

---

# Dynamical Linear Bandits

---

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 does not immediately reflect on the feedback and spreads its effects over a long time frame. For instance, in online advertising, investing in a platform produces an instantaneous increase of *awareness*, but the actual reward, i.e., a *conversion*, might occur far in the future. Furthermore, whether a conversion takes place depends on: how fast the awareness grows, its vanishing effects, and the 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 in the future, disregarding possible dynamical effects. In this paper, we introduce a novel setting, the Dynamical Linear Bandits (DLB), an extension of the 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 of the action. Then, the hidden state evolves according to linear dynamics, affected by the performed action too. We start by introducing the setting, discussing the notion of optimal policy, and deriving an expected regret lower bound. Then, we provide an optimistic regret minimization algorithm, Dynamical Linear Upper Confidence Bound (DynLin-UCB), that suffers an expected regret of order  $\tilde{\mathcal{O}}\left(\frac{d\sqrt{T}}{(1-\bar{\rho})^{3/2}}\right)$ , where  $\bar{\rho}$  is a measure of the stability of the system, and  $d$  is the dimension of the action vector. Finally, we conduct a numerical validation on a synthetic environment and on real-world data to show the effectiveness of DynLin-UCB in comparison with several baselines.

## 1. Introduction

In a large variety of sequential decision-making problems, a learner must choose 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 observes a reward (i.e., feedback) representing the combined effect of multiple actions played in the past. For instance, in online advertising campaigns, the process that leads to a *conversion*, i.e., *marketing funnel* (Court et al., 2009), is characterized by complex dynamics and comprises several phases. When heterogeneous campaigns/platforms are involved, a profitable budget investment policy has to account for the interplay between campaigns/platforms. In this scenario, a conversion (e.g., a user’s purchase of a promoted product) should be attributed not only to the latest ad the user was exposed to, but also to previous ones (Berman, 2018).

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 brand awareness, while the latter aims at creating the actual conversion. If we evaluate the performances in terms of conversions only, we will discover that impression ads are not instantaneously effective in creating conversions, so we will be tempted to reduce the budget invested in such a campaign. However, this approach is sub-optimal because impression ads increase the chance to convert when a conversion ad is shown after the impression (e.g., Hoban & Bucklin, 2015). 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 consequences that non-dynamical models cannot capture. This kind of interplay comprises more general scenarios than the simple reward delay, including the case where the interaction is governed by a dynamics *hidden* to the observer.

While this scenario can be indubitably modeled as a Partially Observable Markov Decision Process (POMDP, Åström, 1965), the complexity of the framework and its general-

---

<sup>1</sup>Politecnico di Milano, Milan, Italy.

Correspondence to: Marco Mussi <marco.mussi@polimi.it>.

Proceedings of the 40<sup>th</sup> International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

ity 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 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 to 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 the 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 a hidden internal state  $\mathbf{x}_t$  (e.g., awareness) which evolves via linear dynamics fed by the action  $\mathbf{u}_t$  (e.g., amount invested) and 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, reduces to linear bandits (Abbasi-Yadkori et al., 2011) when no dynamics are involved. Because of the dynamics, the effect of each action persists over time indefinitely but, under stability conditions, it vanishes asymptotically. This allows representing interference and synergy between platforms, thanks to the dynamic nature of the system.

**Contributions** 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. We show that, under stability conditions, the optimal policy corresponds to playing the *constant action* that leads the system to the most profitable steady state. Then, we derive an expected regret lower bound of order  $\left(\frac{d\sqrt{T}}{(1-\bar{\rho})^{1/2}}\right)$ , being  $d$  the dimensionality of the action space and  $\bar{\rho} < 1$  the spectral radius of the dynamical matrix of the system evolution law.<sup>1</sup> In Section 3, we propose a novel optimistic regret minimization algorithm, *Dynamical Linear Upper Confidence Bound* (DynLin-UCB), for the DLB setting. DynLin-UCB takes inspiration from Lin-UCB but subdivides the optimization horizon  $T$  into increasing-length epochs. In each epoch, an action is selected optimistically and kept constant (i.e., persisted) so that the system approximately reaches the steady state. We provide a regret analysis for DynLin-UCB showing that, under certain assumptions,

<sup>1</sup>The smaller  $\bar{\rho}$ , the faster the system reaches its steady state.

it enjoys  $\tilde{O}\left(\frac{d\sqrt{T}}{(1-\bar{\rho})^{3/2}}\right)$  expected regret. In Section 5, we provide a numerical validation, with both synthetic and real-world data, compared with bandit baselines. The proofs of all the results are reported in Appendix B.

**Notation** Let  $a, b \in \mathbb{N}$  with  $a \leq b$ , we introduce the symbols:  $\mathcal{J}_{a,b} : \{a, \dots, b\}$ ,  $\mathcal{I}_{1,b} : \{1, \dots, b\}$ , and  $\mathcal{J}_a, \mathcal{I}_a : \{a, \dots, a\}$ . Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , we denote with  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i$  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  $\mathbf{A}^\top \mathbf{A}$ . We introduce the maximum spectral norm to spectral radius ratio of the powers of  $\mathbf{A}$  defined as  $\|\mathbf{A}\|_{\infty} = \sup_{\tau \geq 0} \|\mathbf{A}^\tau\|_2 / \rho(\mathbf{A})^\tau$  (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  $\zeta \in \mathbb{R}^n$  it holds that  $\mathbb{E} \exp\{\langle \mathbf{x}, \zeta \rangle\} \leq \exp\{\frac{1}{2} \zeta^\top \zeta\}$ .

## 2. Setting

In this section, we introduce the *Dynamical Linear Bandits* (DLBs), the learner-environment interaction, assumptions, and regret (Section 2.1). Then, we derive a closed-form expression for the optimal policy for DLBs (Section 2.2). Finally, we derive a lower bound to the regret, highlighting the intrinsic complexities of the DLB setting (Section 2.3).

### 2.1. Problem Formulation

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  $y_t = \langle \mathbf{x}_t, \boldsymbol{\omega} \rangle + \langle \mathbf{u}_t, \boldsymbol{\theta} \rangle + \eta_t \in \mathcal{Y}$ , where  $\mathcal{Y} \subseteq \mathbb{R}$  is the reward space,  $\boldsymbol{\omega} \in \mathbb{R}^n$ ,  $\boldsymbol{\theta} \in \mathbb{R}^d$  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  $\mathbf{x}_{t+1} = \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t + \boldsymbol{\epsilon}_t$ , 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  $\boldsymbol{\epsilon}_t$  is a zero-mean  $\sigma^2$ -subgaussian random noise, conditioned to the past, independent of  $\eta_t$ .<sup>2</sup>

**Remark 2.1.** *The setting proposed above is a particular case of a POMDP (Åström, 1965), in which the state  $\mathbf{x}_t$  is non-observable, while the learner accesses the noisy*

<sup>2</sup> $n$  is the order of the LTI system (Kalman, 1963). We make no assumption on the value of  $n$  and on its knowledge.

observation  $y_t$  that 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 (Abbasi-Yadkori et al., 2011) when the hidden state does not affect the reward, i.e., when  $\omega = \mathbf{0}$ .

**Markov Parameters** We revise a useful representation, that for every  $H \in \mathbb{N}$  allows expressing  $y_t$  in terms of the sequence of the most recent  $H - 1$  actions  $\mathbf{u}_{t-s}$ , reward noise  $\eta_t$ ,  $H$  state noises  $\mathbf{e}_{s \in [t-H, t-1]}$ , and starting state  $\mathbf{x}_{t-H}$  (Ho & Kalman, 1966; Oymak & Ozay, 2019; Tsiamis & Pappas, 2019; Sarkar et al., 2021):

$$y_t = \underbrace{\sum_{s=0}^{H-1} \lambda \mathbf{h}^{(s)} \mathbf{u}_{t-s}}_{\text{action effect}} \underbrace{\omega^\top \mathbf{A}^H \mathbf{x}_{t-H}}_{\text{starting state}} + \underbrace{\eta_t \sum_{s=1}^H \omega^\top \mathbf{A}^{s-1} \mathbf{e}_{t-s}}_{\text{noise}}, \quad (1)$$

where the sequence of vectors  $\mathbf{h}^{(s)} \in \mathbb{R}^d$  for every  $s \in \mathbb{N}$  are called *Markov parameters* and are defined as:  $\mathbf{h}^{(0)} = \boldsymbol{\theta}$  and  $\mathbf{h}^{(s)} = \mathbf{B}^\top \rho \mathbf{A}^{s-1} \mathbf{q}^\top \omega$  if  $s \neq 1$ . Furthermore, we introduce the *cumulative Markov parameters*, defined for every  $s, s' \in \mathbb{N}$  with  $s \leq s'$  as  $\mathbf{h}^{[s, s']} = \sum_{l=s}^{s'} \mathbf{h}^{(l)}$  and the corresponding limit as  $s' \rightarrow \infty$ , i.e.,  $\mathbf{h}^{[s, +\infty]} = \sum_{l=s}^{+\infty} \mathbf{h}^{(l)}$ . Finally, we use the abbreviation  $\mathbf{h} = \mathbf{h}^{[0, +\infty]} = \boldsymbol{\theta} \mathbf{B}^\top \rho \mathbf{I}_n \mathbf{A} \mathbf{q}^{-\top} \omega$ .

We will make use of the following standard assumption related to the *stability* of the dynamic matrix  $\mathbf{A}$ , widely employed in discrete-time LTI literature (Oymak & Ozay, 2019; Lale et al., 2020a,b).

**Assumption 2.1 (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.,  $\rho(\mathbf{A}^k) \leq \mathcal{S}^k$ .*

**Policies and Performance** The learner’s behavior is modeled via a deterministic *policy*  $\underline{\pi} = \{\pi_t\}_{t \in \mathbb{N}}$  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  $\mathcal{H}_{t-1} \subseteq \mathcal{U}^t$  is the set of histories of length  $t - 1$ . The performance of a policy  $\underline{\pi}$  is evaluated in terms of the (*infinite-horizon*) *expected average reward*:

$$J(\underline{\pi}) : \liminf_{H \rightarrow +\infty} \mathbb{E} \left[ \frac{1}{H} \sum_{t=1}^H y_t \right], \quad (2)$$

$$\text{where } \begin{cases} \mathbf{x}_{t+1} = \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t + \mathbf{e}_t \\ y_t = \lambda \omega^\top \mathbf{x}_t + \mathbf{h}^\top \mathbf{u}_t + \eta_t, \quad @t \in \mathbb{N}, \\ \mathbf{u}_t = \pi_t(H_{t-1}) \end{cases}$$

where the expectation is taken w.r.t. the randomness of the state noise  $\mathbf{e}_t$  and reward noise  $\eta_t$ . If a policy  $\underline{\pi}$  is *constant*, i.e.,  $\pi_t(H_{t-1}) = \mathbf{u}^*$  for every  $t \in \mathbb{N}$ , we abbreviate  $J(\mathbf{u}^*)$

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

$J(\mathbf{u}^*)$ . A policy  $\underline{\pi}^*$  is an *optimal policy* if it maximizes the expected average reward, i.e.,  $\underline{\pi}^* = \arg \max_{\underline{\pi}} J(\underline{\pi})$ , and its performance is denoted by  $J^* = J(\underline{\pi}^*)$ .

We further introduce the following assumption that requires the boundedness of the norms of the relevant quantities.

**Assumption 2.2 (Boundedness).** *There exist  $\lambda, \rho, B, U \in \mathbb{R}$  s.t.:  $\|\boldsymbol{\theta}\|_2 \leq \lambda$ ,  $\|\omega\|_2 \leq \rho$ ,  $\|\mathbf{B}\|_2 \leq B$ ,  $\sup_{\mathbf{u} \in \mathcal{U}} \|\mathbf{u}\|_2 \leq U$ , and  $\sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_2 \leq X$ ,  $\sup_{\mathbf{u} \in \mathcal{U}} |J(\mathbf{u})| \leq 1$ .*<sup>4</sup>

**Regret** The *regret* suffered by playing a policy  $\underline{\pi}$ , competing against the optimal infinite-horizon policy  $\underline{\pi}^*$  over a *learning horizon*  $T \in \mathbb{N}$  is given by:

$$R(\underline{\pi}, T) = T J^* - \sum_{t=1}^T y_t, \quad (3)$$

where  $y_t$  is the sequence of rewards collected by playing  $\underline{\pi}$  as in Equation (2). The goal of the learner consists in minimizing the *expected regret*  $\mathbb{E} R(\underline{\pi}, T)$ , where the expectation is taken w.r.t. the randomness of the reward.

## 2.2. Optimal Policy

In this section, we derive a closed-form expression for the optimal policy  $\underline{\pi}^*$  for the infinite-horizon objective function, as introduced in Equation (2).

**Theorem 2.1 (Optimal Policy).** *Under Assumptions 2.1 and 2.2, an optimal policy  $\underline{\pi}^*$  maximizing the (infinite-horizon) expected average reward  $J(\underline{\pi}^*)$  (Equation 2), for every round  $t \in \mathbb{N}$  and history  $H_{t-1} \in \mathcal{H}_{t-1}$  is given by:*

$$\pi_t^*(H_{t-1}) = \mathbf{u}^* \text{ where } \mathbf{u}^* = \arg \max_{\mathbf{u} \in \mathcal{U}} \lambda \mathbf{h}^\top \mathbf{u}. \quad (4)$$

Some remarks are in order. The optimal policy plays the *constant* action  $\mathbf{u}^* \in \mathcal{U}$  which brings the system in the “most profitable” steady-state.<sup>5</sup> Indeed, the expression  $\lambda \mathbf{h}^\top \mathbf{u}$  can be rewritten expanding the cumulative Markov parameter as  $\rho \boldsymbol{\theta}^\top \omega^\top \rho \mathbf{I}_n \mathbf{A} \mathbf{q}^{-1} \mathbf{B} \mathbf{u}^*$  and  $\bar{\mathbf{x}}^* = \rho \mathbf{I}_n \mathbf{A} \mathbf{q}^{-1} \mathbf{B} \mathbf{u}^*$  is the expression of the steady state  $\bar{\mathbf{x}}^* = \mathbf{A} \bar{\mathbf{x}}^* + \mathbf{B} \mathbf{u}^*$ , when applying action  $\mathbf{u}^*$ . It is worth noting the role of Assumption 2.1 which guarantees the existence of the inverse  $\rho \mathbf{I}_n \mathbf{A} \mathbf{q}^{-1}$ . In this sense, our problem shares the constant nature of the optimal policy with the linear bandit setting (Abbasi-Yadkori et al., 2011), although ours 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 restrict to *open-loop stationary* policies. The reason why DLBs do not benefit from *closed-loop* policies, differently from other classical problems, such as

<sup>4</sup>The assumption of the bounded state norm  $\|\mathbf{x}\|_2 \leq X$  holds whenever the state noise  $\mathbf{e}$  is bounded. As shown by Agarwal et al. (2019), this assumption can be relaxed, for unbounded subgaussian noise, 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>5</sup>In Appendix C, we show that the optimal policy is non-stationary for the finite-horizon case.

the LQG (Abbasi-Yadkori & Szepesvári, 2011), lies in the linearity of the reward  $y_t$  and in the additive noise  $\eta_t$  and  $\epsilon_t$ , making their presence irrelevant (in expectation) 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  $\mathbf{u} \in \mathcal{U}$ , instantaneous rewards are not reliable, and we need to let the system evolve to the steady state and, only then, observe the reward.

### 2.3. Regret Lower Bound

In this section, we provide a lower bound to the expected regret that any learning algorithm suffers when addressing the learning problem in a DLB.

**Theorem 2.2 (Lower Bound).** *For any policy  $\pi$  (even stochastic), there exists a DLB fulfilling Assumptions 2.1 and 2.2, such that for sufficiently large  $T \asymp \mathcal{O}\left(\frac{d^2}{1-\rho(\mathbf{A})}\right)$ , policy  $\pi$  suffers an expected regret lower bounded by:*

$$\mathbb{E}R_{\pi}(T) \asymp \left( \frac{d^2 T}{\rho(\mathbf{A})} \right)^{\frac{1}{2}}.$$

The lower bound highlights the main challenges of the DLB learning problem. First of all, we observe a dependence on  $\frac{1}{\rho(\mathbf{A})}$ , being  $\rho(\mathbf{A})$  the spectral radius of the matrix  $\mathbf{A}$ . This is in line with the intuition that, as  $\rho(\mathbf{A})$  approaches 1, the problem becomes more challenging. Furthermore, we note that when  $\rho(\mathbf{A}) = 0$ , i.e., the problem has no dynamical effects, the lower bound matches the one of linear bandits (Lattimore & Szepesvári, 2020). It is worth noting that, for technical reasons, the result of Theorem 2.2 is derived under the assumption that, at every round  $t \in [T]$ , the agent observes *both* the state  $\mathbf{x}_t$  and the reward  $y_t$  (see Appendix B). Clearly, this represents a simpler setting w.r.t. DLBs (in which  $\mathbf{x}_t$  is hidden) and, consequently, Theorem 2.2 is a viable lower bound for DLBs too.

## 3. Algorithm

In this section, we present an *optimistic* regret minimization algorithm for the DLB setting. *Dynamical Linear Upper Confidence Bound* (DynLin-UCB), whose pseudocode is reported in Algorithm 1, requires the knowledge of an upper-bound  $\bar{\rho} \geq 1$  on the spectral radius of the dynamic matrix  $\mathbf{A}$  (i.e.,  $\rho(\mathbf{A}) \leq \bar{\rho}$ ) and on the maximum spectral norm to spectral radius ratio  $\bar{\beta} \geq 1$  (i.e.,  $\rho(\mathbf{A}) \leq \bar{\beta}$ ), as well as the bounds on the relevant quantities of Assumption 2.2.<sup>6</sup>

<sup>6</sup>As an alternative, one can consider a more demanding requirement of the knowledge of a bound on the spectral norm  $\|\mathbf{A}\|_2$  of  $\mathbf{A}$ . Similar assumptions regarding the knowledge of analogous quantities are considered in the literature, e.g., *decay of Markov operator norms* (Simchowitz et al., 2020) and *strong stability* (Plevrakis & Hazan, 2020), spectral norm bound (Lale et al., 2020a). As a side note, the knowledge of  $\bar{\rho} \asymp \rho(\mathbf{A})$  (or an equivalent quantity) is proved to be unavoidable by Theorem 2.2. Indeed, if no restriction

DynLin-UCB is based on the following simple observation. To assess the quality of action  $\mathbf{u} \in \mathcal{U}$ , we need to *persist* in applying it so that the system approximately reaches the corresponding steady state and, then, observe the reward  $y_t$ , representing a reliable estimate of  $\mathbf{u}^\top \mathbf{h}$ . We shall show that, under Assumption 2.1, 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  $\bar{\rho}$ . After initializing the Gram matrix  $\mathbf{V}_0 = \lambda \mathbf{I}_d$  and the vectors  $\mathbf{b}_0$  and  $\hat{\mathbf{h}}_0$  both to  $\mathbf{0}_d$  (line 1), DynLin-UCB subdivides the learning horizon  $T$  in  $M \asymp T$  epochs. Each epoch  $m \in [M]$  is composed of  $H_m \geq 1$  rounds, where  $H_m \asymp \lceil \log_{\rho} \frac{1}{\beta} \rceil$  is logarithmic in the epoch index  $m$ . At the beginning of each epoch,  $m \in [M]$ , DynLin-UCB computes the upper confidence bound (UCB) index (line 4) defined for every  $\mathbf{u} \in \mathcal{U}$  as:

$$\text{UCB}_t(\mathbf{u}) := \mathbf{x}_{t-1}^\top \hat{\mathbf{h}}_{t-1} + \beta_{t-1} \|\mathbf{u}\|_{\mathbf{V}_{t-1}^{-1}}, \quad (5)$$

where  $\hat{\mathbf{h}}_{t-1} = \mathbf{V}_{t-1}^{-1} \mathbf{b}_{t-1}$  is the Ridge regression estimator of the cumulative Markov parameter  $\mathbf{h}$ , as in Equation (4) and  $\beta_{t-1} \geq 0$  is an exploration coefficient to be defined later. Similar to Lin-UCB (Abbasi-Yadkori et al., 2011), the index  $\text{UCB}_t(\mathbf{u})$  is designed to be optimistic, i.e.,  $\text{UCB}_t(\mathbf{u}) \geq \mathbf{u}^\top \mathbf{h}$  in high-probability for all  $\mathbf{u} \in \mathcal{U}$ . Then, the optimistic action  $\mathbf{u}_t \in \mathcal{U}$  is selected such that  $\text{UCB}_t(\mathbf{u}_t) = \max_{\mathbf{u} \in \mathcal{U}} \text{UCB}_t(\mathbf{u})$  (line 6) and persisted for the next  $H_m$  rounds (lines 8-11). The length of the epoch  $H_m$  is selected such that, under Assumption 2.1, the system has approximately reached the steady state after  $H_m \geq 1$  rounds. In this way, at the end of epoch  $m$ , the reward  $y_t$  is an almost-unbiased sample of the steady-state performance  $\mathbf{u}_t^\top \mathbf{h}$ . This sample is employed to update the Gram matrix estimate  $\mathbf{V}_t$  and the vector  $\mathbf{b}_t$  (line 13), while the samples collected in the previous  $H_m$  rounds are discarded (line 9). It is worth noting that by setting  $H_m = 0$  for all  $m \in [M]$ , DynLin-UCB reduces to Lin-UCB. The following sections provide the concentration of the estimator  $\hat{\mathbf{h}}_{t-1}$  of  $\mathbf{h}$  (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 estimate  $\hat{\mathbf{h}}_t$  of the cumulative Markov parameter  $\mathbf{h}$ . For every epoch  $m \in [M]$ , we denote with  $t_m$  the last round of epoch  $m$ :  $t_0 = 0$  and  $t_m = t_{m-1} + H_m$ . At the end of each epoch  $m$ , we solve the Ridge regression problem, defined for every round  $t \in [T]$  as:

$$\hat{\mathbf{h}}_t = \arg \min_{\mathbf{h} \in \mathbb{R}^d} \sum_{l \in [M]: t_l \leq t} \rho y_{t_l} \|\mathbf{x}_{t_l} - \mathbf{u}_{t_l}\|_{\mathbf{V}_t}^2 + \lambda \|\hat{\mathbf{h}}_t\|_2^2 + \mathbf{b}_t^\top \mathbf{V}_t^{-1} \mathbf{b}_t.$$

on  $\rho(\mathbf{A})$  is enforced (i.e., just  $\rho(\mathbf{A}) \leq 1$ ), one can always consider the DLB in which  $\rho(\mathbf{A}) \leq 1$  making the regret lower bound degenerate to linear.



**Algorithm 1:** DynLin-UCB.

---

**Input :** Regularization parameter  $\lambda \succ 0$ , exploration coefficients  $\rho, \beta_t \in [0, 1]$ , spectral radius upper bound  $0 \preceq \bar{\rho} \preceq 1$

- 1 Initialize  $t \in \{1, \dots, T\}$ ,  $\mathbf{V}_0 = \lambda \mathbf{I}_d$ ,  $\mathbf{b}_0 = \mathbf{0}_d$ ,  $\hat{\mathbf{h}}_0 = \mathbf{0}_d$ ,
- 2 Define  $M = \min\{M', \lfloor \frac{T}{\log \frac{1}{\bar{\rho}}} \rfloor\}$ ,  $\sum_{m=1}^M \frac{1}{\log \frac{1}{\bar{\rho}}}$
- 3 **for**  $m \in \{1, \dots, M\}$  **do**
- 4     Compute  $\mathbf{u}_t = \arg \max_{\mathbf{u} \in \mathcal{U}} \text{UCB}_t(\rho, \mathbf{u})$
- 5     where  $\text{UCB}_t(\rho, \mathbf{u}) = \mathbf{x}_t^\top \hat{\mathbf{h}}_{t-1} + \beta_t \|\mathbf{u}\|_{\mathbf{V}_{t-1}^{-1}}$
- 6     Play arm  $\mathbf{u}_t$  and observe  $y_t$
- 7     Define  $H_m = \lfloor \frac{\log m}{\log \frac{1}{\bar{\rho}}} \rfloor$
- 8     **for**  $j \in \{1, \dots, H_m\}$  **do**
- 9         Update  $\mathbf{V}_t = \mathbf{V}_{t-1} + \mathbf{b}_t \mathbf{b}_t^\top$
- 10          $t \in \{t+1, \dots, t+H_m\}$
- 11         Play arm  $\mathbf{u}_t = \mathbf{u}_{t-1}$  and observe  $y_t$
- 12     **end**
- 13     Update  $\mathbf{V}_t = \mathbf{V}_{t-1} + \mathbf{u}_t \mathbf{u}_t^\top$ ,  $\mathbf{b}_t = \mathbf{b}_{t-1} + \mathbf{u}_t y_t$
- 14     Compute  $\hat{\mathbf{h}}_t = \mathbf{V}_t^{-1} \mathbf{b}_t$
- 15      $t \in \{t+1, \dots, t+H_m\}$
- 16 **end**

---

We now present the following self-normalized maximal concentration inequality and, then, we compare it with the existing results in the literature.

**Theorem 3.1** (Self-Normalized Concentration). *Let  $\{\hat{\mathbf{h}}_t\}_{t \in \mathbb{N}}$  be the sequence of solutions of the Ridge regression problems of Algorithm 1. Then, under Assumption 2.1 and 2.2, for every  $\lambda \succ 0$  and  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , simultaneously for all rounds  $t \in \mathbb{N}$ , it holds that:*

$$\|\hat{\mathbf{h}}_t - \mathbf{h}\|_{\mathbf{V}_t} \preceq \frac{c_1}{\lambda} \log \rho e \rho t + 1 + c_2 \sqrt{\frac{1}{\lambda}}$$

$$\sqrt{2\tilde{\sigma}^2 \left( \log \left( \frac{1}{\delta} \right) + \frac{1}{2} \log \left( \frac{\det \rho \mathbf{V}_t \mathbf{Q}}{\lambda^d} \right) \right)},$$

where  $c_1 = U \left( \frac{UB}{1-\rho(\mathbf{A})} + X \right)$ ,  $c_2 = \frac{\Omega B \Phi(\mathbf{A})}{1-\rho(\mathbf{A})}$ , and  $\tilde{\sigma}^2 = \sigma^2 \left( 1 + \frac{\Omega^2 \Phi(\mathbf{A})^2}{1-\rho(\mathbf{A})^2} \right)$ .

First, we note that when  $\mathbf{w} = \mathbf{0}_n$ , i.e., the state does not affect the reward, the bound perfectly reduces to the self-normalized concentration used in linear bandits (Abbasi-Yadkori et al., 2011, Theorem 1). In particular, we recognize the second term due to the regularization parameter  $\lambda \succ 0$  and the third one, which involves the subgaussianity parameter  $\tilde{\sigma}^2$ , related to the joint contribution of the state and reward noises. Furthermore, the first term is an additional bias that derives from the epochs of length  $H_m = 1$ . The choice of the value  $H_m$  represents one of the main technical novelties that, on the one hand, leads to a bias that conveniently grows logarithmically with  $t$  and, on the other hand, can be computed without the knowledge of  $T$ .

It is worth looking at our result from the perspective of learning the LTI system parameters. We can compare our

Theorem 3.1 with the concentration presented in (Lale et al., 2020a, Appendix C), which represents, to the best of our knowledge, the only result for the closed-loop identification of LTI systems with non-observable states. First, note that, although we focus on a MISO system ( $y_t$  is a scalar, being our reward), extending our estimator to multiple-outputs (MIMO) is straightforward. Second, the approach of (Lale et al., 2020a) employs the *predictive form* of the LTI system to cope with the correlation introduced by closed-loop control. This choice allows for convenient analysis of the estimated Markov parameters of the predictive form. However, recovering the parameters of the original system requires an application of the Ho-Kalman method (Ho & Kalman, 1966) which, unfortunately, does not preserve the concentration properties in general, but only for *persistently exciting* actions. Our method, instead, forces to play an open-loop policy within a single epoch (each with logarithmic duration), while the overall behavior is closed-loop, as the next action depends on the previous-epoch estimates. In this way, we are able to provide a concentration guarantee on the parameters of the original system without assuming additional properties on the action signal.

### 3.2. Regret Analysis

In this section, we provide the analysis of the regret of DynLin-UCB, when we select the exploration coefficient  $\beta_t$  based on the knowledge of the upper bounds  $\bar{\rho} \preceq 1$ ,  $\bar{\mathcal{S}}$ , and those specified in Assumption 2.2, defined for every round  $t \in \{1, \dots, T\}$  as:

$$\beta_t = \frac{\bar{c}_1}{\lambda} \log \rho e \rho t + 1 + \bar{c}_2 \sqrt{\frac{1}{\lambda}}$$

$$\sqrt{2\bar{\sigma}^2 \left( \log \left( \frac{1}{\delta} \right) + \frac{d}{2} \log \left( 1 + \frac{tU^2}{d\lambda} \right) \right)},$$

where  $\bar{c}_1 = U \left( \frac{UB}{1-\bar{\rho}} + X \right)$ ,  $\bar{c}_2 = \frac{\Omega B \bar{\Phi}}{1-\bar{\rho}}$ , and  $\bar{\sigma}^2 = \sigma^2 \left( 1 + \frac{\Omega^2 \bar{\Phi}^2}{1-\bar{\rho}^2} \right)$ . The following result provides the bound on the expected regret of DynLin-UCB.

**Theorem 3.2** (Upper Bound). *Under Assumptions 2.1 and 2.2, selecting  $\beta_t$  as in Equation (6) and  $\delta = 1/\{T, \text{DynLin-UCB}\}$  suffers an expected regret bounded as (highlighting the dependencies on  $T$ ,  $\bar{\rho}$ ,  $d$ , and  $\sigma$  only):*

$$\mathbb{E} \text{R} \rho \pi^{\text{DynLin-UCB}}, T \preceq \mathcal{O} \left( \frac{d \sigma^2 T \log T q^{\frac{3}{2}}}{1 - \bar{\rho}} + \frac{d T \log T q^2}{\rho \bar{\rho} q^{\frac{3}{2}}} + \frac{1}{\rho \bar{\rho} \mathbf{A} q q^2} \right).$$

*Proof Sketch.* The analysis of DynLin-UCB poses additional challenges compared to that of Lin-UCB (Abbasi-Yadkori et al., 2011) because of the dynamic effects of the hidden state. The idea behind the proof is to first derive

a bound on a different notion of regret, i.e., the *offline regret*:  $R^{\text{off}}(\underline{\rho}, T) = T J^* - \sum_{t=1}^T J(\mathbf{u}_t)$ , that compares  $J^*$  with the steady-state performance  $J(\mathbf{u}_t)$  of the action  $\mathbf{u}_t$  (Theorem B.2). This analysis of  $R^{\text{off}}(\underline{\rho}, T)$  can be comfortably carried out, by adopting a proof strategy similar to that of `Lin-UCB`. However, when applying action  $\mathbf{u}_t$ , the DLB does not immediately reach the performance  $J(\mathbf{u}_t)$  as the expected reward  $\mathbb{E}r_t$  experiences a transitional phase before converging to the steady state. Under stability (Assumption 2.1), it is possible to show that the expected offline regret and the expected regret differ by a constant:  $|R(\underline{\rho}, T) - \mathbb{E}R^{\text{off}}(\underline{\rho}, T)| \leq \mathcal{O}(1)$  (Lemma B.1).  $\square$

Some observations are in order. We first note a dependence on the term  $1/\bar{\rho}$ , which, in turn, depends on the upper bound  $\bar{\rho}$  of the spectral gap  $\rho\mathbf{A}$ . If the system does not display a dynamics, i.e., we can set  $\bar{\rho} = 0$ , we obtain a regret bound that, apart from logarithmic terms, coincides with that of `Lin-UCB`, i.e.,  $\tilde{\mathcal{O}}(d\sqrt{T})$ . Instead, for slow-converging systems, i.e.,  $\bar{\rho} \ll 1$ , the regret bound enlarges, as expected. Clearly, a value of  $\bar{\rho}$  too large compared to the optimization horizon  $T$  (e.g.,  $\bar{\rho} \sim 1/\sqrt{T}$ ) makes the regret bound degenerate to linear. This is a case in which the underlying system is so slow that the whole horizon  $T$  is insufficient to approximately reach the steady state. Third, the regret bound is the sum of three components: the first one depends on the subgaussian proxy  $\sigma$  and is due to the noisy estimation of the relevant quantities; the second one is a bias due to the epoch-based structure of `DynLin-UCB`; finally, the third one is constant (does not depend on  $T$ ) accounts for the time needed to reach the steady state.

**Remark 3.1** (Regret upper bound (Theorem 3.2) and lower bound (Theorem 2.2) Comparison). *Apart from logarithmic terms, we notice a tight dependence on  $d$  and on  $T$ . Instead, concerning the spectral properties of  $\mathbf{A}$ , in the upper bound, we experience a dependence on  $1/\bar{\rho}$  raised to a higher power (either 1 for the term multiplied by  $d$  and 3 for the term multiplied by  $\bar{d}$ ) w.r.t. the exponent appearing in the lower bound (i.e., 1/2). It is currently an open question whether the lower bound is not tight (which is obtained for a simpler setting in which the state is observable  $\mathbf{x}_t$ ) or whether more efficient algorithms for DLBs can be designed. Furthermore, Theorem 3.2 highlights the impact of the upper bound  $\bar{\rho}$  compared with the true  $\rho\mathbf{A}$ .*

## 4. Related Works

In this section, we survey and compare the literature with a particular focus on bandits with delayed, aggregated, and composite feedback (Joulani et al., 2013) and online control for Linear Time-Invariant (LTI) systems (Hespanha, 2018). Additional related works are reported in Appendix A.

**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 by Joulani et al. (2013) and widely exploited in 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., 2021), arm-dependent delays (Manegueu et al., 2020), non-stochastic delays (Ito et al., 2020; Thune et al., 2019; Jin 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  $\tilde{\mathcal{O}}(d\sqrt{T})$  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 a *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.

**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), achieving, in most of the cases  $\tilde{\mathcal{O}}(d\sqrt{T})$  regret for strongly convex functions and  $\tilde{\mathcal{O}}(dT^{2/3})$  for convex functions. Recently,  $\tilde{\mathcal{O}}(d\sqrt{T})$  regret rate has been obtained for convex function too, by means of geometric exploration methods (Plevrakis & Hazan, 2020). Compared to `DynLin-UCB`, the algorithm of Plevrakis & Hazan (2020) considers general convex costs but assumes the observability of the state and limits to the class of disturbance response controllers (Li & Bosch, 1993) that do not include the constant policy. Moreover, the regret bound of Plevrakis & Hazan (2020) differs from Theorem 3.2, as it shows a cubic dependence on the system order<sup>7</sup> and an implicit non-trivial dependence on the dynamic matrix  $\mathbf{A}$ . Instead, our Theorem 3.2 is remarkably independent of the system order  $n$ . Furthermore, Lale et al. (2020a) reach  $\mathcal{O}(\log p T)$  regret in the case of strongly convex cost functions compet-

<sup>7</sup>This holds for *known* cost functions. Instead, for *unknown* costs, the exponent becomes 24 (Plevrakis & Hazan, 2020).

ing against the best *persistently exciting* controller (i.e., a controller implicitly maintaining a non-null exploration). Some approaches are designed to deal with adversarial noise (Simchowitz et al., 2020). All of these solutions, however, look for the best closed-loop controller within a specific class, e.g., disturbance response control (Li & 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, e.g., in advertising domains. DynLin-UCB works with an arbitrary action space and, thanks to the linearity of the reward, does not require complex closed-loop controllers.

## 5. Numerical Simulations

In this section, we provide numerical validations of DynLin-UCB in both a synthetic scenario and a domain obtained from real-world data. The goal of these simulations is to highlight the behavior of DynLin-UCB in comparison with bandit baselines, describing advantages and disadvantages. The first experiment is a synthetic setting in which we can evaluate the performances of all the solutions and the sensitivity of DynLin-UCB w.r.t. the  $\bar{\rho}$  parameter (Section 5.1). Then, we show a comparison in a DLB scenario retrieved from real-world data (Section 5.2). The code of the experiments can be found at <https://github.com/marcomussi/DLB>. Details and additional experiments can be found in Appendix E.

**Baselines** We consider as main baseline Lin-UCB (Abbasi-Yadkori et al., 2011), designed for linear bandits. We include Exp3 (Auer et al., 1995) usually employed in (non-adaptive) adversarial settings, and its extension to  $k$ -length memory (adaptive) adversaries Exp3-k by Dekel et al. (2012).<sup>8</sup> Additionally, we perform a comparison with algorithms for regret minimization in non-stationary environments: D-Lin-UCB (Russac et al., 2019), an extension of Lin-UCB for non-stationary settings, and AR2 (Chen et al., 2021), a bandit algorithm for processes presenting temporal structure. Lastly, in the case of real-world data, we compare our solution with a human-expert policy (Expert). This policy is directly generalized from the original dataset by learning via regression the average budget allocation over all platforms from the available data.

For the baselines which do not support vectorial actions, we perform a discretization of the action space  $\mathcal{U}$  that surely contains optimal action. Concerning the hyperparameters of the baselines, whenever possible, they are selected as in the respective original papers. The experiments are presented with a regularization parameter  $\lambda \in \{1, \log T\}$  for the algorithms which require it (i.e., DynLin-UCB, Lin-UCB, and

<sup>8</sup> $k$  is proportional to  $\lceil \log M \rceil \lceil \log \frac{1}{\bar{\rho} \eta} \rceil$ . In Appendix A.3 we elaborate on the use of adversarial bandit algorithms for DLBs.

D-Lin-UCB).<sup>9</sup> Further information about the hyperparameters of the baselines and the adopted optimistic exploration bounds are presented in Appendix E.1.

### 5.1. Synthetic Data

**Setting** We consider a DLB defined by the following matrices  $\mathbf{A} = \text{diag}(0.2, 0, 0.1) \mathbf{q} \mathbf{q}^\top$ ,  $\mathbf{B} = \text{diag}(0.25, 0, 0.1) \mathbf{q} \mathbf{q}^\top$ ,  $\boldsymbol{\theta} = [0, 0.5, 0.1]^\top$ ,  $\boldsymbol{\omega} = [1, 0, 0.1]^\top$  and a Gaussian noise with  $\sigma = 0.01$  (diagonal covariance matrix for the state noise).<sup>10</sup> This way, the spectral gap of the dynamical matrix is  $\rho \mathbf{p} \mathbf{A} \mathbf{q} = 0.2$  and  $\mathbf{p} \mathbf{A} \mathbf{q} = 1$ . Moreover, the cumulative Markov parameter is given by  $\mathbf{h} = [0.56, 0.5, 0.11]^\top$ . We consider the action space  $\mathcal{U} = \{ \mathbf{u} = [u_1, u_2, u_3]^\top \mid \mathbf{u} \succeq \mathbf{0}, 1s^\top \mathbf{u} = 1.5 \}$  with  $u_1 = u_2 = u_3 \propto 1.5u$  that simulates a total budget of 1.5 to be allocated to the three platforms. Thus, a “myopic” agent would simply look at how the action immediately propagates to the reward through  $\boldsymbol{\theta}$ , and will invest the budget in the second component of the action, which is weighted by 0.5. Instead, a “far-sighted” agent, aware of the system evolution, will look at the cumulative Markov parameter  $\mathbf{h}$ , realizing that the most convenient action is investing in the first component, weighted by 0.56. Therefore, the optimal action is  $\mathbf{u}^* = [1, 0.5, 0]^\top$  leading to  $J^* = 0.81$ .

**Comparison with the bandit baselines** Figure 1 shows the performance in terms of cumulative regret of DynLin-UCB, Lin-UCB, D-Lin-UCB, AR2, Exp3, and Exp3-k. The experiments are conducted over a time horizon of 1 million rounds. For DynLin-UCB, we employed, for the sake of this experiment, the true value of the spectral gap, i.e.,  $\bar{\rho} = \rho \mathbf{p} \mathbf{A} \mathbf{q} = 0.2$ . First of all, we observe that both Exp3 and Exp3-k suffers a significantly large cumulative regret. Similar behavior is displayed by AR2. Moreover, all the versions of Lin-UCB and D-Lin-UCB suffer linear regret. The best performance of D-Lin-UCB is obtained when the discount factor  $\gamma$  is close to 1 (the weights take the form  $w_t = \gamma^{-t}$ ), and the behavior is comparable with the one of Lin-UCB. Even for a quite fast system ( $\rho \mathbf{p} \mathbf{A} \mathbf{q} = 0.2$ ), ignoring the system dynamics, and the presence of the hidden state, has made both Lin-UCB and D-Lin-UCB commit (in their best version, with  $\lambda = \log T$ ) to the sub-optimal (myopic) action  $\mathbf{u}^\circ = [0.5, 1, 0]^\top$  with performance  $J^\circ = 0.78 < J^*$ , with also a relevant variance. On the other hand, DynLin-UCB is able to maintain

<sup>9</sup>For DynLin-UCB,  $\log T$  is a nearly optimal choice for  $\lambda$  as it can be seen by looking at the first two addenda of the exploration factor in Equation (6).

<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 in the cumulative Markov parameter  $\mathbf{h}$  and we could have obtained the same results with an equivalent (non-diagonal) representation, by applying an inevitable transformation  $\mathbf{T}$  as  $\mathbf{A}^1 = \mathbf{T} \mathbf{A} \mathbf{T}^{-1}$ ,  $\boldsymbol{\omega}^1 = \mathbf{T}^\top \boldsymbol{\omega}$ , and  $\mathbf{B}^1 = \mathbf{T} \mathbf{B}$ .

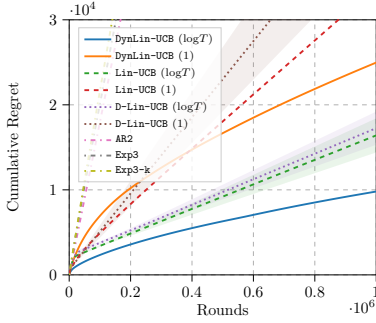


Figure 1. Cumulative regret as a function of the rounds comparing DynLin-UCB and the other bandit baselines (50 runs, mean std).

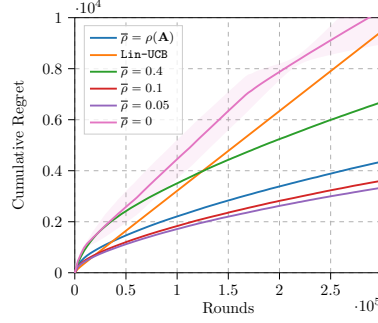


Figure 2. Cumulative regret as a function of the rounds comparing Lin-UCB, and DynLin-UCB with  $\lambda = \log T$ , varying the upper bound on the spectral radius  $\bar{\rho}$  (50 runs, mean std).

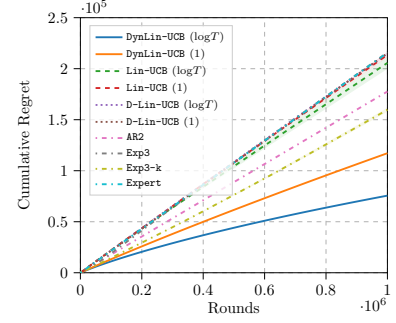


Figure 3. Cumulative regret for DynLin-UCB, the other bandit baselines and the Expert in the system generalized from real-world data (50 runs, mean std).

a smaller and stable (variance is negligible) sublinear regret in both its versions, with a notable advantage when using  $\lambda = \log T$ .

**Sensitivity to the Choice of  $\bar{\rho}$**  The upper bound  $\bar{\rho}$  of the spectral radius  $\rho p A q = 0.2$  represents a crucial parameter of DynLin-UCB. While an overestimation  $\bar{\rho} \gg \rho p A q$  does not compromise the regret rate but tends to slow down the convergence process, a severe underestimation  $\bar{\rho} \ll \rho p A q$  might prevent learning at all. In Figure 2, we test DynLin-UCB against a misspecification of  $\bar{\rho}$ , when  $\lambda = \log T$ . We can see that by considering  $\bar{\rho} = 2\rho p A q$ , DynLin-UCB experiences a larger regret but still sublinear and smaller w.r.t. Lin-UCB with  $\lambda = \log T$ . Even by reducing  $\bar{\rho}$  to 0.1, 0.05, DynLin-UCB is able to keep the regret sublinear, showing remarkable robustness to misspecification. Clearly, setting  $\bar{\rho} = 0$  makes the regret almost degenerate to linear.

## 5.2. Real-world Data

We present an experimental evaluation based on real-world data coming from three web advertising platforms (Facebook, Google, and Bing), related to several campaigns for an invested budget of 5 Million EUR over 2 years. Starting from such data, we learn the best DLB model by means of a specifically designed variant of the Ho-Kalman algorithm (Ho & Kalman, 1966).<sup>11</sup> We used the learned model to build up a simulator. The resulting system has  $\rho p A q = 0.67$ . We evaluate DynLin-UCB against the baselines for  $T = 10^6$  steps over 50 runs.

**Results** Figure 3 shows the results in terms of cumulative regret. It is worth noting that no algorithm, except for DynLin-UCB, is able to converge to the optimal choice. Indeed, they immediately commit to a sub-optimal solution.

<sup>11</sup>See Appendix D.

DynLin-UCB, instead, shows a convergence trend towards the optimal policy over time for both  $\lambda = 1$  and  $\lambda = \log T$ , even if the best-performing version is the one which employs  $\lambda = \log T$ . The Expert, which has a preference towards maximizing the instantaneous effect of the actions only and does not take into account correlations between platforms, displays a sub-optimal performance.

## 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 linear dynamics and by an observable noisy reward that linearly combines the hidden state and the action played. This model accounts for scenarios that cannot be easily represented by existing bandit models that consider delayed and aggregated feedback. We have derived a regret lower bound that highlights the main complexities of the DLB problem. Then, we have proposed a novel optimistic regret minimization approach, DynLin-UCB, that, under stability assumption, is able to achieve sub-linear regret. The numerical simulation in both synthetic and real-world domains succeeded in showing that, in a setting where the baselines mostly suffer linear regret, our algorithm consistently enjoys sublinear regret. Furthermore, DynLin-UCB proved to be robust to misspecification of its most relevant hyper-parameter  $\bar{\rho}$ . To the best of our knowledge, this is the first work addressing this family of problems, characterized by hidden linear dynamics, with a simple, yet effective, bandit-like approach. Short-term future directions include efforts in closing the gap between the regret lower and upper bounds. Long-term future directions should focus on extending the present approach to non-linear system dynamics and embedding in the algorithm additional budget constraints enforced over the optimization horizon.



## Acknowledgements

This paper is supported by PNR-PE-AI FAIR project funded by the NextGeneration EU program.

## References

- Abbasi-Yadkori, Y. and Szepesvári, C. Regret bounds for the adaptive control of linear quadratic systems. In *The 24th Annual Conference on Learning Theory*, pp. 1–26, 2011.
- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Agarwal, N., Hazan, E., and Singh, K. Logarithmic regret for online control. In *Advances in Neural Information Processing Systems*, 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. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of IEEE 36th annual foundations of computer science*, pp. 322–331. IEEE, 1995.
- Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R. E. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, 2002.
- Bacchiocchi, F., Genalti, G., Maran, D., Mussi, M., Restelli, M., Gatti, N., and Metelli, A. M. Autoregressive bandits. *CoRR*, abs/2212.06251, 2022.
- 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. In *Conference On Learning Theory*, 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*, pp. 1097–1105. Association for Computing Machinery, 2014.
- Chen, Q., Golrezaei, N., and Bouneffouf, D. Dynamic bandits with temporal structure. *Available at SSRN 3887608*, 2021.
- Court, D., Elzinga, D., Mulder, S., and Vetvik, O. J. The consumer decision journey. *McKinsey Quarterly*, 3:96–107, 2009.
- Dekel, O., Tewari, A., and Arora, R. Online bandit learning against an adaptive adversary: from regret to policy regret. In *International Conference on Machine Learning*, 2012.
- Garg, S. and Akash, A. K. Stochastic bandits with delayed composite anonymous feedback. *CoRR*, abs/1910.01161, 2019.
- Gur, Y., Zeevi, A., and Besbes, O. Stochastic multi-armed-bandit problem with non-stationary rewards. In *Advances in Neural Information Processing Systems*, pp. 199–207, 2014.
- Hespanha, J. P. *Linear Systems Theory: Second Edition*. Princeton University Press, 2018.
- Ho, B. L. and Kalman, 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.
- Isom, J. D., Meyn, S. P., and Braatz, R. D. Piecewise linear dynamic programming for constrained pomdps. In *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence*, pp. 291–296. AAAI Press, 2008.
- Ito, S., Hatano, D., Sumita, H., Takemura, K., Fukunaga, T., Kakimura, N., and Kawarabayashi, K. Delay and cooperation in nonstochastic linear bandits. In *Advances in Neural Information Processing Systems*, 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. In *Proceedings of the 30th International Conference on Machine Learning*, 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 *Proceedings of the 22nd International Joint Conference on Artificial Intelligence*, pp. 1968–1974, 2011.

- Lale, S., Azizzadenesheli, K., Hassibi, B., and Anandkumar, A. Logarithmic regret bound in partially observable linear dynamical systems. In *Advances in Neural Information Processing Systems*, 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*, pp. 5969–5978, 2021.
- Lattimore, T. and Szepesvári, C. *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*, 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.
- 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*, pp. 3348–3356, 2020.
- Mania, H., Tu, S., and Recht, B. Certainty equivalence is efficient for linear quadratic control. In *Advances in Neural Information Processing Systems*, pp. 10154–10164, 2019.
- Nobari, S. DBA: dynamic multi-armed bandit algorithm. In *The Thirty-Third AAAI Conference on Artificial Intelligence*, pp. 9869–9870, 2019.
- Oymak, S. and Ozay, N. Non-asymptotic identification of LTI systems from a single trajectory. In *2019 American Control Conference*, 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*, pp. 4102–4110, 2018.
- Plevrakis, O. and Hazan, E. Geometric exploration for online control. In *Advances in Neural Information Processing Systems*, 2020.
- Russac, Y., Vernade, C., and Cappé, O. Weighted linear bandits for non-stationary environments. In *Advances in Neural Information Processing Systems*, pp. 12017–12026, 2019.
- 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 *Conference on Learning Theory*, volume 125, pp. 3320–3436. PMLR, 2020.
- Thune, T. S., Cesa-Bianchi, N., and Seldin, Y. Nonstochastic multiarmed bandits with unrestricted delays. In *Advances in Neural Information Processing Systems*, pp. 6538–6547, 2019.
- Tsiamis, A. and Pappas, G. J. Finite sample analysis of stochastic system identification. In *58th IEEE Conference on Decision and Control*, pp. 3648–3654, 2019.
- Undurti, A. and How, J. P. An online algorithm for constrained pomdps. In *IEEE International Conference on Robotics and Automation*, 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*, 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*, 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*, 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*, 10(1):174–205, 1965.

## A. Additional Related Works

In this appendix, we report additional details about the related works.

### A.1. Delayed/Aggregated Feedback with DLBs

In this appendix, we show how we can model *delayed* and *composite* feedback with DLBs. For the delayed feedback, we focus on the case in which either the delay is fixed to the value  $\tau \neq 1$ , i.e., the reward of the pull performed at round  $t$  is experienced at round  $t + \tau$ . For the composite feedback, we assume that the reward of the pull performed at round  $t$  is spread over the next  $\tau \neq 1$  rounds with fixed weights  $\rho w_1, \dots, w_\tau$ . Denoting with  $R_t$  the full reward (not observed) due to the pull performed at round  $t$ , the agent at round  $t$  observes the weighted sum of the rewards reported below:<sup>12</sup>

$$\sum_{l=1}^{\tau} w_l R_{t-l}. \quad (6)$$

These two cases can be modeled as DLBs with a suitable encoding of the arms and choice of matrices. In particular, assuming to have  $K$  arms, we take the arm set  $\mathcal{U}$  to be the canonical basis of  $\mathbb{R}^K$ , and we denote with  $\boldsymbol{\mu}$  the vector of expected rewards. We define  $\boldsymbol{\theta} = \mathbf{0}$  and:

$$\mathbf{A} = \begin{pmatrix} 0 & 0 & \dots & 0 & 0 \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 \end{pmatrix} \in \mathbb{R}^{\tau \times \tau}, \quad \mathbf{B} = \begin{pmatrix} \boldsymbol{\mu}_K^T \\ \mathbf{0}_K^T \\ \mathbf{0}_K^T \\ \vdots \\ \mathbf{0}_K^T \end{pmatrix} \in \mathbb{R}^{\tau \times K},$$

$$\boldsymbol{\omega}_{\text{delay}} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^{\tau}, \quad \boldsymbol{\omega}_{\text{composite}} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_\tau \end{pmatrix} \in \mathbb{R}^{\tau}.$$

However, DLBs cannot model random or adversarial delays. Nevertheless, DLBs can capture scenarios of composite feedback in which the reward is spread over an infinite number of rounds. Keeping the  $K$ -armed case introduced above, we can consider the simplest example of a reward that spreads as an autoregressive process AR(1) with parameter  $\gamma \in [0, 1]$ , that cannot be represented using the standard composite feedback. In such a case, we simply need a system with order  $n = 1$  with matrices (actually scalars):

$$\mathbf{A} = \gamma, \quad \mathbf{B} = \mathbf{u}^T, \quad \boldsymbol{\omega} = 1.$$

Clearly, one can consider AR( $m$ ) processes (Bacchiocchi et al., 2022) by employing systems of order  $n = m + 1$ .

### A.2. 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 hidden state and the action. Nevertheless, several works accounted for the presence of constraints (Isom et al., 2008; Undurti & How, 2010; Kim et al., 2011) without exploiting the linearity and without regret guarantees.

### A.3. Adversarial Bandits

It is worth elaborating on the adaptation of adversarial MAB algorithms to this setting. First, since the reward distribution in DLBs depends at every round  $t$  on the sequence of actions played by the agent prior to  $t$ , we can reduce the DLB setting to an adversarial bandit with an *adaptive* (or non-oblivious) adversary. Second, such an adversary must have *infinite memory* in principle. Third, our regret definition of Section 2 is a *policy regret* (Dekel et al., 2012) that compares the algorithm performance against playing the optimal policy in hindsight from the beginning, as opposed to the *external regret* often

<sup>12</sup>It is worth noting that the fixed-delay case is a particular case of composite feedback, where  $w_1 = \dots = w_{\tau-1} = 0$  and  $w_\tau = 1$ .

employed for non-adaptive adversaries. It is well known that for infinite-memory adaptive adversaries, no algorithm can achieve sublinear policy regret. Nevertheless, for DLB setting, we know that the effect of the past is always vanishing (given Assumption 2.1 enforcing  $\rho \mathbf{A} \mathbf{q} \geq 1$ ), so we can approximate our setting as a *finite-memory* setting, by considering memory length  $k \asymp \frac{\log M}{\log 1/\rho} S$ , where  $M$  is the one defined in Algorithm 1 (line 2), with an additional regret term only logarithmic in the optimization horizon  $T$ . Then, given this approximation, we can make use of an adversarial bandit algorithm (designed for non-adaptive adversaries) in the framework proposed by Dekel et al. (2012) to make it effective for the finite-memory adaptive adversary setting. In the case of an optimal algorithm, such as EXP3 (Auer et al., 2002), suffering an external regret of order  $\tilde{O}(\sqrt{MT})$ , being  $M$  the number of arms, the version to address this finite-memory adaptive adversary setting suffers a regret bounded by  $\tilde{O}(\sqrt{k} \sqrt{1/\rho} M^{1/3} T^{2/3})$ , as shown in Theorem 2 of Dekel et al. (2012).

#### A.4. Other Approaches

Non-stationary bandits (Gur et al., 2014) can be regarded as bandits with a hidden state that evolves through a (possibly non-linear) dynamics. The main difference compared with our DLBs is that the hidden state evolves in an *uncontrollable* way, i.e., it does not depend on the sequence of actions performed so far. Russac et al. (2019) extend the linear bandit setting by considering a non-stationary evolution of the parameter  $\theta_t^*$ . The notion of *dynamic* bandit is further studied by Chen et al. (2021), where an auto-regressive process is considered for the evolution of the reward through time and by Nobari (2019) that propose a practical approach to cope with this setting.

## B. Proofs and Derivations

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

### B.1. Proofs of Section 2

Before we proceed, we introduce a different notion of regret useful for analysis purposes, that we name *offline regret*. This notion of regret compares  $J^*$  with the steady-state performance of the action  $\mathbf{u}_t = \pi_t \rho H_{t-1} \mathbf{q}$  played at each round  $t \in \mathcal{T}^K$  by the agent:

$$R^{\text{off}}(\pi, T) := T J^* - \sum_{t=1}^T J(\rho \mathbf{u}_t). \quad (7)$$

We denote with  $\mathbb{E} R^{\text{off}}(\pi, T)$  the *expected offline regret*, where the expectation is taken w.r.t. the randomness of the reward. Clearly, the two notions of regret coincide when the system has no dynamics.

The following result relates the offline and the (online) expected regret.

**Lemma B.1.** *Under Assumptions 2.1 and 2.2, for any policy  $\pi$ , it holds that:*

$$|\mathbb{E} R^{\text{off}}(\pi, T) - \mathbb{E} R(\pi, T)| \leq \frac{\rho \mathbf{A} \mathbf{q} B U}{\rho \mathbf{A} \mathbf{q}^2} \frac{\rho \mathbf{A} \mathbf{q} X}{1 - \rho \mathbf{A} \mathbf{q}}.$$

*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. For every  $t \in \mathcal{T}^K$ , let us denote with  $\mathbb{E} r_{y_t S}$  the expected reward at time  $t$  and with  $J(\rho \mathbf{u}_t)$  as the steady-state performance when executing action  $\mathbf{u}_t$ :

$$\begin{aligned} \mathbb{E} r_{y_t S} &= \sum_{s=0}^{t-1} \mathbf{x}^{\mathbf{h}\{s\}} \mathbf{u}_{t-s} \mathbf{y} = \boldsymbol{\omega}^T \mathbf{A}^{t-1} \mathbf{x}_1 + \boldsymbol{\theta}^T \mathbf{u}_t + \boldsymbol{\omega}^T \sum_{s=1}^{t-1} \mathbf{A}^{s-1} \mathbf{B} \mathbf{u}_{t-s} + \boldsymbol{\omega}^T \mathbf{A}^{t-1} \mathbf{x}_1, \\ J(\rho \mathbf{u}_t) &= \boldsymbol{\theta}^T \mathbf{u}_t + \boldsymbol{\omega}^T \rho \mathbf{I}_d + \mathbf{A} \mathbf{q}^{-1} \mathbf{u}_t + \boldsymbol{\theta}^T \mathbf{u}_t + \boldsymbol{\omega}^T \sum_{s=0}^{+\infty} \mathbf{A}^s \mathbf{u}_t. \end{aligned}$$

We now proceed by summing over  $t \in \mathcal{T}^K$ . First of all, we consider the following preliminary result involving  $y_t$ , which is obtained by rearranging the summations:

$$\sum_{t=1}^T \mathbb{E} r_{y_t S} = \boldsymbol{\theta}^T \sum_{t=1}^T \mathbf{u}_t + \boldsymbol{\omega}^T \sum_{t=1}^T \sum_{s=1}^{t-1} \mathbf{A}^{s-1} \mathbf{B} \mathbf{u}_{t-s} + \boldsymbol{\omega}^T \sum_{t=1}^T \mathbf{A}^{t-1} \mathbf{x}_1$$



$$\theta^\top \sum_{t=1}^T \mathbf{u}_t - \omega^\top \sum_{t=1}^{T-1} \left( \sum_{s=0}^{T-t-1} \mathbf{A}^s \right) \mathbf{B} \mathbf{u}_t - \omega^\top \sum_{t=1}^T \mathbf{A}^{t-1} \mathbf{x}_1.$$

Thus, we have:

$$\begin{aligned} \left| \sum_{t=1}^T \rho J \rho \mathbf{u}_t^\top \mathbf{q} - \mathbb{E} r_{y_t} s_{\mathbf{q}} \right| & \leq \left| \omega^\top \sum_{t=1}^T \left( \sum_{s=0}^{+\infty} \mathbf{A}^s - \sum_{s=0}^{T-t-1} \mathbf{A}^s \right) \mathbf{B} \mathbf{u}_t - \omega^\top \sum_{t=1}^T \mathbf{A}^{t-1} \mathbf{x}_1 \right| \\ & \leq \left| \omega^\top \sum_{t=1}^T \left( \sum_{s=T-t}^{+\infty} \mathbf{A}^s \right) \mathbf{B} \mathbf{u}_t - \omega^\top \sum_{t=1}^T \mathbf{A}^{t-1} \mathbf{x}_1 \right| \\ & \leq \rho \mathbf{A} \mathbf{q} \mathbf{B} \mathbf{U} \sum_{t=1}^T \sum_{s=T-t}^{+\infty} \rho \rho \mathbf{A} \mathbf{q}^s - \rho \mathbf{A} \mathbf{q} \mathbf{X} \sum_{t=1}^T \rho \rho \mathbf{A} \mathbf{q}^{t-1} \end{aligned} \quad (8)$$

$$\leq \frac{\rho \mathbf{A} \mathbf{q} \mathbf{B} \mathbf{U}}{1 - \rho \rho \mathbf{A} \mathbf{q}} \sum_{t=1}^T \rho \rho \mathbf{A} \mathbf{q}^{T-t} - \frac{\rho \mathbf{A} \mathbf{q} \mathbf{X}}{1 - \rho \rho \mathbf{A} \mathbf{q}} \quad (9)$$

$$\leq \frac{\rho \mathbf{A} \mathbf{q} \mathbf{B} \mathbf{U}}{\rho 1 - \rho \rho \mathbf{A} \mathbf{q} \mathbf{q}^2} - \frac{\rho \mathbf{A} \mathbf{q} \mathbf{X}}{1 - \rho \rho \mathbf{A} \mathbf{q}}, \quad (10)$$

where line (8) follows from Assumptions 2.1 and 2.2, lines (9) and (10) follow from bounding the summations with the series. The result follows by observing that:

$$\mathbb{E} R^{\text{off}}(\rho \underline{\pi}, T, \mathbf{q}) - \mathbb{E} R(\rho \underline{\pi}, T, \mathbf{q}) \leq \sum_{t=1}^T \rho J \rho \mathbf{u}_t^\top \mathbf{q} - \mathbb{E} r_{y_t} s_{\mathbf{q}}.$$

□

**Theorem 2.2 (Lower Bound).** *For any policy  $\underline{\pi}$  (even stochastic), there exists a DLB fulfilling Assumptions 2.1 and 2.2, such that for sufficiently large  $T \asymp \mathcal{O}\left(\frac{d^2}{1-\rho(\mathbf{A})}\right)$ , policy  $\underline{\pi}$  suffers an expected regret lower bounded by:*

$$\mathbb{E} R(\rho \underline{\pi}, T, \mathbf{q}) \asymp \left( \frac{d \sqrt{T}}{\rho 1 - \rho \rho \mathbf{A} \mathbf{q} \mathbf{q}^{\frac{1}{2}}} \right).$$

*Proof.* To derive the lower bound, we take inspiration from the construction of (Lattimore & Szepesvári, 2020) for linear bandits (Theorem 24.1). We consider a class of DLBs defined in terms of fixed  $0 \leq \rho \leq 1$  and  $0 \leq \epsilon \leq \rho$  with  $\omega = \mathbf{1}_d$ ,  $\theta = \frac{2(1-\rho)+\epsilon}{2(1-(\rho-\epsilon))} \mathbf{1}_d$ ,  $\mathbf{B} = \rho 1 - \rho \mathbf{q} \mathbf{I}_d$  and with a diagonal dynamical matrix  $\mathbf{A} = \text{diag}(\rho \mathbf{a}_j)$ , defined in terms of the vector  $\mathbf{a}$  belonging to the set  $\mathcal{A} = \{ \mathbf{t} \in \{1, 1\}^d \}$ . Let us note that  $|\mathcal{A}| = |\mathcal{U}| = 2^d$ . Thus, in our set of DLBs, the vector  $\mathbf{a}$  fully characterizes the problem. Moreover, we observe that, given the diagonal  $\mathbf{a} = \text{diag}(\rho \mathbf{a}_j)$ , we can compute the cumulative Markov parameter  $\mathbf{h}_{\mathbf{a}} = \text{sign}(\rho \mathbf{a}_j) \frac{\epsilon}{2(1-(\rho-\epsilon))}$ .<sup>13</sup> As a consequence the optimal action can be defined as  $\mathbf{u}_{\mathbf{a}}^* = \text{sign}(\rho \mathbf{a}_j)$ , whose performance is given by  $J_{\mathbf{a}}^* = \chi(\mathbf{h}_{\mathbf{a}}, \mathbf{u}_{\mathbf{a}}^*) = \frac{\epsilon d}{2(1-(\rho-\epsilon))}$ .

Let us consider the probability distribution over the canonical bandit model induced by executing a policy  $\underline{\pi}$  in a DLB characterized by the diagonal of the dynamical matrix  $\mathbf{a} \in \mathcal{A}$  and with Gaussian diagonal noise:

$$P_{\mathbf{a}} = \prod_{t=1}^T \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t, \sigma^2 \mathbf{I}_d) \mathcal{N}(y_t | \chi(\theta, \mathbf{u}_t), \chi(\omega, \mathbf{x}_t), \sigma^2) \mathcal{N}(\mathbf{u}_t | H_{t-1}, \mathbf{q}),$$

where  $H_{t-1}$  is the history of observations up to time  $t-1$ . We denote with  $\mathbb{E}_{\mathbf{a}}$  the expectation induced by the distribution  $P_{\mathbf{a}}$ . For every  $i \in \mathcal{P} \setminus \mathcal{J} \setminus \mathcal{K}$ , let us now consider an alternative DLB instance that differs on the dynamical matrix only. Specifically:

$$\mathbf{a}'_j = \begin{cases} \mathbf{a}_j & \text{if } j = i \\ \rho & \text{if } j = i \text{ and } \mathbf{a}_j = \rho - \epsilon, \\ \rho - \epsilon & \text{if } j = i \text{ and } \mathbf{a}_j = \rho \end{cases}, \quad @j \in \mathcal{P} \setminus \mathcal{J} \setminus \mathcal{K}.$$

<sup>13</sup>For a vector  $\mathbf{v} \in \mathbb{R}^d$ , we denote with  $\text{sign}(\mathbf{v}) \in \{1, -1\}^d$  the vector of the signs of the components of  $\mathbf{v}$ . It is irrelevant how we convene to define the sign of 0.

By relative entropy identities (Lattimore & Szepesvári, 2020), let  $\mathbf{A} = \text{diag}(\rho, \rho)$  and  $\mathbf{A}' = \text{diag}(\rho, \rho)$ , we have:

$$D_{\text{KL}}(\mathcal{P}_{\mathbf{a}, \mathbf{P}_{\mathbf{a}'}} | \mathcal{P}_{\mathbf{a}, \mathbf{P}_{\mathbf{a}'}}) = \mathbb{E}_{\mathbf{a}} \left[ \sum_{t=1}^T D_{\text{KL}}(\mathcal{N}(\rho | \mathbf{A} \mathbf{x}_t, \sigma^2 \mathbf{I}_d) | \mathcal{N}(\rho | \mathbf{A}' \mathbf{x}_t, \sigma^2 \mathbf{I}_d)) \right]$$

$$\frac{1}{2\sigma^2} \sum_{t=1}^T \mathbb{E}_{\mathbf{a}} \left[ \left\| (\mathbf{A} - \mathbf{A}') \mathbf{x}_t \right\|_2^2 \right] = \epsilon^2 \mathbb{E}_{\mathbf{a}} [\mathbf{x}_{t,i}^2].$$

We proceed at properly bounding the KL-divergence, letting  $\mathbf{e}_i$  be the  $i$ -th vector of the canonical basis of  $\mathbb{R}^d$  and convening that  $\mathbf{x}_0 = \mathbf{0}_d$ :

$$\mathbb{E}_{\mathbf{a}} [\mathbf{x}_{t,i}^2] = \mathbb{E}_{\mathbf{a}} \left[ \left( \sum_{s=1}^{t-1} \mathbf{e}_i^\top \mathbf{A}^s \mathbf{B} \mathbf{u}_{t-s} + \sum_{s=1}^{t-1} \mathbf{e}_i^\top \mathbf{A}^s \boldsymbol{\epsilon}_{t-s} \right)^2 \right]$$

$$= \mathbb{E}_{\mathbf{a}} \left[ \left( \rho \mathbf{1} + \rho \sum_{s=1}^{t-1} \mathbf{a}_i^s \mathbf{u}_{t-s,i} + \sum_{s=1}^{t-1} \mathbf{a}_i^s \boldsymbol{\epsilon}_{t-s,i} \right)^2 \right]$$

$$= \mathbb{E}_{\mathbf{a}} \left[ \underbrace{\rho^2 \sum_{s=1}^{t-1} \sum_{l=1}^{t-1} \mathbf{a}_i^{s+l} \mathbf{u}_{t-s,i} \mathbf{u}_{t-l,i}}_{(a)} + \underbrace{2\rho \sum_{s=1}^{t-1} \sum_{l=1}^{t-1} \mathbf{a}_i^{s+l} \mathbf{u}_{t-s,i} \boldsymbol{\epsilon}_{t-l,i}}_{(b)} + \underbrace{\sum_{s=1}^{t-1} \sum_{l=1}^{t-1} \mathbf{a}_i^{s+l} \boldsymbol{\epsilon}_{t-s,i} \boldsymbol{\epsilon}_{t-l,i}}_{(c)} \right]$$

Let us start with (a):

$$\rho^2 \sum_{s=1}^{t-1} \sum_{l=1}^{t-1} \mathbf{a}_i^{s+l} \mathbf{u}_{t-s,i} \mathbf{u}_{t-l,i} \preceq \rho^2 \sum_{s=1}^{t-1} \sum_{l=1}^{t-1} \rho^{s+l} \preceq 1,$$

having observed that  $|\mathbf{u}_{t-s,i}|, |\mathbf{u}_{t-l,i}| \preceq 1$ , that  $|\mathbf{a}_i| \preceq \rho$ , and bounding the summations with the series. Let us move to (b):

$$\rho \sum_{s=1}^{t-1} \sum_{l=1}^{t-1} \mathbf{a}_i^{s+l} \mathbf{u}_{t-s,i} \boldsymbol{\epsilon}_{t-l,i} \preceq \rho \sum_{s=1}^{t-1} \sum_{l=s+1}^{t-1} \mathbf{a}_i^{s+l} \mathbf{u}_{t-s,i} \boldsymbol{\epsilon}_{t-l,i}$$

$$\preceq \rho \sum_{s=1}^{t-1} \sum_{l=s+1}^{t-1} \rho^{s+l} \mathbb{E}_{\mathbf{a}} |\boldsymbol{\epsilon}_{t-l,i}|$$

$$\preceq \frac{\sigma}{1-\rho} \sqrt{\frac{2}{\pi}},$$

having observed that  $\mathbf{u}_{t-s,i}$  and  $\boldsymbol{\epsilon}_{t-l,i}$  are independent when  $s \neq l$  and  $\boldsymbol{\epsilon}_{t-l,i}$  has zero mean, that  $|\mathbf{u}_{t-s,i}| \preceq 1$ , that  $\mathbf{a}_i^{s+l} \preceq \rho^{s+l}$ , and that the expectation of the absolute value of random variable normally distributed is given by  $\mathbb{E} |\boldsymbol{\epsilon}_{t-l,i}| = \sigma \sqrt{\frac{2}{\pi}}$ .

Finally, let us consider (c):

$$\sum_{s=1}^{t-1} \sum_{l=1}^{t-1} \mathbf{a}_i^{s+l} \boldsymbol{\epsilon}_{t-s,i} \boldsymbol{\epsilon}_{t-l,i} \preceq \sum_{s=1}^{t-1} \mathbf{a}_i^{2s} \boldsymbol{\epsilon}_{t-s,i} \boldsymbol{\epsilon}_{t-s,i} \preceq 2 \sum_{s=1}^{t-1} \sum_{l=s+1}^{t-1} \mathbf{a}_i^{s+l} \boldsymbol{\epsilon}_{t-s,i} \boldsymbol{\epsilon}_{t-l,i}$$

$$\preceq \sigma^2 \sum_{s=1}^{t-1} \rho^{2s} \preceq \frac{\sigma^2}{1-\rho^2} \preceq \frac{\sigma^2}{1-\rho},$$

having observed that the noise vectors  $\boldsymbol{\epsilon}_{t-l,i}$  and  $\boldsymbol{\epsilon}_{t-s,i}$  are independent whenever  $s \neq l$ , that  $\mathbb{E}_{\mathbf{a}} \boldsymbol{\epsilon}_{t-s,i}^2 = \sigma^2$ , and having bounded the sum with the series. Coming back to the original bound, we have:

$$\mathbb{E}_{\mathbf{a}} [\mathbf{x}_{t,i}^2] \preceq 1 + \frac{1}{1-\rho} \left( \sigma^2 + 2\sigma \sqrt{\frac{2}{\pi}} \right).$$

For  $i \in [d]$  and  $\mathbf{a} \in \mathcal{A}$ , we introduce the symbol:

$$p_{\mathbf{a},i} = \mathbb{P}_{\mathbf{a}} \left( \sum_{t=1}^T \mathbf{1}_{\text{sign}(\mathbf{p}_{\mathbf{a},i} \mathbf{u}_{t,i} - \text{sign}(\mathbf{h}_{\mathbf{a},i} \mathbf{q})) \neq 0} \right).$$

Thus, for  $\mathbf{a}$  and  $\mathbf{a}'$  defined as above, by the Bretagnolle-Huber inequality (Lattimore & Szepesvári, 2020, Theorem 14.2), we have:

$$\begin{aligned} p_{\mathbf{a},i} - p_{\mathbf{a}',i} &\leq \frac{1}{2} \exp(-D_{\text{KL}}(\mathbb{P}_{\mathbf{a}} \|\mathbb{P}_{\mathbf{a}'}) - \frac{1}{2} \exp \left( \frac{1}{2\sigma^2} \sum_{t=1}^T \mathbb{E}_{\mathbb{P}} [\|\mathbf{A} - \mathbf{A}'\|_2^2] \right) \\ &\leq \frac{1}{2} \exp \left( \frac{T\epsilon^2}{2} \left( \frac{1}{\sigma^2} + \frac{1}{1-\rho} \left( 1 + \frac{2}{\sigma} \sqrt{\frac{2}{\pi}} \right) \right) \right) \\ &\leq \frac{1}{2} \exp \left( \frac{2T\epsilon^2}{1-\rho} \right), \end{aligned}$$

having selected  $\sigma^2 = 1$ . We use the notation  $\sum_{\mathbf{a}_{-i}}$  to denote the multiple summation  $\sum_{\mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \mathbf{a}_{i+1}, \dots, \mathbf{a}_d \in \{\rho, \rho-\epsilon\}^{d-1}}$ :

$$\begin{aligned} \sum_{\mathbf{a} \in \mathcal{A}} 2^{-d} \sum_{i=1}^d p_{\mathbf{a},i} &= \sum_{i=1}^d \sum_{\mathbf{a}_{-i}} 2^{-d} \sum_{\mathbf{a}_i \in \{\rho, \rho-\epsilon\}} p_{\mathbf{a},i} \\ &\leq \sum_{i=1}^d \sum_{\mathbf{a}_{-i}} 2^{-d} \frac{1}{2} \exp \left( \frac{2T\epsilon^2}{1-\rho} \right) \\ &= \frac{d}{4} \exp \left( \frac{2T\epsilon^2}{1-\rho} \right). \end{aligned}$$

Therefore, with this averaging argument, we can conclude that there exists  $\mathbf{a}^* \in \mathcal{A}$  such that  $\sum_{i=1}^d p_{\mathbf{a}^*,i} \leq \frac{d}{4} \exp \left( \frac{2T\epsilon^2}{1-\rho} \right)$ .

For this choice  $\mathbf{a}^*$ , we consider  $\mathbf{u}_{\mathbf{a}^*} = \text{sign}(\mathbf{p}_{\mathbf{a}^*} \mathbf{q}) \in \mathcal{U}$ , we can proceed to the lower bound on the expected offline regret:

$$\begin{aligned} \mathbb{E} R_{\mathbf{p}, \mathbf{q}}^{\text{off}}(T) &= \sum_{t=1}^T \mathbb{E}_{\mathbf{a}^*} [\mathbf{r}(\mathbf{h}_{\mathbf{a}^*}, \mathbf{u}_{\mathbf{a}^*}) - \mathbf{u}_{t,y}] \\ &= \sum_{t=1}^T \mathbb{E}_{\mathbf{a}^*} \left[ \sum_{i=1}^d \mathbf{1}_{\text{sign}(\mathbf{p}_{\mathbf{a}^*,i} \mathbf{u}_{t,i} - \text{sign}(\mathbf{h}_{\mathbf{a}^*,i} \mathbf{q})) \neq 0} \frac{\epsilon}{1-\rho-\epsilon} \right] \\ &= \frac{\epsilon}{1-\rho-\epsilon} \sum_{t=1}^T \sum_{i=1}^d \mathbb{P}_{\mathbf{a}^*} (\text{sign}(\mathbf{p}_{\mathbf{a}^*,i} \mathbf{u}_{t,i} - \text{sign}(\mathbf{h}_{\mathbf{a}^*,i} \mathbf{q})) \neq 0) \\ &\leq \frac{T\epsilon}{2(1-\rho-\epsilon)} \sum_{i=1}^d \mathbb{P}_{\mathbf{a}^*} \left( \sum_{t=1}^T \mathbf{1}_{\text{sign}(\mathbf{p}_{\mathbf{a}^*,i} \mathbf{u}_{t,i} - \text{sign}(\mathbf{h}_{\mathbf{a}^*,i} \mathbf{q})) \neq 0} \right) \\ &= \frac{T\epsilon}{2(1-\rho-\epsilon)} \sum_{i=1}^d p_{\mathbf{a}^*,i} \leq \frac{Td\epsilon}{8(1-\rho-\epsilon)} \exp \left( \frac{2T\epsilon^2}{1-\rho} \right). \end{aligned}$$

We now maximize over  $0 < \epsilon < \rho$ . To this end, we perform the substitution  $\epsilon = \frac{(1-\rho)\tilde{\epsilon}}{1-\tilde{\epsilon}}$ , with  $0 < \tilde{\epsilon} < \rho$ :

$$\frac{Td\epsilon}{8(1-\rho-\epsilon)} \exp \left( \frac{2T\epsilon^2}{1-\rho} \right) = \frac{Td\tilde{\epsilon}}{8} \exp \left( \frac{2\tilde{\epsilon}^2 T(1-\rho)}{(1-\tilde{\epsilon})^2} \right) \leq \frac{Td\tilde{\epsilon}}{8} \exp \left( 8\tilde{\epsilon}^2 T(1-\rho) \right),$$

where the last inequality holds for  $\tilde{\epsilon} \leq \frac{1}{2}$ . We not take  $\tilde{\epsilon} = \frac{1}{\sqrt{8T(1-\rho)}}$  which is smaller than  $\frac{1}{2}$  if  $T \leq \frac{1}{2(1-\rho)}$ , to get:

$$\mathbb{E} R_{\mathbf{p}, \mathbf{q}}^{\text{off}}(T) \leq \frac{d \sqrt{T}}{\sqrt{512e(1-\rho)}}.$$

Notice that with this choice of  $\tilde{\epsilon}$  (and, consequently, of  $\epsilon$ ), for sufficiently large  $T$ , we fulfill Assumption 2.2. Indeed:

$$\theta = 1 - \frac{1}{\sqrt{32T(1-\rho)}}, \quad J_{\mathbf{a}^*} = \frac{d}{\sqrt{32T(1-\rho)}}.$$

Thus, we require  $T \asymp \mathcal{O}\left(\frac{d^2}{1-\rho}\right)$ . Finally, to convert this result to the expected regret, we employ Lemma B.1:

$$\mathbb{E}R^{\text{off}}(\underline{\pi}, T) \asymp \mathbb{E}R^{\text{off}}(\underline{\pi}, T) \frac{d}{1-\rho}.$$

Under the constraint  $T \asymp \mathcal{O}\left(\frac{d^2}{1-\rho}\right)$ , we observe that:

$$\mathbb{E}R^{\text{off}}(\underline{\pi}, T) \asymp \left(\frac{d \sqrt{T}}{\rho \sqrt{1-\rho}}\right).$$

□

**Theorem 2.1** (Optimal Policy). *Under Assumptions 2.1 and 2.2, an optimal policy  $\underline{\pi}^*$  maximizing the (infinite-horizon) expected average reward  $J(\underline{\pi})$  (Equation 2), for every round  $t \in \mathbb{N}$  and history  $H_{t-1} \in \mathcal{H}_{t-1}$  is given by:*

$$\pi_t^* \in \arg\max_{\mathbf{u} \in \mathcal{U}} \mathbf{u}^\top \mathbf{h}_{t-1} \quad \text{where } \mathbf{u}^* \in \arg\max_{\mathbf{u} \in \mathcal{U}} \mathbf{u}^\top \mathbf{h}_{t-1}. \quad (4)$$

*Proof.* Referring to the notation of Appendix C, we first observe that for every policy  $\underline{\pi}$ , we have  $J(\underline{\pi}) = \liminf_{H \rightarrow +\infty} J_H(\underline{\pi})$ , where  $J_H(\underline{\pi}) = \frac{1}{H} \mathbb{E} \sum_{t=1}^H y_t$ , is the  $H$ -horizon expected average reward. Let us start with Equation (18), a fixed finite  $H \in \mathbb{N}$ , and considering the sequence of actions  $\mathbf{u}_1, \mathbf{u}_2, \dots$  generated by policy  $\underline{\pi}$ :

$$\begin{aligned} J_H(\underline{\pi}) &= \frac{1}{H} \sum_{s=1}^H \mathbf{x}_s^\top \mathbf{h}^{J_0, H-s}, \mathbb{E} \mathbf{u}_s^\top \mathbf{y}_s = \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1 \mathbf{y}_t \\ &= \frac{1}{H} \sum_{s=1}^H \mathbf{x}_s^\top \mathbf{h}, \mathbb{E} \mathbf{u}_s^\top \mathbf{y}_s = \frac{1}{H} \sum_{s=1}^H \mathbf{x}_s^\top \mathbf{h}^{J_{H-s+1, +\infty}}, \mathbb{E} \mathbf{u}_s^\top \mathbf{y}_s = \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1 \mathbf{y}_t. \end{aligned}$$

Now, we consider two bounds on  $J_H(\underline{\pi})$ , obtained by an application of Cauchy-Schwarz inequality on the second addendum:

$$\begin{aligned} J_H(\underline{\pi}) &\leq \frac{1}{H} \sum_{s=1}^H \mathbf{x}_s^\top \mathbf{h}, \mathbb{E} \mathbf{u}_s^\top \mathbf{y}_s = \frac{1}{H} \sum_{s=1}^H \left\| \mathbf{h}^{J_{H-s+1, +\infty}} \right\|_2 \left\| \mathbb{E} \mathbf{u}_s \mathbf{y}_s \right\|_2 \\ &= \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1 \mathbf{y}_t =: J_H^\uparrow(\underline{\pi}), \\ J_H(\underline{\pi}) &\geq \frac{1}{H} \sum_{s=1}^H \mathbf{x}_s^\top \mathbf{h}, \mathbb{E} \mathbf{u}_s^\top \mathbf{y}_s = \frac{1}{H} \sum_{s=1}^H \left\| \mathbf{h}^{J_{H-s+1, +\infty}} \right\|_2 \left\| \mathbb{E} \mathbf{u}_s \mathbf{y}_s \right\|_2 \\ &= \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1 \mathbf{y}_t =: J_H^\downarrow(\underline{\pi}). \end{aligned}$$

Concerning the term  $\left\| \mathbb{E} \mathbf{u}_s \mathbf{y}_s \right\|_2$ , we have that  $\left\| \mathbb{E} \mathbf{u}_s \mathbf{y}_s \right\|_2 \leq U$ , having used Jensen's inequality and under Assumption 2.2. Regarding the second term, using Assumptions 2.1 and 2.2, we obtain:

$$\begin{aligned} \left\| \mathbf{h}^{J_{H-s+1, +\infty}} \right\|_2 &= \left\| \sum_{l=H-s+1}^{+\infty} \mathbf{B}^\top \rho \mathbf{A}^{l-1} \mathbf{q}^\top \boldsymbol{\omega} \right\|_2 \\ &\leq B \sum_{l=H-s+1}^{+\infty} \rho \mathbf{A}^l \rho \mathbf{A}^{l-1} \\ &= B \frac{\rho \mathbf{A}^H \rho \mathbf{A}^{H-s}}{1 - \rho \mathbf{A}^H}. \end{aligned} \quad (11)$$

Plugging this result into the summation over  $s$ , we obtain:

$$\frac{1}{H} \sum_{s=1}^H \frac{B}{1-\rho} \frac{\rho \mathbf{A}^H}{\rho \mathbf{A}^H} \sum_{s=1}^H \rho \mathbf{A}^{H-s} = \frac{B}{H} \frac{\rho \mathbf{A}^H}{1-\rho} \frac{1-\rho \mathbf{A}^H}{\rho \mathbf{A}^H}.$$

It is simple to observe that the last term approaches zero as  $H \rightarrow \infty$ . Moreover, with an analogous argument, it can be proved that  $\left\| \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1 \mathbf{y}_t \right\|_2 \rightarrow 0$  as  $H \rightarrow \infty$ . Thus, we have that  $\liminf_{H \rightarrow +\infty} J_H(\underline{\pi}) = J(\underline{\pi})$ .



$\liminf_{H \rightarrow +\infty} J_H^\uparrow \rho \underline{\pi} \mathbf{q}$ . Consequently, by the squeezing theorem of limits, we have:

$$\begin{aligned} J \rho \underline{\pi} \mathbf{q} &= \liminf_{H \rightarrow +\infty} J_H^\uparrow \rho \underline{\pi} \mathbf{q} = \liminf_{H \rightarrow +\infty} J_H^\downarrow \rho \underline{\pi} \mathbf{q} \\ &= \liminf_{H \rightarrow +\infty} \frac{1}{H} \sum_{s=1}^H \mathbf{x}_s \mathbf{h}_s, \mathbb{E} \mathbf{r}_{\mathbf{u}_s} \mathbf{y}_s = \mathbf{h}^\top \left( \liminf_{H \rightarrow +\infty} \frac{1}{H} \sum_{s=1}^H \mathbb{E} \mathbf{r}_{\mathbf{u}_s} \mathbf{y}_s \right). \end{aligned}$$

It follows that an optimal policy is a policy that plays the constant action  $\mathbf{u}^* = \arg \max_{\mathbf{u} \in \mathcal{U}} \mathbf{x} \mathbf{h}, \mathbf{u} \mathbf{y}$ .  $\square$

## B.2. Proofs of Section 3

**Theorem 3.1** (Self-Normalized Concentration). *Let  $\hat{\mathbf{h}}_t \mathbf{q}_{t \in \mathbb{N}}$  be the sequence of solutions of the Ridge regression problems of Algorithm 1. Then, under Assumption 2.1 and 2.2, for every  $\lambda \neq 0$  and  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , simultaneously for all rounds  $t \in \mathbb{N}$ , it holds that:*

$$\begin{aligned} \|\hat{\mathbf{h}}_t - \mathbf{h}\|_{\mathbf{V}_t} &\leq \frac{c_1}{\lambda} \log \frac{e}{\delta} + c_2 \frac{1}{\lambda} \\ &\quad \sqrt{2\tilde{\sigma}^2 \left( \log \left( \frac{1}{\delta} \right) + \frac{1}{2} \log \left( \frac{\det \mathbf{p} \mathbf{V}_t \mathbf{q}}{\lambda^d} \right) \right)}, \end{aligned}$$

where  $c_1 = \frac{UB}{1-\rho(\mathbf{A})} \|\mathbf{A}\|$ ,  $c_2 = \frac{\Omega_B \Phi(\mathbf{A})}{1-\rho(\mathbf{A})}$ , and  $\tilde{\sigma}^2 = \sigma^2 \left( 1 + \frac{\Omega^2 \Phi(\mathbf{A})^2}{1-\rho(\mathbf{A})^2} \right)$ .

*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 matrices  $\mathbf{V}_t$  and  $\mathbf{b}_t$ ).<sup>14</sup>

$$t_0 = 0, \quad t_m = t_{m-1} + 1 + H_m.$$

We now proceed to define suitable filtrations. Let  $\mathcal{F}_t = \sigma(\mathcal{F}_t \mathbf{q}_{t \in \mathbb{N}})$  such that for every  $t \in \mathbb{N}$ , the random variables  $\mathbf{u}_1, y_1, \dots, \mathbf{u}_{t-1}, y_{t-1}, \mathbf{u}_t \mathbf{u}$  are  $\mathcal{F}_{t-1}$ -measurable, i.e.,  $\mathcal{F}_{t-1} = \sigma(\mathbf{u}_1, y_1, \dots, \mathbf{u}_{t-1}, y_{t-1}, \mathbf{u}_t \mathbf{q})$ . Let us also consider the filtration indexed by  $m$ , denoted with  $\tilde{\mathcal{F}}_m = \sigma(\tilde{\mathcal{F}}_m \mathbf{q}_{m \in \mathbb{N}})$  and defined for all  $m \in \mathbb{N}$  as  $\tilde{\mathcal{F}}_m = \mathcal{F}_{t_{m+1}-1}$ . Thus, the random variables  $\tilde{\mathcal{F}}_{m-1}$ -measurable are those realized until the end of epoch  $m$  except for  $y_{t_m}$ .

Since the estimates  $\hat{\mathbf{h}}_t$  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:

$$\begin{aligned} \tilde{y}_m &= y_{t_m}, \\ \tilde{\mathbf{u}}_m &= \mathbf{u}_{t_m}, \quad (\text{or any } \mathbf{u}_l \text{ with } l \in [t_{m-1} + 1, t_m] \text{ since they are all equal}) \\ \tilde{\xi}_m &= \eta_{t_m} + \sum_{s=1}^{H_m+1} \boldsymbol{\omega}^\top \mathbf{A}^{s-1} \boldsymbol{\epsilon}_{t_m-s}, \\ \tilde{\mathbf{x}}_{m-1} &= \mathbf{x}_{t_{m-1}}, \\ \tilde{\mathbf{h}}_m &= \hat{\mathbf{h}}_{t_m}, \\ \tilde{\mathbf{V}}_m &= \mathbf{V}_{t_m}, \\ \tilde{\mathbf{b}}_m &= \mathbf{b}_{t_m}. \end{aligned}$$

We prove that  $\rho \tilde{\xi}_m \mathbf{q}_{m \in \mathbb{N}}$  is a martingale difference process adapted to the filtration  $\tilde{\mathcal{F}}$ . To this end, we recall that, by construction,  $\rho \eta_t \mathbf{q}_{t \in \mathbb{N}}$  and  $\rho \boldsymbol{\epsilon}_t \mathbf{q}_{t \in \mathbb{N}}$  are martingale difference processes adapted to the filtration  $\mathcal{F}$ . It is clear that  $\tilde{\xi}_m$  is  $\mathcal{F}_{t_m}$ -measurable and, being  $\sigma^2$ -subgaussian it is absolutely integrable. Furthermore, using the tower law of expectation:

$$\begin{aligned} \mathbb{E} \left[ \tilde{\xi}_m | \tilde{\mathcal{F}}_{m-1} \right] &= \mathbb{E} \left[ \eta_{t_m} + \sum_{s=1}^{H_m+1} \boldsymbol{\omega}^\top \mathbf{A}^{s-1} \boldsymbol{\epsilon}_{t_m-s} | \mathcal{F}_{t_m-1} \right] \\ &= \mathbb{E} \left[ \eta_{t_m} | \mathcal{F}_{t_m-1} \right] + \mathbb{E} \left[ \sum_{s=1}^{H_m+1} \boldsymbol{\omega}^\top \mathbf{A}^{s-1} \mathbb{E} \left[ \boldsymbol{\epsilon}_{t_m-s} | \mathcal{F}_{t_m-s-1} \right] | \mathcal{F}_{t_m-1} \right] = 0, \end{aligned}$$

since the system is operating by persisting the action after having decided it at the beginning of the epoch. Thus, by

<sup>14</sup>It is worth noting that the variables  $t_m$  are deterministic.

exploiting the decomposition in Equation (1), we can write:

$$\begin{aligned}
 \tilde{\mathbf{y}}_m &= \mathbf{y}_{t_m} \quad \chi \mathbf{h}^{J_0, H_m+1K}, \tilde{\mathbf{u}}_m \mathbf{y} \quad \boldsymbol{\omega}^\top \mathbf{A}^{H_m+1} \mathbf{x}_{t_{m-1}} \quad \eta_{t_m} \quad \sum_{s=1}^{H_m+1} \boldsymbol{\omega}^\top \mathbf{A}^{s-1} \boldsymbol{\epsilon}_{t_m-s} \\
 & \quad \chi \mathbf{h}^{J_0, H_m+1K}, \tilde{\mathbf{u}}_m \mathbf{y} \quad \boldsymbol{\omega}^\top \mathbf{A}^{H_m+1} \tilde{\mathbf{x}}_{m-1} \quad \tilde{\boldsymbol{\xi}}_m \\
 & \quad \chi \mathbf{h}, \tilde{\mathbf{u}}_m \mathbf{y} \quad \chi \mathbf{h}^{J_{H_m+2, \infty M}}, \tilde{\mathbf{u}}_m \mathbf{y} \quad \boldsymbol{\omega}^\top \mathbf{A}^{H_m+1} \tilde{\mathbf{x}}_{m-1} \quad \tilde{\boldsymbol{\xi}}_m,
 \end{aligned} \tag{12}$$

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

$$\begin{aligned}
 \tilde{\mathbf{U}}_m &= \begin{pmatrix} \tilde{\mathbf{u}}_1^\top \\ \vdots \\ \tilde{\mathbf{u}}_m^\top \end{pmatrix} \in \mathbb{R}^{m \times d}, & \tilde{\mathbf{y}}_m &= \begin{pmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_m \end{pmatrix} \in \mathbb{R}^m, \\
 \tilde{\boldsymbol{\xi}}_m &= \begin{pmatrix} \tilde{\xi}_1 \\ \vdots \\ \tilde{\xi}_m \end{pmatrix} \in \mathbb{R}^m, & \tilde{\boldsymbol{\nu}}_m &= \begin{pmatrix} \boldsymbol{\omega}^\top \mathbf{A}^{H_m+2} \tilde{\mathbf{x}}_0 \\ \vdots \\ \boldsymbol{\omega}^\top \mathbf{A}^{H_m+2} \tilde{\mathbf{x}}_{m-1} \end{pmatrix} \in \mathbb{R}^m, \\
 \tilde{\mathbf{g}}_m &= \begin{pmatrix} \chi \mathbf{h}^{J_{H_m+1, \infty M}}, \tilde{\mathbf{u}}_1 \mathbf{y} \\ \vdots \\ \chi \mathbf{h}^{J_{H_m+1, \infty M}}, \tilde{\mathbf{u}}_m \mathbf{y} \end{pmatrix} \in \mathbb{R}^m.
 \end{aligned}$$

Using the vectors and matrices above, we observe that  $\tilde{\mathbf{V}}_m = \lambda \mathbf{I} + \tilde{\mathbf{U}}_m^\top \tilde{\mathbf{U}}_m$  and  $\tilde{\mathbf{b}}_m = \tilde{\mathbf{U}}_m^\top \tilde{\mathbf{y}}_m$ . Furthermore, by exploiting Equation (12), we can write:

$$\tilde{\mathbf{y}}_m = \tilde{\mathbf{U}}_m \mathbf{h} + \tilde{\mathbf{g}}_m + \tilde{\boldsymbol{\nu}}_m + \tilde{\boldsymbol{\xi}}_m.$$

Let us consider the estimate at  $m$   $\hat{\mathbf{h}}_m$ :

$$\begin{aligned}
 \hat{\mathbf{h}}_m &= \tilde{\mathbf{V}}_m^{-1} \tilde{\mathbf{b}}_m = \left( \lambda \mathbf{I} + \tilde{\mathbf{U}}_m^\top \tilde{\mathbf{U}}_m \right)^{-1} \tilde{\mathbf{U}}_m^\top \tilde{\mathbf{y}}_m \\
 &= \left( \lambda \mathbf{I} + \tilde{\mathbf{U}}_m^\top \tilde{\mathbf{U}}_m \right)^{-1} \tilde{\mathbf{U}}_m^\top \left( \tilde{\mathbf{U}}_m \mathbf{h} + \tilde{\mathbf{g}}_m + \tilde{\boldsymbol{\nu}}_m + \tilde{\boldsymbol{\xi}}_m \right) \\
 &= \mathbf{h} + \left( \lambda \mathbf{I} + \tilde{\mathbf{U}}_m^\top \tilde{\mathbf{U}}_m \right)^{-1} \left( \lambda \tilde{\mathbf{g}}_m + \tilde{\mathbf{U}}_m^\top \tilde{\boldsymbol{\nu}}_m + \tilde{\mathbf{U}}_m^\top \tilde{\boldsymbol{\xi}}_m \right).
 \end{aligned}$$

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

$$\begin{aligned}
 \|\hat{\mathbf{h}}_m - \mathbf{h}\|_{\tilde{\mathbf{V}}_m} &\leq \lambda \|\tilde{\mathbf{V}}_m^{-1} \mathbf{h}\|_{\tilde{\mathbf{V}}_m} + \|\tilde{\mathbf{V}}_m^{-1} \tilde{\mathbf{U}}_m^\top \tilde{\mathbf{g}}_m\|_{\tilde{\mathbf{V}}_m} + \|\tilde{\mathbf{V}}_m^{-1} \tilde{\mathbf{U}}_m^\top \tilde{\boldsymbol{\nu}}_m\|_{\tilde{\mathbf{V}}_m} + \|\tilde{\mathbf{V}}_m^{-1} \tilde{\mathbf{U}}_m^\top \tilde{\boldsymbol{\xi}}_m\|_{\tilde{\mathbf{V}}_m} \\
 &= \underbrace{\lambda \|\mathbf{h}\|_{\tilde{\mathbf{V}}_m^{-1}}}_{(a)} + \underbrace{\|\tilde{\mathbf{U}}_m^\top \tilde{\mathbf{g}}_m\|_{\tilde{\mathbf{V}}_m^{-1}}}_{(b)} + \underbrace{\|\tilde{\mathbf{U}}_m^\top \tilde{\boldsymbol{\nu}}_m\|_{\tilde{\mathbf{V}}_m^{-1}}}_{(c)} + \underbrace{\|\tilde{\mathbf{U}}_m^\top \tilde{\boldsymbol{\xi}}_m\|_{\tilde{\mathbf{V}}_m^{-1}}}_{(d)},
 \end{aligned}$$

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

$$\begin{aligned}
 (a)^2 &= \lambda^2 \|\mathbf{h}\|_{\tilde{\mathbf{V}}_m^{-1}}^2 = \lambda^2 \mathbf{h}^\top \tilde{\mathbf{V}}_m^{-1} \mathbf{h} \\
 &\leq \lambda^2 \|\tilde{\mathbf{V}}_m^{-1}\|_2 \|\mathbf{h}\|_2^2 \\
 &\leq \lambda \|\mathbf{h}\|_2^2 \\
 &\leq \lambda \left( \frac{B \rho \mathbf{A} \mathbf{q}}{1 - \rho \mathbf{A} \mathbf{q}} \right)^2,
 \end{aligned}$$

where we observed that  $\|\tilde{\mathbf{V}}_m^{-1}\|_2 \leq \|\tilde{\mathbf{V}}_m\|_2^{-1} \leq \lambda^{-1}$ . Finally, we have bounded the norm of  $\mathbf{h}$ :

$$\|\mathbf{h}\|_2 = \left\| \sum_{s=0}^{+\infty} \mathbf{h}^{(s)} \right\|_2$$

$$\begin{aligned}
 & \propto \sum_{s=0}^{+\infty} \|\mathbf{h}^{(s)}\|_2 \\
 & \propto \|\boldsymbol{\omega}\|_2 \|\mathbf{B}\|_2 \sum_{s=1}^{+\infty} \|\mathbf{A}\|^{s-1} \\
 & \propto \frac{B \rho \mathbf{A} \mathbf{q}}{1 - \rho \mathbf{A} \mathbf{q}},
 \end{aligned}$$

where we have exploited Assumptions 2.1 and 2.2.

We now move to term (b):

$$\begin{aligned}
 \text{(b)}^2 \quad & \left\| \tilde{\mathbf{U}}_m^T \tilde{\mathbf{g}}_m \right\|_{\tilde{\mathbf{V}}_m^{-1}}^2 = \tilde{\mathbf{g}}_m^T \tilde{\mathbf{U}}_m \tilde{\mathbf{V}}_m^{-1} \tilde{\mathbf{U}}_m^T \tilde{\mathbf{g}}_m \\
 & \propto \frac{1}{\lambda} \left\| \tilde{\mathbf{g}}_m^T \tilde{\mathbf{U}}_m \right\|_2^2 \\
 & \quad \frac{1}{\lambda} \left\| \sum_{l=1}^m \chi \tilde{\mathbf{u}}_l, \mathbf{h}^{J_{H_l+2, \infty M}} \mathbf{y} \tilde{\mathbf{u}}_l \right\|_2^2 \\
 & \propto \frac{1}{\lambda} \left( \sum_{l=1}^m \|\tilde{\mathbf{u}}_l\|_2^2 \left\| \mathbf{h}^{J_{H_l+2, \infty M}} \right\|_2 \right)^2 \\
 & \propto \frac{U^4 \|\mathbf{B}\|_2^2 \rho \mathbf{A} \mathbf{q}^2}{\lambda \rho^2 \mathbf{A} \mathbf{q}^2} \left( \sum_{l=1}^m \rho \mathbf{A} \mathbf{q}^{H_l+1} \right)^2,
 \end{aligned}$$

where we have employed the following inequality:

$$\begin{aligned}
 \left\| \mathbf{h}^{J_{H_l+2, \infty M}} \right\|_2 & \leq \left\| \boldsymbol{\omega}^T \sum_{j=H_l+2}^{+\infty} \mathbf{A}^{j-1} \mathbf{B} \right\|_2 \\
 & \propto \|\boldsymbol{\omega}\|_2 \|\mathbf{B}\|_2 \sum_{j=H_l+2}^{+\infty} \|\mathbf{A}^{j-1}\|_2 \\
 & \propto B \rho \mathbf{A} \mathbf{q} \frac{\rho \mathbf{A} \mathbf{q}^{H_l+1}}{1 - \rho \mathbf{A} \mathbf{q}}.
 \end{aligned}$$

Let us now consider term (c):

$$\begin{aligned}
 \text{(c)}^2 \quad & \left\| \tilde{\mathbf{U}}_m^T \tilde{\boldsymbol{\nu}}_m \right\|_{\tilde{\mathbf{V}}_m^{-1}}^2 = \tilde{\boldsymbol{\nu}}_m^T \tilde{\mathbf{U}}_m \tilde{\mathbf{V}}_m^{-1} \tilde{\mathbf{U}}_m^T \tilde{\boldsymbol{\nu}}_m \\
 & \propto \frac{1}{\lambda} \left\| \tilde{\mathbf{U}}_m^T \tilde{\boldsymbol{\nu}}_m \right\|_2^2 \\
 & \quad \frac{1}{\lambda} \left\| \sum_{s=1}^m \boldsymbol{\omega}^T \mathbf{A}^{H_l+1} \tilde{\mathbf{x}}_{l-1} \tilde{\mathbf{u}}_l \right\|_2^2 \\
 & \propto \frac{1}{\lambda} \left( \sum_{s=1}^m \|\boldsymbol{\omega}\|_2 \|\mathbf{A}^{H_l+1}\|_2 \|\tilde{\mathbf{x}}_{l-1}\|_2 \|\tilde{\mathbf{u}}_l\|_2 \right)^2 \\
 & \propto \frac{X^2 \|\mathbf{U}\|_2^2 \rho \mathbf{A} \mathbf{q}^2}{\lambda} \left( \sum_{l=1}^m \rho \mathbf{A} \mathbf{q}^{H_l+1} \right)^2.
 \end{aligned}$$

We now bound the summations, exploiting the inequality  $\rho \mathbf{A} \mathbf{q} \propto \bar{\rho}$ , holding by assumption:

$$\begin{aligned}
 \sum_{l=1}^m \rho \mathbf{A} \mathbf{q}^{H_l+1} & \leq \sum_{l=1}^m \rho \mathbf{A} \mathbf{q}^{\left\lfloor \frac{\log l}{\log \frac{1}{\bar{\rho}}} \right\rfloor + 1} \\
 & \propto \sum_{l=1}^m \rho \mathbf{A} \mathbf{q}^{\frac{\log l}{\log \frac{1}{\bar{\rho}}}}
 \end{aligned}$$





$$\asymp \chi \tilde{\mathbf{h}}_{m-1}^\uparrow \quad \mathbf{h}, \tilde{\mathbf{u}}_m \mathbf{y} \quad (13)$$

$$\begin{aligned} &\asymp \left\| \tilde{\mathbf{h}}_{m-1}^\uparrow \quad \mathbf{h} \right\|_{\tilde{\mathbf{v}}_{m-1}} \tilde{\mathbf{u}}_m \tilde{\mathbf{v}}_{m-1}^{-1} \\ &\asymp \left( \left\| \tilde{\mathbf{h}}_{m-1}^\uparrow \quad \tilde{\mathbf{h}}_{m-1} \right\|_{\tilde{\mathbf{v}}_{m-1}} \quad \left\| \tilde{\mathbf{h}}_{m-1} \quad \mathbf{h} \right\|_{\tilde{\mathbf{v}}_{m-1}} \right) \tilde{\mathbf{u}}_m \tilde{\mathbf{v}}_{m-1}^{-1} \end{aligned} \quad (14)$$

$$\asymp 2\tilde{\beta}_{m-1} \tilde{\mathbf{u}}_m \tilde{\mathbf{v}}_{m-1}^{-1}. \quad (15)$$

where line (13) follows from the optimism, line (14) derives from triangle inequality, line (15) is obtained by observing that  $\mathbf{h} \in \mathcal{C}_{m-1}$  with probability at least  $1 - \delta$ , simultaneously for all  $m \in [M]$ , thanks to Theorem 3.1, having observed that  $\tilde{\beta}_{m-1}$  is larger than the right hand side of Theorem 3.1.

We now move to the cumulative offline 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:

$$R^{\text{off}}_{\text{pDynLin-UCB}, T} \asymp \sum_{m=1}^M \rho H_m \quad 1 \mathbf{q} \tilde{r}_m \asymp \sqrt{\sum_{m=1}^M \rho H_m \quad 1 \mathbf{q}^2} \sqrt{\sum_{m=1}^M \tilde{r}_m^2}.$$

Concerning the first summation, we proceed as follows, recalling that  $M \asymp T$  and  $H_m \asymp H_M$  for all  $m \in [M]$ :

$$\sum_{m=1}^M \rho H_m \quad 1 \mathbf{q}^2 \asymp T \rho H_M \quad 1 \mathbf{q} \asymp T \left( 1 \quad \frac{\log T}{\log \frac{1}{\rho}} \right).$$

For the second summation, we follow the usual derivation for linear bandits, recalling that  $\tilde{\beta}_{M-1} \asymp \max\{1, \tilde{\beta}_{M-1}\mathbf{u}\}$  for all  $m \in [M]$  and that under Assumption 2.2 we have that  $\tilde{r}_m^2 \asymp 2$ . In particular:

$$\tilde{r}_m^2 \asymp \min \left\{ 2, 2\tilde{\beta}_{M-1} \tilde{\mathbf{u}}_m \tilde{\mathbf{v}}_{m-1}^{-1} \right\} \asymp 2\tilde{\beta}_{M-1} \min \left\{ 1, \tilde{\mathbf{u}}_m \tilde{\mathbf{v}}_{m-1}^{-1} \right\}.$$

Plugging this inequality into the second summation, we obtain:

$$\begin{aligned} \sum_{m=1}^M \tilde{r}_m^2 &\asymp 4\tilde{\beta}_{M-1}^2 \sum_{m=1}^M \min \left\{ 1, \tilde{\mathbf{u}}_m \tilde{\mathbf{v}}_{m-1}^{-1} \right\} \\ &\asymp 8d\tilde{\beta}_{M-1}^2 \log \left( 1 \quad \frac{MU^2}{d\lambda} \right) \asymp 8d\beta_{T-1}^2 \log \left( 1 \quad \frac{TU^2}{d\lambda} \right), \end{aligned}$$

where the last passage follows from the elliptic potential lemma (Lattimore & Szepesvári, 2020, Lemma 19.4). Putting all together, we obtain the inequality holding with probability at least  $1 - \delta$ :

$$R^{\text{off}}_{\text{pDynLin-UCB}, T} \asymp \sqrt{8dT\beta_{T-1}^2 \left( 1 \quad \frac{\log T}{\log \frac{1}{\rho}} \right) \log \left( 1 \quad \frac{TU^2}{d\lambda} \right)},$$

having observed that  $\tilde{\beta}_{M-1} \asymp \beta_{T-1}$ . We can also arrive at a problem-dependent regret bound, by setting  $\delta = \frac{1}{2}$ :  $\inf_{\mathbf{u} \in \mathcal{U}(\mathbf{h}, \mathbf{u}) < \langle \mathbf{h}, \mathbf{u}^* \rangle} \chi \mathbf{h}, \mathbf{u}^* \quad \mathbf{u} \mathbf{y}$  (if it exists  $\neq 0$ ). Since the instantaneous regret is either 0 or at least  $\frac{1}{2}$ , we have:

$$\begin{aligned} R^{\text{off}}_{\text{pDynLin-UCB}, T} &\asymp \sum_{m=1}^M \rho H_m \quad 1 \mathbf{q} \tilde{r}_m^2 \\ &\asymp \frac{H_M}{2} 8d\tilde{\beta}_{M-1}^2 \log \left( 1 \quad \frac{MU^2}{d\lambda} \right) \\ &\asymp \frac{8d}{2} \left( 1 \quad \frac{\log T}{\log \frac{1}{\rho}} \right) \beta_{T-1}^2 \log \left( 1 \quad \frac{TU^2}{d\lambda} \right). \end{aligned}$$

By setting  $\delta = \frac{1}{2}$ , replacing the value of  $\beta_{T-1}$ , we obtain the offline regret in expectation, highlighting the dependence on  $T$ ,  $\bar{\rho}$ ,  $d$ , and  $\sigma$  only:

$$\mathbb{E} R^{\text{off}}_{\text{pDynLin-UCB}, T} \asymp \mathcal{O} \left( \frac{d\sigma^2 \bar{T} \log T \mathbf{q}^{\frac{3}{2}}}{1 \quad \bar{\rho}} \quad \frac{\bar{d} \bar{T} \log T \mathbf{q}^2}{\rho \bar{\rho}^{\frac{3}{2}}} \right),$$

where we used the fact that  $\frac{1}{\log \frac{1}{\rho}} \asymp \frac{1}{1-\rho}$  and  $\rho \rho \mathbf{A} \mathbf{q} \asymp \bar{\rho}$ .  $\square$

The following lemma relates the expected offline regret with the expected online regret.

**Theorem 3.2** (Upper Bound). *Under Assumptions 2.1 and 2.2, selecting  $\beta_t$  as in Equation (6) and  $\delta = \frac{1}{\sqrt{t}}$ , DynLi n-UCB suffers an expected regret bounded as (highlighting the dependencies on  $T$ ,  $\bar{\rho}$ ,  $d$ , and  $\sigma$  only):*

$$\mathbb{E} \text{Rp}\underline{\pi}^{\text{DynLi n-UCB}}, Tq \leq \mathcal{O} \left( \frac{d\sigma \sqrt{T} \log T q^{\frac{3}{2}}}{1 - \bar{\rho}} \frac{1}{\rho \mathbf{A} q q^2} \right).$$

*Proof.* The result is simply obtained by exploiting the offline regret bound of Theorem B.2 and by upper bounding the expected regret using Lemma B.1.  $\square$

### C. 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 2.1, 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:

$$J_H(\underline{\pi}) : \mathbb{E} \left[ \frac{1}{H} \sum_{t=1}^H y_t \right] \quad \text{where} \quad \begin{cases} \mathbf{x}_{t+1} = \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t + \boldsymbol{\epsilon}_t \\ y_t = \boldsymbol{\omega}^\top \mathbf{x}_t + \eta_t \\ \mathbf{u}_t \in \mathcal{H}_{t-1} \end{cases}, \quad t \in \{1, \dots, H\}, \quad (16)$$

where the expectation is taken w.r.t. the randomness of the state noise  $\boldsymbol{\epsilon}_t$  and reward noise  $\eta_t$ . We now show that the optimal policy for the finite-horizon setting is a non-stationary open-loop policy.

**Theorem C.1** (Optimal Policy for the  $H$ -Horizon Setting). *If  $H \in \mathbb{N}$ , an optimal policy  $\underline{\pi}_H^* = \{\pi_{H,t}^*\}_{t \in \{1, \dots, H\}}$  maximizing the  $H$ -horizon expected average reward  $J_H(\underline{\pi})$  as in Equation (16) is given by:*

$$\forall t \in \{1, \dots, H\}, \forall \mathbf{h}_{t-1} \in \mathcal{H}_{t-1} : \quad \pi_{H,t}^* = \arg \max_{\mathbf{u} \in \mathcal{U}} \boldsymbol{\omega}^\top \mathbf{h}^{J_0, H-t}, \mathbf{u},$$

*Proof.* We start by expressing for every  $t \in \{1, \dots, H\}$  the reward  $y_t$  as a function of the sequence of actions  $\mathbf{u} = \{\mathbf{u}_1, \dots, \mathbf{u}_H\}$  produced by a generic policy  $\underline{\pi}$ . By exploiting Equation (4) instanced with  $H = t - 1$ , we have:

$$y_t = \sum_{s=0}^{t-1} \boldsymbol{\omega}^\top \mathbf{h}^{\{s\}}, \mathbf{u}_{t-s} + \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbf{x}_1 + \eta_t + \sum_{s=1}^{t-1} \boldsymbol{\omega}^\top \mathbf{A}^{s-1} \boldsymbol{\epsilon}_{t-s}.$$

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

$$\mathbb{E} y_t = \sum_{s=0}^{t-1} \boldsymbol{\omega}^\top \mathbf{h}^{\{s\}}, \mathbb{E} \mathbf{u}_{t-s} + \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1.$$

By averaging over  $t \in \{1, \dots, H\}$ , we obtain the  $H$ -horizon expected average reward:

$$\begin{aligned} J_H(\underline{\pi}) &= \frac{1}{H} \sum_{t=1}^H \mathbb{E} y_t \\ &= \frac{1}{H} \sum_{t=1}^H \sum_{s=0}^{t-1} \boldsymbol{\omega}^\top \mathbf{h}^{\{s\}}, \mathbb{E} \mathbf{u}_{t-s} + \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1 \\ &= \frac{1}{H} \sum_{s=1}^H \left( \sum_{t=s}^H \mathbf{h}^{\{t-s\}} \right)^\top \mathbb{E} \mathbf{u}_s + \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1 \end{aligned} \quad (17)$$

$$\frac{1}{H} \sum_{s=1}^H \boldsymbol{\omega}^\top \mathbf{h}^{J_0, H-s}, \mathbb{E} \mathbf{u}_s + \frac{1}{H} \sum_{t=1}^H \boldsymbol{\omega}^\top \mathbf{A}^{t-1} \mathbb{E} \mathbf{x}_1. \quad (18)$$

where line (17) is obtained by renaming the indexes of the summations, and line (18) comes from the definition of cumulative Markov parameter  $\mathbf{h}^{J_0, H-s}$ . It is now simple to see, as no noise is present in the expression, that the performance  $J_H(\underline{\pi})$  is maximized by taking at each round  $s \in \{1, \dots, H\}$  an action  $\mathbf{u}_s^* = \pi_s^*(\mathbf{h}_{s-1})$  such that whose expectation satisfies  $\mathbb{E} \mathbf{u}_s^* = \arg \max_{\mathbf{u} \in \mathcal{U}} \boldsymbol{\omega}^\top \mathbf{h}^{J_0, H-s}, \mathbf{u}$ . Clearly, we can take the deterministic action such that  $\mathbf{u}_s^* = \mathbb{E} \mathbf{u}_s^*$ .  $\square$

We now show that for sufficiently large  $H$ , the  $H$ -horizon expected average reward  $J_H$  tends to coincide with the infinite-horizon expected average reward.

**Proposition C.2.** *Let  $H \in \mathbb{N}$ . Then, for every policy  $\pi$  it holds that:*

$$|J_H(\rho, \pi) - J(\rho, \pi)| \leq \frac{BU}{H} \frac{\rho \mathbf{A} \mathbf{q} \mathbf{1}}{\rho \mathbf{A} \mathbf{q} \mathbf{q}}.$$

*Proof.* Consider two horizons  $H, H' \in \mathbb{N}$ , and let  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_q$  be the sequence of actions played by policy  $\pi$ . Using Equation (18), we have:

$$J_H(\rho, \pi) - J_{H'}(\rho, \pi) = \frac{1}{H} \sum_{s=1}^H \mathbf{x}^T \mathbf{h}^{J_0, H-s}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}} - \frac{1}{H'} \sum_{s=1}^{H'} \mathbf{x}^T \mathbf{h}^{J_0, H'-s}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}} \quad (19)$$

$$\frac{1}{H} \sum_{s=1}^H \mathbf{x}^T \mathbf{h}^{J_0, H-s} - \mathbf{h}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}} - \frac{1}{H'} \sum_{s=1}^{H'} \mathbf{x}^T \mathbf{h}^{J_0, H'-s} - \mathbf{h}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}} \quad (20)$$

$$\frac{1}{H} \sum_{s=1}^H \mathbf{x}^T \mathbf{h}^{J^{H-s+1, +\infty}}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}} - \frac{1}{H'} \sum_{s=1}^{H'} \mathbf{x}^T \mathbf{h}^{J^{H'-s+1, +\infty}}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}}. \quad (21)$$

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

$$\frac{1}{H'} \left| \sum_{s=1}^{H'} \mathbf{x}^T \mathbf{h}^{J^{H'-s+1, +\infty}}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}} \right| \xrightarrow{H' \rightarrow \infty} 0 \quad \text{when} \quad H' \rightarrow \infty.$$

Concerning the first addendum, we have:

$$\begin{aligned} \frac{1}{H} \left| \sum_{s=1}^H \mathbf{x}^T \mathbf{h}^{J^{H-s+1, +\infty}}, \mathbb{E} \mathbf{r}_{s, \mathbf{y}} \right| &\leq \frac{U}{H} \sum_{s=1}^H \left\| \mathbf{h}^{J^{H-s+1, +\infty}} \right\|_2 \\ &\leq \frac{BU}{H} \frac{\rho \mathbf{A} \mathbf{q}}{\rho \mathbf{A} \mathbf{q} \mathbf{q}} \sum_{s=1}^H \rho \mathbf{A} \mathbf{q}^{H-s} \\ &= \frac{BU}{H} \frac{\rho \mathbf{A} \mathbf{q} \mathbf{1}}{\rho \mathbf{A} \mathbf{q} \mathbf{q}}. \end{aligned}$$

□

## D. System Identification

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

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t + \boldsymbol{\epsilon}_t, \\ \tilde{\mathbf{y}}_t &= \mathbf{C} \mathbf{x}_t + \mathbf{z}_t. \end{aligned} \quad (22)$$

Our setting can be seen as (not strictly proper):

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t + \boldsymbol{\epsilon}_t, \\ \mathbf{y}_t &= \mathbf{C} \mathbf{x}_t + \mathbf{D} \mathbf{u}_t + \mathbf{z}_t, \end{aligned} \quad (23)$$

with  $\mathbf{x}_t, \boldsymbol{\epsilon}_t \in \mathbb{R}^n$ ,  $\mathbf{u}_t \in \mathbb{R}^p$ , and  $\mathbf{y}_t, \mathbf{z}_t \in \mathbb{R}^m$ . The noise over state transition model  $\boldsymbol{\epsilon}_t$  and output  $\mathbf{z}_t$  are  $\sigma^2$ -subgaussian random variables. We consider in this part the standard control problem notation adopted for LTI systems. The mapping to our problem is straightforward by considering  $\mathbf{C} = \boldsymbol{\omega}^T$  and  $\mathbf{D} = \boldsymbol{\theta}^T$ . In predictive form, the system described in Equation (22) is:

$$\begin{aligned} \hat{\mathbf{x}}_{t+1} &= \mathbf{A} \hat{\mathbf{x}}_t + \mathbf{B} \mathbf{u}_t + \mathbf{F} \tilde{\mathbf{y}}_t, \\ \tilde{\mathbf{y}}_t &= \mathbf{C} \hat{\mathbf{x}}_t + \mathbf{e}_t, \end{aligned}$$

where:

$$\mathbf{A} = \mathbf{A} - \mathbf{F} \mathbf{C},$$

$$\mathbf{F} = \mathbf{A}\Sigma\mathbf{C}^T\rho\mathbf{C}\Sigma\mathbf{C}^T - \sigma^2\mathbf{I}q^{-1},$$

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

$$\Sigma = \mathbf{A}\Sigma\mathbf{A}^T - \mathbf{A}\Sigma\mathbf{C}^T\rho\mathbf{C}\Sigma\mathbf{C}^T - \sigma^2\mathbf{I}q^{-1}\mathbf{C}\Sigma\mathbf{A}^T - \sigma^2\mathbf{I}.$$

In order to identify this LTI system, we want to detect a matrix  $\tilde{\mathcal{G}}_y$ :

$$\tilde{\mathcal{G}}_y = [\mathbf{C}\mathbf{F} \quad \mathbf{C}\mathbf{A}\mathbf{F} \quad \dots \quad \mathbf{C}\mathbf{A}^{H-1}\mathbf{F} \quad \mathbf{C}\mathbf{B} \quad \mathbf{C}\mathbf{A}\mathbf{B} \quad \dots \quad \mathbf{C}\mathbf{A}^{H-1}\mathbf{B}]. \quad (24)$$

To identify through least squares method matrix  $\tilde{\mathcal{G}}_y$ , we construct for each  $t$ , a vector  $\tilde{\phi}_t$ :

$$\tilde{\phi}_t = [\mathbf{y}_{t-1}^T \quad \dots \quad \mathbf{y}_{t-H}^T \quad \mathbf{u}_{t-1}^T \quad \dots \quad \mathbf{u}_{t-H}^T]^T \in \mathbb{R}^{(m+p)H}. \quad (25)$$

The system output  $\tilde{\mathbf{y}}_t$  can be rewritten as:

$$\tilde{\mathbf{y}}_t = \tilde{\mathcal{G}}_y \tilde{\phi}_t + \mathbf{e}_t - \mathbf{C}\mathbf{A}^H \mathbf{x}_{t-H}.$$

The output of the system under analysis (Equation 23) is:

$$\mathbf{y}_{t+q} = \tilde{\mathbf{y}}_t + \mathbf{D}\mathbf{u}_t + \tilde{\mathcal{G}}_y \tilde{\phi}_t + \mathbf{D}\mathbf{u}_t + \mathbf{e}_t - \mathbf{C}\mathbf{A}^H \mathbf{x}_{t-H}$$

We can incorporate the contribution of  $\mathbf{D}\mathbf{u}_t$  in  $\tilde{\mathcal{G}}_y$  obtaining  $\mathcal{G}_y$ :

$$\mathcal{G}_y = [\mathbf{C}\mathbf{F} \quad \mathbf{C}\mathbf{A}\mathbf{F} \quad \dots \quad \mathbf{C}\mathbf{A}^{H-1}\mathbf{F} \quad \mathbf{D} \quad \mathbf{C}\mathbf{B} \quad \mathbf{C}\mathbf{A}\mathbf{B} \quad \dots \quad \mathbf{C}\mathbf{A}^{H-1}\mathbf{B}].$$

The related vector  $\phi_t$  is:

$$\phi_t = [\mathbf{y}_{t-1}^T \quad \dots \quad \mathbf{y}_{t-H}^T \quad \mathbf{u}_t^T \quad \mathbf{u}_{t-1}^T \quad \dots \quad \mathbf{u}_{t-H}^T]^T \in \mathbb{R}^{(m+p)H+p}. \quad (26)$$

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

$$\hat{\mathcal{G}}_y = \arg \min_{\mathbf{X}} \lambda \|\mathbf{X}\|_F^2 + \sum_{\tau=t-H}^t \|\mathbf{y}_\tau - \mathbf{X}\phi_\tau\|_2^2, \quad (27)$$

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

The matrix  $\mathbf{D}$  can be directly retrieved from  $\hat{\mathcal{G}}_y$ . In order to get matrices  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$ , we remove the values related to  $\mathbf{D}$  from  $\hat{\mathcal{G}}_y$  and we retrieve  $\tilde{\mathcal{G}}_y$ . From now on, we refer to the algorithm proposed in Lale et al. (2020a, Appendix B).

## E. Integration on Numerical Simulations

This section is divided in three parts. First, in Section E.1, we provide additional information about the baselines, their hyperparameters and the optimistic bounds. Second, in Section E.2, we provide all the matrices and vectors generalized to run the real-world experiment. Third, in Section E.3, we provide further results for the simulations presented in Section 5.1.

### E.1. Additional Notes on the Baselines

As mentioned in Section 5, the chosen baselines are Lin-UCB (Abbasi-Yadkori et al., 2011), D-Lin-UCB (Russac et al., 2019), AR2 (Chen et al., 2021), Exp3 (Auer et al., 1995), Exp3-k (Dekel et al., 2012; Auer et al., 1995) and the Expert (the latter available only in the case of real-world data). All the hyperparameters, whenever possible, are set as prescribed in the original papers. The bounds used for the exploration are adjusted in order to be able to fairly compete in this setting, and are considered as follows:

$$\beta_t^{\text{Lin-UCB}} : \bar{c}_2 \sqrt{\lambda} \sqrt{2\bar{\sigma}^2 \left( \log\left(\frac{1}{\delta}\right) + \frac{d}{2} \log\left(1 + \frac{tU^2}{d\lambda}\right) \right)},$$

$$\beta_t^{\text{D-Lin-UCB}} : \bar{c}_2 \sqrt{\lambda} \sqrt{2\bar{\sigma}^2 \left( \log\left(\frac{1}{\delta}\right) + \frac{d}{2} \log\left(1 + \frac{tU^2}{d\lambda} \left(\frac{1}{1-\gamma^2}\right)\right) \right)},$$

where  $\bar{c}_2$  and  $\bar{\sigma}^2$  are as prescribed in Section 3.2, and the hyperparameter  $\gamma$  of D-Lin-UCB is tuned.

For AR2, the hyperparameter  $\alpha$ , describing the correlation over time is considered equal to  $\rho\mathbf{p}\mathbf{A}q$ .

In the case of Exp3, the rewards are rescaled in order to make them range in  $[0, 1]$  with high probability, as follows:

$$\bar{r}_t = \frac{r_t - 2\xi}{4\xi}, \quad \text{where} \quad \xi = \left( \frac{B}{1 - \rho\mathbf{p}\mathbf{A}q} \right) U.$$



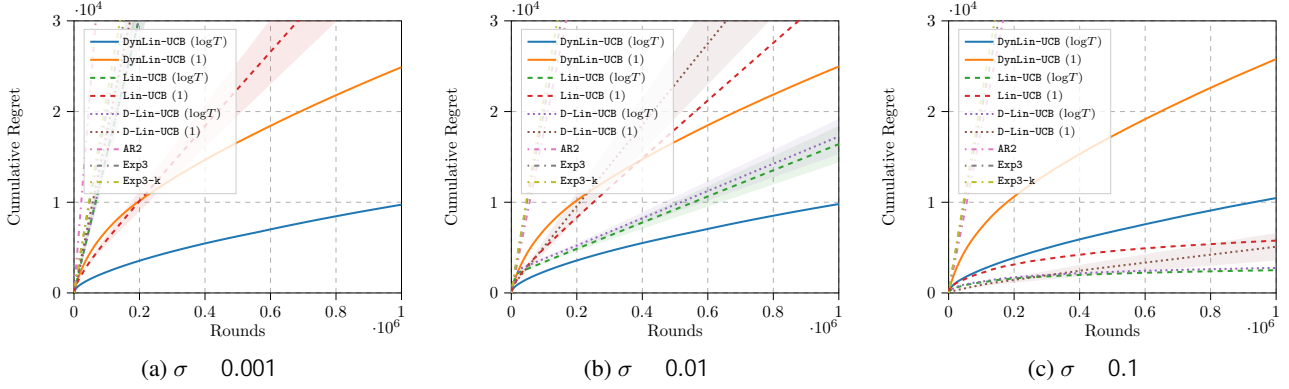


Figure 4. Performance of DynLin-UCB, Lin-UCB, D-Lin-UCB, AR2, Exp3 and Exp3-k at different values of  $\sigma$ . (50 runs, mean std)

Furthermore, in the case of Exp3-k, the batch dimension  $k$  is considered as:

$$k = \left\lceil \frac{\log M}{\log 1/\bar{\rho}} \right\rceil,$$

where  $M$  is the one defined in Algorithm 1 (line 2). This batch size  $k$  ensures that, at each time  $t$ , the contribution of actions  $\mathbf{u}_s$  is negligible, with  $s \in \{1, \dots, k\}$ . The rewards collected in the same batch are averaged and transformed as in Exp3.

## E.2. Further Information on the Real-world Setting

The real-world setting is generalized through a dataset containing real data related to the budgets invested in each advertising platform (i.e., the  $\mathbf{u}_t$ ) and the overall generated conversions (i.e., the  $y_t$ ) collected 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 5 Million USD over 2 years. Starting from such data, we generalized the best model by means of a specifically designed variant of the Ho-Kalman algorithm (Ho & Kalman, 1966), as described in Appendix D. We used the matrices estimated with Ho-Kalman to build up a simulator. The resulting system has  $\rho \mathbf{A} \mathbf{q} = 0.67$ , and is characterized as follows:

$$\mathbf{A} = \begin{pmatrix} 0.38 & 0.33 & 0.6 \\ 0.07 & 0.76 & 0.54 \\ 0.18 & 0.34 & 0.05 \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} 0.14 & 0.34 & 0.05 \\ 0.17 & 0.03 & 0.01 \\ 0.04 & 0.09 & 0.17 \end{pmatrix}, \quad \boldsymbol{\omega} = \begin{pmatrix} 0.61 \\ 0.04 \\ 0.13 \end{pmatrix}, \quad \boldsymbol{\theta} = \begin{pmatrix} 0.13 \\ 0.41 \\ 0.02 \end{pmatrix}.$$

## E.3. Additional Numerical Simulations

These additional results are obtained in the setting presented in Section 5.1. However, here, we want to analyze the behavior of DynLin-UCB and the other bandit baselines at different magnitudes of noise in both the state transition model and the output. The noise in this simulation is a zero-mean Gaussian noise with  $\sigma \in \{0.001, 0.01, 0.1\}$ .

**Results** Figure 4 shows the results of the experiment for the different values of  $\sigma$ . It is clearly visible how DynLin-UCB performs in almost the same way no matter the noise to which the system is subject, always leading to sub-linear regret. On the other hand, the cumulative regret of both Lin-UCB and D-Lin-UCB is different in every simulation we perform. Indeed, with a low level of noise (Figure 4a) reaches linear regret and does not converge, while for large values of noise, it converges very quickly (Figure 4c). This is due to the nature of the confidence bound of linear bandits, which is not able to take into account such a complex scenario and leads to no guarantees in this setting. Exp3, Exp3-k, and AR2 are not able to reach the optimum in this scenario, independently from the noise magnitude  $\sigma$ , and provide large values of (linear) regret.

## E.4. Computational Time

The code used for the results provided in this section has been run on an Intel(R) i5 8259U @ 2.30GHz CPU with 8 GB of LPDDR3 system memory. The operating system was macOS 12.2.1, and the experiments have been run on Python 3.9.7. A single run of DynLin-UCB takes 110 seconds to run. It is worth noting that the time complexity of DynLin-UCB is upper-bounded by the one of Lin-UCB.