
Adapting to a Changing Environment: the Brownian Restless Bandits

Aleksandrs Slivkins* and Eli Upfal†

Abstract

In the multi-armed bandit (MAB) problem there are k distributions associated with the rewards of playing each of k strategies (slot machine arms). The reward distributions are initially unknown to the player. The player iteratively plays one strategy per round, observes the associated reward, and decides on the strategy for the next iteration. The goal is to maximize the reward by balancing *exploitation*: the use of acquired information, with *exploration*: learning new information.

We introduce and study a *dynamic* MAB problem in which the reward functions stochastically and gradually change in time. Specifically, the expected reward of each arm follows a Brownian motion, a discrete random walk, or similar processes. In this setting a player has to continuously keep exploring in order to adapt to the changing environment. Our formulation is (roughly) a special case of the notoriously intractable *restless MAB problem*.

Our goal here is to characterize the cost of learning and adapting to the changing environment, in terms of the stochastic rate of the change. We consider an infinite time horizon, and strive to minimize the average cost per step which we define with respect to a hypothetical algorithm that at every step plays the arm with the maximum expected reward at this step. A related line of work on the *adversarial MAB problem* used a significantly weaker benchmark, the best *time-invariant* policy.

The dynamic MAB problem models a variety of practical online, game-against-nature type optimization settings. While building on prior work, algorithms and steady-state analysis for the dynamic setting require a novel approach based on different stochastic tools.

1 Introduction

The multi-armed bandit (MAB) problem [27, 5, 12] has been studied extensively for over 50 years in Operations Research, Economics and Computer Science literature, modeling online decisions under uncertainty in a setting in which an agent simultaneously attempts to acquire new knowledge and to optimize its decisions based on the existing knowledge. In the basic MAB setting, which we term the *static MAB problem*, there are k time-invariant probability distributions associated with the rewards of playing each of the k strategies (slot machine arms). The distributions are initially unknown to the player. The player iteratively plays one strategy per round, observes the associated reward, and decides on the strategy for the next iteration. The goal of a MAB algorithm is to optimize the total reward¹ by balancing *exploitation*: the use of acquired information, with *exploration*: learning new information. For several algorithms in the literature (e.g. see [5, 2]) as the number of rounds goes to infinity the expected total reward asymptotically approaches that of playing a strategy with the highest expected reward. The quality of an algorithm for the static MAB problem is therefore measured by the expected cost, or *regret*, incurred during an initial finite time interval. The regret in the first t steps is defined as the expected gap between the total reward collected by the algorithm and that collected by playing an optimal strategy in these t steps.

The MAB problem models a variety of practical online optimization problems. As an example consider a packet routing network where a router learns about delays on routes by measuring the time to receive an acknowledgment for a packet sent on that route [4, 16]. The delay for one packet on a given route is a random value drawn from some distribution. The router must try various routes in order to learn about the delays. Trying a loaded route adds unnecessary delay to the routing of one packet, while discovering a route with low delay can improve the routing of the future packets.

Another application is in marketing and advertising. A store would like to display and advertise the products that sell best, but it needs to display and advertise various products to learn how good they sell. Similarly, a web search engine tries to optimize its revenue by displaying advertise-

*Microsoft Research, Mountain View CA. E-mail: slivkins at microsoft.com. Parts of this work has been completed while A. Slivkins was a postdoc at Brown University.

†Computer Science Department, Brown University, Providence RI. E-mail: eli at cs.brown.edu. Supported in part by NSF awards CCR-0121154 and DMI-0600384, and ONR Award N000140610607.

¹In this paper the *total reward* is simply the sum of the rewards, following the line of work in [21, 2, 3] and many other papers. Alternatively, many papers consider the *time-discounted* sum of rewards, e.g. see [5, 12, 29] and references therein.

ments that would bring the largest number of clicks for a given web content. The company needs to experiment with various combinations of advertisements and page contents in order to find the best matches. The cost of these experiments is the loss of advertisement clicks when trying unsuccessful matches [25].

The above examples demonstrate the practical applications of the "explore and exploit" paradigm captured in the MAB model. These examples also point out the limitation of the static approach to the problem. The delay on a route is gradually changing over time, and the router needs to continuously adapt its routing strategy to the changes in route delays. Taste and fashion change over time. A store cannot completely rely on information collected in the previous season to optimize for the next one. Similarly, a web search engine continually updates their content matching strategies to account for the changing customers' response.

A number of models have been proposed for capturing the dynamic aspect of the MAB problem. Motivated by task scheduling, Gittins [13] considered the case where only the state of the active arm (the arm currently being played) can change in a given step, giving an optimal policy for the Bayesian formulation with time discounting. This seminal result gave rise to a rich line of work (e.g. [11, 12, 32, 31, 30, 6, 29]), a proper review of which is beyond the scope of this paper. In particular, Whittle [33] introduced an extension termed *restless bandits* [33, 7, 24], where the states of all arms can change in each step according to a known (but arbitrary) stochastic transition function. Restless bandits are notoriously intractable: e.g. even with deterministic transitions the problem of computing an (approximately) optimal strategy is PSPACE-hard [26]. Guha et al. [14, 15] have recently made a progress on some tractable special cases of the restless MAB problem.² Their motivations, the actual problems they considered, and the techniques they used, are very different from ours. In [14] they gave a constant-factor approximation for the special case of the problem in which arms move stochastically between two possible states. This result was improved to a 2-approximation in [15], and extended to arms assuming a number of possible states, but with a very strict set of transition probabilities that are not compatible with the stochastic processes discussed here.

Auer et al. [3] adopted an adversarial approach: they defined the *adversarial MAB problem* where the reward distributions are allowed to change arbitrarily in time, and the goal is to approach the performance of the best *time-invariant* policy. This formulation has been further studied in [1, 20, 17, 22, 10, 9, 19, 8]. Auer et al. [3, 1] also considered a more general definition of regret, where the comparison is to the best policy that can change arms a limited number of times. Due to the overwhelming strength of the adversary, the guarantees obtained in this line of work are relatively weak when applied to the setting that we consider in this paper.

We propose and study here a somewhat different approach to addressing the dynamic nature of the MAB problem. We note that in a variety of practical applications the time evolution of the system, in particular of the reward functions, is *gradual*. Obvious examples are price, supply and demand

²These papers were published after the initial technical report version of this paper appeared.

in economics, load and delay in networks, etc. A gradual stochastic evolution is traditionally modeled via a random walk or a Brownian motion; for instance, in Mathematical Finance the (geometric) Brownian motion (Wiener process) is the standard model for continuous-time evolution of a stock price. In line with this approach, we describe the *state* of each arm – its expected reward at time t – via a Brownian motion.³ The actual reward at a given time is an independent random sample from the reward distribution parameterized by the current state of this arm, e.g. a 0-1 random variable with an expectation given by the state of the arm (in the web advertising setting this corresponds to a user clicking or not clicking on an ad).

We are interested in systems that exhibit a stationary, steady-state behavior. For this reason instead of the usual Brownian motion on a real line (which diverges to infinity) we consider a Brownian motion on an interval with reflecting bounds. Following the bulk of the stochastic MAB literature, we assume that the evolution of each arm is independent (in fact, we conjecture that regret is maximized in the case of independently evolving arms).

Our goal here is to characterize the long-term average cost of adapting to such changing environment in terms of the stochastic rate of change – the *volatility* of Brownian motion. The paradigmatic setting for us is one in which each arm's state has the same stationary distribution and, therefore, all arms are essentially equivalent in the long term. In such setting the standard benchmark – the *best* time-invariant policy – is uninformative. Instead, we optimize with respect to a more demanding (and also more natural) benchmark – a policy that at each step plays an arm with the currently maximal expected reward.

We consider two versions of the *dynamic MAB problem* described above. In the *state-informed* version an algorithm not only receives a reward of the chosen arm but also finds out the current state of this arm. This is the setting in the restless MAB problem as defined in Whittle [33] and the follow-up literature. In the second, *state-oblivious*, version an algorithm receives its reward and no other information. This formulation generalizes the static MAB problem to stochastically changing expected rewards.

1.1 The Dynamic MAB problem

Let $\{\mathcal{D}(\mu) : \mu \in [0; 1]\}$ be a fixed family of probability distributions on $[0; 1]$ such that $\mathcal{D}(\mu)$ has expectation μ . Time proceeds in rounds. Each arm i at each round t has a *state* $\mu_i(t) \in [0; 1]$ such that the reward from playing arm i in round t is an independent random sample from $\mathcal{D}(\mu_i(t))$. At each round t an algorithm chooses one of the k alternative strategies ("arms") and receives a reward. In the *state-oblivious* version, the reward is the only information that the algorithm receives in a given round. In the *state-informed* version, the algorithm also finds out the current state of the arm that it has chosen. The distributions $\mathcal{D}(\cdot)$ are not revealed to the algorithm (and are not essential to the analysis).

³As we only sample arms at integer time points, we can equivalently describe the state as a sum of t i.i.d. normal increments. In fact, we allow the increments to come from a somewhat more general class of distributions.

The state $\mu_i(\cdot)$ varies in an interval with reflecting boundaries. To clarify the concept of reflecting boundaries, consider an object that starts moving on an interval $I = [0; 1]$, reversing direction every time it hits a boundary. If the object starts at 0 and traverses distance $x \geq 0$, its position is

$$f_I(x) = \begin{cases} x', & x' \leq 1 \\ 1 - (x' - 1), & x' > 1, \end{cases} \quad (1)$$

where $x' = x \pmod{2} = x - 2 \lfloor x/2 \rfloor$. Similarly, we define $f_I(x)$, $x < 0$ as the position of an object that starts moving from 1 and traverses distance $|x|$.

For concreteness we focus here on the case when each arm's state follows a Brownian motion. Similar results hold for related stochastic processes such as discrete random walks (see the Extensions Section).

The state of each arm i undergoes an independent Brownian motion on an interval with reflecting boundaries. Specifically, we define $\mu_i(t) = f_I(B_i(t))$ where $I = [0; 1]$ is the *fundamental interval* and B_i is an independent Brownian motion with volatility σ_i . Since we only sample $\mu_i(\cdot)$ at integer times, we can also define it as a Markov chain:

$$\mu_i(t) = f_I(\mu_i(t-1) + X_i(t)), \quad (2)$$

where each $X_i(t)$ is an i.i.d. sample from $\mathcal{N}(0, \sigma_i)$. The stochastic rate of change is thus given by σ_i , which we term the *volatility* of arm i .

We assume that for each arm i the initial state $\mu_i(0)$ is an independent uniformly random sample from I . This is a reasonable assumption given our goal to study the stationary behavior of the system. Indeed, the uniform distribution on I is the stationary distribution of the Markov chain 2 to which this Markov chain eventually converges.⁴

In the dynamic MAB problem, we measure the performance of a MAB algorithm with respect to a policy that at every step chooses a strategy with the highest expected reward. This policy changes in time, and thus it is a more demanding benchmark than the *time-invariant regret* that is often used in the MAB literature.

Definition 1.1. Consider an instance of the dynamic MAB problem. For a given MAB algorithm \mathcal{A} , let $W_{\mathcal{A}}(t)$ be the reward received by algorithm \mathcal{A} in round t . Let \emptyset be an algorithm that in every round chooses a strategy with the highest expected reward. The *dynamic regret* in round t is

$$R_{\mathcal{A}}(t) = W_{\emptyset}(t) - W_{\mathcal{A}}(t).$$

Define the *steady-state regret* as

$$\bar{R}_{\mathcal{A}} = \limsup_t \sup_{t_0} E \left[\frac{1}{t} \sum_{s=t_0+1}^{t_0+t} R_{\mathcal{A}}(s) \right]. \quad (3)$$

⁴The convergence follows from the ergodic theorem. It should be noted that the *rate* of convergence for Markov chains with infinite state spaces is a rather delicate matter, e.g. see Rosenthal [28]. In this paper the rate of convergence is non-essential. Moreover, the convergence itself does not appear in the proofs: it is used only as intuition and an (additional) justification for assuming the uniform distribution of the initial state.

Thus, for any fixed $R > \bar{R}_{\mathcal{A}}$ the expected average dynamic regret of algorithm \mathcal{A} over any sufficiently large interval is at most R , and it is the best possible upper bound of this form. Our goal is to bound $\bar{R}_{\mathcal{A}}$ in terms of the arms' volatility.

We use the following notation throughout the paper. The state of arm i at time t is $\mu_i(t)$. The maximal state at time t is $\mu^*(t) = \max_{i \in [k]} \mu_i(t)$. An arm i is *maximal* in round t if $\mu_i(t) = \mu^*(t)$.

1.2 Results: the state-informed case

We present an algorithm whose steady-state regret is optimal up to a poly-log factor.

Theorem 1.2. Consider the state-informed dynamic MAB problem with k arms, each with volatility at most σ . Assume that $k < \sigma^{-\gamma}$ for some $\gamma < \frac{1}{2}$. Then there exists a MAB algorithm whose steady-state regret is at most $\tilde{O}(k\sigma^2)$.

The algorithm is very intuitive. An arm with the highest last-observed state is called a *leader* and is played often, e.g. at least every other round. Suppose the last time some other arm i was observed was t rounds ago. By Azuma inequality the state of this arm changed by at most $\Delta\mu = \tilde{O}(\sigma\sqrt{t})$ since then, with high probability. If $\mu_i(t) + \Delta\mu$ is smaller than the state of the leader, then there is no point yet in trying arm i again. Else, we mark this arm *suspicious* and enqueue it to be played soon.

The main technical contribution here is the analysis, which is quite delicate since we need to deal with the complicated dependencies in the algorithm's behavior induced by the stochastically changing environment. Essentially, we manage to reduce the stochastic aspect of the problem to simple events in the state space. We achieve it as follows. Every time each arm is played, we spread the corresponding dynamic regret evenly over the corresponding idle time. This way we express the cumulative dynamic regret as a sum over the contributions of each arm in each round. We prove a uniform bound on the expectation of each such contribution. To this end, we identify a useful high-probability behavior of the system, derive deterministic guarantees conditional on this behavior (which is the tricky part), and then argue in terms of the corresponding conditional expectations.

Surprisingly, the steady-state regret of our algorithm essentially matches a lower bound based on a very simple idea: if in a given round the states of the best two arms are within $\frac{\sigma}{4}$ from one another, then in the next round with constant probability either one of them can be $\frac{\sigma}{4}$ above another, so any algorithm incurs expected dynamic regret $\Omega(\sigma)$.⁵

Theorem 1.3. Consider the state-informed dynamic MAB problem with k arms of volatility σ . Then the steady-state regret of any MAB algorithm is at least $\Omega(k\sigma^2)$.

1.3 Results: the state-oblivious case

Our algorithm for the state-oblivious case builds on an algorithm from [2] for the static MAB problem. That algorithm implicitly uses a simple "padding" function that for a given

⁵The former event happens with probability $\Omega(k\sigma)$, so the steady-state regret is $\Omega(k\sigma^2)$. This is the entire proof!

arm bounds the drift of an average reward from its (static) expected value. We design a new algorithm UCB_f which relies on a novel "padding" function f that accounts for the changing expected rewards. The analysis is quite technical: the specific results from [2] do not directly apply to our setting; instead, we need to "open up the hood" and combine the technique from [2] with some new ideas.

Theorem 1.4. *Consider the state-oblivious dynamic MAB problem with k arms such that each arm i has volatility at most σ_i . Then there exists a MAB algorithm whose steady-state regret is $\tilde{O}(k\sigma_{av})$, where $\sigma_{av}^2 = \frac{1}{k} \sum_{i=1}^k \sigma_i^2$.*

Note that (unlike the guarantee in Theorem 1.2), the guarantee here is in terms of an average volatility rather than the maximal one.

1.4 Using off-the-shelf MAB algorithms?

We ask whether similar results can be obtained using off-the-shelf MAB algorithms. Specifically, we investigate the following idea: take an off-the-shelf algorithm, run it and restart it every fixed number of rounds.

For the state-informed version we consider the obvious "greedy" approach: probe each arm, choose the best one, play it for a fixed number m of rounds, restart. The greedy algorithm is parameterized by the *phase length* m which can be tuned depending on the number of arms and their volatility. We show that the greedy algorithm is indeed suboptimal as compared to Theorem 1.2: the dependence on volatility (which is smaller than one) is linear rather than quadratic; we provide both upper and lower bounds.

For the state-oblivious version one can leverage on the existing work for the adversarial MAB problem [3]. This work assumes no restrictions on the state evolution, but provides guarantees only with respect to the best time-invariant policy, or a policy that switches arms a bounded number of times. We consider the following algorithm: run a fresh instance of algorithm EXP3 from [3] for a fixed number m of rounds, then restart. Using the off-the-shelf performance guarantees for EXP3 and fine-tuning m , one can (only) bound the steady-state regret by $\tilde{O}((k\sigma_{av})^{2/3})$, which is inferior to the result in Theorem 1.4. It is an open question whether one can obtain improved guarantees by tailoring the analysis in [3] to our setting.

1.5 Extensions and open questions

We extend our results in several directions. First, we generalize the Markov-chain formulation (2) to allow the random increments $X_i(t)$ to come from other distributions which has a certain "light-tailed" property, such as the discrete random walk. Second, we consider the setting in which each arm has a distinct fundamental interval. Third, we relax the assumption that the upper bound(s) on volatilities are known to the algorithm.

The main question left open by this paper is to close the gap between the upper and lower bounds for the state-oblivious dynamic MAB problem. The only lower bound we have is Theorem 1.3. We conjecture that one may obtain a better bound based on the relative entropy-based technique from [3]. It is also possible that the algorithmic result can

be improved, possibly via a more refined mechanism for discounting information with time.

Another open question is whether one can obtain the optimal $\tilde{O}(k\sigma^2)$ steady-state regret for the state-informed version in the case when $k \geq \sigma^{-1/2}$. Note that the greedy algorithm mentioned in Section 1.4 achieves steady-state regret $\tilde{O}(k\sigma)$ which is non-trivial for any $k \leq \sigma^{-1}$.

1.6 Organization of the paper

In Sections 2 and 3 we present our main results for the state-informed and the state-oblivious versions, respectively. Section 4 discusses using off-the-shelf MAB algorithms. Section 5 covers the extensions.

2 The state-informed dynamic MAB problem

We consider the state-informed dynamic MAB problem where the volatility of each arm is at most σ . Recall that the state of arm i at time t is denoted $\mu_i(t)$.

For arm i and time t , the *last-seen time* $\tau_i(t)$ is the last time this arm has been played strictly before time t ; the *last-seen state* $\nu_i(t) = \mu_i(\tau_i(t))$ is the corresponding state.

Definition 2.1. The *leader* in round t is the arm with a larger last-seen state, among the arms played in rounds $t-1$ and $t-2$; break ties in favor of the arm played in round $t-1$.

In our algorithm, the leader is our running estimate for an arm with the maximal state. We alternate rounds in which we always *exploit* – play the leader, with rounds in which we may explore other options. Since we define the leader in terms of the last two rounds only, our knowledge of its state is essentially up-to-date.

Let $\nu^*(t)$ be the last-seen state of the leader in round t . Let $c_{\text{susp}} = \Theta(\log \frac{1}{\sigma})^{1/2}$ be the factor to be defined later.

Definition 2.2. An arm i is called *suspicious* at time t if

$$\nu^*(t) - \nu_i(t) \leq c_{\text{susp}} \sigma \sqrt{t - \tau_i(t)}. \quad (4)$$

If an arm i is not suspicious at time t , then with high probability its current reward is less than $\nu^*(t)$. If no arm is suspicious then, intuitively, the best bet is to play the leader. Roughly, our algorithm behaves as follows: if the time is even it plays the current leader, and if the time is odd it plays a suspicious arm if one exists, and the leader otherwise. To complete the description of the algorithm, we need to specify what it does when there are multiple suspicious arms. In particular, we need to guarantee that after an arm becomes suspicious, it is played eventually (and preferably soon).

Definition 2.3. An arm i is *active* at time t if it is not the leader and it has been suspicious at some time $t' > \tau_i(t)$. The *activation time* $\tau_i^{\text{act}}(t)$ is the earliest such time t' .

An arm becomes active when it becomes suspicious. It stays active until it is played. The idea is to play an active arm with the earliest activation time.

Algorithm 2.4. *For bootstrapping, each arm is played once. At any later time t do the following. If t is even, play the current leader. If t is odd play an active arm (with the earliest activation time) if one exists, else play the leader.*

We will use a slightly more refined algorithm which allows for a more efficient analysis. Essentially, we give priority to arms whose state is close to the leader's.

Definition 2.5. Arm i is *high-priority* at time t if it is active at this time and moreover $\tau_i^{\text{act}}(t) - \tau_i(t) \leq 4k$.

Algorithm 2.6. For bootstrapping, each arm is played once. At any later time t do the following. If t is even, play the current leader. If t is odd play an active arm if one exists, else play the leader. If there are multiple active arms:

- if $t \equiv 1 \pmod{4}$ then play an active arm with the earliest activation time; break ties arbitrarily
- if $t \equiv 3 \pmod{4}$ then play a high-priority arm with the earliest activation time if one exists; else, play any active arm; break ties arbitrarily.

The analysis of these two algorithms are very similar, except that Algorithm 2.4 has inefficiencies which lead to an extra k^2 factor in its regret. We focus on Algorithm 2.6.

Theorem 2.7. Consider the state-informed dynamic MAB problem with k arms, each with volatility at most σ . Assume that $k < \sigma^{-\gamma}$ for some $\gamma < \frac{1}{2}$. Then Algorithm 2.6 achieves steady-state regret $O(k \sigma^2 \log^2 1/\sigma)$.

In the rest of this section we prove Theorem 2.7.

Let $\bar{R}_A(t)$ be the average dynamic regret up to time t . Then, letting $T_i(t)$ be the set of times arm i was played before and including time t , we have

$$E[\bar{R}_A(t)] = \frac{1}{t} \sum_{i \in [k]} \sum_{t' \in T_i(t)} E[\mu^*(t') - \mu_i(t')]. \quad (5)$$

Let us spread contributions of individual arms evenly over the corresponding idle time. Specifically, let us define

$$\begin{aligned} \Delta\mu_i(t) &= \mu^*(t) - \mu_i(t), \\ \Delta\tau_i(t) &= \tau_i^+(t) - \tau_i(t), \end{aligned}$$

where $\tau_i^+(t)$ is the next time arm i is played after time $\tau_i(t)$.⁶ Then we can re-write (5) as follows:

$$E[\bar{R}_A(t)] = \frac{1}{t} \sum_{i \in [k]} \sum_{t' \in [t]} E\left[\frac{\Delta\mu_i(\tau_i(t'))}{\Delta\tau_i(t')}\right]. \quad (6)$$

We define the *contribution* of arm i in round t as

$$C_i(t) = \frac{\Delta\mu_i(\tau_i(t))}{\Delta\tau_i(t)}.$$

A crucial idea is that we upper-bound $E[C_i(t)]$ for each round t separately. Namely, we will prove that

$$E[C_i(t)] < O(\sigma^2 \log^2 \frac{1}{\sigma}). \quad (7)$$

We identify the high-probability behavior of the processes $\{\mu_i(t)\}_{i \in [k]}$. Specifically, we consider the $\tilde{O}(\sqrt{t})$ bound on deviations, and an $O(1)$ bound on the number of near-optimal arms. A large portion of our analysis is deterministic conditional on such behavior.

⁶In other words, $\tau_i(t)$ and $\tau_i^+(t)$ are the two consecutive times arm i is played such that $\tau_i(t) < t \leq \tau_i^+(t)$.

Definition 2.8. A real-valued function f is *well-behaved* on an interval $[t_1; t_2]$ if for any $t, t + \Delta t \in [t_1; t_2]$ we have

$$|f(t + \Delta t) - f(t)| < c_{\text{well}} \sigma \sqrt{\Delta t}. \quad (8)$$

where $c_{\text{well}} = \Theta(\log \frac{1}{\sigma})^{1/2}$ will be chosen later.

Definition 2.9. An instance of the dynamic MAB problem is *well-behaved* on a time interval I if

- functions $\mu_1(t), \dots, \mu_k(t)$ are well-behaved on I ;
- at each time $t \in I$ there are at most $c_{\text{near}} = O(1)$ arms i such that $\Delta\mu_i(t) < (8k + 15\sqrt{k}) c_{\text{well}} \sigma$.⁷

A problem instance is *well-behaved near time t* if it is well-behaved on the time interval $[t - 3\sigma^{-2}; t + \sigma^{-2}]$.

Choosing the parameters. The factors c_{well} and c_{near} are chosen so that for any fixed t a problem instance is well-behaved near time t with probability at least $1 - \sigma^{-3}$. In Definition 2.1, define $c_{\text{susp}} = 5 c_{\text{well}}$.

Our conditionally deterministic guarantees (conditional on the problem instance being well-behaved) are expressed by the following lemma.

Lemma 2.10 (The Deterministic Lemma). *Suppose a problem instance is well-behaved near time t . Fix arm i and let $\delta = \Delta\mu_i(t)$. Then:*

- If $\delta = 0$ and $C_i(t) > 0$ then

$$C_i(t) \leq O(\sigma \log \frac{1}{\sigma}) / \sqrt{t - \tau_i(t)}, \quad (9)$$

and moreover for some arm $j \neq i$ we have

$$\Delta\mu_j(t) < O(\sigma \log \frac{1}{\sigma}) \sqrt{t - \tau_i(t)}. \quad (10)$$

- If $\delta > 0$ then $C_i(t) \leq O(\sigma^2/\delta) \log^2 \frac{1}{\sigma}$.

Let us use Lemma 2.10 to derive the main result.

Proof of Theorem 2.7: It suffices to prove (7). Let $\mathcal{E}(t)$ denote the event that the problem instance is well-behaved near time t . By Lemma 2.10(a), letting $x = \sqrt{t - \tau_i(t)}$ and suppressing the $\log \frac{1}{\sigma}$ factors under the $\tilde{O}(\cdot)$ notation,

$$\begin{aligned} E[C_i(t) \mid \Delta\mu_i(t) = 0, \mathcal{E}(t)] \\ \leq \tilde{O}(\sigma/x) \Pr[\exists j \neq i : \Delta\mu_j(t) < \tilde{O}(\sigma x)] \\ \leq \tilde{O}(\sigma^2). \end{aligned} \quad (11)$$

By Lemma 2.10(b) for any $\delta > 0$ we have

$$E[C_i(t) \mid \Delta\mu_i(t) \geq \delta, \mathcal{E}(t)] \leq \tilde{O}(\sigma^2/\delta) \quad (12)$$

$$\Pr[\Delta\mu_i(t) \leq \delta \mid \Delta\mu_i(t) > 0, \mathcal{E}(t)] \leq \tilde{O}(\delta). \quad (13)$$

Now (7) follows from (11-13) via a simple computation. \square

In the rest of this section we prove Lemma 2.10.

⁷This expression is tailored to (16) in the subsequent analysis. The event in question happens with probability at least $1 - O(c_{\text{well}} k^2 \sigma)^{c_{\text{near}}}$. Recall that we assume $k < \sigma^{-\gamma}$ for some $\gamma < \frac{1}{2}$. Thus, a (large enough) constant c_{near} suffices to guarantee a sufficiently low failure probability.

2.1 Deterministic bounds for the leader

We will argue deterministically assuming that the problem instance is well-behaved. We split our argument into a chain of claims and lemmas. The proofs are quite detailed; one can skip them for the first reading. For shorthand, let $\mathcal{E}[t_1; t_2]$ denote the event that the (fixed) problem instance is well-behaved on the time interval $[t_1; t_2]$.

First, we argue that the leader's last-seen state, $\nu^*(\cdot)$, does not decrease too much in one round.

Claim 2.11. *If $\mathcal{E}[t-2; t]$ then*

$$\nu^*(t+1) \geq \nu^*(t) - 2c_{\text{well}}\sigma.$$

Proof. Assume that t is even (if t is odd the proof proceeds similarly). Recall that the leader in round t is some arm i played in one of the previous two rounds. It follows that

$$\nu^*(t) = \nu_i(t) \leq \mu_i(t) + 2c_{\text{well}}\sigma.$$

Moreover, the leader (i.e. arm i) is played in round t and therefore $\nu^*(t+1) \geq \mu_i(t)$, claim proved. \square

Second, each arm becomes active eventually.

Claim 2.12. *Any arm i becomes active at most σ^{-2} rounds after it is played: $\tau_i^{\text{act}}(t) - \tau_i(t) \leq \sigma^{-2}$ for any time t .*

Proof. If $t - \tau_i(t) \geq \sigma^{-2}$ then (4) is trivially true. \square

Third, we show that a currently maximal arm has been activated within $4k$ rounds from its last-seen time, and therefore it has been played in the previous $8k$ rounds. The proof of this lemma is one of the crucial arguments in our analysis.

Lemma 2.13. *Suppose $\mathcal{E}[t - \sigma^{-2}; t]$ and arm i is maximal at time t . Then*

$$\tau_i^{\text{act}}(t) - \tau_i(t) \leq t - \tau_i^{\text{act}}(t) \leq 4k.$$

Proof. Note that $t - \tau_i^{\text{act}}(t) \leq 4k$, since otherwise after becoming active at time $\tau_i^{\text{act}}(t)$ arm i would have been played strictly before round t , contradiction.

Let $\tau = \tau_i(t)$. For the sake of contradiction assume that

$$\tau_i^{\text{act}}(t) - \tau > t - \tau_i^{\text{act}}(t). \quad (14)$$

Since arm i is not suspicious at time $t' = \tau_i^{\text{act}}(t) - 1$, by Definition 2.2 we have

$$\nu^*(t') - \nu_i(t') \geq c_{\text{susp}}\sigma\sqrt{t' - \tau}. \quad (15)$$

By Claim 2.12 the problem instance is well-behaved on $[\tau; t]$. It follows that

$$\begin{aligned} \nu_i(t') &= \mu_i(\tau) \geq \mu_i(t) - c_{\text{well}}\sigma\sqrt{t - \tau} \\ \nu^*(t') &= \mu_j(t'') \leq \mu_j(t) + c_{\text{well}}\sigma\sqrt{t - t''}, \end{aligned}$$

where arm j is the leader in round t' , and t'' is one of the two rounds preceding t' . Plugging this into (15) and using (14), we see that $\mu_j(t) > \mu_i(t)$, contradiction. \square

Fourth, we show that the leader's last-seen state is not much worse than the maximal state.

Claim 2.14. *If $\mathcal{E}[t - \sigma^{-2}; t]$ then*

$$\mu^*(t) - \nu^*(t) \leq (8k + \sqrt{8k})c_{\text{well}}\sigma.$$

Proof. Let $\mu^*(t) = \mu_i(t)$ for some arm i , and let $\tau = \tau_i(t)$ be the last time this arm was played. By Lemma 2.13 we have $t - \tau \leq 8k$. Therefore

$$\nu^*(\tau + 1) \geq \mu_i(\tau) \geq \mu_i(t) - c_{\text{well}}\sigma\sqrt{8k},$$

and the claim follows by Claim 2.11. \square

Fifth, we show that high-priority arms are played very soon after they become active.

Claim 2.15. *Suppose arm i is a high-priority active arm at time t . Assume $\mathcal{E}[t - \sigma^{-2}; t]$. Then $t - \tau_i^{\text{act}}(t) \leq 4c_{\text{near}}$.*

Proof. Fix time t and let $t' = \tau_i^{\text{act}}(t)$ be the activation time of arm i . Then by Definition 2.4 and Definition 2.3

$$\nu^*(t') - \nu_i(t') \leq c_{\text{susp}}\sigma\sqrt{t - t'} \leq c_{\text{susp}}\sigma\sqrt{4k}.$$

Using Claim 2.14 to relate $\nu^*(t')$ and $\mu^*(t')$, and using the fact that $\nu_i(t') = \mu_i(\tau)$ and that $\mu_i(\cdot)$ is well-behaved, we obtain

$$\Delta\mu_i(t') \leq (8k + 15\sqrt{k})c_{\text{well}}\sigma. \quad (16)$$

Lemma follows by Definition 2.9(ii) which is, in fact, tailored to (16). \square

Now we have the tools needed to prove a stronger version of Claim 2.14: $\mu^*(t) - \nu^*(t) \leq \tilde{O}(\sigma)$.

Lemma 2.16. *If the problem instance is well-behaved on $[t - \sigma^{-2}; t]$ then $\mu^*(t) - \nu^*(t) \leq O(c_{\text{well}}\sigma)$.*

Proof. Let i be an active arm at time t . By Lemma 2.13 $\tau_i^{\text{act}}(t) - \tau_i(t) \leq 4k$, so at time $\tau_i^{\text{act}}(t)$ arm i is a high-priority active arm. By Claim 2.15 $t - \tau_i^{\text{act}}(t) \leq 4c_{\text{near}} = O(1)$. By Lemma 2.13 it follows that $t - \tau_i(t) \leq O(1)$.

Now $\nu^*(\tau + 1) \geq \mu_i(\tau)$ by definition of the leader; $\nu^*(t) \geq \nu^*(\tau + 1) - O(c_{\text{well}}\sigma)$ by Claim 2.11; and also $\mu^*(t) \leq \mu^*(\tau) + O(c_{\text{well}}\sigma)$ since the problem instance is well-behaved. Putting it together, we obtain the lemma. \square

2.2 Proof of The Deterministic Lemma

Let $\tau = \tau_i(t)$ and recall that we denote $\delta = \Delta\mu_i(t)$.

By Lemma 2.13 we have $t - \tau \leq 8k$. Since the problem instance is well-behaved on $[t - 8k; t]$, it follows that $\mu^*(\cdot)$ is well-behaved, too, and therefore

$$|\Delta\mu_i(t) - \Delta\mu_i(\tau)| \leq 2c_{\text{well}}\sigma\sqrt{t - \tau}, \quad (17)$$

which immediately implies (9). To obtain (10) note that (17) in fact applies to any arm j , in particular to an arm j that is maximal at time τ .

To prove Lemma 2.10(b), it suffices to prove the following two inequalities:

$$\Delta\tau_i(t) \geq \Omega(\delta/\sigma)^2 / \log \frac{1}{\sigma}, \quad (18)$$

$$\Delta\mu_i(\tau) \leq O(\delta + \sigma \log \frac{1}{\sigma}). \quad (19)$$

Proof of (18): We consider two cases.

First, if we have $\Delta\mu_i(\tau) < \delta/2$ then by (17) we obtain

$$2c_{\text{well}}\sigma\sqrt{t - \tau} \geq |\Delta\mu_i(t) - \Delta\mu_i(\tau)| \geq \delta/2,$$

and (18) follows since $\Delta\tau_i(t) \geq t - \tau$.

Second, assume $\Delta\mu_i(\tau) \geq \delta/2$. Then by Lemma 2.16 for any time $t' \in (\tau; t + \sigma^{-2})$ we have

$$\begin{aligned} \nu^*(t') - \mu_i(\tau) &\geq \mu^*(t') - O(c_{\text{well}} \sigma) + \Delta\mu_i(\tau) - \mu^*(\tau) \\ &\geq \delta/2 - c_{\text{well}} \sigma \sqrt{t' - \tau} + O(1). \end{aligned}$$

This is at least $\geq c_{\text{susp}} \sigma \sqrt{t' - \tau}$ as long as it is the case that $t' - \tau \leq (12 c_{\text{well}} \sigma / \delta)^{-2}$. So for any such t' arm i is not suspicious, proving (18). \square

Proof of (19): First, note that if $\tau_i^{\text{act}}(t) - \tau_i(t) \leq 4k$ then by Definition 2.5 arm i is a high-priority active arm at time $\tau_i^{\text{act}}(t)$, so by Claim 2.15 we have $t - \tau_i^{\text{act}}(t) \leq O(1)$ and so $t - \tau_i(t) \leq O(1)$ by Lemma 2.13. It follows by (17) that

$$\Delta\mu_i(\tau) \leq \Delta\mu_i(t) + O(\sigma),$$

and we are done. In what follows we will assume that

$$\tau_i^{\text{act}}(t) - \tau_i(t) > 4k. \quad (20)$$

Note that for any time t' we have

$$\begin{aligned} \nu^*(t') &\leq \max(\mu^*(t' - 1), \mu^*(t' - 2)) \\ &\leq \mu^*(t') + 2 c_{\text{well}} \sigma. \end{aligned}$$

Let $t' = \tau_i^{\text{act}}(t) - 1$ be the round immediately preceding the activation time. Since arm i is not suspicious at time t' ,

$$\begin{aligned} c_{\text{susp}} \sigma \sqrt{t' - \tau} &\leq \nu^*(t') - \mu_i(\tau) \\ &\leq \mu^*(t') - \mu_i(\tau) + 2 c_{\text{well}} \sigma \\ &\leq \Delta\mu_i(t') + c_{\text{well}} \sigma (2 + \sqrt{t' - \tau}). \end{aligned}$$

Since $c_{\text{susp}} = 5 c_{\text{well}}$, it follows that

$$\Delta\mu_i(t') + 2 c_{\text{well}} \sigma \geq 4 c_{\text{well}} \sigma \sqrt{t' - \tau}. \quad (21)$$

Combining (17) and (21), we obtain

$$\begin{aligned} \Delta\mu_i(\tau) &\leq \Delta\mu_i(t') + 2 c_{\text{well}} \sigma \sqrt{t' - \tau} \\ &\leq \frac{3}{2} \Delta\mu_i(t') + 2 c_{\text{well}} \sigma. \end{aligned}$$

Finally, by (17), (20) and (21) we obtain

$$\begin{aligned} \Delta\mu_i(t') &\leq \Delta\mu_i(t) + 2 c_{\text{well}} \sigma \sqrt{t - t'} \\ &\leq \Delta\mu_i(t) + \frac{1}{2} \Delta\mu_i(t') + 2 c_{\text{well}} \sigma \\ \Delta\mu_i(t') &\leq 2 \Delta\mu_i(t) + 4 c_{\text{well}} \sigma \\ \Delta\mu_i(\tau') &\leq 3 \Delta\mu_i(t) + 6 c_{\text{well}} \sigma. \end{aligned} \quad \square$$

3 The state-oblivious dynamic MAB problem

We consider the state-oblivious dynamic MAB problem with k arms where the volatility of each arm i is at most σ_i .

Definition 3.1. For each arm i , $N_i(t)$ is the number of times it has been played in the first $t - 1$ rounds, and $\bar{W}_i(t)$ is the corresponding average reward. Let $\bar{W}_i(0) = 0$ if $N_i(t) = 0$. For shorthand, let $\mu_i = \mu_i(0)$ be the initial state.

Definition 3.2. Consider an instance of the state-oblivious dynamic MAB problem. A function $f_i : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{R}_+$ is a *padding* for arm i if the following two properties hold:

- $f_i(t, t_i)$ is increasing in t and decreasing in t_i ,

- for any time t , letting $t_i = N_i(t)$ we have

$$\Pr [|\bar{W}_i(t) - \mu_i(0)| > f_i(t, t_i)] < O(t^{-4}). \quad (22)$$

The family $\{f_i\}_{i \in [k]}$ is a *padding* for the problem instance.

We build on an algorithm UCB1 from [2] for the static MAB problem. We define a generalization of UCB1, which we call UCB_f , which is parameterized by a padding $f = \{f_i\}_{i \in [k]}$.

Algorithm 3.3 (UCB_f). *In each round t play any arm*

$$i \in \operatorname{argmax}_{i \in [k]} [\bar{W}_i(t) + f_i(t, N_i(t))].$$

The original UCB1 algorithm is defined for a specific padding f , and in fact does not explicitly uses the notion of a padding. We introduce this notion here in order to extend the ideas from [2] to our setting.

We incorporate the analysis from [2] via the following lemma which, essentially, bounds the number of times a sub-optimal arm is played by the algorithm.

Lemma 3.4 (Auer et al. [2]). *Consider an instance of the state-oblivious MAB problem with a padding $f = \{f_i\}_{i \in [k]}$. Consider the behavior of algorithm UCB_f in the first t rounds. Then for each arm i and any $t_i < t$ we have*

$$f_i(t, t_i) \leq \frac{1}{2} \Delta\mu_i(t) \Rightarrow E[N_i(t)] \leq t_i + O(1). \quad (23)$$

This lemma is implicit in Auer et al. [2], where it is the crux of the main proof. That proof considers the static MAB problem and (implicitly) a specific padding f .

We will use UCB_f where $f = \{f_i\}_{i \in [k]}$ is defined by

$$f_i(t, t_i) = \sqrt{2 \ln(t)/t_i} + \sigma_i \sqrt{8t \log t}. \quad (24)$$

Define the *average dynamic regret* of an algorithm \mathcal{A} $\bar{R}_{\mathcal{A}}(t) = \frac{1}{t} \sum_{s \in [t]} R_{\mathcal{A}}(s)$. We prove the following guarantee for algorithm UCB_f :

Theorem 3.5. *Consider the state-oblivious dynamic MAB problem with k arms. Suppose the volatility of each arm i is at most σ_i . Then there exists time t_0 such that*

$$E[\bar{R}_{\text{UCB}_f}(t_0)] \leq O(k \sigma_{\text{av}}) \log^{3/2}(\sigma_{\text{av}}^{-1}), \quad (25)$$

where $\sigma_{\text{av}}^2 = \frac{1}{k} \sum_{i=1}^k \sigma_i^2$.

To obtain Theorem 1.4 from Theorem 3.5, we start a fresh instance of algorithm UCB_f after every t_0 steps. We take advantage of the facts that (i) the "restarting times" are deterministic and, in particular, independent of the past history, and (ii) in any fixed round each $\mu_i(t)$ is distributed independently and uniformly on $[0; 1]$.

In the rest of this section we prove Theorem 3.5. We start with a very useful fact about the state evolution $\mu_i(t)$. In general, if $\mu_i(0) > \frac{1}{2}$ then due to the influence of the upper boundary the expected state $E[\mu_i(\cdot)]$ drifts down from its initial value. The following claim upper-bounds such drift.

Let us use a shorthand for the second summand in (24):

$$\delta_i(t) = \sigma_i \sqrt{8t \log t}.$$

Claim 3.6. Fix arm i and integer times $t \leq t_*$. Then

$$\Pr[|\mu_i(t) - \mu_i| > \delta_i] < t_*^{-3} \quad (26)$$

where $\mu_i = \mu_i(0)$ and $\delta_i = \delta_i(t_*)$, and therefore

$$E[\mu_i(t) | \mu_i] \geq \min(\mu_i, 1 - \delta_i) - t_*^{-2}. \quad (27)$$

Proof. Recall that the state $\mu_i(t)$ is defined as $f_I(B_i(t))$ where B_i is a Brownian motion with volatility σ_i , and f_I is the "projection" (1) into the interval $I = [0; 1]$ with reflective boundaries. Note that $\mu_i = \mathcal{B}_i(0)$.

It follows that $|\mu_i(t) - \mu_i| > \delta_i$ only if $|\mathcal{B}_i(t) - \mu_i| > \delta_i$. We know that for any $c > 1$ we have

$$\Pr[|\mathcal{B}_i(t) - \mu_i| > c \sigma_i \sqrt{t}] < 2e^{-c^2/2}.$$

We obtain (26) setting $c = \sqrt{6 \log t_*}$.

Now let us prove (27). Define

$$f(\mu) = E[\mu_i(t) | \mu_i].$$

Note that if $\mu < \frac{1}{2}$ then $f(\mu) > \mu$. Also, note that $f(\mu)$ is increasing and $f(\frac{1}{2}) = \frac{1}{2}$ by symmetry. Therefore, it suffices to prove (27) under the assumption that $\frac{1}{2} < \mu_i \leq 1 - \delta_i$.

Consider $T = \min(t, T_B)$, where

$$T_B = \min\{s \in \mathbb{N} : B_i(s) \notin (0; 1)\}.$$

Then $Z_s = \mu_i(\min(s, T))$ is a martingale such that $Z_0 = \mu$ and T is a bounded stopping time. By the Optional Stopping Theorem it follows that $E[Z_T] = \mu$. By (26) we have $T_B \geq t$ with probability at least $1 - t_*^{-2}$, in which case $T = t$ and $\mu = Z_T = \mu_i(t)$. Thus (27) follows. \square

Using Claim 3.6, let us argue that (24) is indeed a padding. Essentially, the first summand in (24) is tuned for an application of Chernoff-Hoeffding Bounds, whereas the second one corrects for the drift.

Lemma 3.7. The family f defined by (24) is a padding.

Proof. We need to prove (22). Fix arm i and time t . Let $\{t_j\}_{j=1}^\infty$ be the enumeration of all times when arm i is played. Let $X_j = \mu_i(t_j)$ be the state of arm i in round t . Let \hat{X}_j be the actual reward collected by the algorithm from arm i in round t_j . Let us define the sums $S = \sum_{j \in [n]} X_j$ and $S^* = \sum_{j \in [n]} \hat{X}_j$, where $n = N_i(t)$ is the number of times arm i is played before time t . Let $\mu = \mu_i(0)$ and $\delta = \delta_i(t)$.

We can rewrite (22) as follows:

$$\Pr[|S^* - \mu n| > \sqrt{2n \ln t} + \delta n] < O(t^{-4}). \quad (28)$$

Let F be the failure event when $|\mu_i(s) - \mu| > \delta$ for some $s \in [t]$. Recall that by Claim 3.6 the probability of F is at most t^{-4} . In the probability space induced by conditioning on $\hat{X}_1, \dots, \hat{X}_{j-1}$ and the event \bar{F} , we have

$$\begin{aligned} E[\hat{X}_j] &= E[E[\hat{X}_j | t_j, X_j]] = E[E[X_j | t_j]] \\ &= E[X_j] \in [\mu - \delta, \mu + \delta]. \end{aligned}$$

Going back to the original probability space,

$$E[\hat{X}_j | \hat{X}_1 \dots \hat{X}_{j-1}, \bar{F}] \in [\mu - \delta, \mu + \delta]. \quad (29)$$

The Chernoff-Hoeffding bounds (applied to the probability space induced by conditioning on \bar{F}) say precisely that the condition (29) implies the following tail inequality:

$$\Pr[|\hat{S} - \mu n| > \delta m + a | \bar{F}] \leq 2e^{-2a^2/m}$$

for any $a \leq 0$. We obtain (28) by taking $a = \sqrt{2m \ln T}$. \square

To argue about algorithm UCB_f , we will use the following notation:

Definition 3.8. We will use the following notation:

$$\begin{cases} \rho_i(t) &= \min(\mu_i, 1 - \delta_i(t)), & \mu_i &= \mu_i(0), \\ \Delta_i &= \mu^* - \mu_i, & \mu^* &= \mu^*(0) \\ S(t) &= \{\text{arms } i : \Delta_i \geq 4\delta_i(t)\}. \end{cases}$$

Lemma 3.9. Consider any algorithm for the state-oblivious dynamic MAB problem. Then for each arm i and time $t \geq k$

$$E[N_i(t) \bar{W}_i(t) | \mu_i] \geq \rho_i(t) E[N_i(t)] - t^{-2}. \quad (30)$$

The left-hand side of (30) is the total winnings collected by arm i up to time t . If the bandit algorithm always plays arm i , then $N_i(t) = t$ and the left-hand side of (30) is simply equal to $\sum_s E[\mu_i(s)]$, so the lemma follows from Claim 3.6. In this sense, Lemma 3.9 is an extension of Claim 3.6. The proof of (30) is a rather intricate exercise in conditional expectations and martingales. We defer it to Section 3.1.

We combine Lemma 3.9 and Lemma 3.4 to derive a conditional bound on $\bar{R}_{\text{UCB}_f}(t)$:

Corollary 3.10. For any time t we have

$$\begin{aligned} E[\bar{R}_{\text{UCB}_f}(t) | \mu_1, \dots, \mu_k] &\leq \frac{k}{t^2} + O(1) \left[\sum_{i \notin S(t)} \mu^* - \rho_i(t) \right] \\ &\quad + O\left(\frac{1}{t} \log t\right) \left[\sum_{i \in S(t)} \frac{1}{\Delta_i} \right]. \quad (31) \end{aligned}$$

Proof. Fix time t and let $\bar{W}_i = \bar{W}_i(t)$, $\rho_i = \rho_i(t)$ and $N_i = N_i(t)$. Let $R(t)$ be the left-hand side of (31). Using (30),

$$\begin{aligned} t R(t) &= \sum_{i \in [k]} E[(\mu^* - \bar{W}_i) N_i] \\ &\leq \sum_{i \in [k]} E[N_i] (\mu^* - \rho_i) + t^{-2}. \end{aligned}$$

For each $i \in S(t)$ we have $\mu^* - \rho_i \leq 2\Delta_i$ and by Lemma 3.4

$$E[N_i(t)] \leq 32 \ln(m) / \Delta_i^2 + O(1). \quad \square$$

We obtain Theorem 3.5 by integrating both sides of (31) with respect to $\mu_1 \dots \mu_k$.

Proof of Theorem 3.5: Fix time t and let $\delta_i = \delta_i(t)$ and $\rho_i = \rho_i(t)$. Note that (31) is, essentially, the sum over all arms. We partition the arms into three sets and bound the three corresponding sums separately.

Note that the following holds for any fixed $\gamma > 0$: given μ^* and the event $\{\Delta_i > \gamma\}$, the random variable μ_i is distributed uniformly on the interval $[0; \mu^* - \gamma]$. We will use this property in the forthcoming integrations.

First, we consider the set $S = S(t)$. Conditional on μ^* ,

$$\begin{aligned} E \left[\sum_{i \in S} \Delta_i^{-1} \right] &= \sum_{i \in [k]} E [\Delta_i^{-1} | \Delta_i > 4\delta_i] \Pr[\Delta_i > 4\delta_i] \\ &\leq \sum_{i \in [k]} \ln \sigma_i^{-1} \leq O(k \ln \sigma_{\text{av}}^{-1}). \end{aligned} \quad (32)$$

Second, let us consider the set S^+ of all arms i such that $0 < \Delta_i < 4\delta_i$. Conditional on μ^* , we obtain

$$\begin{aligned} E \left[\sum_{i \in S^+} \mu^* - \rho_i \right] &\leq \sum_{i \in [k]} O(\delta_i) \Pr[\Delta_i < 4\delta_i | \Delta_i > 0] \\ &\leq \sum_{i \in [k]} O(\delta_i) \min(1, \delta_i / \mu^*). \end{aligned}$$

Integrating over μ^* , we obtain

$$\begin{aligned} E \left[\sum_{i \in S^+} \mu^* - \rho_i \right] &\leq \sum_{i \in [k]} O(\delta_i^2) \\ &\leq O(k \sigma_{\text{av}}^2 t \log t). \end{aligned} \quad (33)$$

Third, we consider the set S^* of all maximal arms, i.e. the set of all arms i such that $\Delta_i = 0$. We show the main steps of the argument, omitting the details of some straightforward integrations:

$$\begin{aligned} Z_i &:= \mathbb{I}_{\{\Delta_i=0\}} (\mu^* - \rho_i) \\ E[Z_i] &= E[E[Z_i | \mu^*]] = \frac{1}{k} E[\mu^* - \rho_i] = O(\delta_i^2) \\ E \left[\sum_{i \in S^*} \mu^* - \rho_i \right] &= \sum_{i \in [k]} E[Z_i] \leq O(k \sigma_{\text{av}}^2) (t \log t). \end{aligned} \quad (34)$$

Finally, using (32-34), we take expectations in (31):

$$E[\bar{R}_{\text{UCB}_f}(t)] = O\left(\frac{k}{t} \log t\right) ((\sigma_{\text{av}} t)^2 + \log \sigma_{\text{av}}^{-1}).$$

The theorem follows if we take $t_0 = \sigma_{\text{av}} \sqrt{\log \sigma_{\text{av}}^{-1}}$. \square

3.1 Proof of Lemma 3.9: conditional expectations

Fix arm i and time t . Let us introduce a more concise notation which gets rid of the subscript i . Let $\mu = \mu_i(0)$ and $\delta = \delta_i(t)$, and denote $N = N_i(t)$. For every time s , let $Y_s = \mu_i(s)$, and let X_s be the winnings from arm i at time s if it is played by the algorithm.⁸ Let ζ_s be equal to 1 if arm i is played at time s , and 0 otherwise.

To prove (30), we will show that

$$\begin{aligned} E \left[\sum_{s \in [t]} \zeta_s X_s \right] &= E \left[\sum_{s \in [t]} \zeta_s Y_s \right] \\ &\geq \min(\mu, 1 - \delta) E[N] + t^{-2}. \end{aligned} \quad (35)$$

Note that ζ_s and X_s are conditionally independent given Y_s . It follows that

$$\begin{aligned} E[\zeta_s X_s | Y_s] &= E[\zeta_s | Y_s] E[X_s | Y_s] = E[\zeta_s | Y_s] Y_s \\ &= E[\zeta_s Y_s | Y_s]. \end{aligned}$$

⁸That is, X_s is an independent random sample from distribution $\mathcal{D}(Y_s)$, as defined in Section 1.1.

Taking expectations on both sides, we obtain

$$E[\zeta_s X_s] = E[\zeta_s Y_s],$$

which proves (35).

Going from (35) to (36) is somewhat more complicated. In what follows we denote $S = \sum_{t \in [m]} \zeta_s Y_s$.

Claim 3.11. *If $\mu \leq 1 - \delta$ then $E[S] \geq \mu E[N] - t^{-2}$.*

Proof. As in Claim 3.6, we recall the definition $\mu_i(s) = f_I(B_i(s))$ where B_i is a Brownian motion with volatility σ_i , and f_I is the "projection" (1) into the interval $I = [0; 1]$ with reflective boundaries. Note that $\mu_i = B_i(0)$.

For brevity, denote $\hat{Y}_s = B_i(s)$, and define the corresponding shorthand $\hat{S} = \sum_{s \in [t]} \zeta_s \hat{Y}_s$. Let F be the *failure event* when $\hat{Y}_s \geq 1$ for some $t \leq m$. Note that if this event does not occur, then $Y_s \geq \hat{Y}_s$ for every time $t \in [m]$ and therefore $S \geq \hat{S}$. We use this observation to express $E[S]$ in terms of $E[\hat{S}]$. Let $p := \Pr[F]$ and note that it is at most m^{-4} . Then:

$$\begin{aligned} E[\hat{S}] &= (1 - p) E[\hat{S} | \text{not } F] + p E[\hat{S} | F] \\ &\leq (1 - p) E[\hat{S} | \text{not } F] + p(\mu + t\sigma_i) \\ E[S] &\geq (1 - p) E[S | \text{not } F] + p E[S | F] \\ &\geq (1 - p) E[\hat{S} | \text{not } F] \\ &\geq E[\hat{S}] - pt\sigma_i - p. \end{aligned}$$

To prove the claim, it remains to bound $E[\hat{S}]$.

Let $\{s_j\}_{j=1}^\infty$ be the enumeration of all times when arm i is played. Note that $N = \max\{j : s_j \leq t\}$. Define $\hat{Z}_j = \hat{Y}_{s_j}$ for each j . We would like to argue that $\{\hat{Z}_j\}_{j=1}^\infty$ is a martingale and N is a stopping time. More precisely, claim that this is true for some common filtration. Indeed, one way to define such filtration $\{\mathcal{F}_j\}_{j=1}^\infty$ is to define \mathcal{F}_j as the σ -algebra generated by s_{j+1} and all tuples $(s_l, Z_l, Z_l^*, \hat{Z}_l)$ such that $l \leq j$. Now using the Optional Stopping Theorem one can show that

$$E[\hat{S}] = \sum_{j \in [N]} Z_j = E[N] E[\hat{Z}_0],$$

which proves the claim since $\hat{Z}_0 = \mu$. \square

To prove (36), it remains to consider the case $\mu > 1 - \delta$.

Claim 3.12. *if $\mu > 1 - \delta$ then*

$$E[S] \geq (1 - \delta) E[N] - t^{-2}.$$

Proof. Let T be the smallest time s such that $Y_s \leq 1 - \delta$. Let $\{s_j\}_{j=1}^\infty$ be the enumeration of all times when arm i is played, and let $J = \max j : t_j \leq T$. Conditioning on T and J , consider the entire problem starting from time $T+1$. Then by Claim 3.11 we have:

$$E \left[\sum_{s=T+1}^m \zeta_s Y_s | T, J \right] \geq (1 - \delta)(E[N] - J) - t^{-2}.$$

Let $S_T = \sum_{s=T+1}^t \zeta_s Y_s$. It follows that

$$\begin{aligned} S &= S_T + \sum_{t \in [T]} \zeta_s Y_s \geq S_T + (1 - \delta) J \\ E[S] &= E[S_T] + (1 - \delta) E[J] \\ &\geq (1 - \delta) E[N - J] - t^{-2} + (1 - \delta) E[J] \\ &\geq (1 - \delta) E[N] - t^{-2}. \quad \square \end{aligned}$$

4 Using off-the-shelf algorithms

In this section we investigate the following idea: take an off-the-shelf MAB algorithm, run it, and restart it every fixed number of rounds. We consider both the state-informed and state-oblivious versions of the dynamic MAB problem.

We use the following notation: there are k arms, each arm i has volatility σ_i , and the average volatility σ_{av} is defined by $\sigma_{\text{av}}^2 = \frac{1}{k} \sum_{i=1}^k \sigma_i^2$. We rely on the following lemma:

Lemma 4.1. *Let $\mu^* = \mu^*(0)$ and let $i^* \in \operatorname{argmax} \mu_i(0)$, ties broken arbitrarily. Then for any times $t \leq t_*$*

$$E[\mu^* - \mu_{i^*}(t)] \leq O(k)(t_*^{-4} + \sigma_{\text{av}}^2 t_* \log t_*). \quad (37)$$

More generally, we can consider arbitrary fixed times

$$0 \leq t_1 \leq t_2 \leq \dots \leq t_k \leq t \leq t_*$$

and define $\mu^* = \max \mu_i(t_i)$ and $i^* \in \operatorname{argmax} \mu_i(t_i)$.

The lemma is obtained, essentially, by combining Claim 3.6 and (34); we omit the details of the proof.

Remark. The intuition is that each arm i is probed in round t_i , so that $\mu_i(t_i)$ is the expected value of the corresponding probe. This lemma is similar to Claim 3.6 in that it bounds the downwards drift of $E[\mu_i(\cdot)]$ which is caused by the proximity of the upper boundary. The difference is that here we specifically consider a "maximal" arm, e.g. when $t_i \equiv 0$ we consider an arm which is maximal at time 0.

4.1 State-informed version: greedy algorithm

For the state-informed version we consider a very simple, "greedy" approach: probe each arm once, choose one with the largest state, play it for a fixed number $m - k$ of rounds, restart. Call this a *greedy algorithm* with phase length m .

Theorem 4.2. *Consider the state-informed dynamic MAB problem with k arms such that the volatility of each arm i is σ_i . With phase length $m = \sigma_{\text{av}} \sqrt{\log \sigma_{\text{av}}^{-1}}$, the steady-state regret of the greedy algorithm is at most*

$$O(k \sigma_{\text{av}} \log \sigma_{\text{av}}^{-1}), \text{ where } \sigma_{\text{av}}^2 = \frac{1}{k} \sum_{i=1}^k \sigma_i^2.$$

Proof. For the algorithmic result, fix phase length $m > k$ and consider a single phase of the greedy algorithm. Assume without loss of generality that in the first k rounds of the phase our algorithm plays arm i in step i . Let $\mu_i = \mu_i(i)$ be the corresponding rewards, and let μ^* be the largest of them. Then the greedy algorithm chooses arm $i^* \in \operatorname{argmax}_{i \in [k]} \mu_i$ and plays it for $m - k$ rounds. Consider the t -th of these $m - k$ rounds and let $Y_t = \mu_{i^*}(t + k)$ be the state of arm i^* in this round. By Lemma 4.1 we have $E[Y_t] \geq E[\mu^*] - z$, where z is the right-hand side of (37). Therefore, letting \bar{W} be the per-round average reward in this phase, we have

$$\begin{aligned} E[\bar{W}] &\geq \frac{1}{m} \sum_{t=1}^{m-k} Y_t \geq \frac{m-k}{m} (E[\mu^*] - z) \\ E[\mu^* - \bar{W}] &\leq z + \frac{k}{m} E[\mu^*] \\ &\leq O(km \sigma_{\text{av}}^2 \log m) + \frac{k}{m} (1 + \frac{1}{m}) \\ &= O(k \sigma_{\text{av}}) \sqrt{\log \sigma_{\text{av}}^{-1}} \end{aligned}$$

for $m = \sigma_{\text{av}} \sqrt{\log \sigma_{\text{av}}^{-1}}$. \square

We provide a matching lower bound.

Theorem 4.3. *Consider the setting in Theorem 4.2. Then the steady-state regret of the greedy algorithm is $\tilde{\Omega}(k \sigma_{\text{av}})$.*

Proof Sketch. For simplicity assume $\sigma_i \equiv \sigma$. It is known that in time t a Brownian motion with volatility σ drifts by at least $\Delta = \tilde{\Omega}(\sigma \sqrt{t})$ with high probability. Thus for each arm i with high probability $\mu_i(t) \leq 1 - \Delta/2$, regardless of the initial value $\mu_i(0)$. Now we can obtain a lower bound that corresponds to Lemma 4.1: letting $\mu^* = \max \mu_i(i)$ and $i^* \in \operatorname{argmax} \mu_i(i)$ be the arm chosen by the greedy algorithm,

$$E[\mu^* - \mu_{i^*}(t)] \geq \tilde{\Omega}(k \sigma^2 t), \quad (38)$$

for any $t > k$. Now consider a given phase of the greedy algorithm. In the first k rounds the algorithm accumulates regret $\Omega(k)$, and in each subsequent round t the regret is the left-hand side of (38). The theorem follows easily. \square

4.2 State-oblivious version via adversarial MAB

For the state-oblivious dynamic MAB problem, we use a very general result of Auer et al. [3] for the adversarial MAB problem. For simplicity, here we only state this result in terms of the present setting.

Let $\bar{W}_{\mathcal{A}}(t)$ be the average reward collected by algorithm \mathcal{A} during the time interval $[1; t]$.

Theorem 4.4 (Auer et al.[3]). *Consider the state-oblivious dynamic MAB problem with k arms. Let \mathcal{A}_i be an algorithm that plays arm i at every step. Then there exists an algorithm, call it EXP3, such that for any arm i and any time t*

$$E[\bar{W}_{\text{EXP3}}(t)] \geq E[\bar{W}_{\mathcal{A}_i}(t)] - O(\frac{k}{t} \log t)^{1/2}.$$

For our problem, we restart EXP3 every m steps, for some fixed m ; call this algorithm EXP3(m).

Theorem 4.5. *Consider the state-informed dynamic MAB problem with k arms such that the volatility of each arm i is at most σ_i . Then there exists m such that algorithm EXP3(m) has steady-state regret*

$$O(k \sigma_{\text{av}} \log \sigma_{\text{av}}^{-1})^{2/3}, \text{ where } \sigma_{\text{av}}^2 = \frac{1}{k} \sum_{i=1}^k \sigma_i^2.$$

Proof. Let use shorthand $\mathcal{A} = \text{EXP3}(m)$. Let μ^* be the maximal expected reward at time 0, and suppose it is achieved by some arm i^* . Let \mathcal{A}^* be the algorithm that plays this arm at every step. Let $Y_t = \mu_{i^*}(t)$ the state of arm i^* in round t . Then by Lemma 4.1 we have $E[Y_t] \geq E[\mu^*] - z(m)$, where $z(m)$ is the right-hand side of (37). Therefore:

$$\begin{aligned} E[\bar{W}_{\mathcal{A}^*}(m)] &= E[E[\bar{W}_{\mathcal{A}^*}(m) | Y_1, \dots, Y_m]] \\ &= \frac{1}{m} E[\sum_{t=1}^m Y_t] \\ &\geq \mu^* - z(m) \\ E[\bar{R}_{\mathcal{A}}(m)] &= E[\mu^* - \bar{W}_{\mathcal{A}^*}(m)] \\ &\quad + E[\bar{W}_{\mathcal{A}^*}(m) - \bar{W}_{\mathcal{A}}(m)] \end{aligned} \quad (39)$$

Now using (39) and Theorem 4.4 we obtain

$$E[\bar{R}_{\mathcal{A}}(m)] \leq z(m) + O(\frac{k}{m} \log m)^{1/2}. \quad (40)$$

We choose m that minimizes the right-hand side of (40). \square

We note in passing that we can also get non-trivial (but worse) guarantees for the state-oblivious dynamic MAB problem using two other off-the-shelf approaches:

- a version of the greedy algorithm which probes each arm a few times in the beginning of each phase,
- a version of Theorem 4.4 in which the benchmark algorithm is allowed to switch arms a few times [3].

Essentially, the first approach is too primitive, while the second one makes overly pessimistic assumptions about the environment. In both cases we obtain guarantees of the form $\tilde{O}(k\sigma_{\text{av}})^\gamma$, $\gamma < \frac{2}{3}$, which are inferior to Theorem 4.5.

5 Extensions

Recall that the state evolution of arm i in the dynamic MAB problem is described by (2), where the i.i.d. increments $\mu_i(t)$ are distributed with respect to some fixed distribution \mathcal{X}_i . Can we relax the assumption that \mathcal{X}_i is normal?

Definition 5.1. Random variable X is *stochastically* (ρ, σ) -bounded if its moment-generating function satisfies

$$E[e^{r(X-E[x])}] \leq e^{r^2\sigma^2/2} \text{ for } |r| \leq \rho.$$

This is precisely the condition needed to establish an Azuma-type inequality: if S is the sum of t independent stochastically (ρ, σ) -bounded random variables with zero mean, then with high probability $S \leq \tilde{O}(\sigma\sqrt{t})$. Specifically, for any $\lambda \leq \frac{1}{2}\rho\sigma\sqrt{t}$ we have

$$\Pr[S > \lambda\sigma\sqrt{t}] \leq \exp(-\lambda^2/2). \quad (41)$$

Note that a normal distribution $\mathcal{N}(0, \sigma)$ is (∞, σ) -bounded, and any distribution with support $[-\sigma, \sigma]$ is $(1, \sigma)$ -bounded.

We can recover all of our algorithmic results if we assume that each distribution \mathcal{X}_i has zero mean and is stochastically (ρ, σ_i) -bounded for some σ_i , where $\rho > 0$ is a fixed absolute constant. We re-define the *volatility* of arm i as the infimum of all σ such that \mathcal{X}_i is (ρ, σ) -bounded.

It is appealing to tackle a more general setting when the only restriction on each distribution \mathcal{X}_i is that it has mean 0 and variance σ_i^2 . We can extend our analysis (at the cost of somewhat weaker guarantees) if we further assume that, essentially, the absolute third moment of \mathcal{X}_i is comparable to σ_i^3 . Then instead of (41) we can use a weaker inequality called the *non-uniform Berry-Esseen theorem* [23]:

$$\Pr\left[\sum_{s=1}^t \mu_i(s) > \sigma_i t^\gamma\right] \leq O\left(\left(\frac{\rho_i}{\sigma_i}\right)^3 t^{1-3\gamma}\right), \quad (42)$$

for any $\gamma > 1/2$, where $\rho_i^3 = E[|\mu_i(s)|^3]$. We omit further discussion of this extension from the present version.

Let us discuss one other direction in which our setting can be generalized. Recall that in the dynamic MAB problem the state of each arm evolves on the same interval $I = [0; 1]$ (see Section 1.1) which we term the *fundamental interval*. What if we allow each arm to have a distinct fundamental interval? All our algorithms fit this extended setting with little or no modification. The performance guarantees

should look like a weighted sum of contributions from different arms, where the weights depend (perhaps in rather complicated way) on the respective fundamental intervals. To illustrate this point, we worked out the guarantees for the two algorithms discussed in Section 4, see Appendix A for details. It is an open question to derive similar closed-form guarantees for the other algorithms in this paper.

Recall that in all our results we assumed that the volatilities are known to the algorithm. In fact, this assumption is not necessary: we are interested in the stationary performance of our algorithms and, as it turns out, we can afford to learn the static parameters of the model. Roughly, the argument goes as follows. It suffices for our analysis if for each arm an algorithm knows a 2-approximate upper bound on volatility σ_i , rather than the exact value. One can learn such bound by playing arm i for $O(\log^2 \sigma_i)$ rounds, with failure probability as low as $O(\sigma_i^{-10})$, and repeat this learning phase every σ_i^{-1} rounds (we omit the details).

Acknowledgments. The first author would like to thank Bobby Kleinberg for many stimulating conversations about multi-armed bandits.

References

- [1] P. Auer. Using confidence bounds for exploitation-exploration trade-offs. *J. Machine Learning Research*, 3:397–422, 2002. Preliminary version in *41st IEEE FOCS*, 2000.
- [2] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002. Preliminary version in *15th ICML*, 1998.
- [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, 2002. Preliminary version in *36th IEEE FOCS*, 1995.
- [4] B. Awerbuch and R. D. Kleinberg. Adaptive routing with end-to-end feedback: distributed learning and geometric approaches. In *36th ACM Symp. on Theory of Computing (STOC)*, pages 45–53, 2004.
- [5] D. A. Berry and B. Fristedt. *Bandit problems: sequential allocation of experiments*. Chapman and Hall, 1985.
- [6] D. Bertsimas and J. Nino-Mora. Conservation laws, extended polymatroids and multi-armed bandit problems: A unified polyhedral approach. *Math. of Oper. Res.*, 21(2):257–306, 1996.
- [7] D. Bertsimas and J. Nino-Mora. Restless bandits, linear programming relaxations, and a primal-dual index heuristic. *Operations Research*, 48(1):80–90, 2000.
- [8] V. Dani and T. P. Hayes. How to beat the adaptive multi-armed bandit. Technical report. Available from arXiv at <http://arxiv.org/cs.DS/0602053>, 2006.
- [9] V. Dani and T. P. Hayes. Robbing the bandit: less regret in online geometric optimization against an adaptive adversary. In *17th ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 937–943, 2006.
- [10] A. Flaxman, A. Kalai, and H. B. McMahan. Online Convex Optimization in the Bandit Setting: Gradient Descent without a Gradient. In *16th ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 385–394, 2005.
- [11] J. C. Gittins. Bandit processes and dynamic allocation indices (with discussion). *J. Roy. Statist. Soc. Ser. B*, 41:148–177, 1979.
- [12] J. C. Gittins. *Multi-Armed Bandit Allocation Indices*. John Wiley & Sons, 1989.

- [13] J. C. Gittins and D. M. Jones. A dynamic allocation index for the sequential design of experiments. In J. G. et al., editor, *Progress in Statistics*, pages 241–266. North-Holland, 1974.
- [14] S. Guha and K. Munagala. Approximation algorithms for partial-information based stochastic control with Markovian rewards. In *48th Symp. on Foundations of Computer Science (FOCS)*, 2007.
- [15] S. Guha, K. Munagala, and P. Shi. On Index Policies for Restless Bandit Problems. arXiv:0711.3861v1 [cs.DS], 2007.
- [16] F. Heidari, S. Mannor, and L. Mason. Reinforcement learning-based load shared sequential routing. In *IFIP Networking*, 2007.
- [17] R. D. Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. In *18th Advances in Neural Information Processing Systems (NIPS)*, 2004. Full version appeared as Chapters 4-5 in [18].
- [18] R. D. Kleinberg. *Online Decision Problems with Large Strategy Sets*. PhD thesis, MIT, Boston, MA, 2005.
- [19] R. D. Kleinberg. Anytime algorithms for multi-armed bandit problems. In *17th ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 928–936, 2006. Full version appeared as Chapter 6 in [18].
- [20] R. D. Kleinberg and F. T. Leighton. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *44th Symp. on Foundations of Computer Science (FOCS)*, pages 594–605, 2003.
- [21] T. Lai and H. Robbins. Asymptotically efficient Adaptive Allocation Rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- [22] H. B. McMahan and A. Blum. Online Geometric Optimization in the Bandit Setting Against an Adaptive Adversary. In *17th Conference on Learning Theory (COLT)*, pages 109–123, 2004.
- [23] K. Neammanee. On the constant in the nonuniform version of the Berry-Esseen theorem. *Intl. J. of Mathematics and Mathematical Sciences*, 2005:12:1951–1967, 2005.
- [24] J. Nino-Mora. Restless bandits, partial conservation laws and indexability. *Advances in Applied Probability*, 33:76–98, 2001.
- [25] S. Pandey, D. Agarwal, D. Chakrabarti, and V. Josifovski. Bandits for Taxonomies: A Model-based Approach. In *SIAM Intl. Conf. on Data Mining (SDM)*, 2007.
- [26] C. H. Papadimitriou and J. N. Tsitsiklis. The complexity of optimal queueing network control. In *Structure in Complexity Theory Conference*, pages 318–322, 1994.
- [27] H. Robbins. Some Aspects of the Sequential Design of Experiments. *Bull. Amer. Math. Soc.*, 58:527–535, 1952.
- [28] J. S. Rosenthal. Markov chain convergence: From finite to infinite. *Stochastic Processes Appl.*, 62(1):55–72, 1996.
- [29] R. K. Sundaram. Generalized Bandit Problems. In D. Austen-Smith and J. Duggan, editors, *Social Choice and Strategic Decisions: Essays in Honor of Jeffrey S. Banks (Studies in Choice and Welfare)*, pages 131–162. Springer, 2005. First appeared as *Working Paper, Stern School of Business*, 2003.
- [30] J. N. Tsitsiklis. A short proof of the Gittins index theorem. *Annals of Applied Probability*, 4(1):194–199, 1994.
- [31] G. Weiss. Branching bandit processes. *Probab. Engng. Inform. Sci.*, 2:269–278, 1988.
- [32] P. Whittle. Arm acquiring bandits. *Ann. Probab.*, 9:284–292, 1981.
- [33] P. Whittle. Restless bandits: Activity allocation in a changing world. *J. of Appl. Prob.*, 25A:287–298, 1988.

A Distinct fundamental intervals

Recall that in the dynamic MAB problem the state of each arm evolves on the same interval $I = [0; 1]$ (see Section 1.1)

which we term the *fundamental interval*. In this section we consider a generalization in which we allow each arm to have a distinct fundamental interval. We work out the guarantees for the two algorithms discussed in Section 1.4.

The main contribution of this appendix is that we find a way to upper-bound the steady-state regret of the respective algorithms in terms of reasonably defined averages of the arms’ properties. The actual derivations are rather tedious but not that illuminating; we omit them from this version.

A.1 The setting and notation

We consider the following setting. There are k arms. Each arm has volatility σ_i and fundamental interval $[a_i; b_i]$. Without loss of generality we assume that $b_1 \leq \dots \leq b_k$ and that $\max a_i < \min b_i$. (If the latter fails then we can always ignore the arm with the smallest upper boundary b_i .) To simplify the derivation we assume that $\max \sigma_i \leq \frac{1}{3}$.

Define the *weight* of arm i as

$$w_i = \prod_{l=1}^k \frac{b_l - a_l}{b_l - a_l},$$

Define the *average volatility* σ_{av} by

$$\sigma_{av}^2 = \frac{\sum_{i \in [k]} w_i (b_i - a_i) \sigma_i^2}{\sum_{i \in [k]} w_i (b_i - a_i)}$$

Define the *average length* as

$$d_{av} = \frac{1}{k} \sum_{i \in [k]} w_i (b_i - a_i).$$

To see that the quantities we defined above are reasonable as *averages*, note that if all arms have the same fundamental interval $[a; b]$ then all weights are 1 and $d_{av} = b - a$ and, moreover, the average volatility σ_{av} coincides with the one defined in the body of the paper.

A.2 Results

We present two results that extend, respectively, Theorem 4.2 and Theorem 4.5 to the setting from Section A.1. In both cases the algorithms are exactly the same. The main tool is a version of Lemma 37, where the guarantee (37) looks exactly the same in our notation, except the right-hand side is multiplied by d_{av} .

Theorem A.1. *Consider the deterministic dynamic MAB problem in the setting from Section A.1. Let $a_{min} = \min a_i$. Then for phase length*

$$m = \sigma_{av}^{-1} \sqrt{(b_k - a_{min}) / \log \sigma_{av}^{-1}}$$

the greedy algorithm has steady-state regret

$$O(k \sigma_{av}) \sqrt{(b_k - a_{min}) d_{av} \log \sigma_{av}^{-1}}.$$

Theorem A.2. *Consider the state-informed dynamic MAB problem in the setting from Section A.1. Then there exists m such that algorithm EXP3(m) has steady-state regret*

$$O(d_{av})^{1/3} (k \sigma_{av} \log \sigma_{av}^{-1})^{2/3}.$$