
Long-Memory AutoRegressive Bandits

Uladzimir Charniauski

Department of Statistics
University of Connecticut
Storrs, CT 06286
uladzimir.charniauski@uconn.edu

Yao Zheng

Department of Statistics
University of Connecticut
Storrs, CT 06286
yao.zheng@uconn.edu

Abstract

The Autoregressive Fractionally Integrated (ARFIMA) processes naturally occur in the context of real-world scenarios that exhibit long memory properties. The construct of ARFIMA process allows us to easily formulate a separate class of stochastic multi-armed bandits problem (Abbasi-Yadkori et al. [2011]) by modifying its general mechanism. In this work, we introduce a novel setting named Long-Memory AutoRegressive Bandits (LM-ARBs), where the environment-generated reward evolves according to the fractionally integrated autoregressive process of the autoregressive order p and the memory parameter d , an extension of the AutoRegressive Bandits (Bacchiocchi et al. [2024]) characterized only by the autoregressive process. Then, we provide an optimistic regret-minimization algorithm Long-Memory AutoRegressive Upper Confidence Bound (ARLM-UCB) that suffers a sublinear regret of order $\mathcal{O}\left(\frac{(p+1)^2 \sqrt{n} T^{2d+0.5} \log^2(T) \log \log(T)}{(1-\Gamma)^2}\right)$, where n is the number of actions, T is the optimization horizon, and $\Gamma < 1$ is a stability index of the fractionally integrated process. Finally, we conduct numerical experiments in synthetic environments to validate our algorithm effectiveness w.r.t. several bandit baselines.

1 Introduction

Many real-world sequential decision-making problems require the learner to select an action that determines a long-term reward evolution, creating a temporal dependence for future rewards over a long-range horizon. When analyzing this reward, the agent must account for a much slower decay of temporal dependence between the current reward and the sequence of past observations. Autoregressive Fractionally Integrated Moving-Average (ARFIMA) (Hosking [1981], McLeod and Hipel [1986], Granger and Joyeux [1980]) is widely used to model the long memory in the real-world phenomena, such as stock market volatility, temperature variations, earthquake magnitude sequences, and traffic data (Bakar and Hafner [2019], Yuan et al. [2014], Kondo Lembang et al. [2021], Doodipala [2020]). In the context of Reinforcement Learning (Sutton and Barto [2020]), this method flexibly allows one to model the long-term behavior of reward sequences. For example, the learner’s ability to capture long-range dependence in stock price return dynamics enables more accurate forecasting of future trends, volatility clustering, and regime shifts that are not evident in short memory models (Liu [2000]).

In this paper, we model the reward of a decision-making process as an ARFIMA process, whose parameters depend on the action selected by the agent at every round. This scenario can be viewed as a separate class of stochastic bandit algorithms (Abbasi-Yadkori et al. [2011]), where the temporal structure of the reward is governed by the fractionally integrated autoregressive process whose action-dependent parameters are unknown to the agent. In this setting, the agent faces a multi-staged challenge of estimating the ARFIMA parameters responsible for generating the reward in the given environment. Given such a complex learning process, this scenario displays remarkable

differences to more traditional non-stationary learning problems. Indeed, it is worth mentioning that the environment does not change by any exogenous sources of non-stationarity, which is often represented by a smooth changing in the mean of rewards for each arm over time, studied by [Trella et al. \[2024\]](#) in the bandit context.

1.1 Original Contribution

In this work, we propose a novel setting named Long-Memory AutoRegressive Bandits (LM-ARBs), in which the reward follows an ARFIMA($p, d, 0$) process, where p stands for AR order and the memory parameter d . We then devise a new optimistic algorithm AutoRegressive Long-Memory Upper Confidence Bound (ARLM-UCB) to learn an optimal policy in this online settings and show that it suffers a sublinear regret of order $\mathcal{O}\left(\frac{(p+1)^2 \sqrt{n} T^{2d+0.5} \log^2(T) \log \log(T)}{(1-\Gamma)^2}\right)$, where n is the number of actions, T is the optimization horizon and $\Gamma < 1$ is a stability index of the autoregressive fractionally integrated process. In the end, we empirically evaluate ARLM-UCB and compare its performance with other bandit baselines in our setting, illustrating that our proposed algorithm outperforms a number of popular multi-armed bandit (MAB) benchmarks that do and do not account for the temporal dynamics in the reward evolution process.

1.2 Related Work

This work proposes a modification of a traditional MAB learning problem by incorporating long-range temporal dependence into the reward evolution process. Many related work on creating temporal dynamics in bandit settings focused on addressing challenges like delayed feedback ([Tang et al. \[2024\]](#)), influence of past actions on future rewards ([Tang et al. \[2021\]](#)), or the non-stationarity of AR dynamics ([Chen et al. \[2023\]](#)). [Bacchiocchi et al. \[2024\]](#) presented AutoRegressive Bandits (ARBs) setting. This setting explicitly models autoregressive (AR) sequences, where the reward depends on several past observations of the length of a finite AR order p , which eventually makes this method neglect the long-range dependence of rewards over extended horizons.

Limited studies have addressed MAB settings with long-range temporal dependence between rewards. A recent study by [Qin et al. \[2023\]](#) introduces a framework for contextual bandits, where rewards depend on long-range temporal dependence between past actions and contexts. The major limitation of this method hinders in the assumption of sparsity in the reward structure, where only a finite number of contexts, much smaller from their total number, influence the current reward, which may not be realistic for many real-world settings. For instance, in scenarios where rewards depend on dense or complex long-range temporal patterns (e.g., cumulative effects across many past contexts), this assumption may fail to capture the full dependency structure, loosening the ability to model richer temporal dynamics.

The ARFIMA process has been well studied in the classical time series literature. Asymptotic theory for ARFIMA estimation was developed by [Dahlhaus \[1989\]](#) for maximum likelihood methods and [Fox and Taqqu \[1986\]](#) for general long memory processes. In the frequency domain, [Robinson \[1995\]](#) made seminal contributions by proposing Gaussian semiparametric estimators of the memory parameter and establishing their asymptotic properties, while [Beran \[1995\]](#) advanced time domain approaches including maximum likelihood estimation with rigorous asymptotic theory. However, only a few studies have addressed MAB settings with long-range temporal dependence between rewards. The key limitation of many existing machine learning methods is their inability to analyze temporal dependence on a long-range horizon of past observations. This notion was particularly validated by different studies ([Al-Selwi et al. \[2023\]](#), [Zucchet et al. \[2023\]](#)), demonstrating that traditional sequential modeling algorithms do not learn long-range dependence.

[Gupta et al. \[2021\]](#) proposed fractional dynamical systems with long-range filtering operations on state vectors, which closely relate to ARFIMA processes. In the context of deep learning, [Zhao et al. \[2020\]](#) have proposed a modification to a traditional recurrent neural network (RNN) architecture that enables capturing long memory from a time series perspective via new memory filter component directly incorporating ARFIMA process. Motivated by successes in extending machine learning frameworks with long memory processes, we directly implement ARFIMA modeling in our bandit setting.

2 Problem Formulation

In this section, we briefly discuss the formal representation of the ARFIMA(p, d, q) process in terms of the parameters used throughout this document and introduce the LM-ARB setting (Section 2.1). Subsequently, we establish several essential assumptions for convergent evolution of the utilized ARFIMA process (Section 2.2), present the closed-form solution for the optimal policy in our setting (Section 2.3), and describe the stochastic properties of the AR reward process (Section 2.4).

Notation. We will employ the following notation across the paper. Let $a, b \in \mathbb{N}$, with $a \leq b$, we introduce the symbols: $\llbracket a, b \rrbracket := \{a, \dots, b\}$ and $\llbracket b \rrbracket := \{1, \dots, b\}$. Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ be real-valued vectors, we denote with $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i$ the inner product. For a positive semi-definite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, we denote $\|\mathbf{x}\|_{\mathbf{A}}^2 = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ the weighted 2-norm. We define a zero-mean random variable ε is σ^2 -subgaussian if $\mathbb{E}[e^{\lambda \varepsilon}] \leq e^{\frac{\lambda^2 \sigma^2}{2}}$.

2.1 Long-Memory Autoregressive Bandits

The AutoRegressive Fractionally Integrated Moving-Average (ARFIMA) process characterizes long-range dependence in time series data. The ARFIMA(p, d, q) process $\{X_t, t \in \mathbb{Z}\}$ is represented through the following form:

$$(1 - B)^d (1 - \sum_{i=1}^p \phi_i B^i) X_t = (1 + \sum_{i=1}^q \theta_i B^i) \varepsilon_t, \quad (1)$$

where ϕ_i, θ_j , ($i \in \llbracket p \rrbracket$ and $j \in \llbracket q \rrbracket$) are the coefficients of the model, B is the backshift operator defined as $B^j X_t = X_{t-j}$ for $j \in \mathbb{N}$, ε_t is the zero-mean *i.i.d.* σ^2 -subgaussian error term of the system, and $(1 - B)^d = \sum_{j=0}^{\infty} \psi_j(d)$, where $\psi_j(d) = \prod_{i=1}^j \frac{i-1-d}{i}$ for $j \geq 1$, with $\psi_0(d) \equiv 1$, are the *decay coefficients* defined for $j \in \mathbb{N}$. We employ this process representation across this paper to create a long memory reward evolution dynamics within our environment. For simplicity, we consider $q = 0$ in our proposed LM-ARB setting, with the potential to be further generalized.

We denote x_t to be an AR process of order p , which is represented in the form of [Bacchiocchi et al. \[2024\]](#):

$$x_t = \phi_0(u_t) + \sum_{i=1}^p \phi_i(u_t) x_{t-i} + \varepsilon_t = \langle \phi(u_t), \mathbf{z}_{t-1} \rangle + \varepsilon_t, \quad (2)$$

where $\phi(u) := (\phi_0(u), \dots, \phi_p(u))^\top \in \mathbb{R}^{p+1}$ is the *parameter vector* containing the the *unknown parameters* $(\phi_i(u_t))_{i \in \llbracket p \rrbracket}$ depending on the choice of an action u_t , $\mathbf{z}_{t-1} = (1, x_{t-1}, \dots, x_{t-p})^\top \in \mathcal{Z} := \{1\} \times \mathcal{X}^p$ ($\mathcal{X} \subseteq \mathbb{R}$ is the reward space) is the *vector of past estimated rewards* expressing the past history of observed rewards, and ε_t is a random *i.i.d.* zero mean σ^2 sub-Gaussian noise, conditioned to the past.

We introduce a novel Long-Memory Autoregressive Bandits (LM-ARB) setting, where at each round $t \in \mathbb{N}$, the environment generates a noisy LM reward y_t that evolves according to the ARFIMA($p, d, 0$) process of the following form:

$$y_t = (1 - B)^d x_t = (1 - B)^d (\phi_0(u_t) + \sum_{i=1}^p \phi_i(u_t) x_{t-i} + \varepsilon_t). \quad (3)$$

where $y_t \in \mathcal{Y}$, with the reward space $\mathcal{Y} \subseteq \mathbb{R}$, and d is an unknown memory parameter, whose true value is stored within the environment for reward generation. After observing the reward y_t , the learner estimates a memory parameter (learning rate) \hat{d}_t to convert y_t to an approximate short memory AR(p) reward in the following way:

$$\tilde{x}_t = (1 - B)^{-\hat{d}_t} y_t = \sum_{j=0}^{\infty} \psi_j(-\hat{d}_t) y_{t-j}. \quad (4)$$

The value of \hat{d}_t is estimated by the agent through a defined loss-minimization mechanism, which we will later introduce in Section 3.

In this way, the agent approximates the *short memory* autoregressive reward \tilde{x}_t from an infinite sequence of past observations of long memory rewards (y_t, y_{t-1}, \dots) , reducing the setting to AR(p) to analyze the behavior of the system. In practice, the infinite-sum approximation must be reduced to

a finite window size for the computational efficiency. Furthermore, in this setting, we presuppose that all the rewards played prior to the first round $t = 1$ are zero. Therefore, in our study, we truncate the infinite summation with the number of rounds played $t \in \mathbb{N}$, which gives us the following reward approximation we use across our learning horizon:

$$\hat{x}_t = \sum_{j=0}^t \psi_j(-\hat{d}_t) y_{t-j}. \quad (5)$$

2.2 Assumptions

We introduce the following assumptions, which we will utilize across the paper, and comment on their roles.

Assumption 2.1. (Non-negativity). $c \leq \phi_i(u)$ for every $u \in \mathcal{U}, i \in \llbracket p \rrbracket$ and $c \in (0, 1)$.

Assumption 2.2. (Stability). $\max_{u \in \mathcal{U}} \sum_{i=1}^p \phi_i(u) \leq \Gamma$ for $\Gamma < 1$.

Assumption 2.3. (Boundedness). $\max_{u \in \mathcal{U}} \phi_0(u) \leq m$ for $m \in (0, \infty)$

Assumption 2.4. (Long-memoryness). $0 < d < 0.5$

The Assumption 2.1 enforces the non-negativity of AR coefficients. Many real-world processes (i.e., pricing, stock markets, temperature anomalies etc.) are characterized by this assumption, where processes violating such will generate unrealistic and counterintuitive behaviors. Assumption 2.2 ensures that the sum of $(\phi_i(u))_{i \in \llbracket p \rrbracket}$ is bounded to a value $\Gamma \in [0, 1)$, and Assumption 2.3 enforces the boundedness on $\phi_0(u)$ and the sequence of rewards of the environment. These latter assumptions guarantee the stability of the considered ARFIMA process, preventing it from diverging for any action sequence played. Finally, Assumption 2.4 is a necessary requirement that enables the process to model long memory temporal sequences (Box et al. [2015]).

2.3 Policy and Regret:

We model the learner's behavior by a deterministic policy $\pi = (\pi_t)_{t \in \mathbb{N}}$, defined for every round $t \in \mathbb{N}$ as $\pi : \mathcal{H}_{t-1} \rightarrow \mathcal{U}$ that maps the history of observations $H_{t-1} := (\hat{x}_0, u_1, \hat{x}_1, \dots, u_{t-1}, \hat{x}_{t-1}) \in \mathcal{H}_{t-1}$ to an action $u_t = \pi(H_{t-1}) \in \mathcal{U}$, where $\mathcal{H}_{t-1} := \mathcal{X} \times (\mathcal{U} \times \mathcal{X})^{-1}$ is the set of length histories $t - 1$. The performance of a policy is evaluated in terms of the expected cumulative estimated reward over the finite horizon $T \in \mathbb{N}$:

$$J_T(\pi) = \mathbb{E}[\sum_{t=1}^T \hat{x}_t], \quad (6)$$

where the expectation is taken w.r.t. the randomness of the reward noise ε_t . The goal of the learner is to *minimize the expected cumulative (policy) regret* by playing a policy π , competing against the optimal policy π^* on a learning horizon $T \in \mathbb{N}$:

$$\hat{R}(\pi, T) = J^* - \mathbb{E}[\sum_{t=1}^T \hat{x}_t] = \mathbb{E}[\sum_{t=1}^T \hat{r}_t], \quad (7)$$

where $\hat{r}_t := x_t^* - \hat{x}_t$ is the instantaneous policy regret and $(x_t^*)_{t \in \llbracket T \rrbracket}$ is the sequence of short memory AR rewards observed playing the optimal policy π^* .

In the LM-ARB setting, because the converted analyzed reward \hat{x}_t is of AR(p) process exhibiting short memory behavior, we employ the definition of the optimal policy of Bacchiocchi et al. [2024] expressed as follows:

Theorem 2.5. (Optimal Policy) Under Assumptions 2.1, an optimal policy π^* maximizing the expected reward $J_T(\pi)$, for every round $t \in \mathbb{N}$ and history $H_{t-1} \in \mathcal{H}_{t-1}$ is given by:

$$\pi_t^*(H_{t-1}) \in \arg \max_{u \in \mathcal{U}} \langle \phi(u), \hat{\mathbf{z}}_{t-1} \rangle. \quad (8)$$

Some important comments on the implementation of this theorem in our setting come in order. First, the optimal action depends on the vector past reconstructed AR rewards $\hat{\mathbf{z}}_{t-1}$ only. Thus, we preserve the Markovian property of the reward policy π^* by representing LM-ARB as a Markov Decision Process (MDP) with the state representation $\hat{\mathbf{z}}_{t-1} = (1, \hat{x}_{t-1}, \dots, \hat{x}_{t-p})$ (Puterman [1994]), defined in the same way as in the ARB setting with true AR reward process $\{x_t\}$. On the other hand, defining the optimal policy in terms of a sequence of environment-generated rewards $(y_t)_{t \in \mathbb{N}}$ will create additional challenges in policy evaluation due to the complex structure of the reward-generating ARFIMA(p, d, 0) process. Second, in every round $t \in \mathbb{N}$, the optimal action maximizes the instantaneous reward expected $\mathbb{E}[\hat{x}_t | H_{t-1}] = \langle \phi(u), \hat{\mathbf{z}}_{t-1} \rangle$. This is because the Assumption 2.1 establishes the non-negativity of parameters, ensuring the meaningful evolution of the ARFIMA process, compatible with real-world settings. In this way, the optimal action maximizes both the expected immediate reward (i.e., myopic policy) and the expected cumulative reward.

2.4 On the Stochastic Properties of the AR Reward Process x_t

Estimating the memory parameter \hat{d}_t by the agent requires strong theoretical guarantees for the asymptotic convergence of this rate to the true parameter. For such convergence guarantees to hold, the true short-memory reward process must satisfy all necessary stochastic properties. We encapsulate the important properties of x_t in the following theorem:

Theorem 2.6. (*Geometric Ergodicity*) Under Assumptions 2.2 and 2.3, the AR reward process x_t is strictly stationary and geometrically ergodic.

This theorem presents a required condition about x_t when showing the asymptotic convergence of the differencing rate estimates:

Theorem 2.7. (*Asymptotic Convergence*) Provided that Theorem 2.6 and Assumptions 2.1-2.4 hold, the asymptotic convergence of the estimated learning rate \hat{d}_t to the true rate d is guaranteed by McAleer and Ling [2008] at the following rate:

$$\hat{d}_t = d + \mathcal{O}\left(\sqrt{\frac{\log \log(t)}{t}}\right) \quad a.s., \quad (9)$$

where *a.s.* represents the convergences in an almost sure sense.

A particular comment deserves the relevance of a non-zero lower bound for AR coefficients made in Assumption 2.1 to 2.7. The reward evolution process in the addressed setting is conditioned by persistent time-varying AR (TVAR) coefficients (Jiang [2023]), whose configuration depends on the arm played in the current round $t \in \mathbb{N}$. The nonzero lower bound condition guaranteed by Assumption 2.1 for every $(\phi_i(u))_{i \in [p]}$ preserves the invertibility and regularity conditions required for consistent long-run inference. This property allows the estimated memory parameter \hat{d}_t to converge asymptotically to its true value d of the environment over the entire horizon T . The complete proofs for Theorems 2.6 and 2.7 are outlined in Appendix B.

3 Algorithm

In this section, we present the algorithm AutoRegressive Long-Memory Upper Confidence Bound (ARLM-UCB), an optimistic regret-minimizing algorithm for the LM-ARB setting whose pseudocode is presented in Algorithm 1. ARLM-UCB leverages the myopic optimal policy defined in Theorem 2.5 implements an incremental regularized least squares procedure across the given grid of memory parameter values $G_{\hat{d}}$ to estimate the unknown parameters d and then $\phi(u)$ for every action $u \in \mathcal{U}$ independently. The algorithm requires knowledge of the autoregressive order p of the reward-governing process.

The algorithm is based on the following procedure. ARLM-UCB starts by initializing for every action $u \in \mathcal{U}$ the Gram matrix $\hat{\mathbf{V}}_0(u) = \lambda \mathbf{I}_{p+1}$, where $\lambda > 0$ is the Ridge regularization parameter, the vectors $\hat{\mathbf{b}}_0(u) = \hat{\phi}_0(u) = \mathbf{0}_{p+1}$ and the vector of estimated short memory observations $\hat{\mathbf{z}}_0 = (1, 0, \dots, 0)^\top$ to store estimated short memory rewards \hat{x}_t converted from the long memory (LM) rewards y_t generated by the environment. The algorithm also initializes the grid of long memory parameters $G_{\hat{d}} := \{0.01, 0.02, \dots, 0.49\}$, s.t. the true $d \in G_{\hat{d}}$, and the decay coefficients as $\hat{\psi}_j(-\hat{d}) = \prod_{i=1}^j \frac{i-1+\hat{d}}{i}$ and $\psi_0 = 1$ for every $\hat{d} \in G_{\hat{d}}$ utilized for the reward conversion process. Additionally, the algorithm initializes for every $u \in \mathcal{U}$ and $\hat{d} \in G_{\hat{d}}$ the internal Gram matrix $\hat{\mathbf{V}}_0(u, \hat{d}) = \lambda \mathbf{I}_{p+1}$ and internal vectors $\hat{\mathbf{b}}_0(u, \hat{d}) = \hat{\phi}_0(u, \hat{d}) = \mathbf{0}_{p+1}$ and $\hat{\mathbf{z}}_0(\hat{d}) = (1, 0, \dots, 0)^\top$ used for a parameter search procedure on $G_{\hat{d}}$.

Then, for every round $t \in \mathbb{N}$, the algorithm computes the *Upper Confidence Bound* index (line 4) defined for every $u \in \mathcal{U}$ as follows:

$$\begin{aligned} u_t \in \arg \max_{u \in \mathcal{U}} \text{UCB}_t(u) &:= \langle \hat{\phi}_{t-1}(u), \hat{\mathbf{z}}_{t-1} \rangle + \mathcal{B}_{\hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1}} \|\hat{\phi}_{t-1}(u)\|_{\hat{\mathbf{V}}_{t-1}^{-1}(u)} \\ &+ \beta_{t-1}(u) \left(\|\hat{\mathbf{z}}_{t-1}\|_{\hat{\mathbf{V}}_{t-1}^{-1}(u)} + \mathcal{B}_{\hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1}} \right). \end{aligned} \quad (10)$$

where $\mathcal{B}_{\hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1}} = \sqrt{p} t^{\hat{d}_{t-1}-0.5} \log(t) \sqrt{\log(\log(t) + 1)}$ arises from the asymptotic boundedness of $\hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1}$. Similarly to Lin-UCB (Abbasi-Yadkori et al. [2011]), the index $\text{UCB}_t(u)$ is designed

Algorithm ARLM-UCB

- 1: **Input:** regularization parameter $\lambda > 0$, grid of memory parameters $G_{\hat{d}}$, exploration coefficient $(\beta_{t-1})_{t \in \mathbb{N}}$.
 - 2: **Initialize:** $t \leftarrow 1$, $\hat{\mathbf{V}}_0(u) = \mathbf{V}_0(u, \hat{d}) = \lambda \mathbf{I}_{p+1}$, $\hat{\mathbf{b}}_0(u) = \mathbf{b}_0(u, \hat{d}) = \mathbf{0}_{p+1}$, $\hat{\phi}_0(u) = \phi_0(u, \hat{d}) = \mathbf{0}_{p+1}$, $\hat{\mathbf{z}}_0 = \hat{\mathbf{z}}_0(\hat{d}) = (1, 0, \dots, 0)^\top$, $\hat{\psi}_j(-\hat{d}) = \prod_{i=1}^j \frac{i-1+\hat{d}}{i}$
 - 3: **for** $t \in \llbracket T \rrbracket$ **do**
 - 4: Compute $u_t \in \arg \max_{u \in \mathcal{U}} \text{UCB}_t(u) := \langle \hat{\phi}_{t-1}(u), \hat{\mathbf{z}}_{t-1} \rangle + \mathcal{B}_{\hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1}} \|\hat{\phi}_{t-1}(u)\|_{\hat{\mathbf{V}}_{t-1}^{-1}(u)} + \beta_{t-1}(u) \left(\|\hat{\mathbf{z}}_{t-1}\|_{\hat{\mathbf{V}}_{t-1}^{-1}(u)} + \mathcal{B}_{\hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1}} \right)$,
 where $\mathcal{B}_{\hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1}} = \sqrt{p} t^{\hat{d}_{t-1} - 0.5} \log(t) \sqrt{\log(\log(t) + 1)}$
 - 5: Play action u_t and observe the LM reward y_t
 - 6: **for** $\hat{d} \in G_{\hat{d}}$ **do**
 - 7: Compute $x_t(\hat{d}) = \sum_{j=0}^t \psi_j(-\hat{d}) y_{t-j}$ and $\mathbf{z}_{t-1}(\hat{d}) = (1, x_{t-1}(\hat{d}), \dots, x_{t-p}(\hat{d}))^\top$
 - 8: **for** $u \in \mathcal{U}$ **do**
 - 9: $\hat{\mathbf{V}}_t(u, \hat{d}) = \hat{\mathbf{V}}_{t-1}(u, \hat{d}) + \mathbf{z}_{t-1}(\hat{d}) \mathbf{z}_{t-1}^\top(\hat{d}) \mathbb{1}_{\{u=u_t\}}$
 - 10: $\hat{\mathbf{b}}_t(u, \hat{d}) = \hat{\mathbf{b}}_{t-1}(u, \hat{d}) + \mathbf{z}_{t-1}(\hat{d})^\top x_t(\hat{d}) \mathbb{1}_{\{u=u_t\}}$
 - 11: $\hat{\phi}_t(u, \hat{d}) = \hat{\mathbf{V}}_t^{-1}(u, \hat{d}) \hat{\mathbf{b}}_t(u, \hat{d})$
 - 12: Compute loss $L_t(\hat{d}) = \sum_{u \in \mathcal{U}} \sum_{i=1}^t \{x_{t-i+1}(\hat{d}) - \langle \phi_{t-i+1}(u, \hat{d}), \mathbf{z}_{t-i}(\hat{d}) \rangle\}^2 \mathbb{1}_{\{u=u_t\}}$
 - 13: Choose $\hat{d}_t = \arg \min_{\hat{d} \in G_{\hat{d}}} L_t(\hat{d})$ and $\hat{x}_t = x_t(\hat{d}_t)$
 - 14: **for** $u \in \mathcal{U}$ **do**
 - 15: $\hat{\mathbf{V}}_t(u) = \hat{\mathbf{V}}_{t-1}(u) + \hat{\mathbf{z}}_{t-1} \hat{\mathbf{z}}_{t-1}^\top \mathbb{1}_{\{u=u_t\}}$
 - 16: $\hat{\mathbf{b}}_t(u) = \hat{\mathbf{b}}_{t-1}(u) + \hat{\mathbf{z}}_{t-1} \hat{x}_t \mathbb{1}_{\{u=u_t\}}$
 - 17: $\hat{\phi}_t(u) = \hat{\mathbf{V}}_t^{-1}(u) \hat{\mathbf{b}}_t(u)$
 - 18: Update $\hat{\mathbf{z}}_t = (1, \hat{x}_t, \dots, \hat{x}_{t-p+1})^\top$
 - 19: $t \leftarrow t + 1$
-

to be optimistic, i.e., $\langle \hat{\phi}_{t-1}(u), \hat{\mathbf{z}}_{t-1} \rangle \leq \text{UCB}_t(u)$ in high probability for every $u \in \mathcal{U}$. The agent then plays an optimistic action $u_t \in \arg \max_{u \in \mathcal{U}} \text{UCB}_t(u)$ and observes the reward y_t (line 5) of the form as in Equation 3.

Once the LM reward y_t is observed, the agent calculates a short memory converted reward \hat{x}_t and later the parameter vector $\hat{\phi}_t(u)$ on a grid of values of memory parameters $G_{\hat{d}}$ using the following minimalized squared distance loss procedure. For each $\hat{d} \in G_{\hat{d}}$, the agent estimates the short memory reward using a truncated infinite series ARFIMA sum $\hat{x}_t = (1 - B)^{-\hat{d}} y_t = \sum_{j=0}^t \psi_j(-\hat{d}) y_{t-j}$ (line 7). The agent then continues with updating the Gram matrix estimate $\mathbf{V}_t(u, \hat{d})$, the vector $\mathbf{b}_t(u, \hat{d})$, and the estimate $\phi_t(u, \hat{d})$ (lines 9-11). Finally, using all previous short memory samples up to round $t \in \mathbb{N}$ ($\hat{x}_i(\hat{d})_{i \in \mathbb{N}}$), and previous short memory observation vectors ($\hat{\mathbf{z}}_i(\hat{d})_{i \in \mathbb{N}}$) and estimates ($\hat{\phi}_i(u, \hat{d})_{i \in \mathbb{N}}$), the agent estimates the loss function (line 12) for the selected \hat{d} and all $u \in \mathcal{U}$ defined as follows:

$$L_t(\hat{d}) = \sum_{u \in \mathcal{U}} \sum_{i=1}^t \{x_{t-i+1}(\hat{d}) - \langle \phi_{t-i+1}(u, \hat{d}), \mathbf{z}_{t-i}(\hat{d}) \rangle\}^2 \mathbb{1}_{\{u=u_t\}} \quad (11)$$

The estimate of the memory parameter \hat{d}_t and the estimated short memory reward sample \hat{x}_t that minimize the loss function $L_t(\hat{d})$ are selected at the end of a search over a grid of memory parameters $G_{\hat{d}}$ in line 13. Using \hat{x}_t , the agents proceeds to compute $\hat{\mathbf{V}}_t(u)$, the vector $\hat{\mathbf{b}}_t(u)$, and the estimate $\hat{\phi}_t(u)$ for the current round (lines 15-17). The observation vector is then updated with a new reward sample as $\hat{\mathbf{z}}_t = (1, \hat{x}_t, \dots, \hat{x}_{t-p+1})^\top$ (line 18), so that it can be used in the next round along with the estimate $\hat{\phi}_t(u)$.

4 Regret Analysis

In this section, we conduct analysis of the regret of ARLM-UCB. We first present the formal self-normalized concentration inequality and compare it with existing results in the literature (Section 4.1). Then, we provide the bound on the expected cumulative (policy) regret (Section 4.2). The complete proofs of the theorems stated in the section are presented in Appendices B.3 and B.4.

4.1 Concentration Inequality for the Parameter Vectors

We first present the concentration result for the estimates $\hat{\phi}_t(u)$ of the true parameters $\phi(u)$, for every action $u \in \mathcal{U}$. At each round $t \in \mathbb{N}$, for the chosen action $u_t \in \mathcal{U}$, we solve the Ridge regression problem as:

$$\hat{\phi}_t(u) := \arg \min_{\tilde{\phi} \in \mathbb{R}^{p+1}} \sum_{l \in \mathcal{O}_t(u_t)} (x_l - \langle \tilde{\phi}, \mathbf{z}_{l-1} \rangle)^2 + \lambda \|\tilde{\phi}\|_2^2 = \hat{\mathbf{V}}_t(u_t)^{-1} \hat{\mathbf{b}}_t(u_t), \quad (12)$$

where $\mathcal{O}_t(u)$ is the set of rounds, where the action u was played, i.e., $\mathcal{O}_t(u) := \{\tau \in \mathbb{N} : u_\tau = u\}$. The following theorem formulates the concentration of $\hat{\phi}_t(u)$ around $\phi(u)$ over the rounds:

Theorem 4.1. (Self-normalized concentration) *Let $u \in \mathcal{U}$ be an action and $(\hat{\phi}_t(u))_{t \in \mathcal{O}_\infty(u)}$ be the sequence of solutions of the Ridge regression problems of Algorithm 1. Then, under Assumption 2.1 and 2.2, for every $\lambda \geq 0$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$, and for all rounds $t \in \mathbb{N}$, we have the following.*

$$\|\hat{\phi}_t(u) - \phi(u)\|_{\hat{\mathbf{V}}_t(u)} \leq \sqrt{\lambda} \|\phi(u)\|_2 + c_1(t) \|\phi(u)\|_2 + c_2(t) + \sigma \sqrt{2 \log \left(\frac{1}{\delta} \right) + \log \left(\frac{\det(\hat{\mathbf{V}}_t(u))}{\lambda^{p+1}} \right)},$$

where $c_1(t) := \mathcal{O} \left(\sqrt{p} t^d \log(t) \sqrt{\log \log(t)} \right)$ and $c_2(t) := \mathcal{O} \left(\sqrt{p+1} t^d \log(t) \sqrt{\log \log(t)} \right)$

Theorem 4.1 resembles the self-normalized concentration inequality of Bacchiocchi et al. [2024], whose idea originates from Theorem 1 of Abbasi-Yadkori et al. [2011]. Likewise, in case of AR-UCB, the exploration coefficients $\beta_t(u)$ are different for every action $u \in \mathcal{U}$. However, the important novelty that distinguishes the algorithm ARLM-UCB from the former is that the agent estimates $\hat{\phi}_t(u)$ using estimated AR rewards $(\hat{x}_t)_{t \in \mathbb{N}}$ calculated from past t environment-generated ARFIMA rewards $(y_t)_{t \in \mathbb{N}}$. This notion results in the derivation of the term $c_1(t) \|\phi(u)\|_2 + c_2(t)$ that emphasizes the convergence of the estimated reward \hat{x}_t with the coefficient \hat{d}_t at every round $t \in \mathbb{N}$ to the true AR reward $x_t = (1 - B)^{-d} y_t$ obtained through direct undifferencing the environment reward y_t .

Using the results in Theorem 4.1, we select the coefficient β_t based on the knowledge of the upper bounds specified in Assumption 2.2 and Assumption 2.3 for every $t \in \mathbb{N}$:

$$\beta_t(u) := \sqrt{\lambda(m^2 + 1)} + c_1(t) \sqrt{m^2 + 1} + c_2(t) + \sigma \sqrt{2 \log \left(\frac{1}{\delta} \right) + \log \left(\frac{\det(\hat{\mathbf{V}}_t(u))}{\lambda^{p+1}} \right)}, \quad (13)$$

where $c_1(t) = \sqrt{p} t^d \log(t) \sqrt{\log(\log(t) + 1)}$ and $c_2(t) = \sqrt{p+1} t^d \log(t) \sqrt{\log(\log(t) + 1)}$. This formula is constructed with three terms. The first term is a *bias* term, the second one is a *estimation error* term, and the third one is a *concentration* term. The *bias* term is derived by utilizing Assumptions 2.2 and 2.3, which guarantee that $\|\phi(u)\|_2 \leq \sqrt{m^2 + \Gamma^2} \leq \sqrt{m^2 + 1}$. The *estimation error* term arises from the fact that our learner is required learn the true memory parameter d throughout the learning interval T . Two components of this term, $c_1(t)$ and $c_2(t)$, arise from their respective bounds derived in Equation 4.1. In this way, the exploration coefficient $\beta_t(u)$ ensures that, with probability $1 - \delta$, the following inequality holds universally for every action $u \in \mathcal{U}$:

$$\|\hat{\phi}_t(u) - \phi(u)\|_{\hat{\mathbf{V}}_t(u)} \leq \beta_t(u) \quad (14)$$

The *bias* and *concentration* terms mimic those for $\beta_t(u)$ of Bacchiocchi et al. [2024], highlighting the independence of the simultaneous knowledge of Γ and c (Assumptions 2.1 and 2.2) introduced in our setting. This feature for ARLM-UCB is plausible for learning, as the true values of these parameters are unknown in practice.

4.2 Regret Bound

In this section, we derive a bound on the expected policy regret bound for ARLM-UCB:

Theorem 4.2. *Let $\delta = (2T)^{-1}$. Under Assumptions 2.1-2.4, ARLM-UCB suffers a cumulative expected (policy) regret bounded by (highlighting the dependence on Γ, p, m, σ, n , and T):*

$$\mathbb{E}[R(\text{ARLM-UCB}, T)] \leq \mathcal{O} \left(\frac{(p+1)^2(m+\sigma)\sqrt{n}T^{2d+0.5}\log^2(T)\log\log(T)}{(1-\Gamma)^2} \right).$$

The regret bound in Theorem 4.2 expands the one for AR-UCB. It’s worth noting that, unlike in the case of AR-UCB, with $p = 0$ and $\Gamma = 0$, our problem becomes ARFIMA(0, d , 0), where we obtain the regret rate $\mathcal{O}((m+\sigma)\sqrt{n}T^{2d+0.5}\log^2(T)\log\log(T))$. This rate does not reduce ARLM-UCB to standard MAB problem, unlike in the case of Bacchiocchi et al. [2024], where $p = 0$ gives the regret $\mathcal{O}((m+\sigma)\sqrt{nT})$ tight to standard MABs. This notion suggests that introducing the fractional integration in the reward-evolution process generally makes the learning problem more complex and difficult to handle by regular multi-armed bandits. All the theoretical derivations allowing us to achieve this upper bound are proved in Appendix B.

5 Experiments

This section presents numerical experiments on the ARLM-UCB, highlighting how this algorithm can outperform various competing bandit baselines in synthetically-generated domains. Appendix A features the bandit comparison on a dynamic strategy optimization task for the real-world stock volatility data. All algorithms were implemented in Python 3.12, and run over an Apple M1 with 8 GB RAM, in no more than a couple of hours.

We compare ARLM-UCB with the following baselines: (a) UCB1 (Auer et al. [1995]), an algorithm developed for stochastic MABs, (b) EXP3 (Auer et al. [2002]), an algorithm designed for adversarial MABs, (c) its finite-memory adaptive adversaries B-EXP3 (Dekel et al. [2012]), (d) AR2 (Chen et al. [2023]) designed to manage non-stationary AR(1) processes, and (e) AR-UCB (Bacchiocchi et al. [2024]) designed to operate in stationary AR(p) environments.

We evaluate the selected bandits in three synthetic scenarios with different properties that govern the reward evolution processes. For all synthetic experiments, we set the number of rounds to be $T = 10000$, the true memory parameter $d = 0.35$, and the grid of memory parameters $G_d = \{0.01, 0.02, \dots, 0.48, 0.49\}$, s.t. $d \in G_d$. The three settings have their ARFIMA($p, d, 0$) process order $p \in \{0, 1, 2\}$, number of actions $n \in \{5, 7\}$, number of initial exploration rounds $K \in \{5, 75, 550\}$ during which the agent pulls arms at random, and scale $m \in \{1.6, 7.4, 920\}$. The values of AR coefficients $\phi(u)$ are drawn randomly from a uniform distribution for each action $u \in \mathcal{U}$ and for each setting. The standard deviations of the noise in three environments are $\sigma \in \{0.9, 1.25, 10\}$. The selected hyperparameters of AR-UCB and ARLM-UCB are $\lambda = 1$ and $\bar{m} \in \{1.6, 7.5, 1000\}$. Table 1 summarizes the parameter settings for each experimental scenario.

Setting	p	n	m	\bar{m}	σ	K
A	0	5	7.4	7.5	1.25	5
B	1	5	1.6	1.6	0.9	75
C	2	7	920	1000	10	550

Table 1: Setting description

5.1 Results

Figure 1 shows the average cumulative regrets for three settings. We observe that ARLM-UCB consistently outperforms all competing bandit algorithms, always demonstrating sublinear behavior. On the other hand, all other bandits exhibit linear regret, since they are unable to process the long memory rewards and converge to sublinear regret over the learning horizon. That is because only ARLM-UCB has a specific mechanism for modeling long memory dynamics, allowing for the precise estimation and analysis of the underlying long-range dependence between environment-generated rewards.

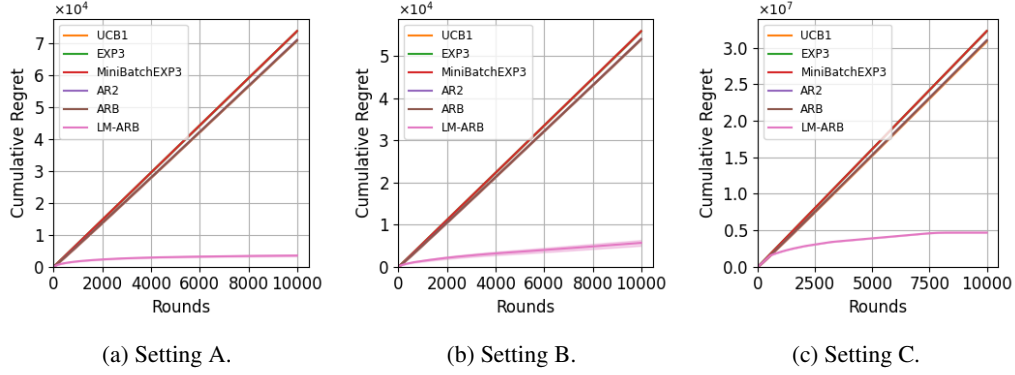


Figure 1: Cumulative Regret of ARLM-UCB and multiple baselines (100 runs, mean \pm std).

5.2 On the Convergence of the Memory Parameter Estimate

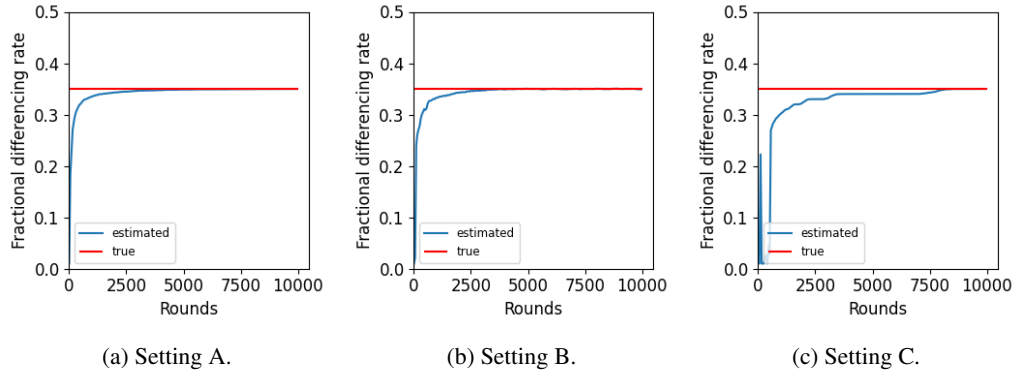


Figure 2: The convergence dynamics of selected memory parameters \hat{d} .

Figures 2 display the average estimates of the sequence of memory parameters $(\hat{d}_t)_{t \in \mathbb{N}}$ by ARLM-UCB on the optimization horizon T (blue) and a straight vertical line (red) representing the true memory parameter $d = 0.35$ as a benchmark. We immediately observe that the learner-selected rates always start at near-zero values in earlier round. This is due to the fact that our agent did not sufficiently explore the environment, which loosens his sense of present long-range dependence. However, the agent quickly regains the perception of long memory over more rounds and eventually converges in his estimates of $(\hat{d}_t)_{t \in \mathbb{N}}$ to the true rate d .

6 Conclusion

In this work, we addressed the online sequential decision-making problem, where the ARFIMA long memory temporal structure between rewards is present. We first formulate a LM-ARB setting by introducing a range of necessary assumptions and a selection of the optimal policy applicable in the context of our problem. We then propose a new online algorithm ARLM-UCB that learns the parameters for each action and the true memory parameter d corresponding to the long memory reward y_t . ARLM-UCB employs the idea that this bandit setting can be reduced to ARBs using an estimated memory parameter \hat{d}_t and solved using a linear contextual bandit. We also present a novel regret bound for this algorithm, accounting for the need to estimate the memory parameter in the addressed setting. Finally, we provided a variety of numerical experiments to assess the performance of ARLM-UCB based on cumulative regret and validate our solution. Future research directions should focus on extending the presented approach by designing a similar setting incorporating moving average part in the reward-generating mechanism. It is also promising in the long-term to derive a policy evaluation directly on a sequence of ARFIMA rewards.

References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems 24 (NeurIPS 2011)*, pages 2312–2320, 2011.
- Safwan Mahmood Al-Selwi, Mohd Fadzil Hassan, Said Jadid Abdulkadir, and Amgad Muneer. Lstm inefficiency in long-term dependencies regression problems. *Journal of Advanced Research in Applied Sciences and Engineering Technology*, 30(3):16–31, 2023.
- Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of the 36th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 322–331. IEEE, 1995.
- Peter Auer, Nicolò Cesa-Bianchi, and Gábor Lugosi. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 31(1):145–157, 2002.
- Francesco Bacchiocchi, Gianmarco Genalti, Davide Maran, Marco Mussi, Marcello Restelli, Nicola Gatti, and Alberto Maria Metelli. Autoregressive bandits. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*. PMLR, 2024.
- S. H. A. Bakar and C. M. Hafner. Forecasting realised volatility using arfima and har models. *Quantitative Finance*, 19(12):1925–1934, 2019.
- Jan Beran. Maximum likelihood estimation of the differencing parameter for invertible short and long memory ARIMA models. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(4):659–672, 1995.
- George E. P. Box, Gwilym M. Jenkins, Gregory C. Reinsel, and Greta M. Ljung. *Time Series Analysis: Forecasting and Control*. John Wiley & Sons, 5th edition, 2015. ISBN 978-1-118-67502-1.
- Uladzimir Charniauskii and Yao Zheng. Autoregressive bandits in near-unstable or unstable environment. *American Journal of Undergraduate Research*, 21(2):15–26, 2024. doi: 10.33697/ajur.2024.116.
- Qinyi Chen, Negin Golrezaei, and Djallel Bouneffouf. Non-stationary bandits with auto-regressive temporal dependency. In *Advances in Neural Information Processing Systems (NeurIPS)*. Curran Associates, Inc., 2023.
- Rainer Dahlhaus. Efficient parameter estimation for self-similar processes. *The Annals of Statistics*, 17(4):1749–1766, 1989.
- Ofer Dekel, Ambuj Tewari, and Raman Arora. Online bandit learning against an adaptive adversary: From regret to policy regret. In *Proceedings of the 29th International Conference on Machine Learning (ICML)*, pages 27–34. PMLR, 2012.
- Mallikarjuna Doodipala. Time series analysis for long memory process of air traffic using arfima. *International Journal of Scientific & Technology Research*, 9(3):6268–6272, 2020.
- Robert Fox and Murad S. Taqqu. Large-sample properties of parameter estimates for strongly dependent stationary gaussian time series. *The Annals of Statistics*, 14(2):517–532, 1986.
- C. W. J. Granger and R. Joyeux. An introduction to long-memory time series models and fractional differencing. *Journal of Time Series Analysis*, 1(1):15–29, 1980.
- Gaurav Gupta, Chenzhong Yin, Jyotirmoy V. Deshmukh, and Paul Bogdan. Non-markovian reinforcement learning using fractional dynamics. *arXiv preprint arXiv:2107.13790*, 2021.
- J. R. M. Hosking. Fractional differencing. *Biometrika*, 68(1):165–176, 1981.
- Xing-Qi Jiang. Time varying coefficient AR and VAR models. In *The Practice of Time Series Analysis*. Springer, 2023.
- Ferry Kondo Lembang, Lexy Janzen Sinay, and Asrul Irfanullah. Arfima modelling for tectonic earthquakes in the maluku region. *Indonesian Journal of Statistics and Its Applications*, 5(1):39–49, 2021.
- Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020. ISBN 9781108486828.
- Ming Liu. Modeling long memory in stock market volatility. *Journal of Econometrics*, 99(1):139–171, 2000.
- Michael McAleer and Shiqing Ling. A general asymptotic theory for time-series models. Technical report, Hong Kong University of Science and Technology, 2008.

- A. I. McLeod and K. Hipel. Fractional time series modelling. *Technometrics*, 28(2):101–111, 1986.
- D. Feigin Paul and L. Tweedie Richard. Random coefficient autoregressive processes: a markov chain analysis of stationarity and finiteness of moments. *Journal of Time Series Analysis*, 6(1):1–14, 1985. doi: 10.1111/j.1467-9892.1985.tb00394.x.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. John Wiley Sons, 1994. ISBN 978-0-471-61977-2.
- Yuzhen Qin, Yingcong Li, Fabio Pasqualetti, Maryam Fazel, and Samet Oymak. Stochastic contextual bandits with long horizon rewards. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence*, pages 8611–8619. Association for the Advancement of Artificial Intelligence, 2023.
- L. Tweedie Richard. Criteria for rates of convergence of markov chains, with application to queueing and storage theory. In C. Kingman J., F. and E. H. Reuter G., editors, *Probability, Statistics and Analysis*, London Mathematical Society Lecture Note Series, page 260–276. Cambridge University Press, 1983. doi: 10.1017/CBO9780511662430.016.
- P. M. Robinson. Gaussian semiparametric estimation of long range dependence. *The Annals of Statistics*, 23(5): 1630–1661, 1995.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2nd edition, 2020.
- Wei Tang, Chien-Ju Ho, and Yang Liu. Bandit learning with delayed impact of actions. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Yifu Tang, Yingfei Wang, and Zeyu Zheng. Stochastic multi-armed bandits with strongly reward-dependent delays. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 3043–3051. PMLR, 02–04 May 2024.
- Anna L. Trella, Walter Dempsey, Finale Doshi-Velez, and Susan A. Murphy. Non-stationary latent auto-regressive bandits. *arXiv preprint arXiv:2402.03110*, 2024.
- JQ (Justin) Veenstra and A. Ian McLeod. *arfima: Fractional ARIMA (and Other Long Memory) Time Series Modeling*, 2022. R package version 1.8-1.
- Naiming Yuan, Zuntao Fu, and Shida Liu. Extracting climate memory using fractional integrated statistical model: A new perspective on climate prediction. *Scientific Reports*, 4:6577, 2014.
- Jingyu Zhao, Feiqing Huang, Jia Lv, Yanjie Duan, Zhen Qin, Guodong Li, and Guangjian Tian. Do RNN and LSTM have long memory? In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 11365–11375. PMLR, 13–18 Jul 2020.
- Nicolas Zucchet, Robert Meier, Simon Schug, Asier Mujika, and João Sacramento. Online learning of long-range dependencies. In *Advances in Neural Information Processing Systems 36 (NeurIPS 2023)*, pages 10477–10493. Curran Associates, Inc., 2023.

A Real-world Data Experiments

In this section, we compare ARLM-UCB with the same baselines as in Section 5 on the dynamic strategy optimization task in an online fashion, for which the real-world stock volatility data. We select stock prices of three technological companies: Apple, Netflix, and Google. For each company listed, the data is obtained from Kaggle and contains daily closing prices and daily trading volumes from 2014 to 2023, 2002 to 2022, and 2004 to 2024, respectively. We convert the closing prices to absolute log-returns and parallelly shift this series by 0.01 numeric increment to prepare our series of log-returns for further processing. We then discretize the prices into $n = 4$ price bands (i.e., our actions) based on whether the magnitude of a return is positive/negative before being converted to an absolute value and whether the volatility exceeds the value in the third quartile of the volatility distribution of each stock. Table 2 summarizes the action selection for our real-world experiment.

Log-Return	Trading Volume	Arm Label
Positive	Exceeds 3rd quartile	1
Negative	Exceeds 3rd quartile	2
Positive	Below 3rd quartile	3
Negative	Below 3rd quartile	4

Table 2: A summary of an action selection method

A.1 Setting Configuration

We construct the simulation environment for each stock using the following methodology. First, we find the true value of a memory parameter d corresponding to each data set using the R package `arfima` (Veenstra and McLeod [2022]). We fit the ARFIMA($p, d, 0$) model on the absolute log-return series to estimate the global memory parameter d . Then we convert the observations of the original log-returns to short memory with d . The earlier parallel shift of our log-return series reinforces the positivity of the reconstructed price series, which is a necessary condition for our bandit analysis. Finally, we estimate the hidden autoregressive parameters for each action using standard regression methods, taking into account the adjacent values of the past p at each time point. Each selected value of p represents the number of significant AR lags in reward-governing ARFIMA model fit to each data. To determine the noise standard deviation σ for each dataset, we compute the square root of the weighted mean of all variances produced from estimating each arm.

Table 3 summarizes the parameter settings for each dataset considered. We globally set the number of rounds $T = 10000$, the grid of memory parameters $G_{\hat{d}} = \{0.01, 0.02, \dots, 0.48, 0.49\}$, s.t. $d \in G_{\hat{d}}$. Each setting has estimated AR parameters $m \in \{0.012, 0.135, 0.007\}$ and $\Gamma \in \{0.75, 0.959, 0.995\}$ with the number of exploration rounds $K \in \{150, 250\}$. The estimated Gaussian noise standard deviations are $\sigma \in \{0.013, 0.014, 0.027\}$.

We select $\lambda \in \{0.001, 0.01, 1\}$, where $\lambda < 1$ values are set for settings with Γ close to 1 numerically. Charniauski and Zheng [2024] showed that setting close-to-zero values for λ could help AR-UCB achieve smaller regrets in such near-unstable (e.g., with $\Gamma \approx 1$) settings. This notion should hold similarly for ARLM-UCB, since we aim to estimate the underlying AR process of converted rewards \hat{x}_t . We also choose $\bar{m} \in \{0.01, 0.015, 0.15\}$ that are set to around 20% greater value than m in each respective environment. Table 3 summarizes the parameter configurations for each scenario.

Stock Data	p	m	\bar{m}	λ	Γ	d	σ	K
Apple	1	0.012	0.015	0.001	0.959	0.25	0.013	150
Netflix	2	0.135	0.15	1	0.75	0.21	0.027	250
Google	3	0.007	0.01	0.01	0.995	0.31	0.014	150

Table 3: Setting description

A.2 Results

Figure 3 illustrates the average cumulative regrets for three datasets. The ARLM-UCB demonstrates the sublinear convergence in all three cases, achieving the smallest cumulative regret. On the other

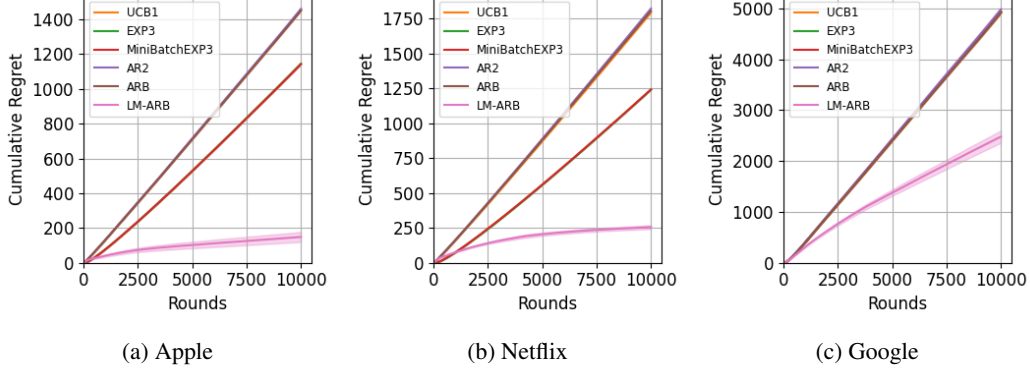


Figure 3: Cumulative regret of ARLM-UCB and baselines on real data (100 runs, mean \pm std).

hand, neither of the other bandits achieve sublinear regret in all considered cases. We also observe that both EXP3 and B-EXP3 suffer the exponential regret in all three scenarios, which is explicitly seen in 3a and 3b. This might be due to the structure of arms and parameter values presented in each of three environments. These observations make ARLM-UCB the algorithm with the best performance over the competitors.

A.3 On the Convergence of the Memory Parameter Estimate

Figure 4 demonstrates the average estimates of the sequence of memory parameters $(\hat{d}_t)_{t \in \mathbb{N}}$ on the optimization horizon T (blue) and a straight vertical line (red) representing the true memory parameter for each dataset. These results replicate the ones achieved on synthetic data cases displayed in Section 5.2. We see that the estimated rates rapidly converge to the true rate, starting at low values. These plots also demonstrate that the estimate \hat{d}_t is converging regardless of the value for a true rate d established in the environment.

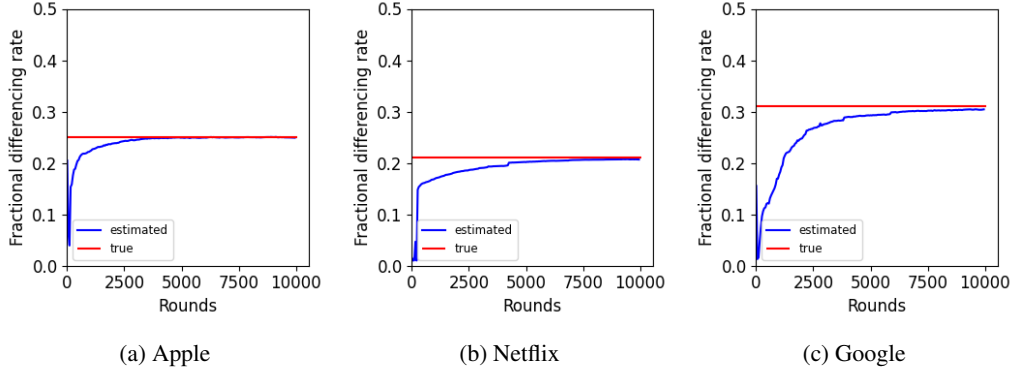


Figure 4: The convergence dynamics of selected memory parameters \hat{d} .

B Omitted Proofs

B.1 Proof of Theorem 2.6

Proof. We prove the important property of geometric ergodicity of the reward process x_t . For the process x_t , we consider the process expressed in Equation 2.

We define the companion vector state $X_t := (x_{t-1}, \dots, x_{t-p})^\top \in \mathbb{R}^p$ for all $t \in \mathbb{N}$. We rewrite the reward evolution from Equation 2 as follows:

$$X_t = A(u_t)X_{t-1} + b(u_t) + \xi_t,$$

where we define:

$$A(u_t) := \begin{pmatrix} \phi_1(u_t) & \phi_2(u_t) & \cdots & \phi_p(u_t) \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{p \times p}, b(u_t) := \begin{pmatrix} \phi_0(u_t) \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \in \mathbb{R}^p, \xi_t := \begin{pmatrix} \varepsilon_t \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \in \mathbb{R}^p,$$

and the transition probability is:

$$P(x, \mathcal{A}) = \int_{\mathcal{A}} \varphi(z - m_t(x)) dz,$$

for $x \in \mathbb{R}^p$, the mean process $m_t(x) = A(u_t)x + b(u_t)$, and $\mathcal{A} \in \mathbb{B}^p$, the class of Borel sets of \mathbb{R}^p .

Because ξ_t is defined in terms of Gaussian noise, $P(x, \mathcal{A}) > 0$ and $\{X_t\}$ is ν_p -irreducible for the Lebesgue measure ν_p on $(\mathbb{R}^p, \mathcal{B}^p)$.

We prove by showing that Tweedie's drift criterion (Richard [1983]) holds, i.e. there is a small set $G \subset \mathbb{R}^p$ with $\nu_p(G) > 0$ and a non-negative continuous function $V(x)$, s.t.

$$\mathbb{E}[V(X_t)|X_{t-1} = x] \leq (1 - \delta)V(x), x \notin G \quad (15)$$

and

$$\mathbb{E}[V(X_t)|X_{t-1} = x] \leq M, x \in G \quad (16)$$

for $0 < \delta < 1$ and $0 < M < \infty$.

By Assumptions 2.2 and 2.3, we observe that $\rho(A(u_t)) = \Gamma < 1$ and $\sup_t \mathbb{E}[\|b(u_t)\|] \leq m$.

We choose Lyapunov criterion as $V(x) = 1 + \|x\|^2$. We are now able to observe the following:

$$\begin{aligned} \mathbb{E}[V(X_t)|X_{t-1} = x] &= \mathbb{E}[1 + \|X_t\|^2|X_{t-1} = x] \\ &\leq 1 + \Gamma^2\|x\|^2 + m^2 + \sigma^2 \leq 1 + \Gamma^2V(x) + m^2 + \sigma^2. \end{aligned}$$

Denote $\delta = 1 - \Gamma^2 - \frac{1 - \Gamma^2 + m^2 + \sigma^2}{V(x)}$ and $G := \{x : \|x\| \leq L\}$, s.t. $V(x) \geq 1 + \frac{m^2 + \sigma^2}{1 - \Gamma^2}$ for every $\|x\| > L$. We obtain that conditions stated in Equations 15 and 16 hold.

Moreover, we observe that $\mathbb{E}[f(X_t)|X_{t-1} = x]$ is continuous w.r.t. x for every bounded function $f(\cdot)$. Thus, $\{X_t\}$ is a Feller chain.

By Paul and Richard [1985], G is a small set. By referring to Theorem 4(ii) in Richard [1983] and Theorem 1 in Paul and Richard [1985], X_t is geometrically ergodic with a unique strictly stationary solution. \square

B.2 Proof of Theorem 2.7

Proof. Let $u \in \mathcal{U}$. For each arm $u \in \mathcal{U}$, the reward y_t at every round $t \in \mathbb{N}$ evolves according to ARFIMA($p, d, 0$) process as stated in Equation 3. Observe that the process for each arm evolves independently, regardless of when or whether the arm is pulled.

Assumptions 2.1-2.4, Theorem 2.6 and the noise $\varepsilon_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ for $t \in \mathbb{N}$ guarantee that x_t is strictly stationary and ergodic, with $\mathbb{E}[x_t^2] < \infty$.

For ARFIMA($p, d, 0$), we also have the following:

$$x_t = \sum_{i=0}^{\infty} \tilde{\psi}_{0i}(d) \varepsilon_{t-i}(d) \text{ and } \varepsilon_t(d) = (1-B)^d x_t = \sum_{i=0}^{\infty} \tilde{\psi}_i(d) x_{t-i} - \sum_{i=0}^{\infty} \left[\sum_{n=0}^{\infty} \tilde{\psi}_n(d) \phi_{i-n}(u_t) \right] x_{t-k},$$

from which we observe that

$$\sum_{i=0}^{\infty} \tilde{\psi}_i(d) x_{t-i} - \sum_{i=0}^{\infty} \left[\sum_{n=0}^{\infty} \tilde{\psi}_n(d) \phi_{i-n}(u_t) \right] x_{t-k} \leq \sum_{i=0}^{\infty} \tilde{\psi}_i(d) x_{t-i} - \sum_{i=0}^{\infty} \left[\sum_{n=0}^{\infty} c \cdot \tilde{\psi}_n(d) \right] x_{t-k} \leq \sum_{i=0}^{\infty} \tilde{\psi}_i(d) x_{t-i}.$$

For the last term of the above inequality, using the conditions of stationarity and ergodicity of the reward x_t , McAleer and Ling [2008] verified that the estimated rates $(\hat{d}_t)_{t \in \mathbb{N}}$ converge to the true rate d as stated in Equation 9 and that concludes the proof. \square

B.3 Proof of Theorem 4.1

Before we develop the self-normalized concentration, we introduce the context vector bound (Lemma B.1) by [Bacchiocchi et al. \[2024\]](#).

Lemma B.1. *Let $(\mathbf{z}_t^*)_{t \in [T]}$ be the sequence of observation vectors observed by executing the learner's policy. If $\mathbf{z}_0 = (1, 0, \dots, 0)^\top$, then, for every $\delta \in (0, 1)$, with probability at least $1 - \delta$, simultaneously for every $t \in [T]$, it holds that:*

$$\|\mathbf{z}_{t-1}\|_2 \leq \sqrt{1 + p \left(\frac{m + \eta}{1 - \Gamma} \right)^2},$$

where $\eta = \sqrt{2\sigma^2 \log(T/\delta)}$.

Proof of Theorem 4.1. We first consider an action at a time; then we obtain the final result with a union bound over $\mathcal{U} := [n]$.

Let $u \in \mathcal{U}$. Observe that the estimates of an action u change only when u is pulled. Let $l \in \mathbb{N}$ be an index and let $t_l(u) \in \mathbb{N}$ be the round in which the action u is pulled for the l -th time, i.e., $\{t_l(u) : l \in \mathbb{N}\} = \mathcal{O}_\infty(u)$. Thus, we have the following:

$$\begin{aligned} \phi_{t_l}(u) &= \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \hat{\mathbf{b}}_{t_l(u)}(u) = \left(\lambda \mathbf{I}_{p+1} + \sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top \right)^{-1} \sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{x}_{t_j} \\ &= \left(\lambda \mathbf{I}_{p+1} + \underbrace{\sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \mathbf{z}_{t_j(u)-1}^\top}_{\Delta_1} + \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \mathbf{z}_{t_j(u)-1}^\top \right)^{-1} \\ &\quad \left(\underbrace{\sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{x}_{t_j} - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} x_{t_j}}_{\Delta_2} + \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} x_{t_j} \right) \\ &= \left(\lambda \mathbf{I}_{p+1} + \Delta_1 + \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \mathbf{z}_{t_j(u)-1}^\top \right)^{-1} \left(\Delta_2 + \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} x_{t_j} \right) \\ &= (\mathbf{V}_{t_l(u)}(u) + \Delta_1)^{-1} (\Delta_2 + \mathbf{b}_{t_l(u)}(u)) = (\mathbf{V}_{t_l(u)}(u) + \Delta_1)^{-1} \Delta_2 + (\mathbf{V}_{t_l(u)}(u) + \Delta_1)^{-1} \mathbf{b}_{t_l(u)}(u) \\ &\stackrel{(a)}{=} \underbrace{\mathbf{V}_{t_l(u)}^{-1}(u) \mathbf{b}_{t_l(u)}(u) - \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \mathbf{b}_{t_l(u)}(u)}_{P_1} + \underbrace{\mathbf{V}_{t_l(u)}^{-1}(u) \Delta_2 - \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \Delta_2}_{P_2}. \end{aligned}$$

where passage (a) arises from the observation that

$$(\mathbf{V}_{t_l(u)}(u) + \Delta_1)^{-1} = \mathbf{V}_{t_l(u)}(u) - \mathbf{V}_{t_l(u)}(u) \Delta_1 (\mathbf{V}_{t_l(u)}(u) + \Delta_1)^{-1} = \mathbf{V}_{t_l(u)}(u) - \mathbf{V}_{t_l(u)}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u)$$

Before we decompose terms P_1 and P_2 , we demonstrate the following decomposition of $\mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u)$, observing that $\Delta_1 = \hat{\mathbf{V}}_{t_l(u)}(u) - \mathbf{V}_{t_l(u)}(u)$:

$$\begin{aligned} \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) &= \mathbf{V}_{t_l(u)}^{-1}(u) (\hat{\mathbf{V}}_{t_l(u)}(u) - \mathbf{V}_{t_l(u)}(u)) \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \\ &= (\mathbf{V}_{t_l(u)}^{-1}(u) \hat{\mathbf{V}}_{t_l(u)}(u) - \mathbf{I}_{p+1}) \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) = \mathbf{V}_{t_l(u)}^{-1}(u) - \hat{\mathbf{V}}_{t_l(u)}^{-1}(u). \end{aligned}$$

We now begin with the term P_1 . The following holds:

$$P_1 = \mathbf{V}_{t_l(u)}^{-1}(u) \mathbf{b}_{t_l(u)}(u) - \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \mathbf{b}_{t_l(u)}(u) =$$

$$\begin{aligned}
&= \mathbf{V}_{t_l(u)}^{-1}(u) \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} x_{t_j} - \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} x_{t_j} \\
&= \mathbf{V}_{t_l(u)}^{-1}(u) \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} (\langle \phi(u), \mathbf{z}_{t_j(u)-1} \rangle + \varepsilon_{t_j}) \\
&\quad + \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} (\langle \phi(u), \mathbf{z}_{t_j(u)-1} \rangle + \varepsilon_{t_j}) \\
&\stackrel{(b)}{=} \phi(u) - \lambda \mathbf{V}_{t_l(u)}^{-1}(u) \phi(u) + \underbrace{\mathbf{V}_{t_l(u)}^{-1}(u) \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \varepsilon_{t_j}}_{\mathbf{s}_{t_j}} + \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) \\
&\quad - \lambda \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) + \underbrace{\mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \varepsilon_{t_j}}_{\mathbf{s}_{t_j}} \\
&= \phi(u) - \lambda \mathbf{V}_{t_l(u)}^{-1}(u) \phi(u) + \mathbf{V}_{t_l(u)}^{-1}(u) \mathbf{s}_{t_j} - \Delta_1 \mathbf{V}_{t_l(u)}^{-1}(u) \phi(u) \\
&\quad + \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1} \phi(u) + \lambda \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1} \phi(u) - \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1} \mathbf{s}_{t_j} \\
&= \phi(u) - \lambda \mathbf{V}_{t_l(u)}^{-1}(u) \phi(u) + \mathbf{V}_{t_l(u)}^{-1}(u) \mathbf{s}_{t_j} - \Delta_1 \mathbf{V}_{t_l(u)}^{-1}(u) \phi(u) \\
&\quad + \Delta_1 \mathbf{V}_{t_l(u)}^{-1}(u) \phi(u) - \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) + \lambda \mathbf{V}_{t_l(u)}^{-1}(u) \phi(u) - \lambda \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) \\
&\quad - \mathbf{V}_{t_l(u)}^{-1}(u) \mathbf{s}_{t_j} + \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \mathbf{s}_{t_j} = \phi(u) - \lambda \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) - \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) + \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \mathbf{s}_{t_j},
\end{aligned}$$

where the passage (b) arises from the observation that $\sum_{j=1}^l \mathbf{z}_{t_j-1} (\langle \phi(u), \mathbf{z}_{t_j-1} \rangle) = \sum_{j=1}^l \mathbf{z}_{t_j-1} \mathbf{z}_{t_j-1}^\top \phi(u)$.

We then proceed with P_2 in a similar fashion as follows:

$$P_2 = \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_2 - \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \Delta_2 = \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_2 - \mathbf{V}_{t_l(u)}^{-1}(u) \Delta_2 + \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \Delta_2 = \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \Delta_2.$$

Thus, we achieve

$$\phi_{t_l}(u) = \phi(u) - \lambda \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) - \Delta_1 \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \phi(u) + \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \mathbf{s}_{t_j} + \hat{\mathbf{V}}_{t_l(u)}^{-1}(u) \Delta_2.$$

Moving $\phi(u)$ on the left side of the equation and taking the Gramian norm of both sides, we have the following inequality

$$\|\phi_{t_l}(u) - \phi(u)\|_{\hat{\mathbf{V}}_{t_l(u)}^{-1}(u)} \leq \sqrt{\lambda} \|\phi(u)\|_2 + \|\Delta_1\|_2 \|\phi(u)\|_2 + \|\Delta_2\|_2 + \|\mathbf{s}_{t_j}\|_{\hat{\mathbf{V}}_{t_l(u)}^{-1}(u)}.$$

We begin deriving the bounds for $\|\Delta_1\|_2$ and $\|\Delta_2\|_2$. First, $\|\Delta_1\|_2$ can be decomposed as

$$\begin{aligned}
\|\Delta_1\|_2^2 &= \left\| \sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \mathbf{z}_{t_j(u)-1}^\top \right\|_2^2 \\
&= \left\| \sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top + \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \mathbf{z}_{t_j(u)-1}^\top \right\|_2^2 \\
&\leq \left\| \sum_{j=1}^l \underbrace{[\hat{\mathbf{z}}_{t_j(u)-1} - \mathbf{z}_{t_j(u)-1}] \hat{\mathbf{z}}_{t_j(u)-1}^\top}_{\mathbf{e}_{t_j(u)-1}} \right\|_2^2 + \left\| \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \underbrace{[\hat{\mathbf{z}}_{t_j(u)-1}^\top - \mathbf{z}_{t_j(u)-1}^\top]}_{\mathbf{e}_{t_j(u)-1}^\top} \right\|_2^2
\end{aligned}$$

$$\leq \sum_{j=1}^l \|\mathbf{e}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top\|_2^2 + \sum_{j=1}^l \|\mathbf{z}_{t_j(u)-1} \mathbf{e}_{t_j(u)-1}^\top\|_2^2,$$

and $\|\Delta_2\|_2$ as

$$\begin{aligned} \|\Delta_2\|_2^2 &= \left\| \sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{x}_{t_j} - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} x_{t_j} \right\|_2^2 \\ &= \left\| \sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{x}_{t_j} - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \hat{x}_{t_j} + \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \hat{x}_{t_j} - \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} x_{t_j} \right\|_2^2 \\ &\leq \left\| \sum_{j=1}^l \hat{\mathbf{z}}_{t_j(u)-1} \hat{x}_{t_j} - \mathbf{z}_{t_j(u)-1} \hat{x}_{t_j} \right\|_2^2 + \left\| \sum_{j=1}^l \mathbf{z}_{t_j(u)-1} \hat{x}_{t_j} - \mathbf{z}_{t_j(u)-1} x_{t_j} \right\|_2^2 \\ &\leq \sum_{j=1}^l \|\mathbf{e}_{t_j(u)-1} \hat{x}_{t_j}\|_2^2 + \sum_{j=1}^l \|\mathbf{z}_{t_j(u)-1} (\hat{x}_{t_j} - x_{t_j})\|_2^2. \end{aligned}$$

To be able to bound the 2-norms of Δ_1 and Δ_2 , we first introduce the following decomposition of the reward difference:

$$|x_t - \hat{x}_t| = \left| \sum_{j=0}^{\infty} \psi_j(-d) y_{t-j} - \sum_{j=0}^t \psi_j(-\hat{d}_t) y_{t-j} \right| \leq \underbrace{\left| \sum_{j=1}^t \psi_j(-d) y_{t-j} - \sum_{j=1}^t \psi_j(-\hat{d}_t) y_{t-j} \right|}_{B_1} + \underbrace{\left| \sum_{j=t+1}^{\infty} \psi_j(-d) y_{t-j} \right|}_{B_2}.$$

We let $x_t = \sum_{j=0}^{\infty} \psi_j(-d) y_{t-j} = y_t + \sum_{j=1}^{\infty} \prod_{i=1}^j \frac{i-1+d}{i} y_{t-j}$, $\tilde{x}_t = \sum_{j=0}^t \psi_j(-d) y_{t-j} = y_t + \sum_{j=1}^t \prod_{i=1}^j \frac{i-1+d}{i} y_{t-j}$, and $\hat{x}_t = \sum_{j=0}^t \psi_j(-\hat{d}_t) y_{t-j} = y_t + \sum_{j=1}^t \prod_{i=1}^j \frac{i-1+\hat{d}_t}{i} y_{t-j}$.

Few important observations about y_t . First, by Theorem 2.6, the process x_t is strictly stationary and ergodic, so the variance of x_t is bounded by a finite constant, i.e., $\text{Var}(x_t) \leq \sigma_{x_t}^2$, for every round $t \in \mathbb{N}$. We also have that, by Lemma B.3, $x_t \leq \frac{m+\eta}{1-\Gamma}$ for every $t \in \mathbb{N}$, where $\eta = \sqrt{2\sigma^2 \log(T/\delta)}$, and so is $y_t = (1-B)^d x_t \leq \frac{m+\eta}{1-\Gamma}$ and with a finite variance, respectively.

With these observations, we begin with the term B_1 . We denote $f(d) = \tilde{x}_t$ and $f(\hat{d}_t) = \hat{x}_t$, for which we observe the following through Taylor decomposition:

$$B_1 = |\hat{x}_t - \tilde{x}_t| = |f(\hat{d}_t) - f(d)| \leq |\hat{d}_t - d| \cdot \left| \frac{\partial f(d)}{\partial d} \right| = |\hat{d}_t - d| \sum_{j=1}^t \psi'_j(d) y_{t-j}.$$

We intermediately provide the decomposition of the term $\psi'_j(d)$ through the log-derivative, which holds for every $j \in \llbracket t \rrbracket$:

$$\psi'_j(-d) = \psi_j(-d) \sum_{i=1}^j \frac{1}{j-1+d} = \psi_j(-d) \cdot \mathcal{O}(\log(j)) = \mathcal{O}(j^{-1+d} \log(j)),$$

where we exploit the notion that $\psi_j(-d) = \mathcal{O}(j^{-1+d})$ from Theorem 3.1 of McAleer and Ling [2008].

Therefore, using Theorem 2.7 and Cauchy-Schwarz, we bound $\mathbb{E}[B_1^2]$ as

$$\mathbb{E}[B_1^2] \leq (\hat{d}_t - d)^2 \left(\sum_{j=1}^t \mathcal{O}(j^{-1+d} \log(j)) y_{t-j} \right)^2 = \mathcal{O} \left(\frac{\log \log(t)}{t} \right) \cdot \mathcal{O}(t^{2d} \log t) = \mathcal{O}(t^{2d-1} \log^2(t) \log \log(t)),$$

so the Root Mean Square of B_1 is bounded by $\mathcal{O}(t^{d-0.5} \log(t) \sqrt{\log \log(t)})$.

We proceed by bounding B_2 in a similar fashion. Given that $\psi_j(-d) = \mathcal{O}(j^{-1+d})$, we have by Cauchy-Schwarz that

$$\mathbb{E}[B_2^2] \leq \left(\sum_{j=t+1}^{\infty} \psi_j(-d) y_{t-j} \right)^2 \leq \mathcal{O}(t^{2d-1}),$$

from which we see that the Root Mean Square of B_2 is bounded by $\mathcal{O}(t^{d-0.5})$.

Thus, the bound for the reward difference is:

$$|x_t - \hat{x}_t| = \mathcal{O}\left(t^{d-0.5} \log(t) \sqrt{\log \log(t)}\right) + \mathcal{O}(t^{d-0.5}) = \mathcal{O}\left(t^{d-0.5} \log(t) \sqrt{\log \log(t)}\right),$$

where the first term dominates the latter by the log-factor.

We now bounding $\|\Delta_1\|_2$. Applying Lemma B.1 and Theorem 2.7, we observe that the following holds:

$$\begin{aligned} \sum_{j=1}^l \|\mathbf{e}_{t_j(u)-1} \hat{\mathbf{z}}_{t_j(u)-1}^\top\|_2^2 &= \sum_{j=1}^l \|\mathbf{e}_{t_j(u)-1}\|_2^2 \|\hat{\mathbf{z}}_{t_j(u)-1}\|_2^2 \\ &\leq \sqrt{1 + p \left(\frac{m+\eta}{1-\Gamma} \right)^2} \cdot \sum_{j=1}^l \sum_{i=1}^p (\hat{x}_{t_j(u)-i-1} - x_{t_j(u)-i-1})^2 \stackrel{(a)}{=} \mathcal{O}(pt^{2d} \log^2(t) \log \log(t)), \end{aligned}$$

where passage (a) follows from $\|\hat{\mathbf{z}}_t\|_2^2 \leq \sqrt{1 + p \left(\frac{m+\eta}{1-\Gamma} \right)^2} = \mathcal{O}(1)$ for every $t \in \llbracket T \rrbracket$. By a similar notion, the same bound holds for the second term $\sum_{j=1}^l \|\mathbf{z}_{t_j(u)-1} \mathbf{e}_{t_j(u)-1}^\top\|_2^2$ as well.

Thus, we have the following bound for the norm of Δ_1 :

$$\|\Delta_1\|_2 = \mathcal{O}\left(\sqrt{p} t^d \log(t) \sqrt{\log \log(t)}\right) := c_1(t)$$

We then bound $\|\Delta_2\|_2$ in a similar fashion. We start with the term $\sum_{j=1}^l \|\mathbf{e}_{t_j(u)-1} \hat{x}_{t_j}\|_2^2$:

$$\sum_{j=1}^l \|\mathbf{e}_{t_j(u)-1} \hat{x}_{t_j}\|_2^2 \leq \left(\frac{m+\eta}{1-\Gamma} \right)^2 \sum_{j=1}^l \sum_{i=1}^p (\hat{x}_{t_j(u)-i-1} - x_{t_j(u)-i-1})^2 \stackrel{(a)}{=} \mathcal{O}(pt^{2d} \log^2(t) \log \log(t)),$$

where in passage (a) we similarly observe that $x_t \leq \frac{m+\eta}{1-\Gamma} = \mathcal{O}(1)$ for every $t \in \llbracket T \rrbracket$ by Lemma B.1.

We finish by bounding the second term $\sum_{j=1}^l \|\mathbf{z}_{t_j(u)-1} (\hat{x}_{t_j} - x_{t_j})\|_2^2$ using all the previous observations as follows:

$$\sum_{j=1}^l \|\mathbf{z}_{t_j(u)-1} (\hat{x}_{t_j} - x_{t_j})\|_2^2 \leq \sqrt{1 + p \left(\frac{m+\eta}{1-\Gamma} \right)^2} \sum_{j=1}^l (\hat{x}_{t_j} - x_{t_j})^2 = \mathcal{O}(t^{2d} \log^2(t) \log \log(t))$$

Thus, we have the following bound for the norm of Δ_2 :

$$\|\Delta_2\|_2 = \mathcal{O}\left(\sqrt{p+1} t^d \log(t) \sqrt{\log \log(t)}\right) := c_2(t).$$

Therefore, the following inequality holds:

$$\|\phi_{t_l(u)} - \phi(u)\|_{\hat{\mathbf{V}}_{t_l(u)}^{-1}} \leq \sqrt{\lambda} \|\phi(u)\|_2 + c_1(t) \|\phi(u)\|_2 + c_2(t) + \|\mathbf{s}_{t_l(u)}\|_{\hat{\mathbf{V}}_{t_l(u)}^{-1}}.$$

Finally, let $\mathcal{F}_{t_l(u)} = \sigma(\mathbf{z}_0, u_1, \mathbf{z}_1, u_2, \dots, \mathbf{z}_{t_l(u)-1}, u_{t_l(u)})$ be the filtration generated by all events realized at round $t_l(u)$. Let us now consider the stochastic processes $(\varepsilon_{t_l(u)})_{l \in \mathbb{N}}$ and $(\mathbf{z}_{t_l(u)-1})_{l \in \mathbb{N}}$. We observe that $\varepsilon_{t_l(u)}$ is $\mathcal{F}_{t_l(u)}$ -measurable and conditionally σ^2 -subgaussian and that $\mathbf{z}_{t_l(u)-1}$ is $\mathcal{F}_{t_l(u)-1}$ -measurable. Thus, by applying Theorem 1 of Abbasi-Yadkori et al. [2011], we have that simultaneously for all $l \in \mathbb{N}$ with probability $1 - \delta$:

$$\|\mathbf{s}_{t_l(u)}\|_{\hat{\mathbf{V}}_{t_l(u)}^{-1}} \leq \sigma \sqrt{2 \log \frac{1}{\delta} + \log \frac{\det \hat{\mathbf{V}}_{t_l(u)}(u)}{\lambda^{p+1}}}$$

for all actions $u \in \mathcal{U}$ and the rounds $t \in \mathbb{N}$.

□

B.4 Proof of the Upper Regret Bound

Before we derive the regret bound, we reconstruct the results of [Bacchiocchi et al. \[2024\]](#) for the External-to-Policy Regret Bound (Lemma B.3) in our setting. We first introduce their results for Policy-Regret-Decomposition (Lemma B.2) to be used for our purposes.

Lemma B.2. (*Policy-Regret-Decomposition*). Let $(x_t^*)_{t \in \mathbb{N}}$ be the sequence of rewards by executing the optimal policy π^* and let $(x_t)_{t \in \mathbb{N}}$ be the sequence of rewards by executing the learner's policy π . Then, for every $t \in \mathbb{N}$, it holds that:

$$\hat{r}_t = r_t + \epsilon_t = \sum_{i=1}^p \phi_i(u_t^*) r_{t-i} + \rho_t + \epsilon_t,$$

where $r_t := x_t^* - x_t$ is the instantaneous policy regret, $\rho_t := \langle \phi(u_t^*) - \phi(u_t), \mathbf{z}_{t-1} \rangle$ is the instantaneous external regret, $\epsilon_t = x_t - \hat{x}_t$ is the error term representing the difference between the converted and the true AR rewards, the $u_t^* = \pi_t(H_{t-1}^*)$, and $r_{t-i} = 0$ if $i \geq t$.

Lemma B.3. (*External-to-Policy Regret Bound*) Let π be the learner's policy and $T \in \mathbb{N}$ be the horizon. Under Assumptions 2.1 and 2.2, it holds that:

$$\begin{aligned} \hat{R}(\pi, T) &= \mathbb{E} \left[\sum_{t=1}^T \left[\sum_{i=1}^p \phi_i(u_t^*) r_{t-i} + \rho_t + \eta_t \right] \right] \\ &\leq \left(\frac{\Gamma p}{1 - \Gamma} + 1 \right) \mathcal{Q}(\pi, T) + \mathcal{O} \left(T^{d+0.5} \log(T) \sqrt{\log \log(T)} \right), \end{aligned}$$

where $\mathcal{Q}(\pi, T) := \mathbb{E}[\sum_{t=1}^T \rho_t]$ is the cumulative expected external regret.

Proof of Lemma B.3. We use the results of the regret decomposition in Lemma B.2. We then express our cumulative regret as the following Triangular Inequality:

$$\hat{R}(\pi, T) \leq \mathbb{E} \left[\left| \sum_{t=1}^T \left[\sum_{i=1}^p \phi_i(u_t^*) r_{t-i} + \rho_t \right] \right| \right] + \left| \sum_{t=1}^T \epsilon_t \right| = \left| \sum_{t=1}^T r_t \right| + \left| \sum_{t=1}^T \epsilon_t \right|.$$

[Bacchiocchi et al. \[2024\]](#) proved that $\sum_{t=1}^T r_t \leq \left(\frac{\Gamma p}{1 - \Gamma} + 1 \right) \mathcal{Q}(\pi, T)$. We now create the bound for the sum of error terms in the following way:

$$\begin{aligned} \left| \sum_{t=1}^T \epsilon_t \right| &= \left| \sum_{t=1}^T (x_t - \hat{x}_t) \right| \leq \sum_{t=1}^T |x_t - \hat{x}_t| \\ &\stackrel{(a)}{=} \sum_{t=1}^T \mathcal{O} \left(t^{d-0.5} \log(t) \sqrt{\log \log(t)} \right) = \mathcal{O} \left(T^{d+0.5} \log(T) \sqrt{\log \log(T)} \right), \end{aligned}$$

where the passage (a) arises from the result for the bound of $\hat{x}_t - x_t$ obtained in Appendix B.3. \square

To derive the upper bound of regret, we also make use of the Elliptic Potential Lemma ([Lattimore and Szepesvári \[2020\]](#), Lemma 19.4) in our derivations of the bound of regret for our setting.

Lemma B.4. (*Elliptic Potential Lemma*). Let $\mathbf{V}_0 \in \mathbb{R}^{b \times b}$ be a positive definite matrix and let $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbb{R}^b$ be a sequence of vectors such that $\|\mathbf{u}_t\|_2 \leq L < +\infty$ for all $t \in \llbracket n \rrbracket$. Let $\mathbf{V}_t = \mathbf{V}_0 + \sum_{s=1}^t \mathbf{u}_s \mathbf{u}_s^\top$, then:

$$\sum_{t=1}^n \min\{1, \|\mathbf{u}_t\|_{\mathbf{V}_t^{-1}}\} \leq 2d \left(\frac{\text{tr}(\mathbf{V}_0) + nL^2}{b \det(\mathbf{V}_0)^{1/b}} \right).$$

Proof of Theorem 4.2. Let $\delta \in (0, 1)$, and define, as in the main paper, for every round $t \in \llbracket T \rrbracket$ and action $u \in \mathcal{U}$:

$$\beta_t(u) := \sqrt{\lambda(m^2 + 1)} + c_1(t) \sqrt{m^2 + 1} + c_2(t) + \sigma \sqrt{2 \log \left(\frac{1}{\delta} \right) + \log \left(\frac{\det(\hat{\mathbf{V}}_t(u))}{\lambda^{p+1}} \right)},$$

where $c_1(t) = \sqrt{pt^d \log(t)} \sqrt{\log(\log(t) + 1)}$ and $c_2(t) = \sqrt{p + 1} t^d \log(t) \sqrt{\log(\log(t) + 1)}$.

Let us define the confidence set $\mathcal{C}_t(u) := \{\phi \in \mathbb{R}^{p+1} : \|\phi - \hat{\phi}_{t-1}(u)\|_{\mathbf{V}_{t-1}(u)} \leq \beta_{t-1}(u)\}$ and the optimistic estimate of the parameter vector $\phi(u)$:

$$\tilde{\phi}_{t-1}(u) = \arg \max_{\phi \in \mathcal{C}_t(u)} \langle \phi, \hat{\mathbf{z}}_{t-1} \rangle.$$

By Theorem 4.1, we have that, for every action $u \in \mathcal{U}$ and round $t \in \llbracket T \rrbracket$, the true parameter vector satisfies $\phi(u) \in \mathcal{C}_t(u)$ with a probability of at least $1 - \delta$. Therefore, with the same probability, we have:

$$\begin{aligned} \langle \phi(u_t^*) - \phi(u_t), \hat{\mathbf{z}}_{t-1} \rangle &= \langle \phi(u_t^*) - \phi(u_t), \mathbf{z}_{t-1} \rangle + \langle \phi(u_t^*) - \phi(u_t), \hat{\mathbf{z}}_{t-1} - \mathbf{z}_{t-1} \rangle \\ &\leq 2\beta_{t-1}(u_t) \left(\|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}} + \mathcal{O} \left(\sqrt{pt^{d-0.5} \log(t)} \sqrt{\log \log(t)} \right) \right), \end{aligned}$$

where the first term is derived by Bacchiocchi et al. [2024] using the Cauchy-Schwartz inequality, for which it holds that $\langle \mathbf{v}, \mathbf{w} \rangle = \|\mathbf{v}\|_{\mathbf{V}_{t-1}(u)^{-1}} \cdot \|\mathbf{w}\|_{\mathbf{V}_{t-1}(u)^{-1}}$ for every couple of vectors \mathbf{v}, \mathbf{w} , and the second term follows the result of Lemma B.3.

We also introduce the following notion about the external regret derived by Bacchiocchi et al. [2024]:

$$\rho_t = \langle \phi(u_t^*) - \phi(u_t), \mathbf{z}_{t-1} \rangle \leq \|\mathbf{z}_{t-1}\|_2 + m.$$

By Lemma B.1 we have:

$$\|\mathbf{z}_{t-1}\|_2 \leq \sqrt{1 + p \left(\frac{m + \eta}{1 - \Gamma} \right)^2} := L,$$

where $\eta = \sqrt{2\sigma^2 \log(T/\delta)}$, and, consequently, we have:

$$\rho_t \leq L + m := C_1.$$

We then proceed as follows:

$$\begin{aligned} \rho_t &\leq 2 \min \left\{ C_1, \beta_{t-1}(u_t) \left(\|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}} + \mathcal{O} \left(\sqrt{pt^{d-0.5} \log(t)} \sqrt{\log \log(t)} \right) \right) \right\} \\ &\leq 2 \max\{C_1, \beta_{t-1}(u_t)\} \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}} + \mathcal{O} \left(\sqrt{pt^{d-0.5} \log(t)} \sqrt{\log \log(t)} \right) \right\} \\ &= 2 \max\{C_1, \beta_{t-1}(u_t)\} \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}} \right\} \\ &\quad + 2 \max\{C_1, \beta_{t-1}(u_t)\} \min \left\{ 1, \mathcal{O} \left(\sqrt{pt^{d-0.5} \log(t)} \sqrt{\log \log(t)} \right) \right\}. \end{aligned}$$

Summing all over $t \in \llbracket T \rrbracket$, we get the following bound on the cumulative external regret:

$$\begin{aligned} \mathcal{Q}(\text{ARLM-UCB}) &= \sum_{t=1}^T \rho_t \leq \sqrt{T \sum_{t=1}^T \rho_t^2} \\ &\leq 2 \max\{C_1, \beta_{T-1}\} \sqrt{T \sum_{t=1}^T \left[\min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}} \right\} \left(1 + \mathcal{O} \left(\sqrt{pt^{d-0.5} \log(t)} \sqrt{\log \log(t)} \right) \right) \right]^2}, \end{aligned}$$

with $\beta_{T-1} := \max_{u \in \mathcal{U}} \beta_{T-1}(u)$ passage about β_{T-1} holds since the sequence $\beta_t(u_t)$ is non-decreasing and thus each term can be bounded with their value at $t = T$. Furthermore, the last inequality follows from an application of Cauchy-Schwartz inequality and the following observation:

$$\begin{aligned} &\left[\min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}} \right\} \left(1 + \mathcal{O} \left(\sqrt{pt^{d-0.5} \log(t)} \sqrt{\log \log(t)} \right) \right) \right]^2 \\ &= \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}}^2 \right\} + \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}}^2 \right\} \mathcal{O} \left(pt^{2d-1} \log^2(t) \log \log(t) \right) \\ &\quad + \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}}^2 \right\} \mathcal{O} \left(\sqrt{pt^{d-0.5} \log(t)} \sqrt{\log \log(t)} \right). \end{aligned}$$

Applying The Elliptic Potential Lemma, [Bacchiocchi et al. \[2024\]](#) proved that

$$\sum_{t=1}^T \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}}^2 \right\} \leq 2n(p+1) \log \left(1 + \frac{TL^2}{n\lambda(p+1)} \right).$$

Similarly, we may show that

$$\begin{aligned} & \sum_{t=1}^T \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}}^2 \right\} \mathcal{O} \left(\sqrt{pt^{d-0.5}} \log(t) \sqrt{\log \log(t)} \right) \\ & \leq n(p+1) \log \left(1 + \frac{TL^2}{n\lambda(p+1)} \right) \mathcal{O} \left(\sqrt{pT^{d+0.5}} \log(T) \sqrt{\log \log(T)} \right), \end{aligned}$$

and

$$\begin{aligned} & \sum_{t=1}^T \min \left\{ 1, \|\mathbf{z}_{t-1}\|_{\mathbf{V}_{t-1}(u)^{-1}}^2 \right\} \mathcal{O} \left(pt^{2d-1} \log^2(t) \log \log(t) \right) \\ & \leq n(p+1) \log \left(1 + \frac{TL^2}{n\lambda(p+1)} \right) \mathcal{O} \left(pT^{2d} \log^2(T) \log \log(T) \right). \end{aligned}$$

Plugging in these results in the bound for cumulative external regret, we obtain the following:

$$\begin{aligned} \mathcal{Q}(\text{ARLM-UCB}) &= \sum_{i=1}^T \rho_t \leq \max\{C_1, \beta_{T-1}\} \sqrt{n(p+1) \log \left(1 + \frac{TL^2}{n\lambda(p+1)} \right)} \\ & \cdot \sqrt{T + \mathcal{O} \left(\sqrt{pT^{d+1.5}} \log(T) \sqrt{\log \log(T)} \right) + \mathcal{O} \left(pT^{2d+1} \log^2(T) \log \log(T) \right)} \\ &= \max\{C_1, \beta_{T-1}\} \sqrt{n(p+1) \log \left(1 + \frac{TL^2}{n\lambda(p+1)} \right) (T + \mathcal{O} \left(pT^{2d+1} \log^2(T) \log \log(T) \right))}. \end{aligned}$$

Finally, we bound the term β_{T-1} :

$$\begin{aligned} \beta_{T-1} &:= \sqrt{\lambda(m^2+1)+c_1(T-1)} \sqrt{m^2+1+c_2(T-1)} + \sigma \max_{u \in \mathcal{U}} \sqrt{2 \log \left(\frac{1}{\delta} \right) + \log \left(\frac{\det(\hat{\mathbf{V}}_{T-1}(u))}{\lambda^{p+1}} \right)}, \\ &\leq \sqrt{\lambda(m^2+1)+c_1(T-1)} \sqrt{m^2+1+c_2(T-1)} + \sigma \sqrt{2 \log \left(\frac{1}{\delta} \right) + (p+1) \log \left(\frac{\lambda(p+1)+TL^2}{\lambda(p+1)} \right)}. \end{aligned}$$

We then set $\delta = (2T)^{-1}$. By highlighting the dependence on m, p, σ, Γ , and T , we have:

$$\beta_{T-1} = \mathcal{O} \left(m \left(1 + T^d \log(T) \sqrt{(p+1) \log \log(T)} \right) + \sigma \sqrt{p+1} \right),$$

and

$$C_1 = \mathcal{O} \left(1 + \sqrt{p} \frac{m + \sigma}{1 - \Gamma} \right).$$

These results hold with probability $1 - 2\delta$:

$$\begin{aligned} \mathcal{Q}(\text{ARLM-UCB}) &= \sum_{i=1}^T \rho_t \leq \mathcal{O} \left(\frac{(m + \sigma) T^d \log(T) \sqrt{n(p+1) \log(\log(T))}}{1 - \Gamma} \right) \\ & \cdot \sqrt{T + \mathcal{O} \left(pT^{2d+1} \log^2(T) \log \log(T) \right)} \\ &= \left(\frac{(p+1)(m + \sigma) \sqrt{n} T^{2d+0.5} \log^2(T) \log \log(T)}{1 - \Gamma} \right). \end{aligned}$$

Finally, applying Lemma [B.3](#), this results in:

$$\hat{R}(\text{ARLM-UCB}, T) \leq \left(\frac{(p+1)^2(m + \sigma) \sqrt{n} T^{2d+0.5} \log^2(T) \log \log(T)}{(1 - \Gamma)^2} \right).$$

□