2), is given by:

$$u^* = \underset{u \in \mathcal{U}}{\operatorname{argmax}} \sum_{q \in \mathcal{Q}} \mathbb{E} \left[ \sum_{t=1}^T R_{t,q} \mathbf{1}_{\{u_t = q\}} \right] \quad (3)$$

where  $u = \underset{u \in \mathcal{U}}{\operatorname{argmax}} \sum_{q \in \mathcal{Q}} \mu_q(x; y)$ .

Some remarks are in order. The optimal policy performs the constant action  $u^* \in \mathcal{U}$  which brings the system in the “most profitable” steady-state. Indeed, the expression  $\sum_{q \in \mathcal{Q}} \mu_q(x; y)$  can be rewritten expanding the cumulative Markov parameter  $\mathbf{a}_t^T = \sum_{q \in \mathcal{Q}} \mathbf{A}^t \mathbf{b}_q$  and  $\mathbf{b}_t = \sum_{q \in \mathcal{Q}} \mathbf{A}^t \mathbf{b}_q$  is the expression of the steady state  $\mathbf{A} \bar{x} = \mathbf{B} \bar{u}$ , when applying action  $\bar{u}$ . It is worth noting the role of Assumption 2.2 which guarantees the existence of the inverse  $\mathbf{A}^{-1}$ . In this sense, our problem shares the constant nature of the optimal policy with the linear bandit setting (Lattimore & Szepesvári, 2020), although our setting is characterized by an evolving state, which introduces a new trade-off in the action selection. From the LTI system perspective, this implies that we can focus on open-loop stationary policies only. The reason why this problem does not benefit from closed-loop policies, differently from other classical problems, such as the LQG (Abbasi-Yadkori & Szepesvári, 2011), lies in the linearity of the reward and in the additive nature of the noise components, which makes its presence irrelevant for control purposes. Nonetheless, as we shall see, our problem poses additional challenges compared to linear bandits since, in order to assess the quality of an action  $u \in \mathcal{U}$ , we need to let the system evolve to the steady-state and, then, observe the reward.

## 3. Algorithm

In this section, we present a anytime optimistic regret minimization algorithm for the DLB setting introduced in Section 2. Dynamical Linear Upper Confidence Bound (DynLin-UCB), whose pseudocode is reported in Algorithm 1, requires the knowledge of an upper-bound  $\bar{1}$  on the spectral radius of the dynamic matrix (i.e.,  $\rho(A) \leq \bar{1}$ )

<sup>5</sup>In Appendix B, we show that the optimal policy is non-stationary for the finite-horizon case.

and on the maximum spectral norm to spectral radius ratio  $\bar{1} \geq \rho(A)$  (i.e.,  $\rho(A) \leq \bar{1}$ ), as well as the bounds on the relevant quantities of Assumption 2. DynLin-UCB is based on the following simple observation. To assess the quality of an action  $u \in \mathcal{U}$ , we need to persist in applying it so that the system approximately reaches the corresponding steady-state and, then, observe the reward representing a reliable estimate of  $\sum_{q \in \mathcal{Q}} \mu_q(x; y)$ . We shall show that, under Assumption 2.2, the number of rounds needed to approximately reach such a steady state is logarithmic in the learning horizon  $T$  and depends on the upper bound of the spectral norm. After initializing the Gram matrix  $V_0 = I_d$  and the vectors  $\mathbf{s}_0$  and  $\mathbf{R}_0$  both to  $\mathbf{0}_d$  (line 1), DynLin-UCB subdivides the learning horizon in  $M = \lceil T/\epsilon \rceil$  epochs. Each epoch  $m \in \{1, \dots, M\}$  is composed of  $H_m = \lceil \log m \rceil$  rounds, where  $H_m$  is logarithmic in the epoch index  $m$ . At the beginning of each epoch  $m \in \{1, \dots, M\}$ , DynLin-UCB computes the upper confidence bound (UCB) index (line 4):

$$u_m = \underset{u \in \mathcal{U}}{\operatorname{argmax}} \sum_{q \in \mathcal{Q}} \left( \mathbf{R}_{t-1} + \frac{\mathbf{b}_t}{V_{t-1}} \right)^T \mathbf{u} \quad (4)$$

where  $\mathbf{R}_{t-1} = \sum_{s=1}^{t-1} \mathbf{b}_s \mathbf{b}_s^T$  is the regression estimator of the cumulative Markov parameter  $\mathbf{b}$  as employed in Equation (3) to define the optimal action and  $\mathbf{b}_t \neq \mathbf{0}$  is an exploration coefficient that will be defined later. Similar to the LinUCB algorithm (Lattimore & Szepesvári, 2020), the index  $u_m$  is designed to be optimistic, i.e.,  $\sum_{q \in \mathcal{Q}} \mu_q(x; y) \leq \sum_{q \in \mathcal{Q}} \mathbf{R}_{t-1} + \frac{\mathbf{b}_t}{V_{t-1}} \mathbf{u}_m$  in high-probability for all  $u_m \in \mathcal{U}$ . Then, the optimistic action  $u_m = \underset{u \in \mathcal{U}}{\operatorname{argmax}} \sum_{q \in \mathcal{Q}} \mathbf{R}_{t-1} + \frac{\mathbf{b}_t}{V_{t-1}} \mathbf{u}$  is executed (line 5) and persisted for the next  $H_m$  rounds (lines 7-10). The length of the epoch  $H_m$  is selected such that, under Assumption 2.2, the system has approximately reached the steady state after  $H_m = \lceil \frac{\log m}{\epsilon} \rceil$  rounds. In this way, at the end of epoch  $m$ , the reward  $r_m$  is an almost-unbiased sample of the steady-state performance  $\sum_{q \in \mathcal{Q}} \mu_q(x; y)$ . This sample is employed to update the Gram matrix estimator  $V_t$  and the vector  $\mathbf{b}_t$  (line 12), while the samples collected in the previous  $H_m$  rounds are discarded (line 8). It is worth noting that by setting  $H_m = 0$  for all  $m \in \{1, \dots, M\}$ , DynLin-UCB reduces to LinUCB. The following sections provide the concentration of the estimator  $\mathbf{R}_{t-1}$  of  $\mathbf{R}$  (Section 3.1) and the regret analysis of DynLin-UCB (Section 3.2).

### 3.1. Self-Normalized Concentration Inequality for the Cumulative Markov Parameter

In this section, we provide a self-normalized concentration result for the estimator  $\mathbf{R}_t$  of the cumulative Markov parameter

<sup>6</sup>As an alternative, one can consider a more demanding requirement of the knowledge of a bound on the spectral norm  $\rho(A) \leq \bar{1}$  of  $A$ . Similar assumptions regarding the knowledge of analogous quantities are considered in the literature, including decay of Markov operator norms (Simchowitz et al., 2020) and strong stability (Plevrakis & Hazan, 2020).





8, and those specified in Assumption 2.1:

@ PJTK:

$$c_1 = \frac{2^{-2} \log \frac{1}{\epsilon} + \frac{d}{2} \log \frac{1}{\epsilon} + \frac{tU^2}{d}}{2^{-2} \log \frac{1}{\epsilon} + \frac{d}{2} \log \frac{1}{\epsilon} + \frac{tU^2}{d}}; \quad (5)$$

where:

$$c_1 = U - \frac{UB}{1 - \epsilon} X; \\ c_2 = \frac{B}{1 - \epsilon}; \\ -2 \quad 2 \quad 1 \quad \frac{2^{-2}}{\rho 1 - \epsilon^2 q};$$

The analysis poses additional challenges than that of Lin-UCB. Indeed, by a straightforward application of the proof strategy for Lin-UCB, it is possible, under Assumptions 2.1 and 2.2, to comfortably obtain an expected regret bound, involving just the steady-state performances (details in Appendix A.2).

$$E R_{\text{off}}^{\text{DynLin-UCB}}; T, q = E \sum_{t=1}^T J_{\mu_t} q \\ \propto \frac{d \bar{T}}{\rho 1 - \epsilon^{3/2}}; \quad (6)$$

However, when applying action  $u_t$ , the DLB does not immediately reach  $J_{\mu_t} q$  as the system needs to converge to the steady-state according to its dynamics (Equation 2). Consequently, the expected online reward experiences a transitional phase. In order to obtain a sublinear expected online regret, we need to guarantee that during the transitional phase, when moving from one epoch  $P_{JM}K$  to the next one  $n-1$ , performance does not degrade too much. This property can be guaranteed by the following assumption.

Assumption 3.1. For every action  $u_1; u_2 \in \mathcal{U}$  it holds that:

$$x; u_1 y \times B^T p \quad A q^T; u_2 y \neq \\ \min_{u \in \mathcal{U}} J_{\mu} q x; u y \times B p \quad A q^T; u y.$$

The rationale behind the assumption is the following. We need to guarantee that when we converge to the steady-state associated with action  $u_2$ , i.e.,  $T p \quad A q^T B u_2$ , and we apply action  $u_1$ , the instantaneous expected reward

cannot take values below the minimum between the steady-state performances of actions  $u_1$  and  $u_2$ , i.e.,  $J_{\mu_1} q$  and  $J_{\mu_2} q$  respectively. In this way, the transitional phase does not have a relevant impact on the regret rate. It is worth noting that the assumption is surely fulfilled for linear bandits (i.e.,  $0_n$ ) and strictly proper LTI systems (i.e.,  $0_d$ ). Under this assumption, we are able to provide a bound on the expected online regret.

Theorem 3.2. Under Assumptions 2.1, 2.2, and 3.1, Algorithm 1 suffers an expected online regret bounded by:

$$E R_{\text{DynLin-UCB}}; T, q = E \sum_{t=1}^T J_{\mu_t} y_t \\ \propto \frac{\rho 1 \} A \} \rho \bar{T}}{\rho 1 - \epsilon^{3/2}};$$

Some remarks are in order. First of all, compared to the expected of the regret bound, we obtain a multiplicative factor that corresponds to the Frobenius norm of the dynamic matrix  $A$ . It is worth noting that the Frobenius norm is related to the spectral gap, since  $\|A\|_F \propto n \rho A q \rho A q$ . Second, if the underlying problem does not have a dynamics, i.e.,  $A = 0$ , and we choose  $\epsilon = 0$ , we obtain a regret bound of order  $\bar{T} q$  which corresponds to the regret bound in UCB. Clearly, the dependence on  $\epsilon$  is relevant and with too large a value of  $\epsilon$  compared to the optimization horizon  $T$  (e.g.,  $\epsilon = 1/T^{1/3}$ ) could make the regret degenerate to linear. This is a case in which the underlying system is as slow that the whole horizon is insufficient to approximately reach steady state.

#### 4. Related Works

Our formulation is placed at the intersection of three research areas: (i) bandits with delayed, aggregated, and composite feedback (Joulani et al., 2013), (ii) Partially-Observable Markov Decision Processes (POMDPs, Åström, 1965), and online control for Linear Time-Invariant (LTI) dynamical systems (Hespanha, 2018). In this section, we survey the related approaches in these areas in comparison with our model and algorithm.

Bandits with Delayed/Aggregated/Composite Feedback The Multi-Armed Bandit setting has been widely employed as a principled approach to address sequential decision-making problems (Lattimore & Szepesvári, 2020). The possibility of experiencing delayed rewards has been introduced in (Joulani et al., 2013) and widely exploited in the advertising applications (Chapelle, 2014; Vernade et al., 2017). A large number of approaches have extended this setting either considering stochastic delays (Vernade et al., 2020), unknown delays (Li et al., 2019; Lancewicki et al.,

<sup>8</sup>For interpretability reasons, we highlight the dependencies on  $T, \epsilon$ , and only and disregard other constant and the logarithmic terms.

2021), arm-dependent delays (Manegueu et al., 2020), no hidden state and the action. Nevertheless, several works accounted for the presence of constraints (Undurti & How, et al., 2022). Some methods relaxed the assumption that the individual reward is revealed after the delay expires, admitting the possibility of receiving anonymous feedback, which can be aggregated (Pike-Burke et al., 2018; Zhang et al., 2021) or composite (Cesa-Bianchi et al., 2018; Garg & Akash, 2019; Wang et al., 2021). Most of these approaches are able to achieve  $\mathcal{O}(\sqrt{T})$  regret, plus additional terms depending on the extent of the delay. In our DLBs, the reward is generated over time as a combined effect of past and present actions through the hidden state while these approaches generate the reward instantaneously and reveal it (individually or in aggregate) to the learner in the future and no underlying state dynamics is present.

2010; Kim et al., 2011; Isom et al., 2008) without exploiting the linearity and with no regret guarantees.

## 5. Numerical Simulations

In this section, we provide a numerical validation of DynLin-UCB in a synthetically generated domain. The goal of this simulation is to highlight the behavior of DynLin-UCB in comparison with bandit baselines, describing advantages and disadvantages. We start by introducing the DLB setting considered and the baselines for comparison, and then discussing the obtained results. The complete experimental results are reported in Appendix E.

### Online Control of Linear Time-Invariant Systems

The particular structure imposed by linear dynamics makes our approach comparable to LTI online control for partially-observable systems (e.g., Lale et al., 2020b; Simchowitz et al., 2020; Plevrakis & Hazan, 2020). While the dynamical model is similar, in online control of LTI systems the perspective is quite different. Most of the works either consider the Linear Quadratic Regulator (Mania et al., 2019; Lale et al., 2020b) or (strongly) convex objective functions (Mania et al., 2019; Simchowitz et al., 2020; Lale et al., 2020a) that simulates a total budget of 15 to be allocated to the achieving, in most cases  $\mathcal{O}(\sqrt{T})$  regret for strongly convex functions and  $\mathcal{O}(T^{2/3})$  for convex functions. Recently,  $\mathcal{O}(\sqrt{T})$  regret rate has been obtained for convex functions too, by means of geometric exploration methods (Plevrakis & Hazan, 2020). Furthermore, (Lale et al., 2020a) reacts  $\mathcal{O}(\log T)$  regret in the case of strongly convex cost functions competing against the best persistently exciting controller (i.e., a controller that implicitly maintains a non-null exploration). Some approaches are designed to deal with adversarial noise (Simchowitz et al., 2020). All of these approaches, however, look for the best closed-loop controller within a specific class (e.g., disturbance response control (Ladakori et al., 2011) designed for linear bandits, and & Bosch, 1993)). These controllers, however, do not allow us to easily incorporate constraints on the action space which could be of crucial importance in practice, especially in advertising domains. Our DynLin-UCB works with an arbitrary action space and, thanks to the linearity of the reward in the hidden state, does not require considering complex closed-loop controllers.

**Setting** We consider a DLB defined by means of the following matrices  $A = \begin{bmatrix} 0.2 & 0 & 0 \\ 0.25 & 0 & 0 \\ 0 & 0.5 & 0.1 \end{bmatrix}$ ,  $B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ ,  $C = \begin{bmatrix} 0.2 & 0 & 0 \\ 0 & 0.5 & 0.1 \end{bmatrix}$ , and a Gaussian noise with  $\Sigma = 0.03$  (diagonal covariance matrix for the state noise). This way the spectral gap of the dynamical matrix is  $\rho(A) = 0.2$  and  $\rho(A) < 1$ . Moreover, the cumulative Markov parameter is given by  $\sum_{t=0}^{\infty} C A^t B^T = \begin{bmatrix} 0.56 & 0.5 & 0.11 \end{bmatrix}$ . We consider the action space  $\mathcal{U} = \{u_1, u_2, u_3\}$  with  $u_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ ,  $u_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ ,  $u_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$  that simulates a total budget of 15 to be allocated to the three platforms. Thus, a “myopic” agent would simply look at how the action immediately propagates to the reward through  $C$ , will invest the budget in the first component of the action, that is weighted by 0.5. Instead, a “far-sighted” agent, aware of the system evolution, will look at the cumulative Markov parameter, realizing that the most convenient action is investing in the first component, weighted by 0.56. Therefore, the optimal action is  $u_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$  leading to  $J = 0.81$ .

**Baselines** We consider as baselines Lin-UCB (Abbasi-Yadkori et al., 2011) designed for linear bandits, and Exp3 (Auer et al., 2002) usually employed in adversarial settings. Concerning the hyperparameters for both DynLin-UCB and Lin-UCB, we tested  $\{P \log T, \log T\}$ .<sup>10</sup> Since Exp3 requires finite actions, we consider a subset of the action space that surely contain the optimal action. Given that  $\mathcal{U}$  is a polytope and since the objective is linear, we enumerate all the vertices of  $\mathcal{U}$  that, in this setting, can-

Partially Observable Markov Decision Processes As already noted, looking at DLBs in their generality, we realize that our model is a particular subclass of the Partially Observable Markov Decision Processes (POMDP) (Åström, 1965). However, in the POMDP literature, no particular structure of the hidden state dynamics is assumed. The specific linear dynamics are rarely considered, as well as the possibility of a reward that is a linear combination of the factor in Equation (5).

<sup>10</sup>It is worth noting that the decision of using diagonal matrices is just for explanation purposes and w.l.o.g. (at least in the class of diagonalizable dynamic matrices). Indeed, we are just interested to the cumulative Markov parameter and we could have obtained the same results with an equivalent (non-diagonal) representation, by applying an inevitable transformation as  $A' = T^{-1} A T$  and  $B' = T^{-1} B$ .

<sup>10</sup>For DynLin-UCB,  $\log T$  is a nearly optimal choice for as it can be seen by looking at the first two addenda of the exploration

Figure 1. Cumulative regret as a function of the rounds comparing DynLin-UCB, Lin-UCB (both with  $P \propto 1/\log T$ ), and Exp3 (3 runs, mean  $\pm$  std).

not be more than 6. The learning rate for Exp3 is set as prescribed in the original paper (Auer et al., 2002).

Comparison with Lin-UCB and Exp3 Figure 1 shows the performance in terms of cumulative regret of DynLin-UCB, Lin-UCB and Exp3. The experiments are conducted over a time horizon of 500k rounds. For DynLin-UCB, we employed, for the sake of this experiment, the true value of the spectral gap, i.e.,  $\rho = 0.2$ . First of all, we observe that Exp3 suffers a significantly large cumulative regret. Moreover, both versions of Lin-UCB suffer linear regret. Indeed, even for a quite fast system ( $\rho = 0.2$ ), ignoring the system dynamics, and the presence of the hidden state, has made Lin-UCB committing to the sub-optimal (myopic) action  $\mathbf{a}^* = [0.5; 1; 0]^T$  with performance  $\approx 0.78 \sqrt{T}$ , with also a relevant variance. On the other hand, DynLin-UCB is able to maintain a smaller and stable (variance is negligible) sublinear regret in both its versions, with a notable advantage when using  $\beta = \log T$ .

Sensitivity to the Choice of  $\beta$  The upper bound of the spectral radius  $\rho = 0.2$  represents a crucial parameter of DynLin-UCB. While an overestimation  $\beta \gg \rho$  does not compromise the regret rate, but tends to slow down the convergence process, a severe underestimation  $\beta \ll \rho$  might prevent learning at all. In Figure 2, we test DynLin-UCB against a misspecification of  $\beta$ , when  $\beta = \log T$ . We can see that by considering  $\beta = 2\rho$ , DynLin-UCB experiences a larger regret but still sublinear and smaller w.r.t. Lin-UCB with  $\beta = \log T$ . Even by reducing  $\beta = 0.1; 0.05$ , DynLin-UCB is able to keep the regret sublinear, showing a remarkable robustness

Figure 2. Cumulative regret as a function of the rounds comparing Lin-UCB and DynLin-UCB with  $\beta = \log T$ , varying the upper bound on the spectral radius (3 runs, mean  $\pm$  std).

misspecification. Clearly, when setting  $\beta$  to 0 makes the regret almost degenerates to linear.

## 6. Discussion and Conclusions

In this paper, we have introduced the Dynamical Linear Bandits (DLBs), a novel model to represent sequential decision-making problems in which the system is characterized by a non-observable hidden state that evolves according to a linear dynamics and by an observable noisy reward that linearly combines the hidden state and the played action. This model accounts for scenarios that cannot be easily represented by existing bandit models that consider delayed and aggregated feedback. On top of it, we have proposed a novel any-time optimistic regret minimization approach, DynLin-UCB, that, under suitable assumptions, is able to achieve sub-linear regret. The numerical simulation in a synthetic domains, succeeded in showing that, in a system where the baselines suffer linear regret, our algorithm enjoys sublinear regret. Furthermore, DynLin-UCB proved to be robust to misspecification of its most relevant hyperparameter. To the best of our knowledge, this is the first work addressing this family of problems, characterized by a hidden linear dynamics, with a bandit-like approach. Short-term future directions include the study of the complexity of the problem by deriving regret lower bounds and understanding whether the considered assumptions (especially Assumption 3.1) and the knowledge of the upper bound are actually unavoidable. Long-term future directions might focus on extending the present approach to non-linear system dynamics (e.g., to model saturation phenomena in the awareness) and non-stationary dynamics (e.g., to account for the natural market evolution).



## References

- Abbasi-Yadkori, Y. and Szepesvári, C. Regret bounds for the adaptive control of linear quadratic systems. *COLT 2011 - The 24th Annual Conference on Learning Theory*, June 9-11, 2011, Budapest, Hungary, pp. 1–26, 2011.
- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In Shawe-Taylor, J., Zemel, R. S., Bartlett, P. L., Pereira, F. C. N., and Weinberger, K. Q. (eds.) *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain* pp. 2312–2320, 2011.
- Agarwal, N., Hazan, E., and Singh, K. Logarithmic regret for online control. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 10175–10184, 2019.
- Åström, K. J. Optimal control of markov processes with incomplete state information. *Journal of mathematical analysis and applications* 10(1):174–205, 1965.
- Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R. E. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.* 32(1):48–77, 2002.
- Berman, R. Beyond the last touch: Attribution in online advertising. *Marketing Science* 37(5):771–792, 2018.
- Cesa-Bianchi, N., Gentile, C., and Mansour, Y. Nonstochastic bandits with composite anonymous feedback. *Conference On Learning Theory, COLT 2018, Stockholm, Sweden, 6-9 July 2018*, pp. 750–773, 2018.
- Chapelle, O. Modeling delayed feedback in display advertising. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining KDD '14*, pp. 1097–1105, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450329569.
- Court, D., Elzinga, D., Mulder, S., and Vetvik, O. J. The consumer decision journey. *McKinsey Quarterly* 3:96–107, 2009.
- Garg, S. and Akash, A. K. Stochastic bandits with delayed composite anonymous feedback. *CoRR abs/1910.01161*, 2019.
- Hespanha, J. P. *Linear Systems Theory: Second Edition* Princeton University Press, 2018.
- Ho, B. L. and Kálmán, R. E. Effective construction of linear state-variable models from input/output functions. *at-Automatisierungstechnik* 14(1-12):545–548, 1966.
- Hoban, P. R. and Bucklin, R. E. Effects of internet display advertising in the purchase funnel: Model-based insights from a randomized field experiment. *Journal of Marketing Research* 52(3):375–393, 2015.
- Hsu, D., Kakade, S., and Zhang, T. A tail inequality for quadratic forms of subgaussian random vectors. *Electronic Communications in Probability* 17:1–6, 2012.
- Som, J. D., Meyn, S. P., and Braatz, R. D. Piecewise linear dynamic programming for constrained pomdps. In Fox, D. and Gomes, C. P. (eds.) *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence, AAAI 2008, Chicago, Illinois, USA, July 13-17, 2008*, pp. 291–296. AAAI Press, 2008.
- Do, S., Hatano, D., Sumita, H., Takemura, K., Fukunaga, T., Kakimura, N., and Kawarabayashi, K. Delay and cooperation in nonstochastic linear bandits. *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, pp. 1020–1030, 2020.
- Jin, T., Lancewicki, T., Luo, H., Mansour, Y., and Rosenberg, A. Near-optimal regret for adversarial MDP with delayed bandit feedback. *CoRR abs/2201.13172*, 2022.
- Joulani, P., György, A., and Szepesvári, C. Online learning under delayed feedback. *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, pp. 1453–1461, 2013.
- Kalman, R. E. Mathematical description of linear dynamical systems. *Journal of the Society for Industrial and Applied Mathematics, Series A: Control* 1(2):152–192, 1963.
- Kim, D., Lee, J., Kim, K., and Poupart, P. Point-based value iteration for constrained pomdps. In Walsh, T. (ed.), *IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011*, pp. 1968–1974. IJCAI/AAAI, 2011.
- Lale, S., Azizzadenesheli, K., Hassibi, B., and Anandkumar, A. Logarithmic regret bound in partially observable linear dynamical systems. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.) *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, pp. 2020a–2020a, 2020a.
- Lale, S., Azizzadenesheli, K., Hassibi, B., and Anandkumar, A. Regret minimization in partially observable linear quadratic control. *CoRR abs/2002.00082*, 2020b.

- Lancewicki, T., Segal, S., Koren, T., and Mansour, Y. Stochastic multi-armed bandits with unrestricted delay distributions. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event* pp. 5969–5978, 2021.
- Lattimore, T. and Szepesvári, G. *Bandit algorithms* Cambridge University Press, 2020.
- Li, B., Chen, T., and Giannakis, G. B. Bandit online learning with unknown delays. In *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan* pp. 993–1002, 2019.
- Li, H. X. and Bosch, P. P. J. V. D. A robust disturbance-based control and its application. *International Journal of Control*, 58(3):537–554, 1993.
- Littman, M. L., Cassandra, A. R., and Kaelbling, L. P. Learning policies for partially observable environments: Scaling up. In *Machine Learning, Proceedings of the Twelfth International Conference on Machine Learning, Tahoe City, California, USA, July 9-12, 1995* pp. 362–370, 1995.
- Manegueu, A. G., Vernade, C., Carpentier, A., and Valko, M. Stochastic bandits with arm-dependent delays. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event* pp. 3348–3356, 2020.
- Mania, H., Tu, S., and Recht, B. Certainty equivalence is efficient for linear quadratic control. In *Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R. (eds) Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada* pp. 10154–10164, 2019.
- Oymak, S. and Ozay, N. Non-asymptotic identification of LTI systems from a single trajectory. *2019 American Control Conference, ACC 2019, Philadelphia, PA, USA, July 10-12, 2019* pp. 5655–5661, 2019.
- Pike-Burke, C., Agrawal, S., Szepesvári, C., and Grünewälder, S. Bandits with delayed, aggregated anonymous feedback. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Stockholm, Sweden, July 10-15, 2018* pp. 4102–4110, 2018.
- Plevrakis, O. and Hazan, E. Geometric exploration for online control. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual* pp. 2020–2020, 2020.
- Sarkar, T., Rakhlin, A., and Dahleh, M. A. Finite time LTI system identification. *J. Mach. Learn. Res.* 22:26:1–26:61, 2021.
- Simchowitz, M., Singh, K., and Hazan, E. Improper learning for non-stochastic control. In *Abernethy, J. D. and Agarwal, S. (eds.) Conference on Learning Theory, COLT 2020, 9-12 July 2020, Virtual Event [Graz, Austria]* pp. 1–11, 2020.
- Thune, T. S., Cesa-Bianchi, N., and Seldin, Y. Nonstochastic multiarmed bandits with unrestricted delays. *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada* pp. 6538–6547, 2019.
- Tsiamis, A. and Pappas, G. J. Finite sample analysis of stochastic system identification. *58th IEEE Conference on Decision and Control, CDC 2019, Nice, France, December 11-13, 2019* pp. 3648–3654, 2019.
- Undurti, A. and How, J. P. An online algorithm for constrained pomdps. *IEEE International Conference on Robotics and Automation, ICRA 2010, Anchorage, Alaska, USA, 3-7 May 2010* pp. 3966–3973. IEEE, 2010.
- Vernade, C., Cappé, O., and Perchet, V. Stochastic bandit models for delayed conversions. In *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence, UAI 2017, Sydney, Australia, August 11-15, 2017* pp. 2017–2017, 2017.
- Vernade, C., Carpentier, A., Lattimore, T., Zappella, G., Ermis, B., and Brückner, M. Linear bandits with stochastic delayed feedback. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event* pp. 9712–9721, 2020.
- Wang, S., Wang, H., and Huang, L. Adaptive algorithms for multi-armed bandit with composite and anonymous feedback. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021* pp. 10210–10217, 2021.
- Zhang, M., Tsuchida, R., and Ong, C. S. Gaussian process bandits with aggregated feedback. *CoRR abs/2112.13029*, 2021.
- Åström, K. Optimal control of markov processes with incomplete state information. *Journal of Mathematical Analysis and Applications* 50(1):174–205, 1965. ISSN 0022-247X.

## A. Proofs and Derivations

In this section, we provide the proofs we have omitted in the main paper.

### A.1. Proofs of Section 2

**Theorem 2.1 (Optimal Policy)** Under Assumptions 2.1 and 2.2, an optimal policy maximizing the in nite-horizon expected average reward  $J_{H,p,q}$  as in Equation (2), is given by:

$$\pi^* = \arg \max_{\pi} J_{H,p,q}(\pi) \quad (3)$$

where  $u = \arg \max_{u \in \mathcal{U}} \sum_{t=1}^H \mathbb{E} [r_t | u_t = u]$

**Proof.** Referring to the notation of Appendix B, we first observe that for every policy we have  $J_{H,p,q} \leq \liminf_{H \rightarrow \infty} J_{H,p,q}$  where  $J_{H,p,q} = \frac{1}{H} \mathbb{E} \sum_{t=1}^H r_t$  is the  $H$ -horizon expected average reward. Let us start with Equation (20), a fixed finite  $H$ , and considering the sequence of actions  $u_1, u_2, \dots, u_H$  generated by policy  $\pi$ :

$$J_{H,p,q} = \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{t=1}^H \mathbb{E} [r_t | u_t]$$

$$= \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{t=1}^H \mathbb{E} [r_t | u_t]$$

Now, we consider two bounds on  $J_{H,p,q}$  obtained by an application of Cauchy-Schwarz inequality on the second addendum:

$$J_{H,p,q} \leq \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{t=1}^H \mathbb{E} [r_t | u_t]$$

$$= \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{t=1}^H \mathbb{E} [r_t | u_t]$$

Concerning the term  $\mathbb{E} [r_s | u_s]$ , we have that  $\mathbb{E} [r_s | u_s] \leq U$ , having used Jensen's inequality and under Assumption 2.1. Regarding the second term, using Assumptions 2.2 and 2.1, we obtain:

$$\mathbb{E} [r_s | u_s] \leq \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s]$$

$$= \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s]$$

Plugging this result into the summation we obtain:

$$\frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s] + \frac{1}{H} \sum_{s=1}^H \mathbb{E} [r_s | u_s]$$

It is simple to observe that the last term approaches zero as  $H \rightarrow \infty$ . Moreover, with analogous argument, it can be proved that  $\frac{1}{H} \sum_{t=1}^H \mathbb{E} [r_t | u_t] \rightarrow 0$  as  $H \rightarrow \infty$ . Thus, we have that  $\liminf_{H \rightarrow \infty} J_{H,p,q} = \liminf_{H \rightarrow \infty} J_{H,p,q}$

Consequently, by the squeezing theorem of limits, we have:

$$\begin{aligned}
 J_{p,q} &= \liminf_{H \rightarrow \infty} J_{H,p,q} = \liminf_{H \rightarrow \infty} J_{H,p,q}^{\circ} \\
 &= \liminf_{H \rightarrow \infty} \frac{1}{H} \sum_{s=1}^H x_{h_s}^T E_{u_s} y_s = h^T \liminf_{H \rightarrow \infty} \frac{1}{H} \sum_{s=1}^H E_{u_s} y_s :
 \end{aligned}$$

It follows that an optimal policy is a policy that plays the constant action  $p \in \arg \max_{u \in U} x_h^T u$ . □

### A.2. Proofs of Section 3

**Theorem 3.1 (Self-Normalized Concentration)** Let  $\{p_t\}_{t \in \mathbb{N}}$  be the sequence of solutions of the regression problems computed by Algorithm 1. Then, under Assumption 2.1 and 2.2, for every  $\epsilon > 0$  and  $P \in (0, 1]$ , with probability at least  $1 - \epsilon$ , it holds that:

$$\begin{aligned}
 \text{① } \forall t \in \mathbb{N} : \\
 R_t &= h^T \frac{\sum_{s=1}^t y_s - \log p_t}{d^{V_t}} - \frac{1}{2} \log \frac{\det p_t^T V_t}{d} ; \\
 & \text{where: } c_1 = U \cdot \frac{UB}{1 - \rho_A}, \quad c_2 = \frac{B \cdot \rho_A}{1 - \rho_A}, \\
 & \quad r^2 = 2 \cdot \frac{\rho_A^2}{1 - \rho_A^2} :
 \end{aligned}$$

where:

$$\begin{aligned}
 c_1 &= U \cdot \frac{UB}{1 - \rho_A} ; \\
 c_2 &= \frac{B \cdot \rho_A}{1 - \rho_A} ; \\
 r^2 &= 2 \cdot \frac{\rho_A^2}{1 - \rho_A^2} :
 \end{aligned}$$

**Proof.** First of all, let us properly relate the round  $t \in \mathbb{N}$  and the index of the epoch  $m \in \mathbb{N}$ . For every epoch  $m \in \mathbb{N}$ , we denote with  $t_m$  the last round of epoch  $m$  (i.e., the one in which we update the relevant matrix  $V_t$ ):<sup>11</sup>

$$t_0 = 0; \quad t_m = t_{m-1} + 1 \quad \forall m \in \mathbb{N}$$

We now proceed at defining suitable filtrations. Let  $\mathcal{F}_t = \sigma(p_{1:t}, y_{1:t}, u_{1:t})$  such that for every  $t \in \mathbb{N}$ , the random variables  $\{u_1, y_1, \dots, u_t, y_t, u_t\}$  are  $\mathcal{F}_t$ -measurable, i.e.  $\mathcal{F}_t = \sigma(p_{1:t}, y_{1:t}, u_{1:t}, y_{1:t}, u_t)$ . Let us also consider the filtration indexed by  $m$ , denoted with  $\mathcal{F}_m = \sigma(p_{1:t_m}, y_{1:t_m}, u_{1:t_m})$  and defined for all  $m \in \mathbb{N}$  as  $\mathcal{F}_m = \mathcal{F}_{t_m}$ . Thus, the random variables  $\mathcal{F}_m$ -measurable are those realized until the end of epoch  $m$ .

Since the estimates do not change within an epoch, we need to guarantee the statement for all rounds  $t \in \mathbb{N}$  only. For these rounds, we define the following quantities:

$$\begin{aligned}
 y_m &= y_{t_m} ; \\
 u_m &= u_{t_m} ; \quad (\text{or any } u_l \text{ with } l \in [t_{m-1} + 1; t_m] \text{ since they are all equal}) \\
 r_m &= \frac{1}{\sum_{s=1}^{t_m} x_{h_s}^T A^{-1} x_{h_s}} ; \\
 x_{m-1} &= x_{t_{m-1}} ; \\
 h_m &= h_{t_m} ; \\
 V_m &= V_{t_m} ; \\
 b_m &= b_{t_m} :
 \end{aligned}$$

<sup>11</sup>It is worth noting that the variables  $t_m$  are deterministic.



We prove that  $\{r_m\}_{m \in \mathcal{P}_{JM}^K}$  is a martingale difference process adapted to the filtration  $\{\mathcal{F}_t\}_{t \in \mathcal{P}_{JM}^K}$ . To this end, we recall that, by construction,  $\{p_t\}_{t \in \mathcal{P}_{JM}^K}$  and  $\{q_t\}_{t \in \mathcal{P}_{JM}^K}$  are martingale difference processes adapted to the filtration  $\{\mathcal{F}_t\}_{t \in \mathcal{P}_{JM}^K}$ . It is clear that  $r_m$  is  $\mathcal{F}_m$ -measurable and, being  $\sigma$ -subgaussian it is absolutely integrable. Furthermore, using the tower law of expectation:

$$\begin{aligned} \mathbb{E}[r_m | \mathcal{F}_{t_m-1}] &= \mathbb{E}\left[\sum_{s=1}^{H_m-1} \mathbb{1}\{A^s = 1\} r_{t_m-s} \mid \mathcal{F}_{t_m-1}\right] \\ &= \mathbb{E}\left[\sum_{s=1}^{H_m-1} \mathbb{1}\{A^s = 1\} \mathbb{E}[r_{t_m-s} | \mathcal{F}_{t_m-s-1}] \mid \mathcal{F}_{t_m-1}\right] = 0; \end{aligned}$$

since the system is operating by persisting the action after having decided it at the beginning of the epoch. Thus, by exploiting the decomposition in Equation (1), we can write:

$$\begin{aligned} \mathbf{y}_m &= \mathbf{y}_{t_m} + \mathbf{x}^{J_0;H_m-1K} \mathbf{a}_m \mathbf{y} + \sum_{s=1}^{H_m-1} \mathbb{1}\{A^s = 1\} \mathbf{x}_{t_m-s} \\ &= \mathbf{x}^{J_0;H_m-1K} \mathbf{a}_m \mathbf{y} + \sum_{s=1}^{H_m-1} \mathbb{1}\{A^s = 1\} \mathbf{x}_{t_m-s} + \mathbf{r}_m \\ &= \mathbf{x}^h; \mathbf{a}_m \mathbf{y} + \mathbf{x}^{J_{H_m-2;8}^M} \mathbf{a}_m \mathbf{y} + \sum_{s=1}^{H_m-1} \mathbb{1}\{A^s = 1\} \mathbf{x}_{t_m-s} + \mathbf{r}_m; \end{aligned} \tag{8}$$

where we simply exploit the identity  $\mathbf{h}^{J_0;H_m-1K} = \mathbf{h}^{J_{H_m-2;8}^M}$ . We now introduce the following vectors and matrices:

$$\begin{aligned} \mathbf{U}_m &= \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix} \in \mathbb{R}^{m \times d}; & \mathbf{y}_m &= \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_m \\ \vdots \\ \mathbf{y}_m \end{bmatrix} \in \mathbb{R}^{m \times d}; \\ \mathbf{r}_m &= \begin{bmatrix} r_1 \\ \vdots \\ r_m \\ \vdots \\ r_m \end{bmatrix} \in \mathbb{R}^{m \times 1}; & \mathbf{r}_m &= \begin{bmatrix} \mathbb{1}\{A^{H_1} = 2\} \mathbf{x}_0 \\ \vdots \\ \mathbb{1}\{A^{H_m} = 2\} \mathbf{x}_{m-1} \end{bmatrix} \in \mathbb{R}^{m \times 1}; \\ \mathbf{g}_m &= \begin{bmatrix} \mathbf{x}^{J_{H_1-1;8}^M} \mathbf{a}_1 \mathbf{y} \\ \vdots \\ \mathbf{x}^{J_{H_m-1;8}^M} \mathbf{a}_m \mathbf{y} \end{bmatrix} \in \mathbb{R}^{m \times d}; \end{aligned}$$

Using the vectors and matrices above, we observe that  $\mathbf{h}_m = \mathbf{U}_m^T \mathbf{U}_m$  and  $\mathbf{b}_m = \mathbf{U}_m^T \mathbf{y}_m$ . Furthermore, by exploiting Equation (8), we can write:

$$\mathbf{y}_m = \mathbf{U}_m \mathbf{h}_m + \mathbf{g}_m + \mathbf{r}_m + \mathbf{r}_m;$$

Let us consider the estimator  $\hat{\mathbf{h}}_m$  at  $\mathcal{P}_{JM}^K$

$$\begin{aligned} \hat{\mathbf{h}}_m &= \mathbf{V}_m^{-1} \mathbf{b}_m = \mathbf{I} \mathbf{U}_m^T \mathbf{U}_m^{-1} \mathbf{U}_m^T \mathbf{y}_m \\ &= \mathbf{I} \mathbf{U}_m^T \mathbf{U}_m^{-1} \mathbf{U}_m^T \mathbf{U}_m \mathbf{h}_m + \mathbf{U}_m^T \mathbf{g}_m + \mathbf{U}_m^T \mathbf{r}_m + \mathbf{U}_m^T \mathbf{r}_m \\ &= \mathbf{h}_m + \mathbf{I} \mathbf{U}_m^T \mathbf{U}_m^{-1} \mathbf{U}_m^T \mathbf{g}_m + \mathbf{U}_m^T \mathbf{r}_m + \mathbf{U}_m^T \mathbf{r}_m; \end{aligned}$$

We now proceed at bounding the  $\|\cdot\|_{\mathbf{V}_m}$ -norm, and exploit the triangle inequality:

$$\begin{aligned} \|\hat{\mathbf{h}}_m - \mathbf{h}_m\|_{\mathbf{V}_m} &\leq \underbrace{\|\mathbf{I} \mathbf{U}_m^T \mathbf{U}_m^{-1} \mathbf{U}_m^T \mathbf{g}_m\|_{\mathbf{V}_m}}_{(a)} + \underbrace{\|\mathbf{U}_m^T \mathbf{r}_m\|_{\mathbf{V}_m}}_{(b)} + \underbrace{\|\mathbf{U}_m^T \mathbf{r}_m\|_{\mathbf{V}_m}}_{(c)} + \underbrace{\|\mathbf{U}_m^T \mathbf{r}_m\|_{\mathbf{V}_m}}_{(d)}; \end{aligned}$$

where we simply exploited the identity  $\|V_m^{-1}x\|_2^2 = x^T V_m^{-1} V_m^{-1} x = x^T V_m^{-1} x$ . We now bound one term at a time. Let us start with (a):

$$\begin{aligned}
 (a)^2 &= \|h\|_{V_m^{-1}}^2 = h^T V_m^{-1} h \\
 &\leq \|V_m^{-1}\|_2 \|h\|_2^2 \\
 &\leq \|h\|_2^2 \\
 &\leq \frac{B \rho A q^2}{1 - \rho A q};
 \end{aligned}$$

where we observed that  $\|V_m^{-1}\|_2 = \|V_m^{-1}\|_2 \leq 1$ . Finally, we have bounded the norm of

$$\begin{aligned}
 \|h\|_2 &\leq \sum_{s=0}^8 h^{tsu} \\
 &\leq \sum_{s=0}^8 h^{tsu} \\
 &\leq \|h\|_2 \sum_{s=1}^8 \|B\|_2^s \|A\|_2^{s-1} \\
 &\leq \frac{B \rho A q}{1 - \rho A q},
 \end{aligned}$$

where we have exploited Assumptions 2.1 and 2.2.

We now move to term (b):

$$\begin{aligned}
 (b)^2 &= \sum_{m=1}^M g_m^T U_m V_m^{-1} U_m^T g_m \\
 &\leq \sum_{m=2}^M g_m^T U_m^2 \\
 &\leq \sum_{l=1}^m \|x_{l1}\|_2 h^{J_{H_l}, 2;8 M} \|y_{l1}\|_2^2 \\
 &\leq \sum_{l=1}^m \|u_{l1}\|_2^2 h^{J_{H_l}, 2;8 M} \\
 &\leq \frac{U^4 B^2 \rho A q^2}{\rho(1 - \rho A q)} \sum_{l=1}^m \rho A q^{H_l - 1};
 \end{aligned}$$

where we have employed the following inequality:

$$\begin{aligned}
 h^{J_{H_l}, 2;8 M} &\leq \sum_{j=H_l-2}^8 \|A^{j-1} B\|_2 \\
 &\leq \|B\|_2 \sum_{j=H_l-2}^8 \|A^{j-1}\|_2 \\
 &\leq B \rho A q \frac{\rho A q^{H_l-1}}{1 - \rho A q};
 \end{aligned}$$

Let us now consider term (c):

$$\begin{aligned}
 (c)^2 &= \mathbb{U}_m^T r_m^2 \mathbb{V}_m^{-1} r_m^T \mathbb{U}_m \mathbb{V}_m^{-1} \mathbb{U}_m^T r_m \\
 &\leq \frac{1}{2} \mathbb{U}_m^T r_m^2 \\
 &= \frac{1}{2} \sum_{s=1}^m \sum_{l=1}^L \mathbb{1}^T A^{H_l-1} x_{l-1} a_l \\
 &\leq \frac{1}{2} \sum_{s=1}^m \sum_{l=2}^L \mathbb{1}^T A^{H_l-1} x_{l-1} a_l \\
 &\leq \frac{\lambda^2 \mathbb{U}^2 \rho A^2}{2} \sum_{l=1}^m \rho A^{H_l-1} :
 \end{aligned}$$

We now bound the summations, exploiting the inequality  $\sum_{l=1}^m \rho A^{H_l-1} \leq \frac{Y}{\log \frac{1}{\rho A}}$ , holding by assumption:

$$\begin{aligned}
 \sum_{l=1}^m \rho A^{H_l-1} &\leq \sum_{l=1}^m \rho A^{\frac{\log l}{\log \frac{1}{\rho A}}} \\
 &\leq \sum_{l=1}^m \rho A^{\frac{\log l}{\log \frac{1}{\rho A}}} \\
 &\leq \sum_{l=1}^m \exp \left( \frac{\log \frac{1}{\rho A} \log l}{\log \frac{1}{\rho A}} \right) \\
 &\leq \frac{1}{\log \frac{1}{\rho A}} \leq \log m \leq \log t \leq \log p \leq \log p \leq \log p
 \end{aligned}$$

having exploited the fact that  $\log t \leq t$  and the bound with the integral to the harmonic sum.

Finally, we consider term (d). In this case, we apply Theorem 1 of (Abbasi-Yadkori et al., 2011), observing that the conditions are satisfied. To this end, we first need to determine the subgaussianity constant for the noise process  $\epsilon_t$  every  $l \in [1, L]$  and  $P \in \mathcal{R}$ , and properly using the tower law of expectation:

$$\begin{aligned}
 \mathbb{E} \exp \left( \sum_{t=1}^T \epsilon_t \right) &= \mathbb{E} \exp \left( \sum_{s=1}^{H_m-1} \sum_{t=s}^T \epsilon_t \mid \mathcal{F}_{t-1} \right) \\
 &= \mathbb{E} \exp \left( \sum_{s=1}^{H_m-1} \sum_{t=s}^T \epsilon_t \mid \mathcal{F}_{t-1} \right) \\
 &\leq \exp \left( \frac{\lambda^2}{2} \sum_{s=1}^{H_m-1} \mathbb{E} \exp \left( \frac{\sum_{t=s}^T \epsilon_t^2}{2} \right) \mid \mathcal{F}_{t-1} \right) \\
 &\leq \exp \left( \frac{\lambda^2}{2} \sum_{s=1}^{H_m-1} \exp \left( \frac{\lambda^2 \rho A^2 \rho A^{2ps-1} q^2}{2} \right) \right) \\
 &\leq \exp \left( \frac{\lambda^2}{2} \sum_{s=1}^{H_m-1} \rho A^2 \rho A^{2ps-1} q^2 \right) \\
 &\leq \exp \left( \frac{\lambda^2}{2} \sum_{s=1}^{H_m-1} \frac{\rho A^2}{\rho A^2 q} \right) :
 \end{aligned}$$

Thus, simultaneously for all  $P, J, M, K$  with probability at least  $1 - \epsilon$ , it holds that:

$$(d)^2 \sum_{m=1}^M \frac{\mathbf{r}_m^T \mathbf{r}_m}{\mathbf{v}_m^T \mathbf{v}_m} \leq 2 \sum_{m=1}^M \frac{1}{\rho} \frac{\rho^2 \mathbf{A} \mathbf{q}^2}{\rho \mathbf{A} \mathbf{q}^2} \log \frac{1}{1 - \frac{1}{2}} \log \frac{\det \mathbf{V}_m}{d} : \quad (9)$$

□

We now proceed at bounding the of ine regret  $R^{\text{off}}$  and, then, relating the of ine regret  $R^{\text{off}}$  with the online regret  $R$ , as defined in the main paper.

Theorem A.1. Under Assumptions 2.1 and 2.2, having selected  $\rho$  in Equation (5), for every  $P \in [0, 1]$ , with probability at least  $1 - \epsilon$  DynLin-UCB suffers an of ine regret  $R^{\text{off}}$  bounded as:

$$R^{\text{off}}_{\text{DynLin-UCB}} ; T, q \leq e^{8dT} \frac{1}{T} \log \frac{T}{\log \frac{1}{1 - \frac{1}{2}}} \log \frac{1}{1 - \frac{1}{2}} \frac{T U^2}{d} : \quad (10)$$

Moreover, the expected of ine regret  $R^{\text{off}}$  is bounded as:

$$\mathbb{E} R^{\text{off}}_{\text{DynLin-UCB}} ; T, q \leq O \left( \frac{d^2 T}{\rho^{1 - \beta/2}} \right) \quad (11)$$

Proof. For every epoch  $m \in P, J, M, K$  let us define  $\mathbf{r}_{m-1} = \mathbf{r}_{t_{m-1}}$  and define the confidence set  $\mathcal{C}_{m-1} = \{ \mathbf{h} \in \mathbb{R}^d : \|\mathbf{h} - \mathbf{h}_{m-1}\|_{\mathbf{V}_{m-1}} \leq \mathbf{r}_{m-1} \}$ . Let us start by considering the instantaneous of ine regret  $r_m$  at epoch  $m \in P, J, M, K$ . Let  $\mathbf{u} = \text{arg max}_{\mathbf{u} \in \mathcal{U}} \mathbf{x}^T \mathbf{h}; \mathbf{u}$  and let  $\mathbf{h}_{m-1}^{\circ} \in \mathcal{C}_{m-1}$  such that  $\mathbf{u}^T \mathbf{h}_{m-1}^{\circ} \geq \mathbf{u}^T \mathbf{h}_{m-1}$ . Thus, with probability at least  $1 - \epsilon$ , we have:

$$r_m = \mathbf{J} \mathbf{u}^T \mathbf{h}_{m-1}^{\circ} - \mathbf{J} \mathbf{u}^T \mathbf{h}_{m-1} \leq \mathbf{J} \mathbf{u}^T (\mathbf{h}_{m-1}^{\circ} - \mathbf{h}_{m-1}) \leq \mathbf{J} \mathbf{u}^T \mathbf{r}_{m-1} \mathbf{u} \quad (9)$$

$$\leq \mathbf{J} \mathbf{u}^T \mathbf{r}_{m-1} \mathbf{u} \leq \mathbf{J} \mathbf{u}^T \mathbf{r}_{m-1} \mathbf{u} \leq \mathbf{J} \mathbf{u}^T \mathbf{r}_{m-1} \mathbf{u} \quad (10)$$

$$\leq 2 \mathbf{r}_{m-1} \mathbf{u}^T \mathbf{u} \leq 2 \mathbf{r}_{m-1} \mathbf{u}^T \mathbf{u} : \quad (11)$$

where line (9) follows from the optimism, line (10) derives from triangle inequality, line (11) is obtained by observing that  $\mathbf{r}_{m-1} \mathbf{u}^T \mathbf{u} \leq \mathbf{r}_{m-1}$  with probability at least  $1 - \epsilon$ , simultaneously for all  $m \in P, J, M, K$  thanks to Theorem 3.1, having observed that  $\mathbf{r}_{m-1}$  is larger than the right hand side of Theorem 3.1.

We now move to the cumulative of ine regret over the whole horizon  $T$ , by decomposing w.r.t. the epochs and recalling that we pay the same instantaneous regret within each epoch:

$$R^{\text{off}}_{\text{DynLin-UCB}} ; T, q \leq \sum_{m=1}^M \mathbf{p} \mathbf{H}_m \mathbf{1} \mathbf{q}^2 \leq e^{8dT} \sum_{m=1}^M \mathbf{p} \mathbf{H}_m \mathbf{1} \mathbf{q}^2 \leq e^{8dT} \sum_{m=1}^M \mathbf{r}_m^2 : \quad (12)$$

Concerning the first summation, we proceed as follows, recalling that  $T$  and  $\mathbf{H}_m = \mathbf{H}_M$  for all  $m \in P, J, M, K$

$$\sum_{m=1}^M \mathbf{p} \mathbf{H}_m \mathbf{1} \mathbf{q}^2 \leq T \mathbf{p} \mathbf{H}_M \mathbf{1} \mathbf{q}^2 \leq T \mathbf{1} \frac{\log T}{\log \frac{1}{1 - \frac{1}{2}}} : \quad (13)$$



For the second summation, we following the usual derivation for linear bandits, recalling that  $r_m \leq 1$  for all  $m \in [M]$ .

$$\begin{aligned} \sum_{m=1}^M r_m^2 &\leq 4 \sum_{m=1}^M \frac{1}{d} \sum_{k=1}^d \mathbb{1}\{a_m(k) \neq \arg \max_{k \in [d]} \mu_k\} \\ &\leq 8d \sum_{m=1}^M \log \frac{1}{\frac{MU^2}{d}} \leq 8d \frac{2}{T} \log \frac{1}{\frac{TU^2}{d}}; \end{aligned}$$

where the last passage follows from the elliptic potential lemma (Lattimore & Szepesvári, 2020, Lemma 19.4). Putting all together, we obtain:

$$R^{\text{off}}(\text{pDynLin-UCB}; T) \leq \frac{g}{\epsilon} \left( 8dT \frac{2}{T} \log \frac{1}{\frac{TU^2}{d}} + \frac{\log T}{\log \frac{1}{\epsilon}} \log \frac{1}{\frac{TU^2}{d}} \right);$$

We can also arrive to a problem-dependent regret bound, by setting  $\sup_{u \in \mathcal{U}} \sum_{x \in \mathcal{X}} u(x) \leq \frac{1}{2}$ . Since the instantaneous regret is either 0 or at least  $\frac{1}{2}$ , we have:

$$\begin{aligned} R^{\text{off}}(\text{pDynLin-UCB}; T) &\leq \sum_{m=1}^M \left( \frac{1}{2} H_m + 4 \sum_{k=1}^d \mathbb{1}\{a_m(k) \neq \arg \max_{k \in [d]} \mu_k\} \right) \\ &\leq \frac{H_M}{2} + 8d \sum_{m=1}^M \log \frac{1}{\frac{MU^2}{d}} \\ &\leq \frac{8d}{2} \log \frac{1}{\frac{MU^2}{d}} + \frac{2}{T} \log \frac{1}{\frac{TU^2}{d}}; \end{aligned}$$

By setting  $\epsilon = \frac{1}{2} \sqrt{\frac{1}{T}}$ , we obtain the of ine regret in expectation, highlighting the dependence on  $\epsilon, \sigma, \eta$  and  $d$  only:

$$E R^{\text{off}}(\text{pDynLin-UCB}; T) \leq \mathcal{O} \left( \frac{d \sqrt{T}}{\epsilon^2} \right);$$

where we used the fact that  $\frac{1}{\log \frac{1}{\epsilon}} \leq \frac{1}{1-\epsilon}$  and  $\frac{1}{\log \frac{1}{\epsilon}} \leq \frac{1}{\epsilon}$ . □

The following lemma relates the expected of ine regret with the expected online regret.

Lemma A.2. Let  $T \in \mathbb{N}$  be the optimization horizon. Then, under Assumptions 2.1, 2.2, 3.1, for Algorithm 1, it holds that:

$$E R(\text{pDynLin-UCB}; T) \leq \mathcal{O} \left( \frac{1}{\epsilon} \sum_{k=1}^d \left( \frac{1}{\sigma_k} + \frac{1}{\eta_k} \right) \right) E R^{\text{off}}(\text{pDynLin-UCB}; T);$$

Proof. First of all, we observe that for any policy, the cumulative effect of the noise components is zero-mean. Thus, it suffices to consider the deterministic evolution of the system. Let  $\{a_m\}_{m=1}^M$  be the sequence of actions played by Algorithm 1 over the  $M$  epochs. We consider one epoch  $J \in [M]$  at a time, starting from the second one. Let us consider the state that is reached at the end of epoch  $J$ , i.e.,  $x_{t_m}$  that can be expressed as:

$$x_{t_m} = A^{H_m-1} x_{t_m-2} + \sum_{l=1}^{H_m-1} A^{l-1} B a_{m-l};$$

We compare this with the steady state that would have been reached by applying the input  $a_m$ , i.e.,  $\bar{x}_{t_m} = (I - A)^{-1} B a_m$ :

$$\begin{aligned} x_{t_m} - \bar{x}_{t_m} &= A^{H_m-1} x_{t_m-2} + \sum_{l=1}^{H_m-1} A^{l-1} B a_{m-l} - \sum_{l=1}^{H_m-1} A^{l-1} B a_m \\ &= \sum_{l=1}^{H_m-1} \left( A^{l-1} B a_{m-l} - A^{l-1} B a_m \right) \\ &= \sum_{l=1}^{H_m-1} \left( \frac{A^l - I}{A - I} B a_{m-l} - \frac{A^l - I}{A - I} B a_m \right) \\ &= \sum_{l=1}^{H_m-1} \frac{A^l - I}{A - I} B (a_{m-l} - a_m); \end{aligned} \tag{12}$$

where we used the fact that  $\mathbb{P} \{ \|\bar{x}_{t_m} - p\| \leq \frac{1}{\sqrt{m}} \} \geq 1 - \frac{1}{m}$ . We now consider the evolution of the state of the system, starting from  $\bar{x}_{t_m}$  and applying action  $a_m$ . To this end, we consider the steady state  $p = (A - \alpha I)^{-1} \alpha B a_m$  that we would obtain by constantly applying action  $a_m$ . Let  $\phi = (A - \alpha I)^{-1} \alpha B a_m$ , we denote with  $\bar{x}_{t_m}$  the state reached after  $s$  rounds starting from state  $\bar{x}_{t_m}$ :

$$\begin{aligned} \bar{x}_{t_m} &= (A - \alpha I)^{-s} \bar{x}_{t_m} + \sum_{s=1}^s (A - \alpha I)^{s-1} \alpha B a_m \\ &= (A - \alpha I)^{-s} \bar{x}_{t_m} + p - (A - \alpha I)^{-s} p + \sum_{s=1}^s (A - \alpha I)^{s-1} \alpha B a_m \\ &= (A - \alpha I)^{-s} \bar{x}_{t_m} + p - (A - \alpha I)^{-s} p + \alpha B a_m \sum_{s=1}^s (A - \alpha I)^{s-1}; \end{aligned}$$

where we exploited Lemma A.3. Now, we sum over  $s$  limited to epoch  $m$ :

$$\sum_{s=1}^{H_m} \bar{x}_{t_m} = \sum_{s=1}^{H_m} (A - \alpha I)^{-s} \bar{x}_{t_m} + \sum_{s=1}^{H_m} p - \sum_{s=1}^{H_m} (A - \alpha I)^{-s} p + \sum_{s=1}^{H_m} \alpha B a_m \sum_{s=1}^s (A - \alpha I)^{s-1}; \quad (13)$$

Thus, passing to the reward and averaging over the rounds of epoch  $m$  introduce the following quantity, representing the average reward in epoch  $m$ :

$$W_{p,a_m} = \frac{1}{H_m} \sum_{s=1}^{H_m} \bar{r}_{t_m} \cdot x_{t_m} \cdot a_m \cdot y_{t_m} = \frac{1}{H_m} \sum_{s=1}^{H_m} x_{t_m}^T \bar{r}_{t_m} \cdot y_{t_m};$$

Let us also consider the sequence of states generated starting from  $\bar{x}_{t_m}$  and applying action  $a_m$ :

$$x_{t_m} = (A - \alpha I)^{-s} \bar{x}_{t_m} + \sum_{s=1}^s (A - \alpha I)^{s-1} \alpha B a_m;$$

and the corresponding outputs, averaged over the rounds of epoch  $m$ :

$$W_{p,a_m} = \frac{1}{H_m} \sum_{s=1}^{H_m} y_{t_m} \cdot x_{t_m} \cdot a_m \cdot y_{t_m} = \frac{1}{H_m} \sum_{s=1}^{H_m} x_{t_m}^T \cdot y_{t_m};$$

Putting all together, and limiting to epoch  $m$ , we obtain:

$$J_{p,a_m} = W_{p,a_m} - \underbrace{\frac{1}{H_m} \sum_{s=1}^{H_m} \bar{r}_{t_m} \cdot x_{t_m} \cdot a_m \cdot y_{t_m}}_{(a)} + \underbrace{\frac{1}{H_m} \sum_{s=1}^{H_m} x_{t_m}^T \cdot y_{t_m}}_{(b)};$$

We now consider one term at a time and start with term (b):

$$\begin{aligned} (b) \quad W_{p,a_m} &= \frac{1}{H_m} \sum_{s=1}^{H_m} x_{t_m}^T \cdot y_{t_m} = \frac{1}{H_m} \sum_{s=1}^{H_m} \left[ (A - \alpha I)^{-s} \bar{x}_{t_m} + \sum_{s=1}^s (A - \alpha I)^{s-1} \alpha B a_m \right]^T \cdot y_{t_m} \\ &= \frac{1}{H_m} \sum_{s=1}^{H_m} \bar{x}_{t_m}^T (A - \alpha I)^{-s} y_{t_m} + \sum_{s=1}^{H_m} \alpha B a_m^T \sum_{s=1}^s (A - \alpha I)^{s-1} y_{t_m} \\ &\leq \frac{1}{H_m} \sum_{s=1}^{H_m} \bar{x}_{t_m}^T (A - \alpha I)^{-s} y_{t_m} + \frac{\alpha B a_m^T y_{t_m}}{1 - \rho(A - \alpha I)} \\ &\leq \frac{\alpha B a_m^T y_{t_m}}{1 - \rho(A - \alpha I)} + \frac{1}{H_m} \sum_{s=1}^{H_m} \bar{x}_{t_m}^T (A - \alpha I)^{-s} y_{t_m}; \end{aligned}$$

having bounded  $\|A^T \bar{x}_{t_m-1} - x_{t_m-1}\| \leq \rho \|A\| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2$  and, further, by setting  $\rho = 1$  and exploiting Equation (12) for bounding the norm. We now move to term (a), and exploit Equation (13):

$$(a) \quad \mathbb{E} \sum_{m=1}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \frac{1}{H_m - 1} \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \\ + \frac{1}{H_m - 1} \sum_{l=1}^{H_m-1} \|A^T p_l - A^T q^T p_l - A^{H_m-1} q\| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2$$

Let us rename  $R = A^T p_l - A^T q^T p_l - A^{H_m-1} q$  and recall that the involved quantity is a scalar:

$$\|R\|_F \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 = \text{tr} \left( R^T (\bar{x}_{t_m-1} - x_{t_m-1}) (\bar{x}_{t_m-1} - x_{t_m-1})^T \right) \\ = \sum_{i,j} R_{ij} (\bar{x}_{t_m-1} - x_{t_m-1})_i (\bar{x}_{t_m-1} - x_{t_m-1})_j$$

having exploited the Cauchy-Schwarz inequality over the inner product space of matrices  $\mathbb{R}^{d \times d}$ . For the first norm, we obtain:

$$\|R\|_F \leq \|A\|_F \rho + \|A^{H_m-1} q\|_2 \|A\|_F \leq \rho \|H_m - 1\| \rho \|A\|_F + \|A\|_F \|A\|_2^{H_m-1} \|q\|_2$$

using the inequality between matrix norms  $\|X\|_F \leq \|X\|_2 \|Y\|_2$  and bounding  $\|A\|_2^{H_m-1} \leq \rho \|A\|_2^{H_m-1} \leq \rho \|H_m - 1\| \rho \|A\|_2$ . For the second term, recalling that  $\bar{x}_{t_m-1} - x_{t_m-1}$  has rank 1, we obtain:

$$\|R\|_F \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \leq \sum_{i,j} |R_{ij}| (\bar{x}_{t_m-1} - x_{t_m-1})_i (\bar{x}_{t_m-1} - x_{t_m-1})_j \\ \leq \sum_{i,j} |R_{ij}| (\bar{x}_{t_m-1} - x_{t_m-1})_i^2 + \sum_{i,j} |R_{ij}| (\bar{x}_{t_m-1} - x_{t_m-1})_j^2 \\ \leq \max_{i,j} |R_{ij}| \sum_{i,j} (\bar{x}_{t_m-1} - x_{t_m-1})_i^2 + \sum_{i,j} (\bar{x}_{t_m-1} - x_{t_m-1})_j^2 \\ \leq 2 \rho \|A\|_F \|q\|_2$$

where we used Assumption 3.1 for the first inequality. Putting all together, we have:

$$\mathbb{E} \text{Reg}_{\text{DynLin-UCB}}(T) \leq \sum_{t=1}^T \sum_{m=2}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \rho \|A\|_F \|q\|_2 \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \tag{14}$$

$$\leq \rho \|A\|_F \|q\|_2 \sum_{m=2}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \tag{15}$$

$$\leq \rho \|A\|_F \|q\|_2 \sum_{m=2}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \\ \leq 2 \rho \|A\|_F \|q\|_2 \sum_{m=2}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \\ \leq 2 \rho \|A\|_F \|q\|_2 \sum_{m=2}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \\ \leq 2 \rho \|A\|_F \|q\|_2 \sum_{m=2}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \tag{16}$$

$$\leq \rho \|A\|_F \|q\|_2 \sum_{m=2}^M \sum_{q=1}^Q \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{x' \in \mathcal{X}} \sum_{y' \in \mathcal{Y}} \sum_{l=1}^{H_m-1} |x_l - y_l| \|\bar{x}_{t_m-1} - x_{t_m-1}\|_2 \tag{17}$$

where in line (15) we exploited the inequality  $\|y\|_2 \leq \|J\|_2 \|U\|_2 \|X\|_2$ , we used the bounds on (a) and (b) in (16), we used the definition of of-line regret and bounded the harmonic summation in line (17).  $\square$

Theorem 3.2. Under Assumptions 2.1, 2.2, and 3.1, Algorithm 1 suffers an expected online regret bounded by:

$$ER_{\text{DynLin-UCB}} \leq \sum_{t=1}^T \sum_{i=1}^m \left( \frac{1}{\alpha} + \frac{\beta}{\alpha^2} \right) \frac{1}{\alpha^{2t}}$$

Proof. We simply combine Theorem A.1 and Lemma A.2. □

### A.3. Technical Lemmas

Lemma A.3. Let  $A \in \mathbb{R}^{n \times n}$ . Then, for every  $m \in \mathbb{N}$ , it holds that:

$$\sum_{i=0}^{m-1} A^i \alpha^i = \alpha^{-1} (I - A^m) (I - A)^{-1}$$

Proof. Let us simply write explicitly the expressions:

$$\sum_{i=0}^{m-1} A^i \alpha^i = \sum_{i=0}^{m-1} \sum_{j=0}^{m-1} A^i \alpha^j \mathbb{1}_{i=j} = \sum_{j=0}^{m-1} \sum_{i=0}^{m-1} A^i \alpha^j \mathbb{1}_{i=j} = \sum_{j=0}^{m-1} \alpha^j \sum_{i=0}^{m-1} A^i \mathbb{1}_{i=j}$$

An analogous derivation holds for the second expression. □

## B. Finite-Horizon Setting

In this section, we compare the finite-horizon setting with the infinite-horizon one presented in the main paper. We shall show that under Assumption 2.2, the two settings tend to coincide when the horizon is sufficiently large. Let us start by introducing the  $H$ -horizon expected average reward with  $H \in \mathbb{N}$  being the optimization horizon:

$$J_H(p, q) = \mathbb{E} \left[ \frac{1}{H} \sum_{t=1}^H y_t \right] \quad \text{where} \quad \begin{cases} x_{t+1} = Ax_t + Bu_t + \xi_t \\ y_t = x_t^\top \theta + \eta_t \\ u_t \in \mathcal{U} \end{cases} \quad (18)$$

where the expectation is taken w.r.t. the randomness of the state noise and reward noise. We now show that the optimal policy for the finite-horizon setting is a non-stationary open-loop policy.

Theorem B.1 (Optimal Policy for the  $H$ -Horizon Setting) If  $H \in \mathbb{N}$ , an optimal policy  $\pi_H^* = \{p_{H,t}^*, q_{H,t}^*\}_{t=1}^H$  maximizing the  $H$ -horizon expected average reward  $J_H(p, q)$  as in Equation (18) is given by:

$$\pi_H^* = \{p_{H,t}^*, q_{H,t}^*\}_{t=1}^H : \quad p_{H,t}^* = \arg \max_{p \in \mathcal{P}} x_t^\top p, \quad q_{H,t}^* = \arg \max_{q \in \mathcal{Q}} q^\top x_t, \quad \text{where} \quad u_{H,t} = \arg \max_{u \in \mathcal{U}} x_t^\top A u + q^\top (Ax_t + Bu_t + \xi_t)$$

Proof. We start by expressing for every  $t \in \{1, \dots, H\}$  the reward  $y_t$  as a function of the sequence of actions  $u_1, \dots, u_H$  produced by a generic policy. By exploiting Equation (3) instantiated with  $t = 1$ , we have:

$$y_t = \sum_{s=0}^{t-1} x_t^\top A^s u_s + \xi_t = \sum_{s=0}^{t-1} (A^s x_1)^\top u_s + \xi_t = \sum_{s=0}^{t-1} (A^s x_1)^\top u_s + \xi_t$$

By computing the expectation, using linearity, and recalling that the noises are zero-mean, we obtain:

$$\mathbb{E} y_t = \sum_{s=0}^{t-1} x_t^\top A^s \mathbb{E} u_s = \sum_{s=0}^{t-1} (A^s x_1)^\top \mathbb{E} u_s$$



By averaging over  $\mathcal{P}_{J,H,K}$ , we obtain the  $H$ -horizon expected average reward:

$$J_{H,p,q} = \frac{1}{H} \sum_{t=1}^H \mathbb{E} r_{t,s} \tag{19}$$

$$= \frac{1}{H} \sum_{t=1}^H \sum_{s=0}^{t-1} x^{t,s;u_s} \mathbb{E} r_{t,s} \mathbb{E} \left[ \frac{1}{H} \sum_{t=1}^H \mathbb{1}^T A^{t-1} \text{Er}_{1,s} \right] \tag{20}$$

where line (19) is obtained by renaming the indexes of the summations, and (20) comes from the definition of cumulative Markov parameter  $x^{J_0;H^s K}$ . It is now simple to see, as no noise is present in the expression, that the performance  $J_{H,p,q}$  is maximized by taking at each round  $t \in \mathcal{P}$  an action  $u_s = p_{s-1,q}$  such that whose expectation satisfies  $\mathbb{E} r_{t,s} = \arg \max_{u_s} x^{J_0;H^s K} \mathbb{E} r_{t,s}$ . Clearly, we can take the deterministic action such that  $\mathbb{E} r_{t,s}$ .

We now show that for sufficiently large  $H$ , the  $H$ -horizon expected average reward  $J_{H,p,q}$  tends to coincide with the infinite-horizon expected average reward.

Proposition B.2. Let  $H \in \mathcal{P}$ . Then, for every policy  $\pi$  it holds that:

$$|J_{H,p,q} - J_{p,q}| \leq \frac{BU - pAq}{H} + \frac{pAq^H}{pAq}.$$

Proof. Consider two horizons  $H, H^1 \in \mathcal{P}$ , and let  $u_1; u_2; \dots; q$  be the sequence of actions played by policy  $\pi$ . Using Equation (20), we have:

$$J_{H,p,q} - J_{H^1,p,q} = \frac{1}{H} \sum_{s=1}^H x^{J_0;H^s K} \mathbb{E} r_{s,s} - \frac{1}{H^1} \sum_{s=1}^{H^1} x^{J_0;H^1 s K} \mathbb{E} r_{s,s} \tag{21}$$

$$= \frac{1}{H} \sum_{s=1}^H x^{J_0;H^s K} \mathbb{E} r_{s,s} - \frac{1}{H^1} \sum_{s=1}^{H^1} x^{J_0;H^1 s K} \mathbb{E} r_{s,s} \tag{22}$$

$$= \frac{1}{H} \sum_{s=1}^H x^{J^H s 1;8^M} \mathbb{E} r_{s,s} - \frac{1}{H^1} \sum_{s=1}^{H^1} x^{J^{H^1} s 1;8^M} \mathbb{E} r_{s,s} \tag{23}$$

As shown in Appendix A.1, we have that the second addendum vanishes as  $H$  approaches  $\infty$ :

$$\frac{1}{H^1} \sum_{s=1}^{H^1} x^{J^{H^1} s 1;8^M} \mathbb{E} r_{s,s} \xrightarrow{H^1 \rightarrow \infty} 0 \quad \text{when} \quad H^1 \xrightarrow{H^1 \rightarrow \infty} \infty$$

Concerning the first addendum, we have:

$$\begin{aligned} \frac{1}{H} \sum_{s=1}^H x^{J^H s 1;8^M} \mathbb{E} r_{s,s} &\leq \frac{U}{H} \sum_{s=1}^H h^{J^H s 1;8^M} \\ &\leq \frac{BU - pAq}{H} + \frac{pAq^H}{pAq} \\ &= \frac{BU - pAq}{H} + \frac{pAq^H}{pAq}. \end{aligned}$$

□

### C. System Identification

This section presents a solution to identify matrices  $A$ ,  $B$ ,  $C$ , and  $D$  characterizing a Linear Dynamical System starting from a single trajectory. We adopt a variant of the Ho-Kalman (Ho & Kálmán, 1966) algorithm. We start from the identification method proposed by Lale et al. (Lale et al., 2020a, Section 3), where authors consider a system of the type:

$$\begin{aligned} x_{t+1} &= Ax_t + Bu_t + \xi_t; \\ y_t &= Cx_t + z_t; \end{aligned} \quad (24)$$

Our setting can be seen as:

$$\begin{aligned} x_{t+1} &= Ax_t + Bu_t + \xi_t; \\ y_t &= Cx_t + Du_t + z_t; \end{aligned} \quad (25)$$

with  $x_t; \xi_t \in \mathbb{R}^n$ ,  $u_t \in \mathbb{R}^p$ , and  $y_t; z_t \in \mathbb{R}^m$ . The noise over state transition model and output  $\xi_t$  and  $z_t$  are  $\sigma^2$ -subgaussian random variables.

We consider in this part the standard control problem notation adopted for Linear Dynamical Systems. The mapping to our problem is straightforward by considering  $G = v_m^{-1} \Gamma$  and  $D = v_p^{-1} (v_m^{-1} \Gamma + v_p^{-1} \Gamma_p)$ .

In predictive form, the system described in Equation (24) is:

$$\begin{aligned} p_{t+1} &= \bar{A} p_t + Bu_t + Fy_t; \\ y_t &= Cp_t + e_t; \end{aligned}$$

where:

$$\begin{aligned} \bar{A} &= A - FC; \\ F &= (A - C^T \rho C C^T)^{-1} C^T \Gamma^{-1}; \end{aligned}$$

and  $\rho$  is the solution to the following DARE (Discrete Algebraic Riccati Equation):

$$A(A - C^T \rho C C^T)^{-1} A^T - A C^T \rho C C^T (A - C^T \rho C C^T)^{-1} A^T - \rho^{-1} I = 0;$$

In order to identify this Linear Dynamical System, we want to detect a matrix

$$G_y = \begin{bmatrix} CF + C^T AF & \dots & CA^H \Gamma^{-1} F & CB + C^T AB & \dots & CA^H \Gamma^{-1} B \end{bmatrix} \quad (26)$$

To identify through least squares method matrix  $G_y$ , we construct for each  $t$  a vector  $r_t$ :

$$r_t = \begin{bmatrix} y_{t-1}^T & \dots & y_{t-H}^T & u_{t-1}^T & \dots & u_{t-H}^T \end{bmatrix}^T \in \mathbb{R}^{pm + pqH}; \quad (27)$$

The system output  $y_t$  can be rewritten as:

$$y_t = G_y r_t + e_t + CA^H x_{t-H};$$

The output of the system under analysis (Equation 25) is:

$$\begin{aligned} y_t &= Du_t + \\ &+ G_y r_t + e_t + CA^H x_{t-H} \end{aligned}$$

We can incorporate the contribution  $Du_t$  in  $G_y$  obtaining  $\tilde{G}_y$ :

$$\tilde{G}_y = \begin{bmatrix} CF + C^T AF & \dots & CA^H \Gamma^{-1} F + D & CB + C^T AB & \dots & CA^H \Gamma^{-1} B \end{bmatrix};$$

The related vector  $r_t$  is:

$$r_t = \begin{bmatrix} y_{t-1}^T & \dots & y_{t-H}^T & u_{t-1}^T & u_{t-1}^T & \dots & u_{t-H}^T \end{bmatrix}^T \in \mathbb{R}^{pm + pqH + p}; \quad (28)$$

The best value of  $\Theta_y$  can be found through regularized least squares as in (Lale et al., 2020a, Equation 10):

$$\Theta_y = \arg \min_X \|X\|_F^2 + \lambda \|y - X\|_2^2; \quad (29)$$

where  $\| \cdot \|_F$  represents the Frobenius norm.

The matrix  $D$  can be directly retrieved from  $\Theta_y$ . In order to get matrices  $A$ ,  $B$ , and  $C$ , we remove the values related to from  $\Theta_y$  and we retrieve  $\Theta_y$ . From now on we can refer to the algorithm proposed in the work of Lale et al. (Lale et al., 2020a, Appendix B).

## E. Experimental Results

In this section, we present two more contributions of the experimental evaluation provided in the main paper. First, we construct a simulator by generalizing real-world data using the variant of the Ho-Kalman method presented in Appendix C, and we test DynLin-UCB compared with Lin-UCB, Exp3 and the Expert Policy (i.e., the policy applied by human experts). Second, we study the behavior of DynLin-UCB, Lin-UCB and Exp3 for different magnitude of noises in both the state transition model and the output for the setting presented in the main paper.

DynLin-UCB is considered in two version, depending on the choice of the hyper-parameter  $\lambda$ . We consider the case with  $\lambda = 1$  (referred as DynLin-UCB (1)), and  $\lambda = \log T$  (referred as DynLin-UCB ( $\log T$ )), where  $T$  is the total number of steps performed in the experiment.

**Baselines** In this section, the quality of DynLin-UCB is compared with several baselines. First, we compare our solution with Lin-UCB algorithm, which is a particular case of our solution, that, given the nature of the problem (output is a linear combination of the actions plus stochastic noise), is a straightforward baseline. As already done in Section 5, we consider two version of Lin-UCB, one with hyper-parameter  $\lambda = 1$  (referred as Lin-UCB (1)), and the other with  $\lambda = \log T$  (referred as Lin-UCB ( $\log T$ )). Second, we compare our algorithm with Exp3 (Auer et al., 2002), an algorithm designed for adversarial bandits<sup>12</sup>. Lastly, only in the case of real-world data, we compare our solution with the human-expert policy (referred as Expert Policy) which does not consider the interactions and delays explicitly. This policy is static, directly generalized from the original dataset by learning the average budget allocation over all platforms from the available data. If such an allocation does not belong to the action space, a projection into  $\mathcal{J}$  is performed.

### E.1. Real-world Data

In this section, we present an experimental evaluation based on real-world data coming from three of the most important advertising platforms of the web (Facebook, Google, and Bing), related to a large number of campaigns for a value of more than 6 Million USD over 2 years. Starting from such data, we generalized the best model by means of a specially designed variant of the Ho-Kalman algorithm (Ho & Kálmán, 1986). We used the matrices estimated with Ho-Kalman to build up a simulator and compare the performance of DynLin-UCB and the baselines in a controlled environment. The resulting system has  $\rho = 0.67$ . To assert the quality of the solution, we evaluate DynLin-UCB in comparison with the baselines presented above for  $10^6$  steps over 10 runs.

**Results** Figure 3 gives an overview of the results of this simulation. More in details, Figure 3a shows the performance of each algorithm in terms of cumulative regret, while Figure 3b shows the instantaneous reward. It is worth noting how neither Lin-UCB nor Exp3 are able to converge to the optimal choice. Indeed, they immediately converge to a sub-optimal solution and persist in it. DynLin-UCB, instead, shows a convergence trend towards the optimal policy over time in both the setting ( $\lambda = 1$  and  $\lambda = \log T$ ), even if the best solution is the one which employs  $\log T$ . The Expert Policy which tends to consider instantaneous effect only and does not take into account correlations between platforms is also sub-optimal.

<sup>12</sup>It is worth noting that in the adversarial setting the regret is defined by comparing the learner's behavior with the best fixed action. However, the regret guarantees of Exp3 are fully meaningful for the non-adaptive adversaries only to which our setting cannot be reduced to, since the hidden state evolution (and so the reward) depends on the actual played actions.

<sup>13</sup>The complete method used to retrieve the best model generalizing data is described in Appendix C.

(a) Cumulative Regret

(b) Immediate Reward

Figure 3. Performance of DynLin-UCB, Lin-UCB and Exp3, and the Expert Policy in the system generalized from real-world data. (10 runs, mean std)

### E.2. Noise Effect

In this experiment, we evaluate the performance of DynLin-UCB and the other bandit baselines described above at different values of noise. The evaluation is performed over 10 runs considering the setting presented in Section 5 and lasts for  $T = 10^6$  time steps. The analysis considers noisy state transition model and output subject to different magnitude of noise. The noise in this simulation is a zero-mean Gaussian noise with variance  $\sigma^2 \in \{10^{-5}; 10^{-4}; 10^{-3}; 10^{-2}; 10^{-1}; 1\}$ .

**Results** Figure 4 shows the results of the experiment for the different values of  $\sigma^2$ . It is clearly visible how DynLin-UCB performs in almost the same way no matter the noise at which the system is subject to, leading always to sub-linear regret. On the other hand, Lin-UCB regret is different in every simulation we perform. Indeed, with the a low level of noise reaches linear regret and does not converge (showing linear regret), while for large values of noise it converges very quickly. This is due to the nature of the confidence bound of Linear bandits that is not able to take into account such a complex scenario and lead to no guarantees in this setting. Exp3 is not able to reach the optimum in this scenario, independently from the value of the noise<sup>2</sup>.

**Infrastructure and Computational Time** The code used for the results provided in this section has been run on a Intel(R) I5(R) 8259U @ 2.30GHz CPU with 8 GB of LPDDR3 system memory. The operating system was macOS, and the experiments have been run on Python 3.7. The experiment performed over real-world data takes over 30 minutes to to run all the algorithms and perform 10 runs over  $10^6$  rounds each. The experiment to evaluate the performance of DynLin-UCB and the other bandit baselines takes over 6 hours to perform the 6 experiments of 10 runs each over  $10^6$  rounds. It is worth noting that the time complexity of DynLin-UCB is upper-bounded by the one of Lin-UCB.

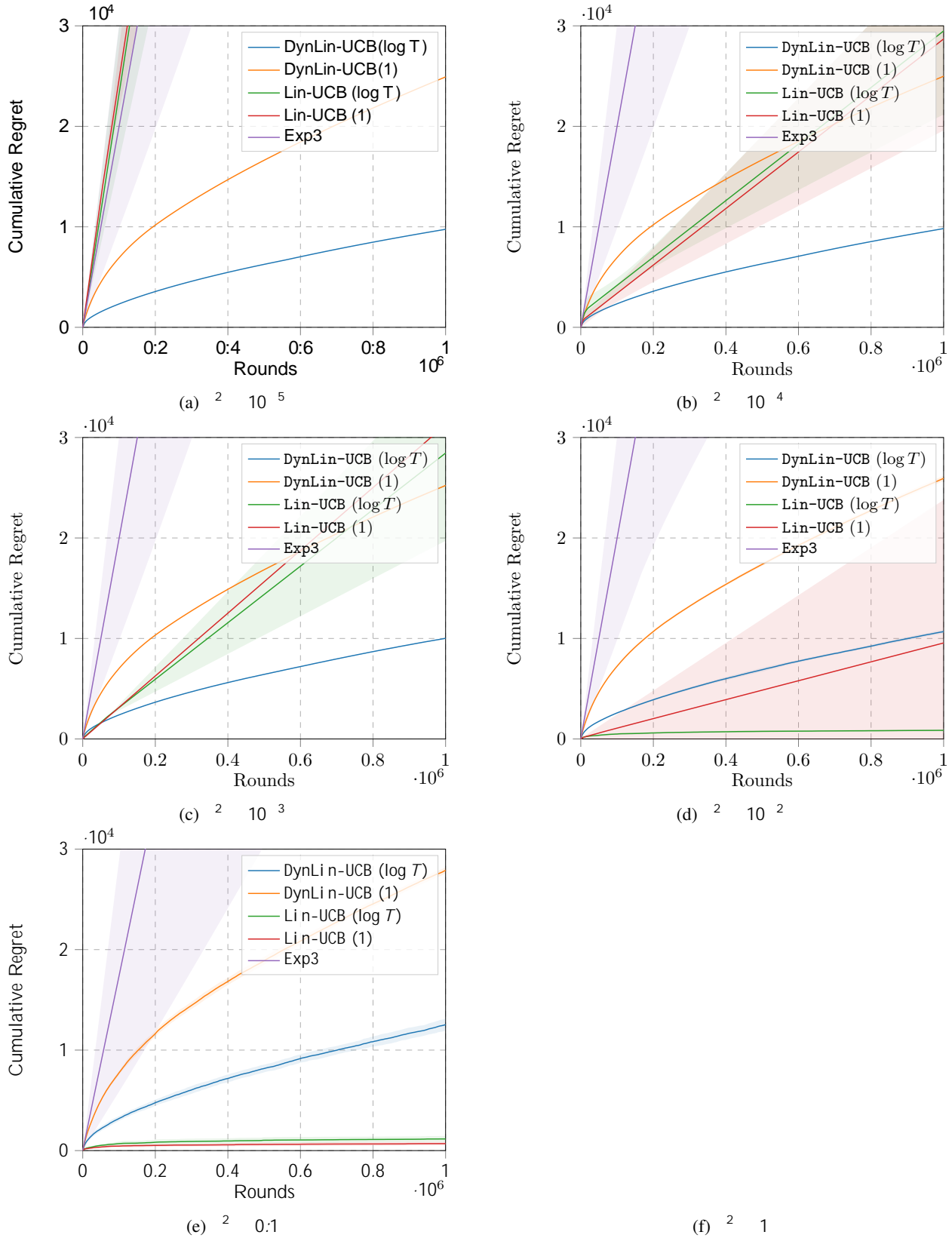


Figure 4. Performance of DynLin-UCB, Lin-UCB and Exp3 at different values of  $\beta$ . (10 runs, mean  $\pm$  std)