

Empirical Bound Information-Directed Sampling

Anonymous authors
Paper under double-blind review

Keywords: bandit algorithms, information-directed sampling, parameter bounds, heteroskedastic noise

Summary

Information-directed sampling (IDS) is a powerful framework for solving bandit problems which has shown strong results in both Bayesian and frequentist settings. However, frequentist IDS, like many other bandit algorithms, requires that one have prior knowledge of a (relatively) tight upper bound on the norm of the true parameter vector governing the reward model in order to achieve good performance. Unfortunately, this requirement is rarely satisfied in practice. As we demonstrate, using a poorly calibrated bound can lead to significant regret accumulation. To address this issue, we introduce a novel frequentist IDS algorithm that iteratively refines a high-probability upper bound on the true parameter norm using accumulating data. We focus on the linear bandit setting with heteroskedastic subgaussian noise. Our method leverages a mixture of relevant information gain criteria to balance exploration aimed at tightening the parameter norm bound and directly searching for the optimal action. We establish regret bounds for our algorithm that do not depend on an initially assumed parameter norm bound and demonstrate that our method outperforms state-of-the-art IDS and UCB algorithms.

Contribution(s)

1. This paper introduces a novel frequentist information-directed sampling (IDS) algorithm that does not require prior knowledge of a tight upper bound of the true parameter norm to achieve good performance. Our method uses accumulating data to generate a sequence of high-probability upper bounds on the parameter norm and accounts for potential heteroskedasticity of the rewards.
Context: The performance of many frequentist bandit algorithms, including various IDS (Kirschner & Krause, 2018; Kirschner et al., 2021) and UCB methods (Auer, 2002; Abbasi-Yadkori et al., 2011), relies heavily on a (at least relatively) tight upper bound on the true parameter norm being available to the algorithm. This is almost never the case in practice which can lead to significant regret accumulation. Recently, some norm-agnostic bandit algorithms have been proposed to address this issue (Gales et al., 2022), however, they do not account for potential heteroskedasticity of the rewards.
2. We introduce a new composite information criterion that balances improving the requisite upper bound on the parameter norm and direct search for the optimal action.
Context: To the best of our knowledge, no other IDS algorithm uses a mixture of information gain criteria to balance acquiring information about different aspects of the environment’s dynamics. We are also not aware of any existing method that uses an information gain criterion aimed at improving the upper bound on the parameter norm.
3. We establish anytime sublinear regret bounds for our algorithm which eventually do not depend on the initially assumed parameter norm bound.
Context: Previously proposed norm-agnostic bandits (Gales et al., 2022) rely on an initial burn-in during which regret accumulation is not controlled, e.g., it need not be sublinear.

Empirical Bound Information-Directed Sampling

Anonymous authors

Paper under double-blind review

Abstract

1 Information-directed sampling (IDS) is a powerful framework for solving bandit prob-
2 lems which has shown strong results in both Bayesian and frequentist settings. How-
3 ever, frequentist IDS, like many other bandit algorithms, requires that one have prior
4 knowledge of a (relatively) tight upper bound on the norm of the true parameter vector
5 governing the reward model in order to achieve good performance. Unfortunately, this
6 requirement is rarely satisfied in practice. As we demonstrate, using a poorly calibrated
7 bound can lead to significant regret accumulation. To address this issue, we introduce a
8 novel frequentist IDS algorithm that iteratively refines a high-probability upper bound on
9 the true parameter norm using accumulating data. We focus on the linear bandit setting
10 with heteroskedastic subgaussian noise. Our method leverages a mixture of relevant
11 information gain criteria to balance exploration aimed at tightening the estimated param-
12 eter norm bound and directly searching for the optimal action. We establish regret bounds
13 for our algorithm that do not depend on an initially assumed parameter norm bound and
14 demonstrate that our method outperforms state-of-the-art IDS and UCB algorithms.

15 1 Introduction

16 We consider linear stochastic bandits (Lattimore & Szepesvári, 2020) with heteroskedastic noise
17 (see Wetz et al., 2023, for applications of such models in marketing and other areas). In this setting,
18 information directed sampling (IDS) and upper confidence bound (UCB) algorithms have been shown
19 to be extremely effective (Auer, 2002; Abbasi-Yadkori et al., 2011; Kirschner & Krause, 2018;
20 Kirschner et al., 2021). However, many of these methods require strong prior information that can be
21 used to inform a high-quality upper bound on the Euclidean norm of the parameter vector indexing
22 the reward model. The choice of this bound is critical to algorithm performance. If the bound is too
23 large, the algorithm risks incurring excess risk due to needless exploration, and if the bound is too
24 small, the algorithm may fail to identify the optimal arm and thus suffer linear regret.

25 To reduce sensitivity on a user-specified bound, we propose a novel version of frequentist IDS that
26 uses accumulating data to generate a sequence of high-probability upper bounds on the norm of
27 the reward model parameters. A key component of our method is a new information gain criterion
28 that balances improving the requisite upper bound and regret minimization. Because improving
29 the bound is critical to avoid over-exploration in early rounds of the bandit process, we develop a
30 two-phase procedure that uses our new information criterion in the first phase and then defaults to a
31 more standard IDS information criterion in the second phase.

32 Unlike other bandit strategies, such as UCB (Auer, 2002; Garivier & Cappé, 2011; Cappé et al., 2013;
33 Zhou et al., 2020) or Thompson sampling (TS) (Thompson, 1933; Agrawal & Goyal, 2013; Phan
34 et al., 2019), which encourage exploration indirectly by leveraging uncertainty about the optimal
35 arm, IDS explicitly balances exploration and exploitation. It selects actions that minimize estimated
36 instantaneous regret while maximizing expected information gain about model parameters. As shown
37 by Russo & Van Roy (2014) and Kirschner & Krause (2018), this approach allows IDS to avoid
38 pitfalls inherent in UCB and TS-based algorithms, particularly in scenarios where certain suboptimal
39 actions provide valuable information about the environment’s dynamics. In such cases, UCB and

40 TS tend to overlook these actions, whereas IDS plays them early on, enabling faster learning of the
41 optimal policy and ultimately achieving superior long-term performance. IDS was first introduced
42 for Bayesian bandits by [Russo & Van Roy \(2014\)](#) and later adapted to the frequentist setting by
43 [Kirschner & Krause \(2018\)](#). Beyond the standard bandit setting, IDS has been applied to problems
44 such as linear partial monitoring ([Kirschner et al., 2020](#)) — a generalization of bandits where the
45 observed signal on the environment model parameters is not necessarily the same as the reward to be
46 optimized — as well as reinforcement learning ([Nikolov et al., 2019](#); [Lindner et al., 2021](#); [Hao &
47 Lattimore, 2022](#)), where the actions taken by the agent influence the state of the environment and the
48 reward dynamics.

49 To the best of our knowledge, no previous work has considered either the strategy of iteratively
50 refining and utilizing a high-probability upper bound on the parameter norm in the heteroskedastic
51 subgaussian linear bandit setting we work with here, or the use of the information gain criterion
52 for tightening the bound on the parameter norm we introduce. We are also not aware of any work
53 utilizing a mixture of information gain criteria to encourage simultaneously obtaining different types
54 of information about the dynamics of the environment. We note that while we introduce this idea in
55 the form of an IDS algorithm, the approach of iteratively refining and utilizing a high-probability
56 upper bound of the true parameter norm can be regarded as a more general design principle beyond
57 its IDS implementation in this setting.

58 The remainder of this manuscript is structured as follows. The next section provides a brief review
59 of related work. Section 3 introduces the problem setup and notation used throughout the paper. In
60 Section 4, we present the necessary background on IDS. Section 5 introduces the novel empirical
61 bound information-directed sampling (EBIDS) algorithm, which removes the need for a tight param-
62 eter norm bound to be known *a priori*. Section 6 establishes regret bound guarantees for EBIDS, and
63 finally, Section 7 evaluates its empirical performance against competitor algorithms in a simulation
64 study.

65 2 Related works

66 The assumption that the norm of the parameter indexing the reward model is known or that one has
67 a (relatively) tight upper bound on this quantity is abundant in the IDS and UCB literature ([Auer,
68 2002](#); [Abbasi-Yadkori et al., 2011](#); [Kirschner & Krause, 2018](#); [Hung et al., 2021](#)); it has also been
69 used in Thompson sampling ([Xu et al., 2023](#)). This assumption commonly arises through the use of
70 self-normalized martingale bounds and related concentration results ([Abbasi-Yadkori et al., 2011](#)).
71 Consequently, algorithms constructed through these concentration results require a user-specified
72 upper bound on the norm or the true parameter vector. Critically, as noted previously, the performance
73 of these algorithms can be highly sensitive to the choice of these bounds. Despite this, only a handful
74 of papers have attempted to alleviate this sensitivity.

75 [Gales et al. \(2022\)](#) propose norm-agnostic linear bandits which construct a series of confidence
76 ellipsoids for the true parameter vector along with a projection interval to construct a UCB-type
77 algorithm. However, their algorithms rely on an initial burn-in during which regret accumulation is
78 not controlled, e.g., it need not be sublinear. In our simulation experiments, we find that the impact of
79 this initial exploration on accumulated regret is not negligible. Furthermore, as UCB algorithms, their
80 methods do not explicitly make use of heteroskedasticity in the reward distributions across arms.

81 The algorithm proposed by [Ghosh et al. \(2021\)](#) shares some underlying ideas with our method in the
82 sense that they use multi-phase exploration to iteratively update the bound on the unknown parameter
83 norm. However, their algorithm is limited to the specialized setting of stochastic linear bandits
84 introduced by [Chatterji et al. \(2020\)](#) with restrictive assumptions on the structure of the rewards
85 which makes their methods generally not applicable to the settings we consider here. Similarly, [Dani
86 et al. \(2008\)](#), [Orabona & Cesa-Bianchi \(2011\)](#), and [Gentile & Orabona \(2014\)](#) do not assume that one
87 has a high-quality (i.e., relatively tight) bound on the norm of the parameter; however, they require
88 bounded rewards for all arms. Other attempts to alleviate the assumption of known parameter norm

bound have been made in spectral bandits (Kocák et al., 2020), and deep active learning (Wang et al., 2021). However, it is not clear how to port these methods to the setup we consider here.

3 Setup and notation

We denote the inner product of two vectors of the same dimension as $\langle \cdot, \cdot \rangle$ so that the squared Euclidean norm of vector \mathbf{v} is $\|\mathbf{v}\|_2^2 = \langle \mathbf{v}, \mathbf{v} \rangle$. For a symmetric positive definite or semi-definite matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, we denote the associated matrix norm (or semi-norm) of a vector $\mathbf{v} \in \mathbb{R}^d$ as $\|\mathbf{v}\|_{\mathbf{A}}^2 = \langle \mathbf{v}, \mathbf{A}\mathbf{v} \rangle$. We let $\lambda_{\max}(\mathbf{A})$ and $\lambda_{\min}(\mathbf{A})$ denote the largest and the smallest eigenvalues of \mathbf{A} . Throughout, $\log(x)$ denotes the natural logarithm of $x \in \mathbb{R}_+$.

At each time step $t \in \{1, \dots, T\}$ the agent selects an action $A_t \in \mathcal{A}$ and observes the outcome $Y_t \in \mathbb{R}$ which is generated from the linear model

$$Y_t(A_t) = \langle \phi(A_t), \boldsymbol{\theta}^* \rangle + \eta_t, \quad (1)$$

where $\boldsymbol{\theta}^* \in \mathbb{R}^d$ is a vector of unknown parameters and $\phi : \mathcal{A} \rightarrow \mathbb{R}^d$ is a feature mapping, such that for any $a \in \mathcal{A}$ we have $\|\phi(a)\|_2 \in [L, U]$ for some positive constants $L \leq U$. The noise term η_t is assumed to be subgaussian and conditionally mean zero, i.e., for every $c \in \mathbb{R}$ we assume that

$$\mathbb{E} \{ \exp(c\eta_t) \mid A_t = a \} \leq \exp \left\{ c^2 \rho(a)^2 / 2 \right\}, \quad (2)$$

where $0 < \rho_{\min} \leq \rho(a) \leq \rho_{\max} < \infty$ for all $a \in \mathcal{A}$ and $\mathbb{E}(\eta_t \mid A_1, \dots, A_t, \eta_1, \dots, \eta_{t-1}) = 0$. Define $B^* := \|\boldsymbol{\theta}^*\|_2$, in some of our theoretical results we assume that one has a conservative upper bound B such that $B^* \leq B$ but that this bound may be quite conservative, i.e., it may be that $B^* \ll B$.

The available history to inform action selection at time t is $\mathbf{H}_t = \{(A_1, Y_1), \dots, (A_{t-1}, Y_{t-1})\}$ of past actions and rewards. A bandit algorithm is thus formalized as a map from histories to distributions over actions $\pi_t(a \mid \mathbf{h}_t) = \mathbb{P}(A_t = a \mid \mathbf{H}_t = \mathbf{h}_t)$. Let

$$\Delta(A_t) = \langle \phi(a^*), \boldsymbol{\theta}^* \rangle - \langle \phi(A_t), \boldsymbol{\theta}^* \rangle$$

be the gap between the action A_t and the optimal action $a^* = \arg \max_{a \in \mathcal{A}} \langle \phi(a), \boldsymbol{\theta}^* \rangle$. Our goal is to design an algorithm $\pi_t(\cdot \mid \mathbf{h}_t)$ which maximizes the cumulative expected reward $\mathbb{E} \left\{ \sum_{t=1}^T Y_t \right\}$, or equivalently, minimizes the regret, defined as $\mathcal{R}_T = \mathbb{E} \left\{ \sum_{t=1}^T \Delta(A_t) \right\}$. While regret is a standard performance metric for bandit algorithms, it involves taking expectation over both the randomness in the policy and the noise in the rewards so it can be a poor indicator of the risk associated with the policy (Lattimore & Szepesvári, 2020). For this reason in this paper we also study the probabilistic bounds on the pseudo-regret defined as $\mathcal{PR}_T = \sum_{t=1}^T \Delta(A_t)$.

4 Review of information-directed sampling

Information-directed sampling (IDS Russo & Van Roy, 2014) is an algorithm design principle that balances minimizing the gap of an action with its potential for information gain. Let $\mathcal{P}(\mathcal{A})$ denote the space of distributions of \mathcal{A} . For any $\mu \in \mathcal{P}(\mathcal{A})$ let $\widehat{\Delta}_t(\mu)$ be an estimator of the expected gap $\mathbb{E}_{\mu} \Delta := \mathbb{E}_{A \sim \mu} \Delta(A)$ constructed from the history \mathbf{H}_t , and, similarly, let $I_t(\mu)$ be a measure of information gain, e.g., the reduction of entropy in the posterior or sampling distribution the parameter indexing the mean reward model (see below for additional details). For any function $f : \mathcal{P}(\mathcal{A}) \rightarrow \mathbb{R}$, if the argument is a point mass at a single action, e.g., where μ is the Dirac delta δ_a , we write $f(a)$ rather than $f(\delta_a)$. The IDS distribution is defined as

$$\mu_t^{\text{IDS}} = \arg \min_{\mu \in \mathcal{P}(\mathcal{A})} \frac{\left\{ \widehat{\Delta}_t(\mu) \right\}^2}{I_t(\mu)}. \quad (3)$$

122 The quantity $\Psi_t(\mu) := \left\{ \widehat{\Delta}_t(\mu) \right\}^2 / I_t(\mu)$ being minimized is known as the *information ratio*. An
 123 IDS algorithm samples the action $A_t \sim \mu_t^{\text{IDS}}$ at each time step t . Note that this results in a randomized
 124 algorithm, which, as shown by [Russo & Van Roy \(2014\)](#) and [Kirschner & Krause \(2018\)](#), always
 125 has at most two actions in its support. However, it is also possible to restrict the optimization in
 126 (3) to Dirac delta functions on the individual actions, thus obtaining what is often referred to as
 127 *deterministic IDS* ([Kirschner & Krause, 2018](#))

$$\widehat{A}_t^{\text{DIDS}} = \arg \min_{a \in \mathcal{A}} \frac{\left\{ \widehat{\Delta}_t(a) \right\}^2}{I_t(a)}. \quad (4)$$

128 Deterministic IDS is typically computationally cheaper, retains the same theoretical regret bounds
 129 as its randomized counterpart, and in simulation experiments was shown to be competitive with or
 130 superior to randomized IDS ([Kirschner & Krause, 2018](#); [Kirschner, 2021](#)). Furthermore, deterministic
 131 IDS may be appealing in settings where randomized policies are unpalatable such as public health
 132 ([Weltz et al., 2022](#)) and site selection ([Ahmadi-Javid et al., 2017](#)).

133 The information ratio provides a natural way of bounding regret within a Bayesian setting ([Russo
 134 & Van Roy, 2014](#)). Notably, the information ratio can also be used to bound the regret under a
 135 frequentist paradigm ([Kirschner & Krause, 2018](#)) as illustrated by the following result based on the
 136 work of [Kirschner \(2021\)](#) which we prove in Section 10.1 of the Supplementary Materials.

Theorem 1 (Kirschner). *For any T let G be a fixed subset of $\{1, \dots, T\}$ and let $\{A_t\}_{t=1}^T$ be an \mathbf{H}_t -adapted sequence in \mathcal{A} . Then*

$$\mathbb{E} \left\{ \sum_{t \in G} \widehat{\Delta}_t(A_t) \right\} \leq \sqrt{\mathbb{E} \left\{ \sum_{t \in G} \Psi_t(A_t) \right\} \mathbb{E} \left\{ \sum_{t \in G} I_t(A_t) \right\}},$$

and if $\widehat{\Delta}_t(A_t) \geq \Delta(A_t)$ for all $t \in G$ then with probability 1 we have

$$\sum_{t \in G} \Delta(A_t) \leq \sqrt{\left\{ \sum_{t \in G} \Psi_t(A_t) \right\} \left\{ \sum_{t \in G} I_t(A_t) \right\}}.$$

137 [Kirschner & Krause \(2018\)](#) used weighted ridge regression to estimate $\boldsymbol{\theta}^*$ at each time step t so that

$$\widehat{\boldsymbol{\theta}}_t^{\text{wls}} = \mathbf{W}_t^{-1} \sum_{s=1}^{t-1} \frac{1}{\rho(A_s)^2} \boldsymbol{\phi}(A_s) Y_s, \quad \text{where} \quad \mathbf{W}_t = \sum_{s=1}^{t-1} \frac{1}{\rho(A_s)^2} \boldsymbol{\phi}(A_s) \boldsymbol{\phi}(A_s)^\top + \gamma \mathbf{I}_d, \quad (5)$$

138 and $\gamma \geq 0$ is a constant chosen by the user. The following result, proposed by [Abbasi-Yadkori et al.
 139 \(2011\)](#) and extended by [Kirschner & Krause \(2018\)](#), provides a means to perform inference using
 140 this estimator.

141 **Theorem 2.** *Suppose that the generative model follows the linear bandit model $Y_t = \langle \boldsymbol{\phi}(A_t), \boldsymbol{\theta}^* \rangle + \eta_t$
 142 given in (1), where the actions A_t are \mathbf{H}_t -adapted and the errors η_t have conditional mean of zero
 143 and satisfy the subgaussian condition in (2). Let $B \geq \|\boldsymbol{\theta}^*\|_2$ be a (potentially conservative) bound
 144 on the norm of the parameters indexing the reward model and define*

$$\mathcal{E}_{t,\delta}^{\text{wls}} := \left\{ \boldsymbol{\theta} \in \mathbb{R}^d : \left\| \boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_{\mathbf{W}_t}^2 \leq \beta_{t,\delta}(B) \right\},$$

145 where

$$\beta_{t,\delta}(B) = \left[\sqrt{2 \log \frac{1}{\delta} + \log \left\{ \frac{\det(\mathbf{W}_t)}{\det(\mathbf{W}_1)} \right\}} + \sqrt{\gamma} B \right]^2. \quad (6)$$

Then

$$\mathbb{P} \left(\bigcap_{t=1}^{\infty} \{\boldsymbol{\theta}^* \in \mathcal{E}_{t,\delta}^{\text{wls}}\} \right) \geq 1 - \delta,$$

146 *i.e.*, $\mathcal{E}_{t,\delta}^{\text{wls}}$ is a $(1 - \delta) \times 100\%$ confidence ellipsoid for $\boldsymbol{\theta}^*$.

147 [Kirschner & Krause \(2018\)](#) use Theorem 2 to formulate a weighted UCB algorithm which at each
148 time step t takes the action

$$A_t^{\text{UCB}(\delta_t)} = \arg \max_{a \in \mathcal{A}} \left\langle \boldsymbol{\phi}(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t,\delta_t}^{1/2}(B) \|\boldsymbol{\phi}(a)\|_{\mathbf{W}_t^{-1}}, \quad (7)$$

149 maximizing the $(1 - \delta_t) \times 100\%$ upper confidence bound on the expected reward based on the $\mathcal{E}_{t,\delta_t}^{\text{wls}}$
150 confidence set. Then they use

$$\check{\Delta}_{t,\delta_t}(a) = \left\langle \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) - \boldsymbol{\phi}(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t,\delta_t}^{1/2}(B) \left(\left\| \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) \right\|_{\mathbf{W}_t^{-1}} + \|\boldsymbol{\phi}(a)\|_{\mathbf{W}_t^{-1}} \right).$$

151 as the gap estimate. This ensures that $\Delta(a) \leq \check{\Delta}_{t,\delta_t}(a)$ for all $a \in \mathcal{A}$ whenever $\boldsymbol{\theta}^* \in \mathcal{E}_{t,\delta_t}^{\text{wls}}$ holds.

152 The choice of the information gain criterion is crucial when designing an IDS algorithm. [Kirschner &](#)
153 [Krause \(2018\)](#) introduce the following criterion

$$I_t^{\text{UCB}(\delta_t)}(a) = \frac{1}{2} \log \left(\frac{\left\| \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) \right\|_{\mathbf{W}_t^{-1}}^2}{\left\| \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) \right\|_{(\mathbf{W}_t + \rho(a)^{-2} \boldsymbol{\phi}(a) \boldsymbol{\phi}(a)^\top)^{-1}}^2} \right),$$

154 for any $a \in \mathcal{A}$. We present the resulting procedure in Algorithm 1, which we hereafter refer to as
IDS-UCB. It can be shown that if one chooses $\delta_t = 1/t^2$, the regret of IDS-UCB satisfies

Algorithm 1 IDS-UCB

Input: Action set \mathcal{A} , penalty parameter $\gamma > 0$, noise function $\rho : \mathcal{A} \rightarrow \mathbb{R}_+$, feature function
 $\boldsymbol{\phi} : \mathcal{A} \rightarrow \mathbb{R}$, sequence of confidence levels $\{\delta_t\}_{t \geq 1} \subset (0, 1)$, assumed true parameter norm bound B .

For $t = 1, 2, \dots, T$:

 Compute \mathbf{W}_t and $\widehat{\boldsymbol{\theta}}_t^{\text{wls}}$ using (5)

$$A_t^{\text{UCB}(\delta_t)} \leftarrow \arg \max_{a \in \mathcal{A}} \left\{ \left\langle \boldsymbol{\phi}(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t,\delta_t}^{1/2}(B) \|\boldsymbol{\phi}(a)\|_{\mathbf{W}_t^{-1}} \right\}$$

$$I_t^{\text{UCB}(\delta_t)}(a) \leftarrow \frac{1}{2} \log \left(\left\| \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) \right\|_{\mathbf{W}_t^{-1}}^2 \right) - \frac{1}{2} \log \left(\left\| \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) \right\|_{(\mathbf{W}_t + \rho(a)^{-2} \boldsymbol{\phi}(a) \boldsymbol{\phi}(a)^\top)^{-1}}^2 \right)$$

$$\check{\Delta}_{t,\delta_t}(a) \leftarrow \left\langle \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) - \boldsymbol{\phi}(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t,\delta_t}^{1/2}(B) \left(\left\| \boldsymbol{\phi} \left(A_t^{\text{UCB}(\delta_t)} \right) \right\|_{\mathbf{W}_t^{-1}} + \|\boldsymbol{\phi}(a)\|_{\mathbf{W}_t^{-1}} \right)$$

$$\mu_t \leftarrow \arg \min_{\mu \in \mathcal{P}(\mathcal{A})} \check{\Delta}_{t,\delta_t}^2(\mu) / I_t^{\text{UCB}(\delta_t)}(\mu)$$

 Sample $A_t \sim \mu_t$

 Play A_t , observe $Y_t = \langle \boldsymbol{\phi}(A_t), \boldsymbol{\theta}^* \rangle + \eta_t$

155

$$\mathcal{R}_T \leq O \left(\max\{U/\sqrt{\gamma}, \rho_{\max}\} \sqrt{\gamma} dB \sqrt{T} \log T \right),$$

while the pseudo-regret \mathcal{PR}_T of IDS-UCB with fixed $\delta_t = \delta$ satisfies with probability at least $1 - \delta$

$$\mathcal{PR}_T \leq O(\max\{U/\sqrt{\gamma}, \rho_{\max}\} \sqrt{\gamma} dB \sqrt{T} \log(T/\delta));$$

156 critically, both regret bounds scale directly with the assumed bound B on the Euclidean norm of the
 157 true parameter (see [Kirschner, 2021](#), for a formal statement of the preceding results and additional
 158 discussion).

159 We now demonstrate via a simple illustrative simulation experiment that the choice of B can have
 160 a significant impact on the finite time performance of IDS-UCB. Large values of B relative to B^*
 161 lead to excess exploration and large regret in early rounds of the algorithm, whereas small values of
 162 B can prevent the algorithm from identifying the optimal arm thus incurring linear regret. In this
 163 experiment we also include the weighted UCB policy given by (7). We evaluate versions of IDS-
 164 UCB and UCB that use a conservative value of $B > B^*$, and those which use an anti-conservative
 165 value $B < B^*$. The parameters indexing the generative model are $\theta^* = [-5, 1, 1, 1.5, 2]^\top$ so
 166 that $B^* = \|\theta^*\|_2 \approx 5.77$. We take $B = 100$ for the conservative bound and $B = 1$ for the anti-
 167 conservative bound. For reference, we also include *oracle* versions of UCB and IDS-UCB that have
 168 access to the true value of B^* . However, we emphasize that these procedures are not generally
 169 possible in practice.

170 We consider a setting with ten arms. Features for each arm are sampled from $\text{Unif}[-1/\sqrt{5}, 1/\sqrt{5}]$.
 171 The error distribution for the first five arms are standard normal and for the remaining five arms
 172 they are normal with mean zero and variance 0.2. Figure 1 shows the mean regret averaged over
 173 200 repeated experiments with $T = 500$ steps along with 95% pointwise confidence bounds. As
 174 anticipated, using a conservative bound of $B = 100$ achieves sublinear regret but pays a strong initial
 175 cost due to excess exploration. Algorithms that used the anti-conservative bound of $B = 1$ fail to
 identify the optimal arm thus sustain linear regret.

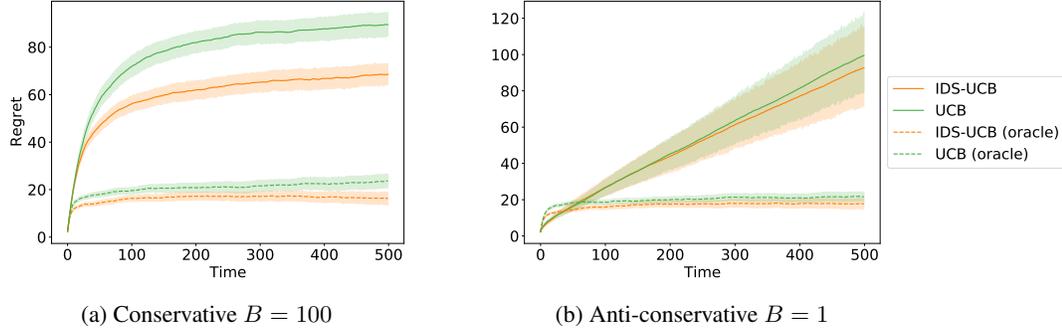


Figure 1: Regret incurred by IDS-UCB and UCB with: (a) conservative $B = 100$; (b) anti-conservative $B = 1$. In both plots we include the oracle versions of IDS-UCB and UCB using $B = B^*$ for reference. However, note that it is not feasible to implement them in most practical settings. The solid and dashes lines represent the regret averaged over 200 repeated experiments, while the shaded bounds are 95% pointwise confidence bands.

176

177 5 Empirical bound information-directed sampling

178 We propose the empirical bound information-directed sampling (EBIDS) algorithm, which, like
 179 existing IDS algorithms, relies on a conservative upper bound B , but, unlike existing algorithms,
 180 EBIDS refines this value with accruing data to obtain a tighter high-probability bound on B^* . Our
 181 algorithm proceeds in two phases. Throughout the first T_B steps, which we will refer to as the *bound*
 182 *exploration phase*, the goal is to gather initial information on the optimal action as well as to improve
 183 the bound on B^* . At each time step t in this first phase, we use

$$\hat{B}_t = \min \left\{ B, \|\hat{\theta}_t^{\text{wls}}\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2}(B) \lambda_{\min}(\mathbf{W}_t)^{-1/2} \right\}. \quad (8)$$

as the upper bound on B^* . The term $\beta_{t, \zeta_t(\delta)}(B)$ is defined in (6) and $\zeta_t(\delta) = \min\{\delta, 1/t^2\}$, where $\delta > 0$ is a user-specified parameter that determines the confidence level for the upper bound on B^* .

The geometric motivation for this estimator stems from the fact that the confidence set $\mathcal{E}_{t,\zeta_t(\delta)}^{\text{wls}}$ is an ellipsoid centered at $\widehat{\boldsymbol{\theta}}_t^{\text{wls}}$ with the longest semi-axis of length $\beta_{t,\zeta_t(\delta)}^{1/2}(B)\lambda_{\min}(\mathbf{W}_t)^{-1/2}$, so by adding it to $\|\widehat{\boldsymbol{\theta}}_t^{\text{wls}}\|_2$, by the triangle inequality, we obtain a conservative upper bound on the distance between the origin and the point of $\mathcal{E}_{t,\zeta_t(\delta)}^{\text{wls}}$ furthest from it. We prove in the Supplementary Materials that

$$\mathbb{P}\left(\bigcap_{t=1}^{\infty}\{\widehat{B}_t \geq B^*\}\right) \geq 1 - \delta.$$

184 Continuing our description of the bound exploration phase, for any $t \leq T_B$ we use \widehat{B}_t to obtain a
 185 UCB algorithm, which we will refer to as empirical bound UCB (EB-UCB) via

$$A_t^{\text{EB-UCB}(\zeta_t(\delta))} = \arg \max_{a \in \mathcal{A}} \left\langle \boldsymbol{\phi}(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \|\boldsymbol{\phi}(a)\|_{\mathbf{W}_t^{-1}}. \quad (9)$$

186 Subsequently, we use

$$\begin{aligned} \widehat{\Delta}_{t,\zeta_t(\delta)}(a) &= \left\langle \boldsymbol{\phi}\left(A_t^{\text{EB-UCB}(\zeta_t(\delta))}\right) - \boldsymbol{\phi}(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle \\ &\quad + \beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \left(\left\| \boldsymbol{\phi}\left(A_t^{\text{EB-UCB}(\zeta_t(\delta))}\right) \right\|_{\mathbf{W}_t^{-1}} + \|\boldsymbol{\phi}(a)\|_{\mathbf{W}_t^{-1}} \right) \end{aligned} \quad (10)$$

187 as the gap estimate for any $a \in \mathcal{A}$. We define a new information gain criterion that combines model
 188 improvement (classic information gain) with bound improvement. The first component of our new
 189 information gain criterion is given by

$$I_t^{\text{EB-UCB}(\zeta_t(\delta))}(a) = \frac{1}{2} \log \left(\frac{\left\| \boldsymbol{\phi}\left(A_t^{\text{EB-UCB}(\zeta_t(\delta))}\right) \right\|_{\mathbf{W}_t^{-1}}^2}{\left\| \boldsymbol{\phi}\left(A_t^{\text{EB-UCB}(\zeta_t(\delta))}\right) \right\|_{(\mathbf{W}_t + \rho(a)^{-2} \boldsymbol{\phi}(a) \boldsymbol{\phi}(a)^\top)^{-1}}^2} \right), \quad (11)$$

for any $a \in \mathcal{A}$. It can be seen that this is analogous to the IDS-UCB information gain criterion considered by [Kirschner & Krause \(2018\)](#). To ensure that improvement in the bound on B^* , we introduce the second component of our information gain criterion I_t^B which is given by

$$I_t^B(a) = \frac{1}{2} \log \left(\|\mathbf{v}_t^{\min}\|_{(\mathbf{W}_t + \rho(a)^{-2} \boldsymbol{\phi}(a) \boldsymbol{\phi}(a)^\top)}^2 \right) - \frac{1}{2} \log \{\lambda_{\min}(\mathbf{W}_t)\},$$

190 where \mathbf{v}_t^{\min} is the unit-length eigenvector of \mathbf{W}_t associated with the smallest eigenvalue $\lambda_{\min}(\mathbf{W}_t)$.
 191 The maximizer of $I_t^B(a)$ corresponds to the feature vector $\boldsymbol{\phi}(a)$ which generates the most (weighted)
 192 information in the direction of the minimum eigenvector of the current information matrix. This
 193 direction corresponds to the longest axis of the confidence ellipsoid defined by the inverse information
 194 and is closely related to E-optimal experimental designs ([Dette & Studden, 1993](#)).

195 In order to balance exploration aimed at reducing the uncertainty about B^* and directly searching for
 196 the optimal arm in the initial phase, we use a mixture of information gain criteria, which we refer to
 197 as the bound-action mixture (BAM) criterion:

$$I_t^{\text{BAM}(\zeta_t(\delta))}(a) = \alpha I_t^B(a) + (1 - \alpha) I_t^{\text{EB-UCB}(\zeta_t(\delta))}(a),$$

198 where $\alpha \in (0, 1)$ is a parameter chosen by the user. Note that while we use the $I_t^{\text{EB-UCB}(\zeta_t(\delta))}$
 199 information gain criterion in this instance, we could use any information gain criterion of choice
 200 instead. For notational convenience we drop the $\zeta_t(\delta)$ term and write $I_t^{\text{EB-UCB}}$ for $I_t^{\text{EB-UCB}(\zeta_t(\delta))}$ and
 201 I_t^{BAM} for $I_t^{\text{BAM}(\zeta_t(\delta))}$ since we will use $\zeta_t(\delta) = \min\{\delta, 1/t^2\}$ in the remainder of this manuscript.

202 Given the advantages of deterministic IDS and its strong performance in various experimental settings,
 203 we focus on this variant of IDS. Hence, we always select the action which minimizes the information

204 ratio on the set \mathcal{A} , as given in (4). So at each time step $t \in \{1, \dots, T_B\}$ of the bound exploration
 205 phase we choose the action

$$A_t^{\text{BAM}} = \arg \min_{a \in \mathcal{A}} \left\{ \Psi_t^{\text{BAM}}(a) := \frac{\widehat{\Delta}_{t, \zeta_t(\delta)}^2(a)}{I_t^{\text{BAM}}(a)} \right\}.$$

206 Throughout the second phase, which we refer to as the *bound exploitation phase*, for any $t \geq T_B + 1$
 207 we use

$$\tilde{B}_t = \min \left\{ B, \min_{\tau \leq t} \left\{ \|\widehat{\boldsymbol{\theta}}_\tau^{\text{wls}}\|_2 + \beta_{\tau, \zeta_\tau(\delta)}^{1/2} (\widehat{B}_\tau) \lambda_{\min}(\mathbf{W}_\tau)^{-1/2} \right\} \right\}$$

208 as the upper bound on B^* , with \widehat{B}_t defined in (8). During this phase we drop the bound information
 209 gain criterion I_t^B from the mixture and use only the $I_t^{\text{EB-UCB}}$ criterion. The quantity \tilde{B}_t is used as
 210 the upper bound for B^* for both the gap estimate $\widehat{\Delta}_{t, \zeta_t(\delta)}$ and $I_t^{\text{EB-UCB}}$, which are defined in the
 211 same way as in equations (9), (10), and (11) with \tilde{B}_t in place of \widehat{B}_t . We summarize this method in
 212 Algorithm 2. Note that in the second phase we could use any algorithm which requires explicit use
 213 of an upper bound on B^* by taking $B = \tilde{B}_t$ as that upper bound. Furthermore, we formulate this
 214 procedure specifically in the context of IDS, however, the approach of estimating a high-probability
 215 upper bound on the true parameter norm and using it to guide decision making can be thought of as a
 216 more general technique, rather than something specific only to IDS.

217 6 Regret analysis of EBIDS algorithm

218 In this section we present the regret and pseudo-regret bounds for both phases of the EBIDS algorithm.
 219 We defer the proofs of these propositions and relevant lemmas to the Supplementary Materials. For
 220 any t and $\xi_t > 0$, let E_{t, ξ_t} be the event

$$E_{t, \xi_t} = \left\{ \left\| \boldsymbol{\theta}^* - \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_{\mathbf{W}_t}^2 \leq \beta_{t, \xi_t}(B^*) \right\}, \quad (12)$$

221 and define $E_\delta = \bigcap_{t=1}^\infty E_{t, \delta}$. Note that by Theorem 2 we have $\mathbb{P}(E_\delta) \geq 1 - \delta$. The following
 222 proposition summarizes the regret and pseudo-regret bounds for EBIDS during the bound exploration
 223 phase.

224 **Proposition 1.** *For any $2 \leq T \leq T_B$ the regret \mathcal{R}_T of Algorithm 2 is bounded above by*

$$\mathcal{R}_T \leq O \left(\frac{d \max\{U/\sqrt{\gamma}, \rho_{\max}\}}{\sqrt{1-\alpha}} \sqrt{T} \log T \sqrt{\log(1/\delta) + \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + \gamma B^2} \right)$$

225 *and whenever event E_δ holds the pseudo-regret \mathcal{PR}_T is bounded above by the same rate.*

226 We also provide guarantees on the estimated upper bound on B^* after the bound exploration phase.
 227 This, in turn, will allow us to obtain an improved bound for the regret and pseudo-regret in the
 228 subsequent phase.

229 **Proposition 2.** *For any constant $g > 0$, with sufficiently large T_B and sufficiently large α , whenever
 230 event E_δ holds we have $B^* \leq \tilde{B}_t \leq (1+g)B^*$ for any $t \geq T_B + 1$.*

231 Please see Section 10.6 in the Supplementary Materials for the exact constants required as lower
 232 bounds for T_B and α depending on g . Finally, using the results of Proposition 2 we are able
 233 to establish a regret bound for the second phase of EBIDS which is independent of the original
 234 conservative bound B .

235 **Proposition 3.** *For any constant $g > 0$, with sufficiently large T_B and sufficiently large α , with
 236 probability at least $1 - \delta$ the regret and pseudo-regret of Algorithm 2 are both bounded above by
 237 $O(dU\rho_{\max}(1+g)B^*\sqrt{T}\log T)$, for any $T \geq T_B + 1$.*

Algorithm 2 EBIDS

Input: Action set \mathcal{A} , penalty parameter $\gamma > 0$, noise function $\rho : \mathcal{A} \rightarrow \mathbb{R}_+$, feature function $\phi : \mathcal{A} \rightarrow \mathbb{R}$, conservative true parameter norm bound B , number of bound exploration steps T_B , information gain mixture parameter $\alpha \in (0, 1)$, error tolerance parameter $\delta \in (0, 1)$.

For $t = 1, 2, \dots, T_B$:

 Compute \mathbf{W}_t and $\widehat{\boldsymbol{\theta}}_t^{\text{wls}}$ using (5)

$$\widehat{B}_t \leftarrow \min \left\{ B, \|\widehat{\boldsymbol{\theta}}_t^{\text{wls}}\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2}(B) \lambda_{\min}(\mathbf{W}_t)^{-1/2} \right\}$$

$$A_t^{\text{EB-UCB}} \leftarrow \arg \max_{a \in \mathcal{A}} \left\langle \phi(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t, \zeta_t(\delta)}^{1/2}(\widehat{B}_t) \|\phi(a)\|_{\mathbf{W}_t^{-1}}$$

$$I_t^{\text{EB-UCB}}(a) \leftarrow \frac{1}{2} \log \left(\|\phi(A_t^{\text{EB-UCB}})\|_{\mathbf{W}_t^{-1}}^2 \right) - \frac{1}{2} \log \left(\|\phi(A_t^{\text{EB-UCB}})\|_{(\mathbf{W}_t + \rho(a)^{-2} \phi(a) \phi(a)^\top)^{-1}}^2 \right)$$

$$I_t^B(a) \leftarrow \frac{1}{2} \log \left(\|\mathbf{v}^{\min}\|_{(\mathbf{W}_t + \rho(a)^{-2} \phi(a) \phi(a)^\top)}^2 \right) - \frac{1}{2} \log \{ \lambda_{\min}(\mathbf{W}_t) \}$$

$$I_t^{\text{BAM}}(a) \leftarrow \alpha I_t^B(a) + (1 - \alpha) I_t^{\text{EB-UCB}}(a)$$

$$\widehat{\Delta}_{t, \zeta_t(\delta)}(a) \leftarrow \left\langle \phi(A_t^{\text{EB-UCB}}) - \phi(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t, \zeta_t(\delta)}^{1/2}(\widehat{B}_t) \left(\|\phi(A_t^{\text{EB-UCB}})\|_{\mathbf{W}_t^{-1}} + \|\phi(a)\|_{\mathbf{W}_t^{-1}} \right)$$

$$A_t \leftarrow \arg \min_{a \in \mathcal{A}} \widehat{\Delta}_{t, \zeta_t(\delta)}^2(a) / I_t^{\text{BAM}}(a)$$

 Play A_t , observe $Y_t = \langle \phi(A_t), \boldsymbol{\theta}^* \rangle + \eta_t$

For $t = T_B + 1, T_B + 2, \dots, T$:

 Compute \mathbf{W}_t and $\widehat{\boldsymbol{\theta}}_t^{\text{wls}}$ using (5)

$$\widehat{B}_t \leftarrow \min \left\{ B, \|\widehat{\boldsymbol{\theta}}_t^{\text{wls}}\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2}(B) \lambda_{\min}(\mathbf{W}_t)^{-1/2} \right\}$$

$$\tilde{B}_t \leftarrow \min \left\{ B, \min_{\tau \leq t} \left\{ \|\widehat{\boldsymbol{\theta}}_\tau^{\text{wls}}\|_2 + \beta_{\tau, \zeta_\tau(\delta)}^{1/2}(\tilde{B}_t) \lambda_{\min}(\mathbf{W}_\tau)^{-1/2} \right\} \right\}$$

$$A_t^{\text{EB-UCB}} \leftarrow \arg \max_{a \in \mathcal{A}} \left\langle \phi(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t, \zeta_t(\delta)}(\tilde{B}_t)^{1/2} \|\phi(a)\|_{\mathbf{W}_t^{-1}}$$

$$I_t^{\text{EB-UCB}}(a) \leftarrow \frac{1}{2} \log \left(\|\phi(A_t^{\text{EB-UCB}})\|_{\mathbf{W}_t^{-1}}^2 \right) - \frac{1}{2} \log \left(\|\phi(A_t^{\text{EB-UCB}})\|_{(\mathbf{W}_t + \rho(a)^{-2} \phi(a) \phi(a)^\top)^{-1}}^2 \right)$$

$$\widehat{\Delta}_{t, \zeta_t(\delta)}(a) \leftarrow \left\langle \phi(A_t^{\text{EB-UCB}}) - \phi(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\rangle + \beta_{t, \zeta_t(\delta)}^{1/2}(\tilde{B}_t) \left(\|\phi(A_t^{\text{EB-UCB}})\|_{\mathbf{W}_t^{-1}} + \|\phi(a)\|_{\mathbf{W}_t^{-1}} \right)$$

$$A_t \leftarrow \arg \min_{a \in \mathcal{A}} \widehat{\Delta}_{t, \zeta_t(\delta)}^2(a) / I_t^{\text{EB-UCB}}(a)$$

 Play A_t , observe $A_t = \langle \phi(A_t), \boldsymbol{\theta}^* \rangle + \eta_t$

238 Similarly, we give the exact constants required as lower bounds for T_B and α in Supplementary
 239 Materials, in Section 10.7. Thus, Propositions 1 and 3 together give us regret and pseudo-regret
 240 guarantees for both bound exploration phase and the subsequent bound exploitation phase of EBIDS.
 241 This is different from Gales et al. (2022) who do not control the regret in the initial stages of their
 242 norm-agnostic algorithms.

243 7 Simulation study

244 We evaluate the performance of EBIDS using simulation studies and compare it with the norm-
 245 agnostic competitor algorithms NAOFUL and OLSOFUL by Gales et al. (2022) which also aim
 246 at alleviating the dependence on access to a high-quality bound on the true parameter norm. We
 247 include the EB-UCB algorithm to demonstrate the advantage of using the IDS strategy in addition
 248 to utilizing the empirical norm bound. We run the comparison also against the oracle versions of
 249 EBIDS, IDS-UCB and UCB with access to the true value of B^* . We use the same setting as in the

250 simulation illustration in Section 4 with $\theta^* = [-5, 1, 1, 1.5, 2]^\top$ as the true parameter and ten arms
 251 with features sampled from $\text{Unif}[-1/\sqrt{5}, 1/\sqrt{5}]$. The error distribution for the first five arms are
 252 standard normal and for the remaining five arms they are normal with mean zero and variance 0.2. We
 253 take the conservative $B = 100$ as the assumed upper bound on B^* . Both the oracle and non-oracle
 254 versions of EBIDS use $\alpha = 0.5$, giving equal weight to both components of the BAM criterion, and
 255 run the bound exploration phase for $T_B = 50$ steps.

256 Figure 2 shows the mean regret averaged over 200 repeated experiments with $T = 500$ steps along
 257 with 95% pointwise confidence bounds. As we can see, EB-UCB is competitive with NAOFUL and
 258 OLSOFUL, while EBIDS performs best among all the algorithms which do not have access to the
 259 true parameter norm. It achieves significantly lower regret than IDS-UCB and UCB. Meanwhile, the
 260 performance of oracle EBIDS is better than that of oracle UCB and almost indistinguishable from the
 261 one achieved by oracle IDS-UCB.

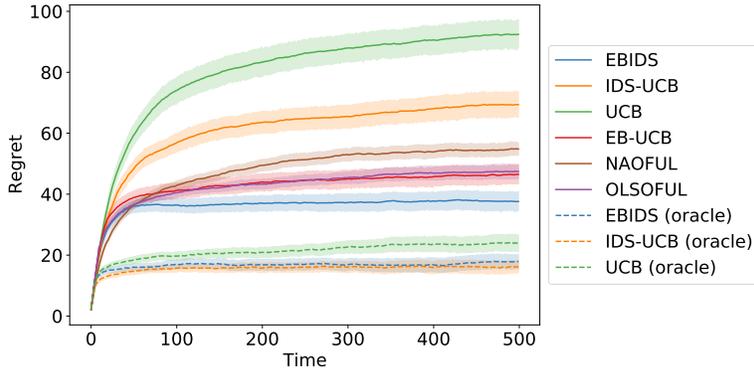


Figure 2: Regret incurred by EBIDS, EB-UCB, NAOFUL, OLSOFUL, IDS-UCB and UCB with conservative $B = 100$. We include the oracle versions of EBIDS, IDS-UCB and UCB using $B = B^*$ for reference. The solid and dashes lines represent the regret averaged over 200 repeated experiments, while the shaded bounds represent 95% pointwise confidence bounds.

262 We also perform an ablation study to determine the sensitivity of EBIDS to the tuning param-
 263 eter α and the length T_B of the bound exploration phase. We consider all combinations of
 264 $\alpha \in \{0.1, 0.3, 0.5, 0.7\}$ and $T_B \in \{50, 100\}$. We use the same setting as above and present the
 265 results for $T = 500$ steps averaged over 200 repeated experiments in Figure 3. Using $T_B = 50$
 266 leads to somewhat better results than $T_B = 100$ and $\alpha = 0.3$ performs best for both values of
 267 T_B . However, the performance is similar for all considered combinations of the tuning parameters,
 268 especially compared to the differences in performance of the competitor algorithms. This shows that
 269 while EBIDS, like most other bandit algorithms, uses tuning parameters, its performance is not very
 270 sensitive to their choice, with several considered combinations of α and T_B achieving practically
 271 indistinguishable regret.

272 8 Discussion

273 Bandit algorithms often require access to a high-quality upper bound on the Euclidean norm of
 274 the true parameter vector in order to achieve good performance. In practice, such information is
 275 rarely available *a priori*, which can lead to significant regret accumulation. Despite its prevalence,
 276 this problem has received relatively little attention in the bandit literature. We introduced the
 277 empirical bound information-directed sampling (EBIDS) algorithm which addresses this challenge by
 278 iteratively refining a high-probability upper bound on the true parameter norm. We developed a novel
 279 information criterion that balances tightening the bound on the true parameter norm and explicitly
 280 searching for the optimal arm. In simulation experiments, EBIDS showed improved performance
 281 compared to the competing norm-agnostic algorithms. Furthermore, we proved regret bounds that

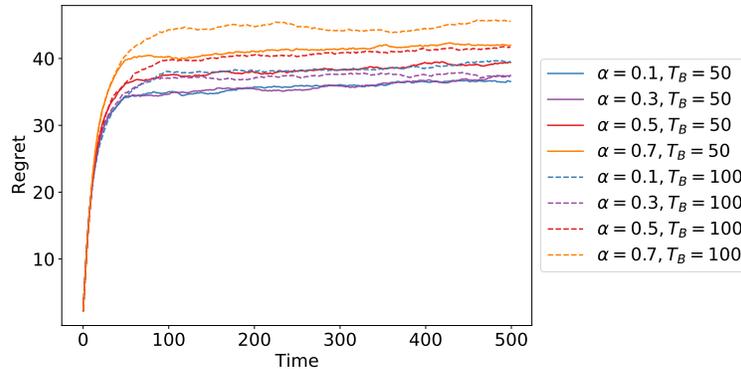


Figure 3: Average regret for EBIDS averaged over 200 repeated experiments with $T = 500$ steps under different values of the tuning parameter α and the length T_B of the bound exploration phase.

282 eventually do not depend on the initially assumed bound for the parameter norm, and unlike prior
 283 regret guarantees, our bounds are anytime in that they apply to all phases of the algorithm.

284 **Broader Impact Statement**

285 This paper introduces novel methodology for frequentist IDS that does not require strong prior
 286 information on the norm of the true parameter indexing the reward model. Our methodology, which
 287 involves a novel information criterion, can be viewed as a general approach to balancing bound
 288 improvement and regret minimization that is applicable in a wide range of UCB and IDS bandit
 289 algorithms.

290 **References**

- 291 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic
 292 bandits. In *Advances in Neural Information Processing Systems*, volume 24, 2011.
- 293 Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In
 294 *International conference on machine learning*, pp. 127–135. PMLR, 2013.
- 295 Amir Ahmadi-Javid, Pardis Seyed, and Siddhartha S Syam. A survey of healthcare facility location.
 296 *Computers & Operations Research*, 79:223–263, 2017.
- 297 Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine
 298 Learning Research*, 3:397–422, 2002.
- 299 Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, and Gilles Stoltz.
 300 Kullback-Leibler upper confidence bounds for optimal sequential allocation. *The Annals of
 301 Statistics*, 41(3):1516–1541, 2013.
- 302 Niladri Chatterji, Vidya Muthukumar, and Peter Bartlett. OSOM: A simultaneously optimal algorithm
 303 for multi-armed and linear contextual bandits. In *Proceedings of the Twenty Third International
 304 Conference on Artificial Intelligence and Statistics*, volume 108, pp. 1844–1854. PMLR, 2020.
- 305 Thomas M Cover and Joy A Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.
- 306 Varsha Dani, Thomas P. Hayes, and Sham M. Kakade. Stochastic linear optimization under bandit
 307 feedback. *21st Annual Conference on Learning Theory*, pp. 355–366, 2008.
- 308 Holger Dette and William J Studden. Geometry of E-optimality. *The Annals of Statistics*, 21(1):
 309 416–433, 1993.

- 310 Joel N. Franklin. *Matrix Theory*. Prentice-Hall, 1968.
- 311 Spencer B. Gales, Sunder Sethuraman, and Kwang-Sung Jun. Norm-agnostic linear bandits. In
312 *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume
313 151, pp. 73–91. PMLR, 2022.
- 314 Aurélien Garivier and Olivier Cappé. The KL-UCB algorithm for bounded stochastic bandits and
315 beyond. In *Proceedings of the 24th Annual Conference on Learning Theory*, volume 19, pp.
316 359–376. PMLR, 2011.
- 317 Claudio Gentile and Francesco Orabona. On multilabel classification and ranking with bandit
318 feedback. *Journal of Machine Learning Research*, 15(70):2451–2487, 2014.
- 319 Avishek Ghosh, Abishek Sankararaman, and Ramchandran Kannan. Problem-complexity adaptive
320 model selection for stochastic linear bandits. In *Proceedings of The 24th International Conference
321 on Artificial Intelligence and Statistics*, volume 130, pp. 1396–1404. PMLR, 2021.
- 322 Botao Hao and Tor Lattimore. Regret bounds for information-directed reinforcement learning. In
323 *Advances in Neural Information Processing Systems*, volume 35, 2022.
- 324 Yu-Heng Hung, Ping-Chun Hsieh, Xi Liu, and P. R. Kumar. Reward-biased maximum likelihood esti-
325 mation for linear stochastic bandits. *Proceedings of the AAAI Conference on Artificial Intelligence*,
326 35(9):7874–7882, 2021.
- 327 Johannes Kirschner. *Information-Directed Sampling-Frequentist Analysis and Applications*. PhD
328 thesis, ETH Zurich, 2021.
- 329 Johannes Kirschner and Andreas Krause. Information directed sampling and bandits with het-
330 eroscedastic noise. In *Proceedings of the 31st Conference On Learning Theory*, volume 75, pp.
331 358–384. PMLR, 2018.
- 332 Johannes Kirschner, Tor Lattimore, and Andreas Krause. Information directed sampling for linear
333 partial monitoring. In *Proceedings of Thirty Third Conference on Learning Theory*, volume 125,
334 pp. 2328–2369. PMLR, 2020.
- 335 Johannes Kirschner, Tor Lattimore, Claire Vernade, and Csaba Szepesvari. Asymptotically optimal
336 information-directed sampling. In *Proceedings of Thirty Fourth Conference on Learning Theory*,
337 volume 134, pp. 2777–2821. PMLR, 2021.
- 338 Tomáš Kocák, Rémi Munos, Branislav Kveton, Shipra Agrawal, and Michal Valko. Spectral bandits.
339 *Journal of Machine Learning Research*, 21(218):1–44, 2020.
- 340 Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- 341 David Lindner, Matteo Turchetta, Sebastian Tschieschek, Kamil Ciosek, and Andreas Krause.
342 Information directed reward learning for reinforcement learning. In *Advances in Neural Information
343 Processing Systems*, volume 34, 2021.
- 344 Nikolay Nikolov, Johannes Kirschner, Felix Berkenkamp, and Andreas Krause. Information-directed
345 exploration for deep reinforcement learning. In *International Conference on Learning Representa-
346 tions*, 2019.
- 347 Francesco Orabona and Nicolo Cesa-Bianchi. Better algorithms for selective sampling. In *Proceed-
348 ings of the 28th International Conference on Machine Learning*, pp. 433–440, 2011.
- 349 My Phan, Yasin Abbasi Yadkori, and Justin Domke. Thompson sampling and approximate inference.
350 In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 351 Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling. In
352 *Advances in Neural Information Processing Systems*, volume 27, 2014.

- 353 William R Thompson. On the likelihood that one unknown probability exceeds another in view of
354 the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 355 Zhilei Wang, Pranjal Awasthi, Christoph Dann, Ayush Sekhari, and Claudio Gentile. Neural active
356 learning with performance guarantees. In *Advances in Neural Information Processing Systems*,
357 volume 34, 2021.
- 358 Justin Wertz, Alex Volfovsky, and Eric B. Laber. Reinforcement learning methods in public health.
359 *Clinical Therapeutics*, 44(1):139–154, 2022.
- 360 Justin Wertz, Tanner Fiez, Alexander Volfovsky, Eric Laber, Blake Mason, Houssam Nassif, and
361 Lalit Jain. Experimental designs for heteroskedastic variance. In *Advances in Neural Information*
362 *Processing Systems*, volume 36, 2023.
- 363 Ruitu Xu, Yifei Min, and Tianhao Wang. Noise-adaptive thompson sampling for linear contextual
364 bandits. In *Advances in Neural Information Processing Systems*, volume 36, 2023.
- 365 Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with UCB-based exploration.
366 In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 11492–
367 11502. PMLR, 2020.

368
369
370

Supplementary Materials

The following content was not necessarily subject to peer review.

371 In these supplementary materials we provide the proofs to the propositions we have stated in the
372 paper.

373 9 Notation and lemmas

We begin by introducing some notation and basic facts. For any unit vector $\mathbf{v} \in \mathbb{R}^d$ and any $a \in \mathcal{A}$, let $\psi_{\mathbf{v}}(\phi(a)), \psi_{\mathbf{v}}^{\perp}(\phi(a)) \in \mathbb{R}$ denote the orthogonal decomposition of $\phi(a)$, i.e.,

$$\phi(a) = \psi_{\mathbf{v}}(\phi(a))\mathbf{v} + \psi_{\mathbf{v}}^{\perp}(\phi(a))\mathbf{v}^{\perp},$$

374 where $\|\mathbf{v}^{\perp}\|_2 = 1$ and $\mathbf{v}^{\perp} \perp \mathbf{v}$. Let

$$\kappa = \min_{\mathbf{v} \in \mathbb{R}^d \text{ s.t. } \|\mathbf{v}\|_2=1} \max_{a \in \mathcal{A}} \{\rho(a)^{-2} \psi_{\mathbf{v}}(\phi(a))^2\}. \quad (13)$$

375 Note that $\kappa > 0$. Let

$$\omega_t(a) = \rho(a)^{-2} \psi_{\mathbf{v}_t^{\min}}(\phi(a))^2. \quad (14)$$

376 Also, note that for any $a \in \mathcal{A}$ we have

$$\|\phi(a)\|_{\mathbf{W}_t^{-1}}^2 = \sum_{i=1}^d \psi_{\mathbf{v}_i}(\phi(a))^2 \lambda_i^{-1}$$

377 where $\{(\lambda_i, \mathbf{v}_i)\}_{i=1}^d$ are the eigenvalue-eigenvector pairs of \mathbf{W}_t . Hence for every $t \geq 1$ and $a \in \mathcal{A}$
378 we have

$$\|\phi(a)\|_2^2 \lambda_{\max}(\mathbf{W}_t)^{-1} \leq \|\phi(a)\|_{\mathbf{W}_t^{-1}}^2 \leq \|\phi(a)\|_2^2 \lambda_{\min}(\mathbf{W}_t)^{-1},$$

379 so

$$L^2 \lambda_{\max}(\mathbf{W}_t)^{-1} \leq \|\phi(a)\|_{\mathbf{W}_t^{-1}}^2 \leq U^2 \lambda_{\min}(\mathbf{W}_t)^{-1}. \quad (15)$$

380 Also from Cauchy-Schwarz inequality

$$\langle \phi(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle^2 \leq \|\phi(a)\|_2^2 \|\widehat{\boldsymbol{\theta}}_t^{\text{wls}}\|_2^2 \leq U^2 \|\widehat{\boldsymbol{\theta}}_t^{\text{wls}}\|_2^2 \quad (16)$$

381 From Weyl's inequality (Franklin, 1968), for any positive semi-definite matrices \mathbf{A}, \mathbf{B} we have

$$\lambda_{\max}(\mathbf{A} + \mathbf{B}) \leq \lambda_{\max}(\mathbf{A}) + \lambda_{\max}(\mathbf{B}).$$

382 Thus, for every $t \geq 1$ we have

$$\lambda_{\max}(\mathbf{W}_t) \leq \lambda_{\max}(\gamma \mathbf{I}_d) + \sum_{\tau=1}^{t-1} \lambda_{\max}(\rho(a_{\tau})^{-2} \phi(a_{\tau}) \phi(a_{\tau})^{\top}) \leq \gamma + (t-1) \rho_{\min}^{-2} U^2, \quad (17)$$

383 so from (15), for any $t \geq 1$ we have

$$\|\phi(a)\|_{\mathbf{W}_t^{-1}}^2 \geq \frac{L^2}{\gamma + (t-1) \rho_{\min}^{-2} U^2} \geq \frac{L^2}{t(\gamma + \rho_{\min}^{-2} U^2)}. \quad (18)$$

384 Also from (17) for $T \geq 2$ we have

$$\begin{aligned} \log \left(\frac{\det(\mathbf{W}_T)}{\det(\mathbf{W}_1)} \right) &= \log(\det(\mathbf{W}_T)) - \log(\det(\gamma \mathbf{I}_d)) \leq d \log(\gamma + (T-1)\rho_{\min}^{-2}U^2) - d \log \gamma \\ &= d \log \left[1 + (T-1) \frac{\rho_{\min}^{-2}U^2}{\gamma} \right] \leq d \log \left[(T-1) \left(1 + \frac{\rho_{\min}^{-2}U^2}{\gamma} \right) \right] \\ &= d \log(T-1) + d \log \left(1 + \frac{\rho_{\min}^{-2}U^2}{\gamma} \right). \end{aligned} \quad (19)$$

385 Applying the data processing inequality (Cover & Thomas, 2012) in an analogous way as Kirschner
386 & Krause (2018), we obtain

$$I_t^{\text{EB-UCB}}(a) \leq \frac{1}{2} \log \left(\frac{\det(\mathbf{W}_t + \rho(a)^{-2} \boldsymbol{\phi}(a) \boldsymbol{\phi}(a)^\top)}{\det(\mathbf{W}_t)} \right) = \frac{1}{2} \log \left(1 + \rho(a)^{-2} \|\boldsymbol{\phi}(a)\|_{\mathbf{W}_t^{-1}}^2 \right) \quad (20)$$

387 for any $a \in \mathcal{A}$. So from (19) we get

$$\sum_{t=1}^T I_t^{\text{EB-UCB}}(a_t) \leq \frac{1}{2} \log \left(\frac{\det(\mathbf{W}_{T+1})}{\det(\mathbf{W}_1)} \right) \leq \frac{1}{2} d \log T + \frac{1}{2} d \log \left(1 + \frac{\rho_{\min}^{-2}U^2}{\gamma} \right) = O(d \log T), \quad (21)$$

388 for any sequence $\{a_t\}_{t=1}^T \subset \mathcal{A}$.

389 We now state and prove some additional lemmas that will be useful throughout the proofs of
390 Propositions 1 - 3.

Lemma 1. Let $\widehat{\Delta}_t : \mathcal{A} \rightarrow \mathbb{R}_+$ be a gap estimate function and let $I_t^X, I_t^Y : \mathcal{A} \rightarrow \mathbb{R}_+$ be two information gain criteria. Let I_t^{XY} be the mixture information gain criterion given by

$$I_t^{XY}(a) = \alpha I_t^X(a) + (1 - \alpha) I_t^Y(a)$$

for some $\alpha \in (0, 1)$. Consider now a deterministic IDS algorithm which at each time step t plays action a_t^{XY} given by

$$a_t^{XY} = \arg \min_{a \in \mathcal{A}} \frac{\widehat{\Delta}_t^2(a)}{I_t^{XY}(a)}$$

Then at each time step t the information gain on to the first criterion I_t^X is lower-bounded by

$$I_t^X(a_t^{XY}) \geq \frac{\widehat{\Delta}_t^2(a_t^{XY})}{\widehat{\Delta}_t^2(a_t^{I,X})} I_t^X(a_t^{I,X}) - \frac{1 - \alpha}{\alpha} I_t^Y(a_t^{XY}),$$

391 where $a_t^{I,X} = \arg \max_{a \in \mathcal{A}} I_t^X(a)$.

Lemma 2. Recall the definition $\omega_t(a) = \rho(a)^{-2} \psi_{\mathbf{v}_{\min}}(\boldsymbol{\phi}(a))^2$. For any $T \geq 1$ and any sequence of actions $\{a_t\}_{t=1}^T$ we have

$$\lambda_{\min}(\mathbf{W}_{T+1}) \geq \gamma - \rho_{\min}^{-2}U^2 + \frac{1}{d} \sum_{t=1}^T \omega_t(a_t).$$

Lemma 3. Let $\{x_t\}_{t=1}^{T+1} \subset [0, U]$ be a bounded sequence for some constant $U > 0$. Then for any constant $c > 0$ we have

$$\sum_{t=1}^T \frac{x_{t+1}}{c + \sum_{\tau=1}^t x_\tau} \leq \log T + \frac{U}{c} + 1.$$

392 10 Proofs of theoretical results

393 In this section we provide the proofs of Theorem 1, Lemmas 1 - 3, and Propositions 1 - 3.

394 **10.1 Proof of Theorem 1**

395 Recall that by Cauchy-Schwarz inequality, for any random variables $\{X_t\}_{t \in G}, \{Y_t\}_{t \in G}$ with non-
 396 negative support, with probability 1 we have

$$\sum_{t \in G} \sqrt{X_t Y_t} \leq \sqrt{\left(\sum_{t \in G} X_t \right) \left(\sum_{t \in G} Y_t \right)},$$

397 and for any random variables X, Y with nonnegative support we have

$$\mathbb{E} \left[\sqrt{XY} \right] \leq \sqrt{\mathbb{E}[X] \mathbb{E}[Y]}.$$

Hence if $\widehat{\Delta}(A_t) \geq \Delta(A_t)$, for all $t \in G$, then with probability 1 we have

$$\sum_{t \in G} \Delta(A_t) \leq \sum_{t \in G} \widehat{\Delta}_t(A_t) = \sum_{t \in G} \sqrt{\Psi_t(A_t) I_t(A_t)} \leq \sqrt{\left[\sum_{t \in G} \Psi_t(A_t) \right] \left[\sum_{t \in G} I_t(A_t) \right]}.$$

398 Also

$$\begin{aligned} \mathbb{E} \left[\sum_{t \in G} \widehat{\Delta}_t(A_t) \right] &= \mathbb{E} \left[\sum_{t \in G} \sqrt{\Psi_t(A_t) I_t(A_t)} \right] \leq \mathbb{E} \left(\sqrt{\left[\sum_{t \in G} \Psi_t(A_t) \right] \left[\sum_{t \in G} I_t(A_t) \right]} \right) \\ &\leq \sqrt{\mathbb{E} \left[\sum_{t \in G} \Psi_t(A_t) \right] \mathbb{E} \left[\sum_{t \in G} I_t(A_t) \right]}. \quad \square \end{aligned}$$

399 **10.2 Proof of Lemma 1**

By the definition of a_t^{XY} we have

$$\frac{\widehat{\Delta}_t^2(a_t^{XY})}{\alpha I_t^X(a_t^{XY}) + (1 - \alpha) I_t^Y(a_t^{XY})} \leq \frac{\widehat{\Delta}_t^2(a_t^{I,X})}{\alpha I_t^X(a_t^{I,X}) + (1 - \alpha) I_t^Y(a_t^{I,X})},$$

hence

$$\alpha I_t^X(a_t^{XY}) + (1 - \alpha) I_t^Y(a_t^{XY}) \geq \frac{\widehat{\Delta}_t^2(a_t^{XY})}{\widehat{\Delta}_t^2(a_t^{I,X})} \left[\alpha I_t^X(a_t^{I,X}) + (1 - \alpha) I_t^Y(a_t^{I,X}) \right],$$

400 and thus

$$\begin{aligned} I_t^X(a_t^{XY}) &\geq \frac{\widehat{\Delta}_t^2(a_t^{XY})}{\widehat{\Delta}_t^2(a_t^{I,X})} I_t^X(a_t^{I,X}) + \frac{(1 - \alpha)}{\alpha} \cdot \frac{\widehat{\Delta}_t^2(a_t^{XY})}{\widehat{\Delta}_t^2(a_t^{I,X})} I_t^Y(a_t^{I,X}) - \frac{1 - \alpha}{\alpha} I_t^Y(a_t^{XY}) \\ &\geq \frac{\widehat{\Delta}_t^2(a_t^{XY})}{\widehat{\Delta}_t^2(a_t^{I,X})} I_t^X(a_t^{I,X}) - \frac{1 - \alpha}{\alpha} I_t^Y(a_t^{XY}). \quad \square \end{aligned}$$

401 **10.3 Proof of Lemma 2**

Recall that we define $\lambda_1^{(t)}, \dots, \lambda_d^{(t)}$ as the (not necessarily ordered) eigenvalues of \mathbf{W}_t . Let

$$i^*(t) = \arg \min_{1 \leq i \leq d} \lambda_i^{(t)}.$$

402 By Weyl's inequality (Franklin, 1968), for any symmetric positive semi-definite matrices $\mathbf{A}, \mathbf{B} \in$
 403 $\mathbb{R}^{m \times m}$ we have

$$\lambda_{(i)}(\mathbf{A} + \mathbf{B}) \geq \lambda_{(i)}(\mathbf{A}) \tag{22}$$

404 where $\lambda_{(i)}(\mathbf{A})$ is the i -th largest eigenvalue of \mathbf{A} for any $1 \leq i \leq m$. Let \mathbf{v}_t^{\min} be the unit eigenvector
 405 corresponding to the smallest eigenvalue of \mathbf{W}_t . Then for any $1 \leq i \leq d$ we have

$$\begin{aligned} \lambda_{(i)}(\mathbf{W}_{t+1}) &= \lambda_{(i)}(\mathbf{W}_t + \rho(a_t)^{-2} \phi(a_t) \phi(a_t)^\top) \\ &= \lambda_{(i)}\left(\mathbf{W}_t + \rho(a_t)^{-2} \psi_{\mathbf{v}_t^{\min}}(\phi(a)) \mathbf{v}_t^{\min} (\mathbf{v}_t^{\min})^\top \right. \\ &\quad \left. + \rho(a_t)^{-2} \psi_{\mathbf{v}_t^{\min\perp}}^\perp(\phi(a)) \mathbf{v}_t^{\min\perp} (\mathbf{v}_t^{\min\perp})^\top\right) \\ &\geq \lambda_{(i)}\left(\mathbf{W}_t + \rho(a_t)^{-2} \psi_{\mathbf{v}_t^{\min}}(\phi(a)) \mathbf{v}_t^{\min} (\mathbf{v}_t^{\min})^\top\right) \\ &= \lambda_{(i)}(\mathbf{W}_t + \omega_t(a_t) \mathbf{v}_t^{\min} (\mathbf{v}_t^{\min})^\top). \end{aligned}$$

Note that the matrix $\mathbf{W}_t + \omega_t(a_t) \mathbf{v}_t^{\min} (\mathbf{v}_t^{\min})^\top$ has the same eigenvectors as \mathbf{W}_t and the smallest eigenvalue of \mathbf{W}_t , i.e., the one corresponding to \mathbf{v}_t^{\min} is increased by $\omega_t(a_t)$. So for any t we can order the eigenvalues $\lambda_1^{(t+1)}, \dots, \lambda_d^{(t+1)}$ in such way that $\lambda_i^{(t+1)} \geq \lambda_i^{(t)}$ and

$$\lambda_{i^*(t)}^{(t+1)} \geq \lambda_{i^*(t)}^{(t)} + \omega_t(a_t).$$

Since we have d eigenvalues and at each time step t we add at least $\omega_t(a_t)$ to the smallest eigenvalue at that time step without reducing the other ones we have

$$\lambda_{i^*(T)}^{(T)} - \lambda_{i^*(1)}^{(1)} + \omega_T(a_T) \geq \frac{1}{d} \sum_{t=1}^T \omega_t(a_t).$$

406 Note that $\lambda_1^{(1)} = \dots = \lambda_d^{(1)} = \gamma$ and $\omega_T(a_T) \leq \rho_{\min}^{-2} U^2$, so

$$\lambda_{\min}(\mathbf{W}_{T+1}) = \lambda_{i^*(T+1)}^{(T+1)} \geq \lambda_{i^*(T)}^{(T)} \geq \gamma - \rho_{\min}^{-2} U^2 + \frac{1}{d} \sum_{t=1}^T \omega_t(a_t). \quad \square$$

407 10.4 Proof of Lemma 3

Let

$$f(x_1, \dots, x_{T+1}) = \sum_{t=1}^T \frac{x_{t+1}}{c + \sum_{\tau=1}^t x_\tau}$$

We use induction to show that the maximum of f is achieved at $x_1 = 0$ and

$$x_2 = x_3 = \dots = x_{T+1} = U.$$

Note that for any $\tilde{x}_1, \dots, \tilde{x}_T \in [0, U]$ we have

$$\arg \max_{x_{T+1} \in [0, U]} f(\tilde{x}_1, \dots, \tilde{x}_T, x_{T+1}) = U.$$

408 Suppose that for any $t \geq 2$ it holds that for any $t \leq k \leq T$ and any $\tilde{x}_1, \dots, \tilde{x}_k \in [0, U]$ we have

$$(x_{k+1}^*, \dots, x_{T+1}^*) := \arg \max_{x_{k+1}, \dots, x_{T+1} \in [0, U]} f(\tilde{x}_1, \dots, \tilde{x}_k, x_{k+1}, \dots, x_{T+1}) = (U, \dots, U) \in \mathbb{R}^{T-k+1}. \quad (23)$$

409 Take any $\tilde{x}_1, \dots, \tilde{x}_{t-1} \in [0, U]$. Then by taking $k = t + 1$ the above statement gives us

$$\max_{x_t, x_{t+1}, \dots, x_{T+1} \in [0, U]} f(\tilde{x}_1, \dots, \tilde{x}_{t-1}, x_t, x_{t+1}, \dots, x_{T+1}) = \max_{x_t, x_{t+1} \in [0, U]} f(\tilde{x}_1, \dots, \tilde{x}_{t-1}, x_t, x_{t+1}, U, \dots, U).$$

410 Let

$$(\tilde{x}_t, \tilde{x}_{t+1}) = \arg \max_{x_t, x_{t+1} \in [0, U]} f(\tilde{x}_1, \dots, \tilde{x}_{t-1}, x_t, x_{t+1}, \dots, x_{T+1}).$$

411 Note that $\tilde{x}_{t+1} = x_{t+1}^* = U$ by taking the induction statement with $k = t$. For notational convenience
 412 let $b = c + \sum_{\tau=1}^{t-1} \tilde{x}_\tau$. Then

$$(\tilde{x}_t, \tilde{x}_{t+1}) = \arg \max_{x_t, x_{t+1} \in [0, U]} \left\{ \frac{x_t}{b} + \frac{x_{t+1}}{b + x_t} + \sum_{\tau=0}^{T-t-1} \frac{U}{b + x_t + x_{t+1} + \tau U} \right\}.$$

413 Let

$$g_t(x_t, x_{t+1}) = \frac{x_t}{b} + \frac{x_{t+1}}{b + x_t} + \sum_{\tau=0}^{T-t-1} \frac{U}{b + x_t + x_{t+1} + \tau U}.$$

414 Suppose that $\tilde{x}_t = x$ for some $0 \leq x < U$. Note that

$$g_t(U, x) - g_t(x, U) = \left(\frac{U}{b} + \frac{x}{b + U} \right) - \left(\frac{x}{b} + \frac{U}{b + x} \right) = \frac{Ux(U - x)}{b(b + U)(b + x)} > 0.$$

So $g_t(U, x) > g_t(x, U) = g_t(\tilde{x}_t, \tilde{x}_{t+1})$ which is a contradiction, since $(\tilde{x}_t, \tilde{x}_{t+1})$ is the maximizer of $g_t(x_t, x_{t+1})$. So $\tilde{x}_t = U$. Thus, we have shown that for any $\tilde{x}_1, \dots, \tilde{x}_{t-1} \in [0, U]$ we have

$$(x_t^*, \dots, x_{T+1}^*) = \arg \max_{x_t, \dots, x_{T+1} \in [0, U]} f(\tilde{x}_1, \dots, \tilde{x}_{t-1}, x_t, \dots, x_{T+1}) = (U, \dots, U) \in \mathbb{R}^{T-t+2}.$$

Hence by induction we get that for any $\tilde{x}_1 \in [0, U]$ we have

$$\arg \max_{x_2, \dots, x_{T+1} \in [0, U]} f(\tilde{x}_1, x_2, \dots, x_{T+1}) = (U, \dots, U) \in \mathbb{R}^T.$$

Clearly

$$\arg \max_{x_1 \in [0, U]} f(x_1, U, \dots, U) = 0,$$

415 so

$$\begin{aligned} \max_{x_1, \dots, x_{T+1} \in [0, U]} f(x_1, \dots, x_{T+1}) &= f(0, U, \dots, U) = \sum_{t=1}^T \frac{U}{c + (t-1)U} \\ &\leq \frac{U}{c} + \sum_{t=2}^T \frac{1}{t-1} < \log T + \frac{U}{c} + 1. \quad \square \end{aligned}$$

416 10.5 Proof of Proposition 1

417 From Theorem 1, for any $T \leq T_B$ we have

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \widehat{\Delta}_{t, \zeta_t(\delta)}(A_t^{\text{BAM}}) \right] &\leq \sqrt{\left(\mathbb{E} \left[\sum_{t=1}^T \Psi_t^{\text{BAM}}(A_t^{\text{BAM}}) \right] \right) \left(\mathbb{E} \left[\sum_{t=1}^T I_t^{\text{BAM}}(A_t^{\text{BAM}}) \right] \right)} \\ &\leq \sqrt{\left(\mathbb{E} \left[\sum_{t=1}^T \frac{\widehat{\Delta}_{t, \zeta_t(\delta)}^2(A_t^{\text{BAM}})}{I_t^{\text{BAM}}(A_t^{\text{BAM}})} \right] \right) \left(\mathbb{E} \left[\sum_{t=1}^T I_t^{\text{BAM}}(A_t^{\text{BAM}}) \right] \right)} \\ &= \sqrt{\mathbb{E} \left[\sum_{t=1}^T \frac{\widehat{\Delta}_{t, \zeta_t(\delta)}^2(A_t^{\text{BAM}})}{\alpha I_t^{\text{B}}(A_t^{\text{BAM}}) + (1 - \alpha) I_t^{\text{EB-UCB}}(A_t^{\text{BAM}})} \right]} \\ &\quad \times \sqrt{\alpha \mathbb{E} \left[\sum_{t=1}^T I_t^{\text{B}}(A_t^{\text{BAM}}) \right] + (1 - \alpha) \mathbb{E} \left[\sum_{t=1}^T I_t^{\text{EB-UCB}}(A_t^{\text{BAM}}) \right]}. \end{aligned}$$

418 By the definition of A_t^{BAM} we have

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2 (A_t^{\text{BAM}})}{\alpha I_t^{\text{B}} (A_t^{\text{BAM}}) + (1-\alpha) I_t^{\text{EB-UCB}} (A_t^{\text{BAM}})} \right] &\leq \mathbb{E} \left[\sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2 (A_t^{\text{EB-UCB}})}{\alpha I_t^{\text{EB-UCB}} (A_t^{\text{EB-UCB}}) + (1-\alpha) I_t^{\text{B}} (A_t^{\text{EB-UCB}})} \right] \\ &\leq \frac{1}{1-\alpha} \mathbb{E} \left[\sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2 (A_t^{\text{EB-UCB}})}{I_t^{\text{EB-UCB}} (A_t^{\text{EB-UCB}})} \right]. \end{aligned} \quad (24)$$

419 The next couple of steps are similar to the analysis by [Kirschner \(2021\)](#). Let $a_t^{\text{EB-UCB}}$ be the realization
420 of $A_t^{\text{EB-UCB}}$. From the Sherman-Morrison formula, we obtain

$$(\mathbf{W}_t + \rho(a_t^{\text{EB-UCB}})^{-2} \phi(a_t^{\text{EB-UCB}}) \phi(a_t^{\text{EB-UCB}})^\top)^{-1} = \mathbf{W}_t^{-1} - \frac{\rho(a_t^{\text{EB-UCB}})^{-2} \mathbf{W}_t^{-1} \phi(a_t^{\text{EB-UCB}}) \phi(a_t^{\text{EB-UCB}})^\top \mathbf{W}_t^{-1}}{1 + \rho(a_t^{\text{EB-UCB}})^{-2} \phi(a_t^{\text{EB-UCB}})^\top \mathbf{W}_t^{-1} \phi(a_t^{\text{EB-UCB}})}$$

421 so

$$\left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{(\mathbf{W}_t + \rho(a_t^{\text{EB-UCB}})^{-2} \phi(a_t^{\text{EB-UCB}}) \phi(a_t^{\text{EB-UCB}})^\top)^{-1}}^2 = \left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}}^2 - \frac{\rho(a_t^{\text{EB-UCB}})^{-2} \left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}}^4}{1 + \rho(a_t^{\text{EB-UCB}})^{-2} \left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}}^2}.$$

422 Thus

$$I_t^{\text{EB-UCB}} (a_t^{\text{EB-UCB}}) = \frac{1}{2} \log \left(1 + \rho(a_t^{\text{EB-UCB}})^{-2} \left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}}^2 \right).$$

423 From (15), we have $\left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}} \leq U^2 \gamma^{-1}$. Thus, using the fact that $\log(1+x) \geq \frac{x}{2q}$ for
424 $q \geq 1$ and $x \in [0, q]$ we get

$$I_t^{\text{EB-UCB}} (a_t^{\text{EB-UCB}}) \geq \frac{1}{4} \min\{U^{-2}\gamma, \rho(a_t^{\text{EB-UCB}})^{-2}\} \left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}}^2.$$

425 So

$$\begin{aligned} \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2 (a_t^{\text{EB-UCB}})}{I_t^{\text{EB-UCB}} (a_t^{\text{EB-UCB}})} &\leq \frac{4\beta_{t,\zeta_t(\delta)}(\widehat{B}_t) \left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}}^2}{\frac{1}{4} \min\{U^{-2}\gamma, \rho(a_t^{\text{EB-UCB}})^{-2}\} \left\| \phi(a_t^{\text{EB-UCB}}) \right\|_{\mathbf{W}_t^{-1}}^2} \\ &= 16\beta_{t,\zeta_t(\delta)}(\widehat{B}_t) \max\{U^2\gamma^{-1}, \rho(a_t^{\text{EB-UCB}})^2\} \\ &\leq 16\beta_{t,\zeta_t(\delta)}(B) \max\{U^2\gamma^{-1}, \rho(a_t^{\text{EB-UCB}})^2\} \\ &\leq 16\beta_{T,\zeta_T(\delta)}(B) \max\{U^2\gamma^{-1}, \rho_{\max}^2\}. \end{aligned} \quad (25)$$

426 So

$$\mathbb{E} \left[\sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2 (A_t^{\text{EB-UCB}})}{I_t^{\text{EB-UCB}} (A_t^{\text{EB-UCB}})} \right] \leq 16T\beta_{T,\zeta_T(\delta)}(B) \max\{U^2\gamma^{-1}, \rho_{\max}^2\}. \quad (26)$$

427 Since $1/\zeta_T(\delta) = \max\{1/\delta, T^2\}$, from (19) we have

$$\begin{aligned} \beta_{T,\zeta_T(\delta)}(B) &= \left(\sqrt{2 \log(1/\zeta_t(\delta)) + \log \left(\frac{\det(\mathbf{W}_t)}{\det(\mathbf{W}_1)} \right)} + \sqrt{\gamma} B \right)^2 \\ &\leq 2 \max\{2 \log T, \log(1/\delta)\} + 2 \log \left(\frac{\det \mathbf{W}_T}{\det \mathbf{W}_1} \right) + 2\gamma B^2 \\ &\leq 2 \max\{2 \log T, \log(1/\delta)\} + 2d \log(T-1) + 2d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + 2\gamma B^2 \\ &< (2d+4) \log T + 2 \log(1/\delta) + 2d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + 2\gamma B^2. \end{aligned} \quad (27)$$

428 So from (24), (26), and (27) we have

$$\mathbb{E} \left[\sum_{t=1}^T \frac{\widehat{\Delta}_{t, \zeta_t(\delta)}^2 (A_t^{\text{BAM}})}{\alpha I_t^B (A_t^{\text{BAM}}) + (1-\alpha) I_t^{\text{EB-UCB}} (A_t^{\text{BAM}})} \right] \leq \frac{16}{1-\alpha} \max\{U^2 \gamma^{-1}, \rho_{\max}^2\} T \\ \times \left[(2d+4) \log T + 2 \log(1/\delta) + 2d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + 2\gamma B^2 \right]. \quad (28)$$

429 Also, for any sequence $\{a_t\}_{t=1}^T \subset \mathcal{A}$ we have

$$I_t^B(a_t) = \frac{1}{2} \log \left(\|\mathbf{v}_t^{\min}\|^2_{(\mathbf{W}_t + \rho(a_t)^{-2} \phi(a_t) \phi(a_t)^\top)} \right) - \frac{1}{2} \log(\lambda_{\min}(\mathbf{W}_t)) \\ = \frac{1}{2} \log \left(\frac{(\mathbf{v}_t^{\min})^\top (\mathbf{W}_t + \rho(a_t)^{-2} \phi(a_t) \phi(a_t)^\top) \mathbf{v}_t^{\min}}{\lambda_{\min}(\mathbf{W}_t)} \right) \\ = \frac{1}{2} \log \left(\frac{\lambda_{\min}(\mathbf{W}_t) + \rho(a_t)^{-2} \mathbf{v}_t^{\min} \phi(a_t) \phi(a_t)^\top \mathbf{v}_t^{\min}}{\lambda_{\min}(\mathbf{W}_t)} \right) \\ = \frac{1}{2} \log \left(1 + \frac{\rho(a_t)^{-2} \psi_{\mathbf{v}_t^{\min}}(\phi(a_t))^2}{\lambda_{\min}(\mathbf{W}_t)} \right) = \frac{1}{2} \log \left(1 + \frac{\omega_t(a_t)}{\lambda_{\min}(\mathbf{W}_t)} \right). \quad (29)$$

Let

$$T_0 = \max \left\{ t : \sum_{\tau=1}^t \omega_\tau(a_\tau) \leq d(\rho_{\min}^{-2} U^2 - \gamma) \right\}.$$

430 Without loss of generality, assume that $T_0 \leq T$. Then using Lemma 2 we get

$$\sum_{t=1}^T I_t^B(a_t) = \sum_{t=1}^T \log \left(1 + \frac{\omega_t(a_t)}{\lambda_{\min}(\mathbf{W}_t)} \right) \\ = \sum_{t=1}^{T_0} \log \left(1 + \frac{\omega_t(a_t)}{\lambda_{\min}(\mathbf{W}_t)} \right) + \sum_{t=T_0+1}^T \log \left(1 + \frac{\omega_t(a_t)}{\lambda_{\min}(\mathbf{W}_t)} \right) \\ \leq \sum_{t=1}^{T_0} \frac{\omega_t(a_t)}{\lambda_{\min}(\mathbf{W}_t)} + \sum_{t=T_0+1}^T \log \left(1 + \frac{\omega_t(a_t)}{\gamma - \rho_{\min}^{-2} U^2 + \frac{1}{d} \sum_{\tau=1}^{t-1} \omega_\tau(a_\tau)} \right) \\ \leq \frac{1}{\gamma} \sum_{t=1}^{T_0} \omega_t(a_t) + \sum_{t=T_0+1}^T \log \left(1 + \frac{d\omega_t(a_t)}{d(\gamma - \rho_{\min}^{-2} U^2) + \sum_{\tau=1}^{t-1} \omega_\tau(a_\tau)} \right) \\ \leq \frac{d(\rho_{\min}^{-2} U^2 - \gamma)}{\gamma} + \sum_{t=T_0+1}^T \frac{d\omega_t(a_t)}{d(\gamma - \rho_{\min}^{-2} U^2) + \sum_{\tau=1}^{T_0} \omega_\tau(a_\tau) + \sum_{\tau=T_0+1}^{t-1} \omega_\tau(a_\tau)}.$$

Let

$$c = d(\gamma - \rho_{\min}^{-2} U^2) + \sum_{\tau=1}^{T_0} \omega_\tau(a_\tau)$$

and

$$x_t = \omega_{T_0+t}(a_{T_0+t}).$$

431 Then from Lemma 3, since $c > 0$ and $x_t \in [0, \rho_{\min}^{-2} U^2]$ for all t we have

$$\sum_{t=1}^T I_t^B(a_t) \leq \frac{d(\rho_{\min}^{-2} U^2 - \gamma)}{\gamma} + \sum_{t=T_0+1}^T \frac{d\omega_t(a_t)}{c + \sum_{\tau=T_0+1}^{t-1} \omega_\tau(a_\tau)} \\ = \frac{d(\rho_{\min}^{-2} U^2 - \gamma)}{\gamma} + d \sum_{t=1}^{T-T_0} \frac{x_t}{c + \sum_{\tau=1}^{t-1} x_\tau} \\ \leq O(d \log(T - T_0)) \leq O(d \log T). \quad (30)$$

432 Thus, from (21) we have

$$\alpha \mathbb{E} \left[\sum_{t=1}^T I_t^B (A_t^{\text{BAM}}) \right] + (1 - \alpha) \mathbb{E} \left[\sum_{t=1}^T I_t^{\text{EB-UCB}} (A_t^{\text{BAM}}) \right] \leq O(d \log T).$$

433 So from Theorem 1 and (28) we have

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \widehat{\Delta}_{t, \zeta_t(\delta)} (A_t^{\text{BAM}}) \right] &\leq O \left(\sqrt{d \frac{16}{1 - \alpha} \max\{U^2 \gamma^{-1}, \rho_{\max}^2\} T \log T} \right. \\ &\quad \times \sqrt{(4 + 2d) \log T + 2 \log(1/\delta) + 2d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + 2\gamma B^2} \left. \right) \\ &\leq O \left(\frac{d \max\{U/\sqrt{\gamma}, \rho_{\max}\}}{\sqrt{1 - \alpha}} \sqrt{T \log T} \right. \\ &\quad \times \sqrt{\log(1/\delta) + \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + \gamma B^2} \left. \right). \end{aligned} \quad (31)$$

Take any $t \geq 1$ and suppose that the event $E_{t, \zeta_t(\delta)}$, as defined in (12), holds. Note that the set

$$\left\{ \boldsymbol{\theta} \in \mathbb{R} : \left\| \boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_{\mathbf{W}_t}^2 \leq \beta_{t, \zeta_t(\delta)}(B^*) \right\}$$

434 is an ellipsoid in \mathbb{R}^d centered at $\widehat{\boldsymbol{\theta}}_t^{\text{wls}}$ with the longest semi-axis of length $\beta_{t, \zeta_t(\delta)}^{1/2}(B^*) \lambda_{\min}(\mathbf{W}_t)^{-1/2}$,
 435 so

$$\left\| \widehat{\boldsymbol{\theta}}_t^{\text{wls}} - \boldsymbol{\theta}^* \right\|_2 \leq \beta_{t, \zeta_t(\delta)}^{1/2}(B^*) \lambda_{\min}(\mathbf{W}_t)^{-1/2}. \quad (32)$$

436 Since $B \geq B^*$ we have $\beta_{t, \zeta_t(\delta)}(B) \geq \beta_{t, \zeta_t(\delta)}(B^*)$, so by the triangle inequality we get

$$B^* = \|\boldsymbol{\theta}^*\|_2 \leq \left\| \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2}(B^*) \lambda_{\min}(\mathbf{W}_t)^{-1/2} \leq \left\| \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2}(B) \lambda_{\min}(\mathbf{W}_t)^{-1/2} = \widehat{B}_t. \quad (33)$$

437 So $B^* \leq \widehat{B}_t$ for all $t \geq 1$ and thus $\beta_{t, \zeta_t(\delta)}(\widehat{B}_t) \geq \beta_{t, \zeta_t(\delta)}(B^*)$, so

$$\boldsymbol{\theta}^* \in \left\{ \boldsymbol{\theta} \in \mathbb{R} : \left\| \boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_{\mathbf{W}_t}^2 \leq \beta_{t, \zeta_t(\delta)}(B^*) \right\} \subseteq \left\{ \boldsymbol{\theta} \in \mathbb{R} : \left\| \boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_{\mathbf{W}_t}^2 \leq \beta_{t, \zeta_t(\delta)}(\widehat{B}_t) \right\}$$

438 Hence $\Delta(a) \leq \widehat{\Delta}_{t, \zeta_t(\delta)}(a)$ for all $a \in \mathcal{A}$. So for any $a \in \mathcal{A}$ we have

$$\mathbb{P} \left(\Delta(a) > \widehat{\Delta}_{t, \zeta_t(\delta)}(a) \right) \leq 1 - \mathbb{P}(E_{t, \zeta_t(\delta)}) \leq \zeta_t(\delta) \leq 1/t^2.$$

439 Thus, letting $\Delta_{\max} = \max_{a \in \mathcal{A}} \Delta(a)$, for any sequence $\{a_t\}_{t=1}^T \subset \mathcal{A}$ we have

$$\mathbb{E} \left[\sum_{t=1}^T \Delta(a_t) - \widehat{\Delta}_{t, \zeta_t(\delta)}(a_t) \right] \leq \Delta_{\max} \sum_{t=1}^T \mathbb{P} \left(\Delta(a_t) > \widehat{\Delta}_{t, 1/t^2}(a_t) \right) \leq \Delta_{\max} \sum_{t=1}^T \frac{1}{t^2} \leq O(\Delta_{\max}). \quad (34)$$

440 So from (31), for $T \leq T_B$ the regret of EBIDS is bounded above by

$$\mathcal{R}_T \leq O \left(\frac{d \max\{U/\sqrt{\gamma}, \rho_{\max}\}}{\sqrt{1 - \alpha}} \sqrt{T \log T} \sqrt{\log(1/\delta) + \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + \gamma B^2} \right).$$

441 From (21) and (30) with probability 1 we have

$$\sum_{t=1}^T I_t^{\text{BAM}}(A_t^{\text{BAM}}) \leq O(d \log T). \quad (35)$$

442 Following the same steps as in (24), using (25) and (28) we have

$$\begin{aligned} \sum_{t=1}^T \Psi_t^{\text{BAM}}(A_t^{\text{BAM}}) &= \sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2(A_t^{\text{BAM}})}{I_t^{\text{BAM}}(A_t^{\text{BAM}})} \leq \sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2(A_t^{\text{BAM}})}{\alpha I_t^B(A_t^{\text{BAM}}) + (1-\alpha)I_t^{\text{EB-UCB}}(A_t^{\text{BAM}})} \\ &\leq \sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2(A_t^{\text{EB-UCB}})}{\alpha I_t^{\text{EB-UCB}}(A_t^{\text{EB-UCB}}) + (1-\alpha)I_t^B(A_t^{\text{EB-UCB}})} \\ &\leq \frac{1}{1-\alpha} \sum_{t=1}^T \frac{\widehat{\Delta}_{t,\zeta_t(\delta)}^2(A_t^{\text{EB-UCB}})}{I_t^{\text{EB-UCB}}(A_t^{\text{EB-UCB}})} \\ &\leq \frac{16}{1-\alpha} \sum_{t=1}^T \beta_{T,\zeta_T(\delta)}(B) \max\{U^2\gamma^{-1}, \rho_{\max}^2\} \\ &\leq \frac{16}{1-\alpha} \max\{U^2\gamma^{-1}, \rho_{\max}^2\} T \\ &\quad \times \left[(4+2d) \log T + 2 \log(1/\delta) + 2d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + 2\gamma B^2 \right]. \quad (36) \end{aligned}$$

443 Following analogous steps as above, since $1/\zeta_t(\delta) \geq 1/\delta$ we have $\beta_{t,\zeta_t(\delta)}(B) \geq \beta_{t,\delta}(B) \geq$
 444 $\beta_{t,\delta}(B^*)$. So for any $t \geq 1$, whenever event $E_{t,\delta}$ holds, the inequality $B^* \leq \widehat{B}_t$ holds as well and
 445 thus $\Delta(a) \leq \widehat{\Delta}_{t,\zeta_t(\delta)}(a)$, for all $a \in \mathcal{A}$. So if $E_\delta = \bigcap_{t=1}^\infty E_{t,\delta}$ holds, then $\Delta(a) \leq \widehat{\Delta}_{t,\zeta_t(\delta)}(a)$, for
 446 all $a \in \mathcal{A}$ and for all $t \geq 1$. So from (36) and (35), by Theorem 1 we have

$$\mathcal{PR}_T \leq O \left(\frac{d \max\{U/\sqrt{\gamma}, \rho_{\max}\}}{\sqrt{1-\alpha}} \sqrt{T} \log T \sqrt{\log(1/\delta) + \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + \gamma B^2} \right). \quad \square$$

447 10.6 Proof of Proposition 2

448 In order to precisely state the conditions on T_B and α , i.e., how large each of them needs to be for
 449 $B^* \leq \widehat{B}_t \leq (1+g)B^*$ to hold for all $t \geq T_B + 1$, we will first define several constants for notational
 450 convenience.

451 Let

$$c_0 = L^2 \left[U^2(\gamma + \rho_{\min}^{-2} U^2) \left(\frac{1}{\kappa} + \frac{1}{\gamma} \right) \right]^{-1} \quad (37)$$

452 and

$$h_0 = 8 \log(5/4) + 4 \log(1/\delta) + 2d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + 2\gamma B^2. \quad (38)$$

453 Then let

$$u_0 = \frac{c_0}{6 + 16g^{-2}} \log 2 + \frac{1-\alpha}{\alpha} d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) \quad (39)$$

$$u_1 = \frac{c_0}{12 + 32g^{-2}} - \frac{1-\alpha}{2\alpha} d \quad (40)$$

$$w_0 = \frac{c_0}{6 + 16g^{-2}} + \frac{1-\alpha}{\alpha} d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) \quad (41)$$

$$w_1 = \frac{c_0}{12 + 32g^{-2}} - \frac{1-\alpha}{\alpha} d. \quad (42)$$

454 and finally let

$$b_0 = \frac{1}{d} \left[w_0 \left(\frac{\gamma}{d} u_0 - \gamma + \rho_{\min}^{-2} U^2 \right) - \gamma u_0 \right] + \gamma - \rho_{\min}^{-2} U^2 \quad (43)$$

$$b_1 = \frac{1}{d} \left(\gamma u_1 - \frac{\gamma}{d} u_1 w_0 - \frac{\gamma}{d} u_0 w_1 + \gamma w_1 - \rho_{\min}^{-2} U^2 w_1 \right) \quad (44)$$

$$b_2 = \frac{\gamma}{d^2} u_1 w_1 \quad (45)$$

455 We make the following assumptions.

456 **Assumption 1.** $B \geq B^*$.

Assumption 2.

$$T_B \geq \max \left\{ 4, \exp \left[\frac{h_0 + 2d + 8}{b_2} \left(4g^{-2} B^{*-2} + \frac{|b_1|}{2d + 8} + \frac{|b_0|}{h_0 + 2d + 8} \right) \right] \right\}.$$

Assumption 3.

$$\alpha \geq \frac{d}{d + \frac{c_0}{12 + 32g^{-2}}}.$$

457 We will now show that if Assumptions 1-3 are satisfied and event E_δ holds then $B^* \leq \tilde{B}_t \leq (1+g)B^*$
 458 for all $t \geq T_B + 1$.

459 *Proof.* Suppose that event E_δ holds. For any t let

$$s(t) = \arg \min_{\tau \leq t} \beta_{\tau, \zeta_\tau(\delta)}^{1/2} (\hat{B}_\tau) \lambda_{\min}(\mathbf{W}_\tau)^{-1/2} \quad (46)$$

460 From (32) in the proof of Proposition 1, using the triangle inequality we get

$$\left\| \hat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_2 \leq \left\| \boldsymbol{\theta}^* \right\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2} (B^*) \lambda_{\min}(\mathbf{W}_t)^{-1/2} = B^* + \beta_{t, \zeta_t(\delta)}^{1/2} (B^*) \lambda_{\min}(\mathbf{W}_t)^{-1/2}. \quad (47)$$

461 From (33) in the proof of Proposition 1, for any t we have $\hat{B}_t \geq B^*$, so

$$\left\| \hat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_2 \leq B^* + \beta_{t, \zeta_t(\delta)}^{1/2} (\hat{B}_t) \lambda_{\min}(\mathbf{W}_t)^{-1/2}.$$

462 Hence

$$\begin{aligned} \tilde{B}_t &= \min_{\tau \leq t} \left\{ \left\| \hat{\boldsymbol{\theta}}_\tau^{\text{wls}} \right\|_2 + \beta_{\tau, \zeta_\tau(\delta)}^{1/2} (\hat{B}_\tau) \lambda_{\min}(\mathbf{W}_\tau)^{-1/2} \right\} \\ &\leq \left\| \hat{\boldsymbol{\theta}}_{s(t)}^{\text{wls}} \right\|_2 + \beta_{s(t), \zeta_{s(t)}(\delta)}^{1/2} (\hat{B}_{s(t)}) \lambda_{\min}(\mathbf{W}_{s(t)})^{-1/2} \\ &\leq B^* + 2\beta_{s(t), \zeta_{s(t)}(\delta)}^{1/2} (\hat{B}_{s(t)}) \lambda_{\min}(\mathbf{W}_{s(t)})^{-1/2}. \end{aligned} \quad (48)$$

463 Also, analogously as in (33), using (32) and the triangle inequality, for any $t \geq 1$ we have

$$B^* = \left\| \boldsymbol{\theta}^* \right\|_2 \leq \left\| \hat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2} (B^*) \lambda_{\min}(\mathbf{W}_t)^{-1/2} \leq \left\| \hat{\boldsymbol{\theta}}_t^{\text{wls}} \right\|_2 + \beta_{t, \zeta_t(\delta)}^{1/2} (\hat{B}_t) \lambda_{\min}(\mathbf{W}_t)^{-1/2}.$$

464 So

$$B^* \leq \tilde{B}_t \quad (49)$$

465 for any $t \geq 1$.

466 From Lemma 1, for any $t \leq T_B$ we have

$$I_t^B(a_t^{\text{BAM}}) \geq \frac{\hat{\Delta}_{t, \zeta_t(\delta)}^2(a_t^{\text{BAM}})}{\hat{\Delta}_{t, \zeta_t(\delta)}^2(a_t^{I, B})} I_t^B(a_t^{I, B}) - \frac{1 - \alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \quad (50)$$

467 where $a_t^{I,B} = \arg \max_{a \in \mathcal{A}} I_t^B(a)$. For any $t \leq T_B$ we have

$$\begin{aligned} \widehat{\Delta}_{t,\zeta_t(\delta)}^2(a_t^{\text{BAM}}) &= \max_{b \in \mathcal{A}} \left\{ \langle \phi(b), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle + \beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \|\phi(b)\|_{\mathbf{W}_t^{-1}} \right\} \\ &\quad - \left(\langle \phi(a_t^{\text{BAM}}), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle - \beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \|\phi(a_t^{\text{BAM}})\|_{\mathbf{W}_t^{-1}} \right) \\ &= \max_{b \in \mathcal{A}} \left\{ \langle \phi(b), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle + \beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \|\phi(b)\|_{\mathbf{W}_t^{-1}} \right\} \\ &\quad - \left(\langle \phi(a_t^{\text{BAM}}), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle + \beta_{t,\zeta_t(\delta)}(\widehat{B}_t)^{1/2} \|\phi(a_t^{\text{BAM}})\|_{\mathbf{W}_t^{-1}} \right) \\ &\quad + 2\beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \|\phi(a_t^{\text{BAM}})\|_{\mathbf{W}_t^{-1}} \\ &\geq 2\beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \|\phi(a_t^{\text{BAM}})\|_{\mathbf{W}_t^{-1}}. \end{aligned}$$

468 So from (18)

$$\widehat{\Delta}_{t,\zeta_t(\delta)}^2(a_t^{\text{BAM}}) \geq 4\beta_{t,\zeta_t(\delta)}(\widehat{B}_t) \|\phi(a_t^{\text{BAM}})\|_{\mathbf{W}_t^{-1}}^2 \geq 4\beta_{t,\zeta_t(\delta)}(\widehat{B}_t) \frac{L^2}{t(\gamma + \rho_{\min}^{-2}U^2)} \quad (51)$$

469 Also

$$\begin{aligned} \widehat{\Delta}_{t,\zeta_t(\delta)}^2(a_t^{I,B}) &= \beta_{t,\zeta_t(\delta)}^{1/2}(\widehat{B}_t) \left(\|\phi(a_t^{\text{EB-UCB}})\|_{\mathbf{W}_t^{-1}} + \|\phi(a_t^{I,B})\|_{\mathbf{W}_t^{-1}} \right) \\ &\quad + \langle \phi(a_t^{\text{EB-UCB}}), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle - \langle \phi(a_t^{I,B}), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle, \end{aligned}$$

470 so

$$\begin{aligned} \widehat{\Delta}_{t,\zeta_t(\delta)}^2(a_t^{I,B}) &\leq 4\beta_{t,\zeta_t(\delta)}(\widehat{B}_t) \left(\|\phi(a_t^{\text{EB-UCB}})\|_{\mathbf{W}_t^{-1}}^2 + \|\phi(a_t^{I,B})\|_{\mathbf{W}_t^{-1}}^2 \right) \\ &\quad + 4 \langle \phi(a_t^{\text{EB-UCB}}), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle^2 + 4 \langle \phi(a_t^{I,B}), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle^2. \end{aligned}$$

471 Since E_δ holds, from (16) and (47) for any t and any $a \in \mathcal{A}$ we have

$$\langle \phi(a), \widehat{\boldsymbol{\theta}}_t^{\text{wls}} \rangle^2 \leq 2U^2(B^{*2} + \beta_{t,\zeta_t(\delta)}(B^*)\lambda_{\min}(\mathbf{W}_t)^{-1})$$

472 so from (15) we have

$$\widehat{\Delta}_{t,\zeta_t(\delta)}^2(a_t^{I,B}) \leq 8\beta_{t,\zeta_t(\delta)}(\widehat{B}_t)U^2\lambda_{\min}(\mathbf{W}_t)^{-1} + 16U^2(B^{*2} + \beta_{t,\zeta_t(\delta)}(B^*)\lambda_{\min}(\mathbf{W}_t)^{-1}). \quad (52)$$

473 From (29) from the proof of Proposition 1, for any $a \in \mathcal{A}$ we have

$$I_t^B(a) = \frac{1}{2} \log \left(1 + \frac{\rho(a)^{-2} \psi_{\mathbf{v}_t^{\min}}(\phi(a))^2}{\lambda_{\min}(\mathbf{W}_t)} \right), \quad (53)$$

474 so

$$\begin{aligned} I_t^B(a_t^{I,B}) &= \max_{a \in \mathcal{A}} I_t^B(a) = \max_{a \in \mathcal{A}} \left\{ \frac{1}{2} \log \left(1 + \frac{\rho(a)^{-2} \psi_{\mathbf{v}_t^{\min}}(\phi(a))^2}{\lambda_{\min}(\mathbf{W}_t)} \right) \right\} \\ &\geq \frac{1}{2} \log \left(1 + \frac{\kappa}{\lambda_{\min}(\mathbf{W}_t)} \right). \end{aligned}$$

475 Thus, since $\log x \geq 1 - \frac{1}{x}$ for all $x > 0$, we have

$$\begin{aligned} I_t^B(a_t^{I,B}) &\geq \frac{\kappa}{2(\lambda_{\min}(\mathbf{W}_t) + \kappa)} = \left[2\lambda_{\min}(\mathbf{W}_t) \left(\frac{1}{\kappa} + \frac{1}{\lambda_{\min}(\mathbf{W}_t)} \right) \right]^{-1} \\ &\geq \left[2\lambda_{\min}(\mathbf{W}_t) \left(\frac{1}{\kappa} + \frac{1}{\gamma} \right) \right]^{-1}. \end{aligned} \quad (54)$$

476 So combining (50), (51), (52), and (54), for any $t \leq T_B$ we have

$$\begin{aligned}
 I_t^B(a_t^{\text{BAM}}) &\geq \frac{L^2 \lambda_{\min}(\mathbf{W}_t)^{-1} \left[2t(\gamma + \rho_{\min}^{-2} U^2) \left(\frac{1}{\kappa} + \frac{1}{\gamma} \right) \right]^{-1}}{2U^2 \lambda_{\min}(\mathbf{W}_t)^{-1} + 4U^2 \beta_{t, \zeta_t(\delta)}(\widehat{B}_t)^{-1} (B^{*2} + \beta_{t, \delta}(B^*) \lambda_{\min}(\mathbf{W}_t)^{-1})} \\
 &\quad - \frac{1 - \alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) = \\
 &= \frac{1}{t} L^2 \left[U^2(\gamma + \rho_{\min}^{-2} U^2) \left(\frac{1}{\kappa} + \frac{1}{\gamma} \right) \left(4 + 8B^{*2} \frac{\lambda_{\min}(\mathbf{W}_t)}{\beta_{t, \zeta_t(\delta)}(\widehat{B}_t)} + 8 \frac{\beta_{t, \delta}(B^*)}{\beta_{t, \zeta_t(\delta)}(\widehat{B}_t)} \right) \right]^{-1} \\
 &\quad - \frac{1 - \alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \geq \\
 &\geq \frac{1}{t} L^2 \left[U^2(\gamma + \rho_{\min}^{-2} U^2) \left(\frac{1}{\kappa} + \frac{1}{\gamma} \right) \left(12 + 8B^{*2} \frac{\lambda_{\min}(\mathbf{W}_t)}{\beta_{t, \zeta_t(\delta)}(\widehat{B}_t)} \right) \right]^{-1} \\
 &\quad - \frac{1 - \alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}),
 \end{aligned}$$

477 where the last inequality follows from the fact that $\widehat{B}_t \geq B^*$ and $1/\zeta_t(\delta) \geq 1/\delta$ which gives us

$$\frac{\beta_{t, \delta}(B^*)}{\beta_{t, \zeta_t(\delta)}(\widehat{B}_t)} \leq 1.$$

478 So from (37) we have

$$I_t^B(a_t^{\text{BAM}}) \geq \frac{1}{t} c_0 \left(12 + 8B^{*2} \frac{\lambda_{\min}(\mathbf{W}_t)}{\beta_{t, \zeta_t(\delta)}(\widehat{B}_t)} \right)^{-1} - \frac{1 - \alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}). \quad (55)$$

479 From (53) we have

$$I_t^B(a_t^{\text{BAM}}) = \frac{1}{2} \log \left(1 + \frac{\omega_t(a_t^{\text{BAM}})}{\lambda_{\min}(\mathbf{W}_t)} \right) \leq \frac{\omega_t(a_t^{\text{BAM}})}{2\lambda_{\min}(\mathbf{W}_t)}.$$

480 So

$$\omega_t(a_t^{\text{BAM}}) \geq 2\lambda_{\min}(\mathbf{W}_t) I_t^B(a_t^{\text{BAM}}). \quad (56)$$

481 If

$$\beta_{t, \zeta_t(\delta)}^{1/2}(\widehat{B}_t) \lambda_{\min}(\mathbf{W}_t)^{-1/2} \leq \frac{1}{2} g B^* \quad (57)$$

482 for some $t \leq T_B + 1$ then

$$\beta_{s(t), \zeta_{s(t)}(\delta)}^{1/2}(\widehat{B}_{s(t)}) \lambda_{\min}(\mathbf{W}_{s(t)})^{-1/2} \leq \frac{1}{2} g B^*,$$

483 so from (48) and (49), since event E_δ holds, for any $t \geq T_B + 1$ we have

$$B^* \leq \tilde{B}_t \leq B^* + 2\beta_{s(t), \zeta_{s(t)}(\delta)}^{1/2}(\widehat{B}_{s(t)}) \lambda_{\min}(\mathbf{W}_{s(t)})^{-1/2} = (1 + g) B^* \quad (58)$$

484 which is what we want to show. We will prove by contradiction that since E_δ holds, (57) holds as

485 well for some $t \leq T_B + 1$. Suppose that (57) does not hold. Then for all $t \leq T_B + 1$ we have

$$\frac{\lambda_{\min}(\mathbf{W}_t)}{\beta_{t, \zeta_t(\delta)}(\widehat{B}_t)} < 4g^{-2} B^{*-2}, \quad (59)$$

486 so from (55) we have

$$\begin{aligned} I_t^B(a_t^{\text{BAM}}) &\geq \frac{1}{t} c_0 \left(12 + 8B^{*2} \frac{\lambda_{\min}(\mathbf{W}_t)}{\beta_{t, \zeta_t(\delta)}(\hat{B}_t)} \right)^{-1} - \frac{1-\alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \\ &> \frac{1}{t} \cdot \frac{c_0}{12 + 32g^{-2}} - \frac{1-\alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}). \end{aligned}$$

487 Hence, from (56) for any $t \leq T_B$ we have

$$\omega_t(a_t^{\text{BAM}}) \geq \lambda_{\min}(\mathbf{W}_t) \left(\frac{1}{t} \cdot \frac{c_0}{6 + 16g^{-2}} - 2 \frac{1-\alpha}{\alpha} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \right).$$

488 Let $\lfloor x \rfloor$ denote the largest integer smaller than or equal to x for any $x \in \mathbb{R}$. From Weyl's inequality

489 (Franklin, 1968)

$$\lambda_{\min}(\mathbf{W}_{t+1}) \geq \lambda_{\min}(\mathbf{W}_t) \geq \gamma \quad (60)$$

490 for any t . Also note that $\omega_t(a) \geq 0$ for any t and any $a \in \mathcal{A}$. So

$$\sum_{t=1}^{\lfloor \sqrt{T_B} \rfloor} \omega_t(a_t^{\text{BAM}}) \geq \gamma \left(\frac{c_0}{6 + 16g^{-2}} \sum_{t=1}^{\lfloor \sqrt{T_B} \rfloor} \frac{1}{t} - 2 \frac{1-\alpha}{\alpha} \sum_{t=1}^{\lfloor \sqrt{T_B} \rfloor} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \right).$$

491 From (21) we have

$$\begin{aligned} \sum_{t=1}^{\lfloor \sqrt{T_B} \rfloor} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) &\leq \frac{1}{2} d \log \lfloor \sqrt{T_B} \rfloor + \frac{1}{2} d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) \\ &\leq \frac{1}{4} d \log T_B + \frac{1}{2} d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right). \end{aligned}$$

492 Also since $T_B \geq 4$ we have $\lfloor \sqrt{T_B} \rfloor \geq \sqrt{T_B} - 1 \geq \sqrt{T_B}/2$, so

$$\sum_{t=1}^{\lfloor \sqrt{T_B} \rfloor} \frac{1}{t} > \log \lfloor \sqrt{T_B} \rfloor \geq \log \left(\frac{1}{2} \sqrt{T_B} \right) = \frac{1}{2} \log T_B - \log 2.$$

493 So

$$\begin{aligned} \sum_{t=1}^{\lfloor \sqrt{T_B} \rfloor} \omega_t(a_t^{\text{BAM}}) &\geq \gamma \left(\frac{c_0}{6 + 16g^{-2}} \left[\frac{1}{2} \log T_B - \log 2 \right] - \frac{1-\alpha}{\alpha} d \left[\frac{1}{2} \log T_B + \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) \right] \right) \\ &\geq \gamma \left(\left[\frac{c_0}{12 + 32g^{-2}} - \frac{1-\alpha}{2\alpha} d \right] \log T_B - \left[\frac{c_0}{6 + 16g^{-2}} \log 2 + \frac{1-\alpha}{\alpha} d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) \right] \right) \\ &= \gamma(u_1 \log T_B - u_0), \end{aligned} \quad (61)$$

494 where the constants u_0 and u_1 were defined in (39) and (40), respectively. Similarly from (60) we

495 have

$$\sum_{t=\lfloor \sqrt{T_B} \rfloor+1}^{T_B} \omega_t(a_t^{\text{BAM}}) \geq \lambda_{\min}(\mathbf{W}_{\lfloor \sqrt{T_B} \rfloor+1}) \left(\frac{c_0}{6 + 16g^{-2}} \sum_{t=\lfloor \sqrt{T_B} \rfloor+1}^{T_B} \frac{1}{t} - 2 \frac{1-\alpha}{\alpha} \sum_{t=\lfloor \sqrt{T_B} \rfloor+1}^{T_B} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \right)$$

496 Note hat

$$\sum_{t=\lfloor \sqrt{T_B} \rfloor+1}^{T_B} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \leq \sum_{t=1}^{T_B} I_t^{\text{EB-UCB}}(a_t^{\text{BAM}}) \leq \frac{1}{2} d \log T_B + \frac{1}{2} d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right)$$

497 and

$$\sum_{t=\lfloor\sqrt{T_B}\rfloor+1}^{T_B} \frac{1}{t} = \sum_{t=1}^{T_B} \frac{1}{t} - \sum_{t=1}^{\lfloor\sqrt{T_B}\rfloor} \frac{1}{t} > \log T_B - (\log \sqrt{T_B} + 1) = \frac{1}{2} \log T_B - 1.$$

498 So

$$\begin{aligned} \sum_{t=\lfloor\sqrt{T_B}\rfloor+1}^{T_B} \omega_t(a_t^{\text{BAM}}) &\geq \lambda_{\min}(\mathbf{W}_{\lfloor\sqrt{T_B}\rfloor+1}) \\ &\quad \times \left(\frac{c_0}{6+16g^{-2}} \left[\frac{1}{2} \log T_B - 1 \right] - \frac{1-\alpha}{\alpha} d \left[\log T_B + \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) \right] \right) \\ &= \lambda_{\min}(\mathbf{W}_{\lfloor\sqrt{T_B}\rfloor+1}) \\ &\quad \times \left(\left[\frac{c_0}{12+32g^{-2}} - \frac{1-\alpha}{\alpha} d \right] \log T_B - \left[\frac{c_0}{6+16g^{-2}} + \frac{1-\alpha}{\alpha} d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) \right] \right) \\ &= \lambda_{\min}(\mathbf{W}_{\lfloor\sqrt{T_B}\rfloor+1}) (w_1 \log T_B + w_0), \end{aligned}$$

499 where the constants w_0 and w_1 were defined in (41) and (42), respectively.

500 From Lemma 2 and (61) we have

$$\lambda_{\min}(\mathbf{W}_{\lfloor\sqrt{T_B}\rfloor+1}) \geq \gamma - \rho_{\min}^{-2} U^2 + \frac{1}{d} \sum_{t=1}^{\lfloor\sqrt{T_B}\rfloor} \omega_t(a_t^{\text{BAM}}) \geq \frac{\gamma}{d} (u_1 \log T_B - u_0) + \gamma - \rho_{\min}^{-2} U^2.$$

501 So

$$\begin{aligned} \sum_{t=1}^{T_B} \omega_t(a_t^{\text{BAM}}) &= \sum_{t=1}^{\lfloor\sqrt{T_B}\rfloor} \omega_t(a_t^{\text{BAM}}) + \sum_{t=\lfloor\sqrt{T_B}\rfloor+1}^{T_B} \omega_t(a_t^{\text{BAM}}) \\ &\geq \gamma (u_1 \log T_B - u_0) + \left(\frac{\gamma}{d} (u_1 \log T_B - u_0) + \gamma - \rho_{\min}^{-2} U^2 \right) (w_1 \log T_B - w_0) \\ &= \frac{\gamma}{d} u_1 w_1 (\log T_B)^2 + \left(\gamma u_1 - \frac{\gamma}{d} u_1 w_0 - \frac{\gamma}{d} u_0 w_1 + \gamma w_1 - \rho_{\min}^{-2} U^2 w_1 \right) \log T_B \\ &\quad + w_0 \left(\frac{\gamma}{d} u_0 - \gamma + \rho_{\min}^{-2} U^2 \right) - \gamma u_0 \\ &= db_2 (\log T_B)^2 + db_1 \log T_B + d(b_0 - \gamma + \rho_{\min}^{-2} U^2), \end{aligned}$$

502 where the constants b_0 , b_1 and b_2 were defined in (43), (44), and (45), respectively.

503 Then, applying Lemma 2 again we get

$$\lambda_{\min}(\mathbf{W}_{T_B+1}) \geq \gamma - \rho_{\min}^{-2} U^2 + \frac{1}{d} \sum_{t=1}^{T_B} \omega_t(a_t^{\text{BAM}}) \geq b_2 (\log T_B)^2 + b_1 \log T_B + b_0. \quad (62)$$

504 Note that by Assumption 3, we have $u_1 > 0$ and $w_1 > 0$, so $b_2 > 0$.

505 From (19) we have

$$\begin{aligned} \beta_{T_B+1, \zeta_{T_B+1}(\delta)}(\widehat{B}_t) &= \left(\sqrt{2 \log(1/\zeta_{T_B+1}(\delta)) + \log \left(\frac{\det \mathbf{W}_{T_B+1}}{\det \mathbf{W}_1} \right)} + \sqrt{\gamma} \widehat{B}_{T_B+1} \right)^2 \leq \\ &\leq 4 \log(1/\zeta_{T_B+1}(\delta)) + 2 \log \left(\frac{\det \mathbf{W}_{T_B+1}}{\det \mathbf{W}_1} \right) + 2\gamma \widehat{B}_{T_B+1}^2 \leq \\ &\leq 4 \max\{\log(1/\delta), 2 \log(T_B + 1)\} + 2d \log T_B + 2d \log \left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma} \right) + 2\gamma B^2. \end{aligned}$$

Since $T_B \geq 4$ we have

$$\log(T_B + 1) \leq \log\left(\frac{5}{4}T_B\right) = \log T_B + \log(5/4),$$

506 so

$$\begin{aligned} \beta_{T_B+1, \zeta_{T_B+1}(\delta)}(\widehat{B}_t) &\leq (2d+8) \log T_B + 8 \log(5/4) + 4 \log(1/\delta) + 2d \log\left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma}\right) + 2\gamma B^2 \\ &= (2d+8) \log T_B + h_0, \end{aligned} \quad (63)$$

507 with h_0 defined in (38). Note that $h_0 > 0$. Also, since $T_B \geq 4$ we have $\log T_B > 1$ so from (62),
508 (63) and the fact that $b_2 > 0$ we get

$$\begin{aligned} \frac{\lambda_{\min}(\mathbf{W}_{T_B+1})}{\beta_{T_B+1, \zeta_{T_B+1}(\delta)}(\widehat{B}_t)} &\geq \frac{b_2(\log T_B)^2 + b_1 \log T_B + b_0}{(2d+8) \log T_B + h_0} \\ &= \frac{b_2}{2d+8 + \frac{h_0}{\log T_B}} \log T_B + \frac{b_1}{(2d+8) + \frac{h_0}{\log T_B}} + \frac{b_0}{(2d+8) \log T_B + h_0} \\ &\geq \frac{b_2}{h_0 + 2d+8} \log T_B - \frac{|b_1|}{2d+8} - \frac{|b_0|}{h_0 + 2d+8}. \end{aligned}$$

509 Note that by Assumption 2 we have

$$T_B \geq \exp\left[\frac{h_0 + 2d+8}{b_2} \left(4g^{-2}B^{*-2} + \frac{|b_1|}{2d+8} + \frac{|b_0|}{h_0 + 2d+8}\right)\right]$$

510 so

$$\frac{\lambda_{\min}(\mathbf{W}_{T_B+1})}{\beta_{T_B+1, \zeta_{T_B+1}(\delta)}(\widehat{B}_t)} \geq 4g^{-2}B^{*-2}$$

511 which is the required contradiction to (59). So there exists $t \leq T_B + 1$ such that

$$\frac{\lambda_{\min}(\mathbf{W}_t)}{\beta_{t, \zeta_t(\delta)}(\widehat{B}_t)} \geq 4g^{-2}B^{*-2}$$

512 and thus, since E_δ holds, from (58) for any $t \geq T_B + 1$ we have

$$B^* \leq \tilde{B}_t \leq (1+g)B^*. \quad \square$$

513 10.7 Proof of Proposition 3

514 The exact assumptions made by Propositions 3 are as follows. We assume that T_B and α are
515 sufficiently large so Assumptions 1 - 3 hold and $(T_B + 1)^2 \geq 1/\delta$. We can now proceed to the proof.

516 *Proof.* Suppose that event E_δ holds.

$$\mathbb{E} \left[\sum_{t=1}^T \widehat{\Delta}_{t, \zeta_t(\delta)}(A_t^{\text{BEIDS}}) \right] = \mathbb{E} \left[\sum_{t=1}^{T_B} \widehat{\Delta}_{t, \zeta_t(\delta)}(A_t^{\text{BAM}}) \right] + \mathbb{E} \left[\sum_{t=T_B+1}^T \widehat{\Delta}_{t, \zeta_t(\delta)}(A_t^{\text{EB-UCB}}) \right].$$

517 From (21) with probability 1 we have

$$\sum_{t=T_B+1}^T I_t^{\text{EB-UCB}}(A_t) \leq O(d \log T). \quad (64)$$

518 Let $a_t^{\text{EB-UCB}}$ be the realization of $A_t^{\text{EB-UCB}}$. Since event E_δ holds and Assumptions 1 - 3 hold, from
 519 Proposition 2 we have $B^* \leq \tilde{B}_t \leq (1+g)B^*$ for all $t \geq T_B + 1$. Also from the assumptions of
 520 this proposition, $2 \log T \geq \log(1/\delta)$, so analogously as in (25) and (27), for any $t \in \{T_B + 1, T_B +$
 521 $2, \dots, T\}$ we have

$$\begin{aligned} \frac{\widehat{\Delta}_{t, \zeta_t(\delta)}^2(a_t^{\text{EB-UCB}})}{I_t^{\text{EB-UCB}}(a_t^{\text{EB-UCB}})} &\leq 16\beta_{T, \zeta_T(\delta)}(\tilde{B}_t) \max\{U^2\gamma^{-1}, \rho_{\max}^2\} \\ &\leq 16 \max\{U^2\gamma^{-1}, \rho_{\max}^2\} \\ &\quad \times \left[2 \max\{2 \log T, \log(1/\delta)\} + 2d \log(T-1) + 2d \log\left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma}\right) + 2\gamma \tilde{B}_t^2 \right] \\ &\leq 16 \max\{U^2\gamma^{-1}, \rho_{\max}^2\} \\ &\quad \times \left[(2d+4) \log T + 2d \log\left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma}\right) + 2\gamma \tilde{B}_t^2 \right] \\ &\leq 16 \max\{U^2\gamma^{-1}, \rho_{\max}^2\} \\ &\quad \times \left[(2d+4) \log T + 2d \log\left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma}\right) + 2\gamma((1+g)B^*)^2 \right]. \end{aligned}$$

522 Hence from Theorem 1 and (64) we have

$$\begin{aligned} \mathbb{E} \left[\sum_{t=T_B+1}^T \widehat{\Delta}_{t, \zeta_t(\delta)}(A_t^{\text{EB-UCB}}) \right] &\leq O \left(d \max\{U/\sqrt{\gamma}, \rho_{\max}\} \sqrt{T} \log T \right. \\ &\quad \times \left. \sqrt{\log\left(1 + \frac{\rho_{\min}^{-2} U^2}{\gamma}\right) + \gamma((1+g)B^*)^2} \right) \\ &\leq O \left(dU\rho_{\max}(1+g)B^* \sqrt{T} \log T \right), \end{aligned}$$

523 and thus from (34) we get that

$$\mathbb{E} \left[\sum_{t=T_B+1}^T \Delta(A_t^{\text{EB-UCB}}) \right] \leq O \left(dU\rho_{\max}(1+g)B^* \sqrt{T} \log T \right).$$

524 and similarly with probability 1 we have

$$\sum_{t=T_B+1}^T \Delta(A_t^{\text{EB-UCB}}) \leq O \left(dU\rho_{\max}(1+g)B^* \sqrt{T} \log T \right).$$

525 Thus, since T_B is fixed with respect to T with probability at least $\mathbb{P}(E_\delta) \geq 1 - \delta$ we have

$$\mathcal{R}_T \leq O \left(dU\rho_{\max}(1+g)B^* \sqrt{T} \log T \right)$$

526 and

$$\mathcal{PR}_T \leq O \left(dU\rho_{\max}(1+g)B^* \sqrt{T} \log T \right). \quad \square$$