
Linear Bandits with Non-i.i.d. Noise

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 We study the linear stochastic bandit problem, relaxing the standard i.i.d. assumption
2 on the observation noise. As an alternative to this restrictive assumption,
3 we allow the noise terms across rounds to be sub-Gaussian but interdependent,
4 with dependencies that decay over time. To address this setting, we develop new
5 confidence sequences using a recently introduced reduction scheme to sequential
6 probability assignment, and use these to derive a bandit algorithm based on the
7 principle of optimism in the face of uncertainty. We provide regret bounds for
8 the resulting algorithm, expressed in terms of the decay rate of the strength of
9 dependence between observations. Among other results, we show that our bounds
10 recover the standard rates up to a factor of the mixing time for geometrically mixing
11 observation noise.

12 1 Introduction

13 The linear bandit problem (Abe and Long, 1999; Auer, 2003) is an instance of a multi-armed bandit
14 framework, where the expected reward is linear in the feature vector representing the chosen arm.
15 More concretely, it is a sequential decision-making problem, where an agent each round picks an arm
16 X_t , and receives a reward $Y_t = \langle \theta^*, X_t \rangle + \varepsilon_t$, with θ^* a fixed parameter unknown to the agent, and
17 ε_t zero-mean random noise. This framework has gained significant attention in the literature as it
18 yields analytic tools that can be applied to several concrete applications, such as online advertising
19 (Abe et al., 2003), recommendation systems (Li et al., 2010; Korkut and Li, 2021), and dynamic
20 pricing (Cohen et al., 2020).

21 A popular strategy to tackle linear bandits leverages the principle of *optimism in the face of uncertainty*,
22 via upper confidence bound (UCB) algorithms. The idea of optimism can be traced back to Lai and
23 Robbins (1985), and its application to linear bandits was already advanced by Auer (2003). Since
24 then, this approach has been improved and analysed by several works (Abbasi-Yadkori et al., 2011;
25 Lattimore and Szepesvári, 2020; Flynn et al., 2023). This class of methods requires constructing an
26 adaptive sequence of confidence sets that, with high probability, contain the true parameter θ^* . Each
27 round, the agent selects the arm maximising the expected reward under the most optimistic parameter
28 (in terms of reward) in the current confidence set. UCB-based algorithms have become popular as
29 they are often easy to implement and come with tight worst-case regret guarantees.

30 For a UCB algorithm to perform well, it is necessary that the confidence sets are tight, which can be
31 ensured by taking advantage of the structure of the problem. In this paper, our focus is on studying
32 various assumptions on the observation noise. A commonly studied situation is when $(\varepsilon_t)_{t \geq 0}$ consists
33 of a sequence of i.i.d. realisations of some bounded or sub-Gaussian random variable (see Lattimore
34 and Szepesvári, 2020, Chapter 20). Often, the standard analysis can be extended to the case in which
35 the realisation are not independent, but conditionally centred and sub-Gaussian (Abbasi-Yadkori
36 et al., 2011). Yet, in real-world settings, this assumption is often unrealistic, as one can expect the
37 presence of interdependencies among the noise at different rounds. For instance, in the context
38 of advertisement selection, the noise models the ensemble of external factors that influence the

39 user's choice on whether to click or not an ad. The i.i.d. assumption implies that across different
40 rounds these external factors are completely independent. In practice, the user choice will be affected
41 by temporally correlated events, such as recent browsing history or exposure to similar content.
42 Therefore, a more realistic assumption is to allow the dependencies to decay with time, rather than
43 being completely absent. This way to model dependencies, often referred to as *mixing*, is common to
44 study concentration for sums of non-i.i.d. random variables, with applications to machine learning
45 (Bradley, 2005; Mohri and Rostamizadeh, 2008; Abélès et al., 2025).

46 In the present paper we relax the assumption that the noise is conditionally zero-mean in the bandit
47 problem, and we allow for the presence of dependencies. Concretely, we replace the standard
48 conditionally sub-Gaussian setting with a more general formulation that accounts for conditional
49 dependence of the noise on the past, by introducing a natural notion of *mixing sub-Gaussianity*. Within
50 this context, we introduce a UCB algorithm for which we rigorously establish regret guarantees.
51 There are two key challenges for our approach: constructing a valid confidence sequence under
52 dependent noise, and deriving a regret upper bound for the UCB algorithm that we propose.

53 We derive the confidence sequence by adapting the *online-to-confidence-sets* technique to accommo-
54 date temporal dependencies in the noise. This approach, originally introduced by Abbasi-Yadkori
55 et al. (2011) and recently extended and improved (Jun et al., 2017; Lee et al., 2024; Clerico et al.,
56 2025), involves constructing an abstract online learning game whose regret guarantees can be turned
57 into a confidence sequence. To deal with the dependencies in the noise, we modify the standard
58 online-to-confidence-sets framework by introducing delays in the feedback received within the ab-
59 abstract online game. This approach is inspired by the recent work of Abélès et al. (2025) on extending
60 online-to-PAC conversions to non-i.i.d. mixing data sets in the context of deriving generalisation
61 bounds for statistical learning. There, a delayed-feedback trick similar to ours is employed to derive
62 statistical guarantees (generalisation bounds) from an abstract online learning game.

63 For the regret analysis of the bandit algorithm, we also need to face some challenges due to the
64 correlated observation noise. We address these by introducing delays into the decision-making policy
65 as well. This makes our approach superficially similar to algorithms used in the rich literature on
66 bandits with delayed feedback (see, e.g., Vernade et al., 2020a; Howson et al., 2023). These works
67 consider delay as part of the problem statement and not part of the solution concept, and are thus
68 orthogonal to our work. In particular, a simple adaptation of results from this literature would not
69 suffice for dealing with dependent observations, which we tackle by developing new concentration
70 inequalities. Another line of work that is conceptually related to ours is that of non-stationary bandits
71 (Garivier and Moulines, 2008; Russac et al., 2019). In that setting, the parameter vector θ_t^* evolves in
72 time according to a nonstationary stochastic process, and the observation noise remains i.i.d., once
73 again making for a rather different problem with its own challenges. Namely, the main obstacle
74 to overcome is that comparing with the optimal sequence of actions becomes impossible unless
75 strong assumptions are made about the sequence of parameter vectors. A typical trick to deal with
76 these nonstationarities is to discard old observations (which may have been generated by a very
77 different reward function), and use only recent rewards for decision-making. This is the polar opposite
78 of our approach that is explicitly *disallowed* to use recent rewards, which clearly highlights how
79 different these problems are. That said, there exists an intersection between the worlds of delayed
80 and nontationary bandits (Vernade et al., 2020b), and thus we would not discard the possibility of
81 eventually building a bridge between bandits with nonstationary reward functions and bandits with
82 nonstationary observation noise. For simplicity, we focus on the second of these two components in
83 this paper.

84 **Notation.** Throughout the paper, we will often use the following notations. For u and v in \mathbb{R}^p , we
85 let $\langle u, v \rangle$ denote their dot product. $\|u\|_2 = \sqrt{\langle u, u \rangle}$ is the Euclidean norm, while for a non-negative
86 definite $(p \times p)$ -matrix A , $\|u\|_A = \sqrt{\langle u, Au \rangle}$ is a semi-norm (a norm if the matrix is strictly positive
87 definite). For $r > 0$, $\mathcal{B}(r)$ denotes the closed centred Euclidean ball in \mathbb{R}^p with radius r . Given a
88 non-empty set $U \subseteq \mathbb{R}^p$, we let Δ_U denote the space of (Borel) probability measures on \mathbb{R}^p whose
89 support in U . Finally, $(u_t)_{t \geq t_0}$ denotes a sequence indexed on the integers, with t_0 its smallest index.

90 2 Preliminaries on linear bandits

91 We consider a version of the classic problem of regret minimisation in stochastic linear bandits, where
92 an agent needs to make a sequence of decisions (or pick an *arm*) from a given contextual decision set

93 that may change over the sequence of rounds. We assume that the environment is oblivious to the
 94 actions of the agent, in the sense that the decision sets are determined in advance, and do not depend
 95 neither on the realisations of the noise nor on the agent's arm-selection strategy.

96 Concretely, we define the problem as follows. Let $\theta^* \in \mathbb{R}^p$ be a parameter vector that is unknown
 97 to the learning agent. We assume as known an upper bound $B > 0$ on its euclidean norm (namely,
 98 $\theta^* \in \mathcal{B}(B)$). Fix a sequence of decision sets $(\mathcal{X}_t)_{t \geq 1}$ in \mathbb{R}^p . We assume that for all t we have
 99 $\mathcal{X}_t \subseteq \mathcal{B}(1)$. At each round t , the agent is required to pick an arm $X_t \in \mathcal{X}_t$, and receives the reward
 100 $Y_t = \langle \theta^*, X_t \rangle + \varepsilon_t$. The sequence $(\varepsilon_t)_{t \geq 1}$ represents the random feedback noise. The noise across
 101 different rounds is typically assumed to be conditionally centred and to have well behaved tails.
 102 For instance, a common assumption is to ask that $\mathbb{E}[\varepsilon_t | \mathcal{F}_{t-1}]$ is centred and sub-Gaussian, where
 103 $\mathcal{F}_t = \sigma(\varepsilon_1, \dots, \varepsilon_t)$ is the σ -field generated by the noise. This is the assumption this work relaxes.
 104 We also remark that, more generally, one can consider the case where the X_t as well are randomised,
 105 namely contain additional randomness that is not included in the noise. To take this into account, one
 106 can add this other source of randomness in the filtration. However, since in our case we will only
 107 consider a non-randomised bandit algorithm, we omit this to simplify our analysis.

108 The agent aims to find a good strategy to pick arms X_t that lead to a high expected T -round reward
 109 $\sum_{t=1}^T \langle X_t, \theta^* \rangle$. To compare their performance to that of an agent playing each round the best available
 110 arm (in expectation), we define the *regret* after T rounds as

$$\text{Reg}(T) = \sum_{t=1}^T \sup_{x \in \mathcal{X}_t} (\langle x, \theta^* \rangle - \langle X_t, \theta^* \rangle).$$

A common approach to tackle the linear bandit problem is to follow an *upper confidence bound* (UCB) strategy. This involves the following protocol. At each round t , we first derive a confidence set \mathcal{C}_{t-1} , based on the arm-reward pairs $(X_s, Y_s)_{s \leq t-1}$. This is a random set (as it depends on the past noise realisations), which must be constructed ensuring that $\theta^* \in \mathcal{C}_{t-1}$ with high probability. More precisely, the regret can be effectively controlled if one can ensure that θ^* uniformly belongs to every set $(\mathcal{C}_t)_{t \geq 1}$, with high probability (a property often referred to as *anytime validity*). Then, for every available arm x , we let

$$\text{UCB}_{\mathcal{C}_{t-1}}(x) = \max_{\theta \in \mathcal{C}_{t-1}} \langle x, \theta \rangle.$$

111 By definition, this is a high-probability upper bound on $\langle x, \theta^* \rangle$, which justifies the name “upper
 112 confidence bound”. The idea is then to *optimistically* pick as $X_t \in \mathcal{X}_t$ the arm maximising $\text{UCB}_{\mathcal{C}_{t-1}}$.

A key technical challenge in designing a UCB algorithm is to construct the anytime valid confidence sequence $(\mathcal{C}_t)_{t \geq 1}$. Typically, under sub-Gaussian assumptions on the noise, these sets take the form of an ellipsoid, centred on a (regularised) maximum likelihood estimator. Explicitly, we often have

$$\mathcal{C}_t = \{ \theta \in \Theta : \|\theta - \hat{\theta}_t\|_{V_t}^2 \leq \beta_t^2 \},$$

113 where $\hat{\theta}_t$ is the least-squares estimator of θ^* , V_t is the *feature-covariance* matrix and β_t is a radius
 114 carefully chosen so that the high-probability coverage requirement is satisfied. In this work, to
 115 construct the confidence sets we will leverage an *online-to-confidence-set-conversion* approach, a
 116 method that reduces the problem of proving statistical concentration bounds to proving existence of
 117 well-performing algorithms for an associated game of *sequential probability assignment*. We refer to
 118 Section 4 for more details on our technique to construct the confidence sequence.

119 3 Linear bandits with non-i.i.d. observation noise

120 We study a variant of the standard linear stochastic bandit problem where the observation-noise
 121 variables feature dependencies across different rounds. We focus on the case of weakly stationary
 122 noise, meaning we assume all the ε_t to have the same marginal distribution. However, the core
 123 assumption we make is what we call *mixing sub-Gaussianity*. This provides a way to control how
 124 dependencies decay as the time between two observations increases. It is defined in terms of a
 125 sequence of mixing coefficients ϕ_d , which quantify this decay.

126 **Assumption 1** (Mixing sub-Gaussianity). *Fix $\sigma > 0$ and let $\phi = (\phi_d)_{d \geq 0}$ be a non-negative and
 127 non-increasing sequence. We say that the random sequence $(\varepsilon_t)_{t \geq 1}$ is (σ, ϕ) -mixing sub-Gaussian if
 128 ε_t is centred and σ -sub-Gaussian for every t , and, for all $d \geq 0$ and all $t > d$, we have*

$$\mathbb{E}[\varepsilon_t | \mathcal{F}_{t-d}] \leq \phi_d \tag{1}$$

129 and

$$\mathbb{E} [\exp \lambda (\epsilon_t - \mathbb{E} [\epsilon_t | \mathcal{F}_{t-d}]) | \mathcal{F}_{t-d}] \leq e^{\frac{\lambda^2 \sigma^2}{2}}, \quad \forall \lambda > 0. \quad (2)$$

130 Clearly, the above assumption generalises the standard conditionally sub-Gaussian assumption (that
131 can be recovered by setting $\phi_d = 0$ for all t), sometimes considered in the bandit literature. Although
132 this might look like an unusual mixing assumption, it is very natural for our problem at hand, and
133 can be weaker than standard mixing hypotheses. For instance, if the noise sequence is φ -mixing
134 (see Bradley, 2005) and each ε_t is centred and bounded in $[-a, b]$, it is straightforward to check that
135 $|\mathbb{E}[\varepsilon_t | \mathcal{F}_{t-d}]| \leq (a+b)\phi_d$, and so Assumption 1 is satisfied since the boundedness automatically
136 implies sub-Gaussianity. In the rest of the paper we assume $\sigma = 1$ for simplicity.

137 Under Assumption 1, we can build the confidence sequence needed for our UCB algorithm. We state
138 this result below, but defer the explicit derivation to Section 4 (see Corollary 1 there).

Proposition 1. *For some given ϕ , let the noise satisfy Assumption 1 with $\sigma = 1$. Fix $\delta \in (0, 1)$,
 $\lambda > 0$, and $d \geq 1$. For $t \geq 1$ let*

$$\mathcal{C}_t = \left\{ \theta \in \mathcal{B}(B) : \frac{1}{2} \|\theta - \hat{\theta}_t\|_{V_t}^2 \leq \frac{dp}{2} \log \frac{(B+1)^2 e \max(dp, t+d)}{dp} + 2\lambda B^2 + t\phi_d(B+1) + d \log \frac{d}{\delta} \right\},$$

where $V_t = \sum_{s=1}^t X_t X_t^\top + \lambda \text{Id}$, and $\hat{\theta}_t = \arg \min_{\theta \in \mathcal{B}(B)} \sum_{s=1}^t (\langle \theta, X_t \rangle - Y_t)^2$. Then, $(\mathcal{C}_t)_{t \geq 1}$ is
an anytime valid confidence sequence, in the sense that

$$\mathbb{P}(\theta^* \in \mathcal{C}_t, \forall t \geq 1) \geq 1 - \delta.$$

139 Leveraging the confidence sequence above, we can define a UCB approach for our problem (Algo-
140 rithm 1). At a high level, the algorithm operates by taking the confidence sets defined in Proposition
141 1, and selecting the arm optimistically, as in the standard UCB. A key point is that a delay d is
142 introduced, which at round t restricts the agent to use only the information available from the first
143 $t - d$ rounds. Although the actual technical reason behind this restriction will become fully clear only
144 with the analysis of the coming sections, one can intuitively think of it as a way to prevent overfitting
145 to recent noise, which might be highly correlated. If d is sufficiently large, the noise observed in
146 each round t will be sufficiently decorrelated from the previous observations, which allows accurate
147 estimation and uncertainty quantification of the true parameter θ^* and the associated rewards.

Algorithm 1 Mixing-LinUCB

```

set  $d > 0$ 
for  $i \in \{1, 2, \dots, d\}$  do
    play an arbitrary  $X_i$  and observe  $Y_i$ 
end for
for  $t \in \{d+1, \dots\}$  do
     $X_t = \arg \max_{x \in \mathcal{X}_t} \text{UCB}_{\mathcal{C}_{t-d}}(x)$ , where  $\mathcal{C}_{t-d}$  is as in Proposition 1
    play  $X_t$  and observe reward  $Y_t$ 
end for

```

148 In Section 5 we provide a detailed analysis of the regret of the algorithm that we proposed. For
149 instance, assuming that the mixing coefficients decay exponentially as $\phi_d = C e^{-d/\tau}$ (geometric
150 mixing), we show that the regret can be upper bounded in high probability as

$$\text{Reg}(T) \leq \mathcal{O} \left(\tau p \sqrt{T} \log(T)^2 + \tau \log T \sqrt{pT \log T} \right).$$

151 We refer to Theorem 2 and Corollary 2 in Section 5 for more details.

152 **4 Constructing the confidence sequence**

153 In this section we derive a confidence sequence for linear models with non-i.i.d. noise. First, we
154 briefly describe the online-to-confidence-set conversion scheme from Clerico et al. (2025), which
155 serves as our starting point. We then extend this technique to handle mixing noise.

156 **4.1 Online-to-confidence set conversion for i.i.d. data**

157 Before proceeding for the analysis of mixing sub-Gaussian noise, which is the focus of this work,
158 we start by describing how to derive a confidence sequence when the noise is independent (or
159 conditionally) centred and sub-Gaussian across different rounds, as in Clerico et al. (2025). The
160 online-to-confidence sets framework that we consider instantiates an abstract game played between
161 an *online learner* and an *environment*. We define the squared loss $\ell_s(\theta) = \frac{1}{2}(\langle \theta, X_s \rangle - Y_s)^2$. For
162 each round $s = 1, \dots, t$, the following steps are repeated:

- 163 1. the environment reveals X_s to the learner;
- 164 2. the learner plays a distribution $Q_s \in \Delta_{\mathbb{R}^p}$;
- 165 3. the environment reveals Y_s to the learner;
- 166 4. the learner suffers the log loss $\mathcal{L}_s(Q_s) = -\log \int_{\mathbb{R}^p} \exp(-\ell_s(\theta)) dQ_s(\theta)$.

167 This game is a special case of a well-studied problem called *sequential probability assignment*
168 (Cesa-Bianchi and Lugosi, 2006). The learner can use any strategy to choose Q_1, \dots, Q_t , as long as
169 each Q_s depends only on $X_1, Y_1, \dots, X_{s-1}, Y_{s-1}, X_s$. We define the *regret* of the learner against a
170 (possibly data-dependent) comparator $\bar{\theta} \in \mathbb{R}^p$ as

$$\text{Regret}_t(\bar{\theta}) = \sum_{s=1}^t \mathcal{L}_s(Q_s) - \sum_{s=1}^t \ell_s(\bar{\theta}).$$

171 Clerico et al. (2025) provide a regret bound upper bound (Proposition 3.1 there) for when the learner's
172 strategy is from an *exponential weighted average* (EWA) forecaster with a centred Gaussian prior
173 Q_1 . However, to account for the presence of dependencies in our analysis, we will need the prior's
174 support to be bounded. We hence state here a regret bound (whose proof is deferred to Appendix
175 A.1) for the regret of an EWA forecaster with a uniform prior.

176 **Proposition 2.** Fix $B > 0$ and consider the EWA forecaster with as prior the uniform distribution on
177 $\mathcal{B}(B + 1)$. Then, for all $\bar{\theta} \in \mathcal{B}(B)$ and any $t \geq 1$,

$$\text{Regret}_t(\bar{\theta}) \leq \frac{p}{2} \log \frac{(B + 1)^2 e \max(p, t)}{p}.$$

178

179 We remark that, by adding and subtracting the total log loss of the learner, the excess loss of θ^*
180 (relative to $\bar{\theta}$) can be rewritten as

$$\sum_{s=1}^t \ell_s(\theta^*) - \sum_{s=1}^t \ell_s(\bar{\theta}) = \text{Regret}_t(\bar{\theta}) + \sum_{s=1}^t \ell_s(\theta^*) - \sum_{s=1}^t \mathcal{L}_s(Q_s). \quad (3)$$

181 This simple decomposition is the key idea in the online-to-confidence sets scheme.

Since the noise is conditionally sub-Gaussian and the distributions played by the online learner
are predictable (Q_s cannot depend on Y_s), $\sum_{s=1}^t \ell_s(\theta^*) - \sum_{s=1}^t \mathcal{L}_s(Q_s)$ is the logarithm of a non-
negative super-martingale (cf. the no-hypercompression inequality in Grünwald, 2007 or Proposition
2.1 in Clerico et al., 2025) with respect to the noise filtration $(\mathcal{F}_t)_{t \geq 1}$. For simplicity, as already
mentioned in Section 2 and since this will be the case for our bandit strategy, we assume throughout
the paper that X_t is fully determined given the past noise. Henceforth, from Ville's inequality (a
classical anytime valid Markov-like inequality that holds for non-negative super-martingales) one can
easily derive that $\theta^* \in \mathcal{C}_t$ (uniformly for all t) with probability at least $1 - \delta$, where

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^p : \sum_{s=1}^t \ell_s(\theta) - \sum_{s=1}^t \ell_s(\bar{\theta}) \leq \text{Regret}_t(\bar{\theta}) + \log \frac{1}{\delta} \right\}.$$

182 This result can be relaxed by replacing $\text{Regret}_t(\bar{\theta})$ by any known regret upper bound for the online
183 algorithm used in the abstract game (e.g., the bound of Proposition 2 for the EWA forecaster).

184 **4.2 Confidence sequence under mixing sub-Gaussian noise**

185 The standard online-to-confidence sets scheme relies on the fact that $\sum_{s=1}^t \ell_s(\theta^*) - \sum_{s=1}^t \mathcal{L}_s(Q_s)$ is
 186 the logarithm of a non-negative super-martingale, whose fluctuations can be controlled uniformly in
 187 time thanks to Ville's inequality. However, this property hinges on the fact that the noise is assumed
 188 to be conditionally centred and sub-Gaussian, which now is not anymore the case. Yet, thanks to
 189 our mixing assumption, if we restrict our focus on rounds that are sufficiently far apart, the mutual
 190 dependencies get weaker, and the exponential of the sum behaves *almost* like a martingale. This
 191 insight suggests to partition the rounds into blocks, whose elements are mutually far apart, then apply
 192 concentration results to each block, and finally use a union bound to recover the desired confidence
 193 sequence spanning all rounds. We remark that this is a classical approach to derive concentration
 194 results for mixing processes, often referred to as the *blocking* technique (Yu, 1994).

195 In order for the online-to-confidence sets scheme to leverage the blocking strategy outlined above,
 196 the abstract online game used for the analysis must be designed in a way that is compatible with
 197 the block structure. To address this point, we adopt an approach inspired by Abélès et al. (2025),
 198 who introduced delays in the feedbacks received by the online learner in order to address a similar
 199 challenge. More precisely, we will now consider the following *delayed-feedback* version of the online
 200 game. Fix a delay $d > 0$. For each round $s = 1, \dots, t$, the following steps are repeated:

- 201 1. the environment reveals to the learner X_s , which is assumed to be \mathcal{F}_{s-d} -measurable;
- 202 2. the learner plays a distribution $Q_s \in \Delta_{\mathbb{R}^p}$;
- 203 3. if $s > d$, the environment reveals Y_{s-d+1} to the learner;
- 204 4. the learner suffers the log loss $\mathcal{L}_s(Q_s) = -\log \int_{\mathbb{R}^p} \exp(-\ell_s(\theta)) dQ_s(\theta)$.

205 Note that the delay d only applies for the rewards, while Q_s can still depend on X_s . Indeed, the choice
 206 of X_s in our mixing UCB algorithm is already “delayed”, as it depends on \mathcal{C}_{t-d} (see Algorithm 1).
 207 Of course, in this setting the decomposition of (3) is still valid. We now want to deal with the
 208 concentration of $\sum_{s=1}^t \ell_s(\theta^*) - \sum_{s=1}^t \mathcal{L}_s(Q_s)$ via the blocking technique. For convenience, let
 209 us write $D_t = \ell_t(\theta^*) - \mathcal{L}_t(Q_t)$. We denote as $S^{(i)} = (S_k^{(i)})_{k \geq 1}$ the subsequence defined as
 210 $S_k^{(i)} = \sum_{j=1}^k D_{i+(j-1)d}$. The key idea is now that each of these $S^{(i)}$ behaves as the log of a
 211 martingale, up to a cumulative remainder that accounts for the conditional mean shift in the mixing
 212 sub-Gaussianity assumption. In particular, Ville's inequality and a union bound yield the following.

Lemma 1. *Fix a delay $d > 0$ and $\delta \in (0, 1)$. We have that*

$$\mathbb{P} \left(\sum_{s=1}^t (\ell_s(\theta^*) - \mathcal{L}_s(Q_s)) \leq t\phi_d B + d \log \frac{d}{\delta}, \forall t \geq 1 \right) \geq 1 - \delta.$$

213 Now that we have a concentration result to control S_t , we only need to be able to upper bound the
 214 regret of an algorithm for the “delayed” online game that we are considering. To this purpose, we
 215 propose the following approach. We run d independent EWA forecaster (with uniform prior), each
 216 one only making prediction and receiving the feedback once every d rounds. More explicitly, the first
 217 forecaster acts at rounds $1, d+1, 2d+1, \dots$, the second at round $2, d+2, 2d+2, \dots$, and so on. As a
 218 direct consequence of Proposition 2, by summing the individual regret upper bounds we get a regret
 219 bound for the joint forecaster, which at each round returns the distribution predicted by the currently
 220 active forecaster. This technique of partitioning rounds into blocks for the regret analysis of online
 221 learning is common in the literature (e.g., see Weinberger and Ordentlich, 2002).

222 **Lemma 2.** *Fix $B > 0$, $d > 0$, and consider a strategy with d independent EWA forecasters outlined
 223 above, all initialised with the uniform distribution on $\mathcal{B}(B+1)$ as prior. For all $\theta \in \mathcal{B}(B)$ and $t \geq 1$,*

$$\text{Regret}_t(\bar{\theta}) \leq \frac{dp}{2} \log \frac{(B+1)^2 e \max(dp, t+d)}{dp}.$$

224 Putting together what we have, we get a confidence sequence suitable for our mixing UCB algorithm.

Theorem 1. *Consider the setting introduced above. Fix $\delta \in (0, 1)$ and a delay $d > 0$. Assume as
 known that $\theta^* \in \mathcal{B}(B)$. Let $\hat{\theta}_t = \arg \min_{\theta \in \mathcal{B}(B)} \{\sum_{s=1}^t \ell_s(\theta)\}$ and $\Lambda_t = \sum_{s=1}^t X_s X_s^\top$. Define*

$$\mathcal{C}_t = \left\{ \theta \in \mathcal{B}(B) : \frac{1}{2} \|\theta - \hat{\theta}_t\|_{\Lambda_t}^2 \leq \frac{dp}{2} \log \frac{(B+1)^2 e \max(dp, t+d)}{dp} + t\phi_d(B+1) + d \log \frac{d}{\delta} \right\}.$$

Then, $(\mathcal{C}_t)_{t \geq 1}$ is an anytime valid confidence sequence for θ^* , namely

$$\mathbb{P}(\theta^* \in \mathcal{C}_t, \forall t \geq 1) \leq 1 - \delta.$$

Proof. The optimality of $\widehat{\theta}_t$ implies $\sum_{s=1}^t \langle \theta - \widehat{\theta}_t, \nabla \ell_s(\widehat{\theta}_t) \rangle \geq 0$, for all $\theta \in \mathcal{B}(B)$. As $\sum_{s=1}^t \ell_s$ is quadratic, it equals its second order Taylor expansion around $\widehat{\theta}_t$ and its Hessian is everywhere Λ_t . So,

$$\frac{1}{2} \|\theta - \widehat{\theta}_t\|_{\Lambda_t}^2 \leq \frac{1}{2} \|\theta - \widehat{\theta}_t\|_{\Lambda_t}^2 + \sum_{s=1}^t \langle \theta - \widehat{\theta}_t, \nabla \ell_s(\widehat{\theta}_t) \rangle = \sum_{s=1}^t (\ell_s(\theta) - \ell_s(\widehat{\theta}_t)),$$

for any $\theta \in \mathcal{B}(B)$. This, together with (3), Lemma 1, and Lemma 2, yields the conclusion. \square

We remark that the confidence sets of Theorem 1 take the form of the intersection between the ball $\mathcal{B}(B)$ and the “ellipsoid” $\{\theta : \|\theta - \widehat{\theta}_t\|_{\Lambda_t} \leq \beta_t\}$, for a suitable radius β_t . In order to implement and analyse the bandit algorithm, it will be more convenient to work with a relaxation of these sets, a pure ellipsoid not intersected with $\mathcal{B}(B)$. We make this explicit in the following corollary.

Corollary 1. *Fix $\lambda > 0$, $d > 0$, and $\delta \in (0, 1)$. For $t \geq 1$, let $V_t = \Lambda_t + \lambda \text{Id}$. Assuming that $\theta^* \in \mathcal{B}(B)$, the following compact ellipsoids define an anytime valid confidence sequence for θ^* :*

$$\mathcal{C}_t = \left\{ \theta \in \mathcal{B}(B) : \frac{1}{2} \|\theta - \widehat{\theta}_t\|_{V_t}^2 \leq \frac{dp}{2} \log \frac{(B+1)^2 e \max(dp, t+d)}{dp} + 2\lambda B^2 + t\phi_d(B+1) + d \log \frac{d}{\delta} \right\}.$$

Proof. Let $\beta_t^2 = dp \log \frac{(B+1)^2 e \max(dp, t+d)}{dp} + 2t\phi_d(B+1) + 2d \log \frac{d}{\delta}$. From Theorem 1, with probability at least $1 - \delta$, uniformly for every t , $\|\theta^* - \widehat{\theta}_t\|_{\Lambda_t}^2 \leq \beta_t^2$. Adding to both sides of this inequality $\frac{\lambda}{2} \|\theta^* - \widehat{\theta}_t\|_2^2$, and relaxing the RHS using that $\|\theta^* - \widehat{\theta}_t\|_2^2 \leq 4B^2$, we conclude. \square

5 Regret bounds for Mixing-LinUCB

In this section, we establish worst-case and gap-dependent cumulative regret bounds for mixing UCB algorithm (Mixing Lin-UCB). However, to account for the fact that Mixing-LinUCB selects actions with delays, the standard elliptical potential arguments must be modified. Throughout this section, we let $R_t = \langle \theta^*, X_t^* - X_t \rangle$ (where $X_t^* = \arg \max_{x \in \mathcal{X}_t} \langle \theta^*, x \rangle$) denote the regret in round t , and $\beta_t^2 = dp \log \frac{(B+1)^2 e \max(dp, t+d)}{dp} + 4\lambda B^2 + 2t\phi_d(B+1) + 2d \log \frac{d}{\delta}$ denote the squared radius of the ellipsoid \mathcal{C}_t in Corollary 1.

5.1 Worst-case regret bounds

First, following the regret analysis in Abbasi-Yadkori et al. (2011) (see also Section 19.3 in Lattimore and Szepesvári, 2020), we upper bound the instantaneous regret. From our boundedness assumptions ($\theta^* \in \mathcal{B}(B)$ and $\mathcal{X}_t \subseteq \mathcal{B}(1)$), we easily deduce that $R_t \leq 2B$. Under the event that our confidence sequence contains θ^* at every step t , we have another bound on R_t . If we define $\widetilde{\theta}_{t-d} \in \mathcal{C}_{t-d}$ to be the point at which $\langle \widetilde{\theta}_{t-d}, X_t \rangle = \text{UCB}_{\mathcal{C}_{t-d}}(X_t)$, then from the definition of X_t we have

$$\langle \theta^*, X_t^* \rangle \leq \max_{x \in \mathcal{X}_t} \max_{\theta \in \mathcal{C}_{t-d}} \langle \theta, x \rangle = \max_{x \in \mathcal{X}_t} \text{UCB}_{\mathcal{C}_{t-d}}(x) = \text{UCB}_{\mathcal{C}_{t-d}}(X_t) = \langle \widetilde{\theta}_{t-d}, X_t \rangle.$$

Recall that, for all s , $V_s = \Lambda_s + \lambda \text{Id}$, which is invertible as $\lambda > 0$. Thus, by Cauchy-Schwarz,

$$R_t \leq \langle \widetilde{\theta}_{t-d} - \theta^*, X_t \rangle \leq \|\widetilde{\theta}_{t-d} - \theta^*\|_{V_{t-d}} \|X_t\|_{V_{t-d}^{-1}} \leq 2\beta_{t-d} \|X_t\|_{V_{t-d}^{-1}}.$$

This means that the instantaneous regret satisfies the bound

$$R_t \leq 2 \max(B, \beta_{t-d}) \min(1, \|X_t\|_{V_{t-d}^{-1}}). \quad (4)$$

248 Next, we separate the regret suffered in the first d rounds and the remaining $T - d$ rounds. We then
249 use Cauchy-Schwarz once more, and the fact that β_t is increasing in t , to obtain

$$\begin{aligned}\text{Reg}(T) &\leq 2dB + \sqrt{(T-d)\sum_{t=d+1}^T R_t^2} \\ &\leq 2dB + \sqrt{4(T-d) \max(B^2, \beta_{T-d}^2) \sum_{t=d+1}^T \min(1, \|X_t\|_{V_{t-d}^{-1}}^2)}.\end{aligned}$$

250 At this point, we must depart from the standard linear UCB analysis (Abbasi-Yadkori et al., 2011; Latti-
251 more and Szepesvári, 2020). We bound the sum of the *elliptical potentials* $\sum_{t=d+1}^T \min(1, \|X_t\|_{V_{t-d}^{-1}}^2)$
252 using the following variant of the well-known “elliptical potential lemma” (see Appendix), which
253 accounts for the fact that the feature covariance matrix V_{t-d} is updated with a delay of d steps.

254 **Lemma 3.** *For all $T \geq d + 1$,*

$$\sum_{t=d+1}^T \min(1, \|X_t\|_{V_{t-d}^{-1}}^2) \leq 2dp \log(1 + \frac{T}{\lambda dp}).$$

255

256 We can now state a worst-case regret upper bound for Mixing-LinUCB.

257 **Theorem 2.** *Fix $\lambda = 1/B^2$, $d > 0$ and $\delta \in (0, 1)$. With probability at least $1 - \delta$, for all $T > d$, the
258 regret of Mixing-LinUCB satisfies*

$$\text{Reg}(T) \leq 2dB + \sqrt{8dpT \max(B^2, \beta_T^2) \log(1 + \frac{B^2 T}{dp})}.$$

259

260 From the definition of β_T , we see that this regret bound is of the order

$$\text{Reg}(T) = \mathcal{O}\left(dB + dp\sqrt{T} \log \frac{TB}{dp} + T\sqrt{Bdp\phi_d \log \frac{TB}{dp}} + d\sqrt{pT \log \frac{TB}{p\delta}}\right).$$

261 For any fixed (*i.e.*, not depending on T) delay d , this regret bound is linear in T . To obtain meaningful
262 regret bounds, it is therefore crucial to set d as a function of T and the rate at which the mixing
263 coefficients decay to zero. We point out that if T is unknown, one could probably use a more
264 general framework where the delay is time dependent which might lead to non-trivial results, but we
265 do not pursue this here. Under the assumption that the noise variables are either geometrically or
266 algebraically mixing, we obtain the following worst-case regret bounds.

267 **Corollary 2.** *Suppose that the noise satisfies Assumption 1 with $\phi_d = Ce^{-\frac{d}{\tau}}$ for some $C, \tau > 0$
268 (geometric mixing), and set $d = \lceil \tau \log \frac{BCT}{p} \rceil$. Then, the regret of Mixing-LinUCB satisfies*

$$\text{Reg}(T) = \mathcal{O}\left(\tau p\sqrt{T} \left(\log \frac{TB \max(1, C)}{p}\right)^2 + p\sqrt{T\tau} \log \frac{TB \max(1, C)}{p} + \tau \log \frac{BCT}{p} \sqrt{pT \log \frac{TB}{p\delta}}\right).$$

269 **Corollary 3.** *Suppose that the noise satisfies Assumption 1 with $\phi_d = Cd^{-r}$ for some $C > 0$ and
270 $r > 0$ (algebraic mixing), and set $d = \lceil CT^{1/(1+r)} \rceil$. Then, the regret of Mixing-LinUCB satisfies*

$$\text{Reg}(T) = \mathcal{O}\left(CBT^{1/(1+r)} + CT^{\frac{3+r}{2(1+r)}} \left(p \log \frac{TB}{p} + \sqrt{Bp \log \frac{T^{r/(1+r)} B}{Cp}} + \sqrt{p \log \frac{TB}{p\delta}}\right)\right).$$

271

272 Up to a factor of $\tau \log T$, the bound for geometrically mixing noise matches the regret bound for
273 linear UCB with i.i.d. noise. This bound is trivial for $r \leq 1$, however for $r > 1$ we get sublinear
274 regret, and in particular we recover standard rates up to logarithmic factors in the limit where $r \rightarrow \infty$.

275 5.2 Gap-dependent regret bounds

276 Under the assumption that, each round, the gap between the expected reward of the optimal arm and
277 the expected reward of any other arm is at least $\Delta > 0$, we get regret bounds with better dependence

278 on T . More precisely, define the *minimum gap* $\Delta = \min_{t \in [T]} \min_{x \in \mathcal{X}_t: x \neq x_t^*} \langle X_t^* - x, \theta^* \rangle$, and
279 assume that $\Delta > 0$. Since we either have $R_t = 0$ or $R_t \geq \Delta > 0$, it follows that

$$R_t \leq R_t^2 / \Delta.$$

280 In our worst-case analysis, we showed that

$$\sum_{t=d+1}^T R_t^2 \leq 8dp \max(B^2, \beta_T^2) \log\left(1 + \frac{T}{\lambda dp}\right).$$

281 Combined with the previous inequality, we obtain the following gap-dependent regret bound.

282 **Theorem 3.** *Fix $\lambda = 1/B^2$, $d > 0$, and $\delta \in (0, 1)$. With probability at least $1 - \delta$, for all $T > d$, the
283 regret of Mixing-LinUCB satisfies*

$$\text{Reg}(T) \leq 2dB + \frac{8dp}{\Delta} \max(B^2, \beta_T^2) \log\left(1 + \frac{B^2 T}{dp}\right).$$

284

285 Similarly to the worst-case bound in Theorem 2, for any fixed $d > 0$, this regret bound is linear in T .
286 By setting d as a suitable function of T , we obtain the following gap-dependent regret bounds under
287 geometrically or algebraically mixing noise.

288 **Corollary 4.** *Suppose that the noise variables are geometrically mixing and set $d = \lceil \tau \log \frac{BCT}{p} \rceil$.
289 Then the regret of Mixing-LinUCB satisfies*

$$\text{Reg}(T) = \mathcal{O}\left(\frac{8\tau p}{\Delta} \left(\log \frac{BCT}{p}\right)^2 \log\left(1 + \frac{B^2 T}{p\tau \log \frac{BCT}{p}}\right) \left(\frac{p}{2} \log \frac{T}{p\tau} + \log \frac{\tau \log \frac{BCT}{p}}{\delta}\right)\right).$$

290

291 **Corollary 5.** *Suppose that the noise variables are algebraically mixing and set $d = \lceil CT^{1/(1+r)} \rceil$.
292 Then the regret of Mixing-LinUCB satisfies*

$$\text{Reg}(T) = \mathcal{O}\left(\frac{8Cp}{\Delta} T^{\frac{2}{1+r}} \log\left(1 + \frac{B^2 T}{pCT^{1/(1+r)}}\right) \left(\frac{p}{2} \log \frac{(B+1)^2 eT}{p} + \log \frac{CT^{1/(1+r)}}{\delta}\right)\right).$$

293

294 6 Conclusion

295 We leave several interesting questions open for future research. Some of these are listed below.

296 An important limitation of our algorithm is that it requires the knowledge of the mixing coefficients
297 (or at least an upper-bound on them). It would be interesting to explore the possibility of relaxing
298 this assumption and to design an algorithm which infers the mixing coefficients while minimizing
299 the regret. We note that the problem of estimating mixing coefficients is already a hard problem on
300 its own right, with tight sample-complexity results only available in special cases such as Markov
301 chains (Hsu et al., 2019; Wolfer, 2020). We also note that in order to recover the standard rate for the
302 regret bound, the delay d introduced in our algorithm need to be chosen as a function of the horizon
303 T . We believe that this could be fixed at little conceptual expense by using time-varying delay in the
304 analysis, but we did not attempt to work out the (potentially non-trivial) details here.

305 Another limitation is that our analysis assumed throughout that the adversary picking the decision sets
306 \mathcal{X}_t is oblivious, which is typically not required in linear bandit problems. For us, this was necessary
307 to avoid potential statistical dependence between decision sets and the nonstationary observations.
308 We believe that this issue can be handled at least for some classes of adversaries. For instance, it
309 is easy to see that our analysis would carry through under the assumption that the decision sets be
310 selected based on delayed information only. We leave the investigation of this question under more
311 realistic assumptions open for future work.

312 **References**

313 Naoki Abe and Philip M. Long. Associative reinforcement learning using linear probabilistic concepts.
314 In *Proceedings of the Sixteenth International Conference on Machine Learning*, 1999.

315 Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *J. Mach. Learn. Res.*, 3:397–422, 2003.

317 Naoki Abe, Alan W. Biermann, and Philip M. Long. Reinforcement learning with immediate rewards
318 and linear hypotheses. *Algorithmica*, 37(4):263–293, 2003.

319 Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
320 personalized news article recommendation. In *Proceedings of the 19th international conference on*
321 *World wide web*, pages 661–670, 2010.

322 Melda Korkut and Andrew Li. Disposable linear bandits for online recommendations. *Proceedings*
323 *of the AAAI Conference on Artificial Intelligence*, 35(5), 2021.

324 Maxime C Cohen, Ilan Lobel, and Renato Paes Leme. Feature-based dynamic pricing. *Management*
325 *Science*, 66(11):4921–4943, 2020.

326 T.L. Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied*
327 *Mathematics*, 6(1):4–22, 1985.

328 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic
329 bandits. *Advances in neural information processing systems*, 24, 2011.

330 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

331 Hamish Flynn, David Reeb, Melih Kandemir, and Jan R Peters. Improved algorithms for stochastic
332 linear bandits using tail bounds for martingale mixtures. *Advances in Neural Information*
333 *Processing Systems*, 36:45102–45136, 2023.

334 Richard C. Bradley. Basic properties of strong mixing conditions: A survey and some open questions.
335 *Probability Surveys*, 2:107–144, 2005.

336 M. Mohri and A. Rostamizadeh. Rademacher complexity bounds for non-i.i.d. processes. *NeurIPS*,
337 2008.

338 Baptiste Abélès, Eugenio Clerico, and Gergely Neu. Generalization bounds for mixing processes
339 via delayed online-to-PAC conversions. In *Proceedings of The 36th International Conference on*
340 *Algorithmic Learning Theory*, 2025.

341 Kwang-Sung Jun, Aniruddha Bhargava, Robert Nowak, and Rebecca Willett. Scalable generalized
342 linear bandits: Online computation and hashing. In *Advances in Neural Information Processing*
343 *Systems*, volume 30, 2017.

344 Junghyun Lee, Se-Young Yun, and Kwang-Sung Jun. Improved regret bounds of (multinomial)
345 logistic bandits via regret-to-confidence-set conversion. In *Proceedings of the 27th International*
346 *Conference on Artificial Intelligence and Statistics*, pages 4474–4482, 2024.

347 Eugenio Clerico, Hamish Flynn, Wojciech Kotłowski, and Gergely Neu. Confidence sequences
348 for generalized linear models via regret analysis, 2025. URL <https://arxiv.org/abs/2504.16555>.

350 Claire Vernade, Alexandra Carpentier, Tor Lattimore, Giovanni Zappella, Beyza Ermis, and Michael
351 Brueckner. Linear bandits with stochastic delayed feedback. In *International Conference on*
352 *Machine Learning*, pages 9712–9721. PMLR, 2020a.

353 Benjamin Howson, Ciara Pike-Burke, and Sarah Filippi. Delayed feedback in generalised linear
354 bandits revisited. In *International Conference on Artificial Intelligence and Statistics*, pages
355 6095–6119. PMLR, 2023.

356 Aurélien Garivier and Eric Moulines. On upper-confidence bound policies for non-stationary bandit
357 problems. *arXiv preprint arXiv:0805.3415*, 2008.

358 Yoan Russac, Claire Vernade, and Olivier Cappé. Weighted linear bandits for non-stationary environments. *Advances in Neural Information Processing Systems*, 32, 2019.

360 Claire Vernade, Andras Gyorgy, and Timothy Mann. Non-stationary delayed bandits with intermediate
361 observations. In *International Conference on Machine Learning*, pages 9722–9732. PMLR, 2020b.

362 Nicolò Cesa-Bianchi and Gabor Lugosi. *Prediction, Learning, and Games*. Cambridge University
363 Press, USA, 2006.

364 Peter D. Grünwald. *The Minimum Description Length Principle (Adaptive Computation and Machine
365 Learning)*. The MIT Press, 2007.

366 Bin Yu. Rates of convergence for empirical processes of stationary mixing sequences. *The Annals of
367 Probability*, 22(1):94–116, 1994.

368 M.J. Weinberger and E. Ordentlich. On delayed prediction of individual sequences. *IEEE Transactions
369 on Information Theory*, 48(7), 2002.

370 Daniel Hsu, Aryeh Kontorovich, David A Levin, Yuval Peres, Csaba Szepesvári, and Geoffrey Wolfer.
371 Mixing time estimation in reversible markov chains from a single sample path. *The Annals of
372 Applied Probability*, 29(4):2439–2480, 2019.

373 Geoffrey Wolfer. Mixing time estimation in ergodic markov chains from a single trajectory with
374 contraction methods. In *Algorithmic Learning Theory*, pages 890–905, 2020.

375 **A Technical Appendices and Supplementary Material**

376 **A.1 Proof of Proposition 2**

For the EWA forecaster with prior Q_1 , we can rewrite the regret via a standard telescoping argument (see Lemma B.1 in Clerico et al., 2025) as

$$\text{Regret}_t(\bar{\theta}) = -\log \int \exp \left(-\sum_{s=1}^t \ell_s(\theta) + \sum_{s=1}^t \ell_s(\bar{\theta}) \right) dQ_1(\theta).$$

377 Using the variational representation of the KL divergence, this can be upper bounded as

$$\begin{aligned} \text{Regret}_t(\bar{\theta}) &= \inf_Q \left\{ \int \sum_{s=1}^t \ell_s(\theta) dQ(\theta) - \sum_{s=1}^t \ell_s(\bar{\theta}) + D_{\text{KL}}(Q||Q_1) \right\} \\ &\leq \inf_{c \in (0,1]} \left\{ \int \sum_{s=1}^t \ell_s(\theta) dP_c(\theta) - \sum_{s=1}^t \ell_s(\bar{\theta}) + D_{\text{KL}}(P_c||Q_1) \right\}, \end{aligned}$$

378 where P_c is the uniform measure on the closed Euclidean ball of radius c in \mathbb{R}^p , centred at $\bar{\theta}$. We
379 remark that for all $c \in (0,1]$, $P_c \ll Q_1$. Therefore, for all $c \in (0,1]$,

$$D_{\text{KL}}(P_c||Q_1) = \int p \log \frac{B+1}{c} dQ_1(\theta) = p \log \frac{B+1}{c}.$$

380 Taking a second-order Taylor expansion of the total squared loss around $\bar{\theta}$, and using the fact that the
381 mean of P_c is $\bar{\theta}$, we obtain

$$\sum_{s=1}^t \int_{\mathbb{R}^p} (\ell_s(\theta) - \ell_s(\bar{\theta})) dP_c(\theta) = \sum_{s=1}^t \int_{\mathbb{R}^p} \left(\langle \theta - \bar{\theta}, \nabla \ell_s(\bar{\theta}) \rangle + \frac{1}{2} \langle \theta - \bar{\theta}, X_s \rangle^2 \right) dP_c(\theta) \leq \frac{tc^2}{2},$$

382 where we used that $\|X_s\|_2 \leq 1$ for all s in the last inequality. Combining everything so far, we obtain

$$\text{Regret}_t(\bar{\theta}) \leq \inf_{c \in (0,1]} \left\{ p \log \frac{B+1}{c} + \frac{tc^2}{2} \right\} \leq \frac{p}{2} \log \frac{(B+1)^2 e \max(p,t)}{p},$$

383 where the last term is obtained taking $c = \min(1, \sqrt{p/t})$.

384 **A.2 Proof of Lemma 1**

Let $D_t = \ell_t(\theta^*) - \mathcal{L}_t(Q_t)$ and $\lambda_t(\theta) = \langle \theta - \theta^*, X_t \rangle$. It is easy to check that

$$D_t = \log \int e^{\lambda_t(\theta) \varepsilon_t - \lambda_t(\theta)^2/2} dQ_t(\theta).$$

385 Fix $i \in \{1, \dots, d\}$. We denote as $S^{(i)} = (S_k^{(i)})_{k \geq 1}$ the subsequence defined as $S_k^{(i)} =$
386 $\sum_{j=1}^k D_{i+(j-1)d}$. We also define $\mathcal{F}_k^{(i)} = \mathcal{F}_{i+(k-1)d}$. It is easy to check that $(S_k^{(i)})_{k \geq 1}$ is adapted
387 with respect to $(\mathcal{F}_k^{(i)})_{k \geq 1}$. Now, let $M_k^{(i)} = \exp(S_k^{(i)} - (k-1)(2B+1)\phi_d)$. We will show that
388 $(M_k^{(i)})_{k \geq 1}$ is a super-martingale with respect to $(\mathcal{F}_k^{(i)})_{k \geq 1}$, with initial expectation bounded by 1.
389 For this it is enough to show that for any $k \geq 1$ we have $\mathbb{E}[e^{D_{i+(k-1)d} - (2B+1)\phi_d} | \mathcal{F}_{k-1}^{(i)}] \leq 1$. This is
390 true for $k = 1$ (where we let $\mathcal{F}_0^{(i)}$ be the trivial σ -field, or more generally a σ -field independent of the
391 noise). Indeed, as $i \leq d$, X_i is \mathcal{F}_0 measurable and hence independent of ε_i . From Assumption 1, we
392 know that ε_i is sub-Gaussian, and so $\mathbb{E}[e^{D_i}] \leq 1$.

393 Let us now check the case $k \geq 2$. For convenience, we define $t_k^{(i)} = i + (k-1)d$. We note that
394 $\mathcal{F}_{t_k^{(i)}} = \mathcal{F}_k^{(i)}$. We have

$$\begin{aligned} \mathbb{E}[e^{D_{i+(k-1)d} - (2B+1)\phi_d} | \mathcal{F}_{k-1}^{(i)}] \\ = \mathbb{E} \left[\int \exp(\lambda_{t_k^{(i)}}(\theta) \varepsilon_{t_k^{(i)}} - \lambda_{t_k^{(i)}}(\theta)^2/2 - (2B+1)\phi_d) dQ_{t_k^{(i)}}(\theta) \middle| \mathcal{F}_{k-1}^{(i)} \right]. \end{aligned}$$

395 Now, $Q_{t_k^{(i)}}$ only depends on the noise up to $\varepsilon_{t_k^{(i)}-d} = \varepsilon_{t_{k-1}^{(i)}}$, thanks to the delayed bandit framework.
396 Henceforth, we can swap the conditional expectation and the integral. In a similar way, we can bring
397 $\exp(-\lambda_{t_k^{(i)}}(\theta)^2/2 - (2B+1)\phi_d)$ outside of the conditional expectation, as it is $\mathcal{F}_{k-1}^{(i)}$ measurable.
398 We get

$$\begin{aligned} & \mathbb{E}[e^{D_{i+(k-1)d} - (2B+1)\phi_d} | \mathcal{F}_{k-1}^{(i)}] \\ &= \int \mathbb{E} \left[\exp(\lambda_{t_k^{(i)}}(\theta) \varepsilon_{t_k^{(i)}}) \middle| \mathcal{F}_{k-1}^{(i)} \right] \exp(-\lambda_{t_k^{(i)}}(\theta)^2/2 - (2B+1)\phi_d) dQ_{t_k^{(i)}}(\theta) \\ &\leq \int \exp(\lambda_{t_k^{(i)}}(\theta)^2/2 + \lambda_{t_k^{(i)}} \mathbb{E}[\varepsilon_{t_k^{(i)}} | \mathcal{F}_{k-1}^{(i)}]) \exp(-\lambda_{t_k^{(i)}}(\theta)^2/2 - (2B+1)\phi_d) dQ_{t_k^{(i)}}(\theta) \\ &\leq \int \exp(|\lambda_{t_k^{(i)}}(\theta)| \phi_d - (2B+1)\phi_d) dQ_{t_k^{(i)}}(\theta), \end{aligned}$$

where the two inequalities use the sub-Gaussianity and mixing properties of Assumption 1. Now, by construction $Q_{t_k^{(i)}}$ has support on $\mathcal{B}(B+1)$, and for every $\theta \in \mathcal{B}(B+1)$

$$|\lambda_{t_k^{(i)}}(\theta)| \leq \|\theta - \theta^*\|_2 \|X_{t_k^{(i)}}\|_2 \leq 2B+1,$$

where we also used that $\|X_{t_k^{(i)}}\|_2 \leq 1$, as for all t we are assuming that $\mathcal{X}_t \subseteq \mathcal{B}(1)$. We thus conclude that $(M_k^{(i)})_{k \geq 1}$ is indeed a super-martingale, non-negative and with initial value bounded by 1. By Ville's inequality it follows that

$$\mathbb{P}(S_k^{(i)} \leq k(2B+1)\phi_d + \log \frac{d}{\delta}, \forall k \geq 1) \geq 1 - \frac{\delta}{d}.$$

399 Now that we have proven that we have a super-martingale for each block, the desired anytime valid
400 concentration result follows directly from a simple union bound.

401 A.3 Proof of Lemma 2

Fix $t \geq 1$, and let $i \in \{1, \dots, d\}$ and $k \geq 1$ be such that $t = i + (k-1)d$. Let $I_j = \{j + d\mathbb{N}\} \cap \{1, \dots, t\}$, for $j \in \{1, \dots, d\}$. We consider d independent EWA forecaster (all initialised with the uniform prior on $\mathcal{B}(B+1)$). The j^{th} forecaster only acts and receive feedback from the rounds in I_j . We note that the j^{th} forecaster acts for t_j rounds, where $t_j = k$ if $j \geq i$, and $t_j = k-1$ otherwise. We denote as $R^{(j)}$ the regret of the j^{th} forecaster (which only takes into account the losses at the rounds in I_j , with comparator $\bar{\theta}$). By Proposition 2 we get

$$\text{Regret}_t(\bar{\theta}) = \sum_{j=1}^d R^{(j)} \leq \sum_{j=1}^d \frac{p}{2} \log \frac{(B+1)^2 e \max(p, t_j)}{p}.$$

402 We conclude by noticing that, for all $j, t_j \leq (t+d)/d$.

403 A.4 Proof of Lemma 3

404 We recall the standard Elliptical Potential Lemma (see e.g. Lemma 11 in Abbasi-Yadkori et al., 2011),
405 which we will use in our proof of Lemma 3.

406 **Lemma 4** (Elliptical Potential Lemma). *Let $(X_t)_t$ be any sequence of vectors in \mathbb{R}^p satisfying
407 $\max_{t \in [T]} \|X_t\|_2 \leq L$ and let $V_T = \sum_{t=1}^T X_t X_t^\top + \lambda I$, for some $\lambda > 0$. Then*

$$\sum_{t=1}^T \min(1, \|X_t\|_{V_{t-1}}^2) \leq 2p \log(1 + \frac{TL^2}{\lambda p}).$$

Next, we introduce some notation. For $t > d$, define $(i(t), k(t)) \in [d] \times [K]$ such that $t = i(t) + k(t)d$ and let

$$V_{k(t)-1}^{i(t)} = \sum_{k=0}^{k(t)-1} X_k^{i(t)} (X_k^{i(t)})^\top + \lambda \text{Id},$$

408 where $X_k^{i(t)} = X_{i(t)+kd}$. With this notation, we can state the following lemma.

409 **Lemma 5.** For any $t > d$, we have

$$V_{t-d} \succcurlyeq V_{k(t)-1}^{i(t)},$$

410 which implies that $\|X_t\|_{V_{t-d}}^2 \leq \|X_t\|_{(V_{k(t)-1}^{i(t)})^{-1}}^2$ for any $t > d$.

411 *Proof.* Notice that we can write $V_{t-d} = \sum_{s=1}^{t-d} X_s X_s^\top + \lambda \text{Id} = V_{k(t)}^{i(t)} + \sum_{s=1, s \notin S_t}^{t-d} X_s X_s^\top$ where
412 $S_t := \{s = i(t) + (k-1)d, k \in [k(t)]\}$ is the set of indices $(i(t), i(t) + d, \dots, i(t) + (k(t) - 1)d)$.
413 The statement now follows from the fact that $\sum_{s=1, s \notin S_t}^{t-d} X_s X_s^\top \succcurlyeq 0$. \square

414 We are now ready to prove Lemma 3. For now, let us assume that $T = Kd$, for some $K > 1$. Using
415 Lemma 5 and then Lemma 4, we have

$$\begin{aligned} \sum_{t=d+1}^T \min(1, \|X_t\|_{V_{t-d}}^2) &\leq \sum_{t=d+1}^T \min(1, \|X_{k(t)}^{i(t)}\|_{(V_{k(t)-1}^{i(t)})^{-1}}^2) \\ &= \sum_{i=1}^d \sum_{k=1}^{K-1} \min(1, \|X_k^i\|_{(V_{k-1}^i)^{-1}}^2) \\ &\leq 2dp \log\left(1 + \frac{(K-1)L^2}{\lambda p}\right). \end{aligned}$$

416 One can verify that if T is not divisible by d , the above inequality still holds if we replace K by $\lceil \frac{T}{d} \rceil$.
417 Therefore, regardless of whether T is divisible by d , we have

$$\sum_{t=d+1}^T \min(1, \|X_t\|_{V_{t-d}}^2) \leq 2dp \log\left(1 + \frac{TL^2}{\lambda dp}\right).$$

418 This concludes the proof of Lemma 3.

419 A.5 Proof of Corollary 2 and Corollary 3

420 We start by recalling the general result

$$\text{Reg}(T) = \mathcal{O}\left(\underbrace{dB}_{(1)} + \underbrace{dp\sqrt{T} \log \frac{TB}{dp}}_{(2)} + \underbrace{T\sqrt{Bdp\phi_d \log \frac{TB}{dp}}}_{(3)} + \underbrace{d\sqrt{pT \log \frac{TB}{p\delta}}}_{(4)}\right). \quad (5)$$

421 To simplify the following calculations, we do not force d to be a positive integer. One can always
422 round d without changing the rates of the regret bounds.

423 Geometric Mixing:

424 Assume $d = \tau \log \frac{BCT}{p}$. We notice that the term (1) is logarithmic in T and thus negligible. From
425 the definition of geometric mixing, it holds that $\phi_d = Ce^{-\frac{d}{\tau}} = \frac{p}{BT}$. Therefore,

$$(3) \leq p\sqrt{\tau T} \log \frac{TB}{p}.$$

426 Substituting the value of d yields the desired bounds for terms (2) and (4) in Equation 5, and hence
427 the desired statement.

428 Algebraic mixing:

429 Assume $d = CT^{\frac{1}{1+r}}$, we notice that in this case since $\phi_d = Cd^{-r}$, we have $d\phi_d = Cd^{1-r}$. In
430 particular this implies that $T\sqrt{d\phi_d} = T^{\frac{3+r}{2(1+r)}}$ and thus

$$(3) \leq C\sqrt{Bp \log \frac{TB}{p}} T^{\frac{3+r}{2(1+r)}}$$

431 The same way (2) and (4) are of order $d\sqrt{T} = T^{\frac{3+r}{2(1+r)}}$ and replacing in Equation 5 yields the desired
432 statement.

433 **A.6 Proof of Corollary 4 and Corollary 5**

434 We start by recalling the general result

$$\text{Reg}(T) \leq 2dB + \frac{8dp}{\Delta} \max(B^2, \beta_T^2) \log \left(1 + \frac{B^2T}{dp} \right),$$

435 where $\beta_T^2 = \underbrace{dp \log \frac{(B+1)^2 e \max(dp, T+d)}{dp}}_{(1)} + \underbrace{2T \phi_d(B+1)}_{(2)} + \underbrace{2d \log \frac{d}{\delta}}_{(3)}.$

436 **Geometric Mixing:**

437 Assume $d = \tau \log \frac{BCT}{p}$, then (2) = $\frac{2p(B+1)}{BC}$ is a constant. Hence we have

$$\text{Reg}(T) \leq 2dB + \frac{8d^2p}{\Delta} \left(p \log \frac{(B+1)^2 e \max(dp, T+d)}{dp} + 2 \log \frac{d}{\delta} + \frac{2p(B+1)}{BC} \right) \log \left(1 + \frac{B^2T}{dp} \right),$$

438 which under the assumption that $\beta_T \geq B$ and replacing d by its definition yields

$$\begin{aligned} \text{Reg}(T) &\leq 2B\tau \log \frac{BCT}{p} \\ &+ \frac{8\tau^2p}{\Delta} \log \left(1 + \frac{B^2T}{p\tau \log \frac{BCT}{p}} \right) \left(\left(\log \frac{BCT}{p} \right)^2 \left(p \log \frac{(B+1)^2 e T}{p} + 2\tau \frac{\log \frac{BCT}{p}}{\delta} \right) + \frac{2p(B+1)}{BC} \right). \end{aligned}$$

439 If Δ is constant, then for large T , the first term and the constant part coming from (2) become
440 negligible. Therefore,

$$\text{Reg}(T) = \mathcal{O} \left(\frac{8\tau^2p}{\Delta} \log \left(1 + \frac{B^2T}{p\tau \log \frac{BCT}{p}} \right) \left(\log \frac{BCT}{p} \right)^2 \left(\frac{p}{2} \log \frac{(B+1)^2 e T}{p} + \tau \frac{\log \frac{BCT}{p}}{\delta} \right) \right)$$

441 **Algebraic mixing:**

442 Assume $d = CT^{\frac{1}{1+r}}$, then we have

$$\beta_T^2 \leq CT^{\frac{1}{1+r}} \log \frac{(B+1)^2 e T}{p} + 2C(B+1)T^{\frac{2}{1+r}} + 2CT^{\frac{1}{1+r}} \log \frac{CT^{\frac{1}{1+r}}}{\delta}.$$

443 Under the regime where $2dB \leq \frac{8dp}{\Delta} \max(B^2, \beta_T^2) \log \left(1 + \frac{B^2T}{dp} \right)$ and $B \leq \beta_T$ this leads to

$$\text{Reg}(T) = \mathcal{O} \left(\frac{8Cp}{\Delta} T^{\frac{2}{1+r}} \log \left(1 + \frac{B^2T}{pCT^{1/(1+r)}} \right) \left(\frac{p}{2} \log \frac{(B+1)^2 e T}{p} + \log \frac{CT^{1/(1+r)}}{\delta} \right) \right).$$