

Optimal Exploration-Exploitation in a Multi-Armed-Bandit Problem with Non-stationary Rewards

Omar Besbes
Columbia University

Yonatan Gur
Columbia University

Assaf Zeevi*
Columbia University

May 13, 2014

Abstract

In a multi-armed bandit (MAB) problem a gambler needs to choose at each round of play one of K arms, each characterized by an unknown reward distribution. Reward realizations are only observed when an arm is selected, and the gambler's objective is to maximize his cumulative expected earnings over some given horizon of play T . To do this, the gambler needs to acquire information about arms (exploration) while simultaneously optimizing immediate rewards (exploitation); the price paid due to this trade off is often referred to as the *regret*, and the main question is how small can this price be as a function of the horizon length T . This problem has been studied extensively when the reward distributions do not change over time; an assumption that supports a sharp characterization of the regret, yet is often violated in practical settings. In this paper, we focus on a MAB formulation which allows for a broad range of temporal uncertainties in the rewards, while still maintaining mathematical tractability. We fully characterize the (regret) complexity of this class of MAB problems by establishing a direct link between the extent of allowable reward "variation" and the minimal achievable regret. Our analysis draws some connections between two rather disparate strands of literature: the adversarial and the stochastic MAB frameworks.

KEYWORDS: Multi-armed bandit, exploration / exploitation, non-stationary, minimax regret

1 Introduction

Background and motivation. In the presence of uncertainty and partial feedback on rewards, an agent that faces a sequence of decisions needs to judiciously use information collected from past observations when trying to optimize future actions. This fundamental paradigm is present in a variety of applications: an internet web site that seeks to customize recommendations to individual users whose tastes are a priori not known; a firm that launches a new product and needs to set a price to maximize profits but does not know the demand curve; a retailer that must select assortments of products among a larger variety of items but does not know the preferences of customers; a firm that selects routes over the internet to efficiently send packets of data to users but does not know the delay along available

*This work is supported by NSF grant 0964170 and BSF grant 2010466. Correspondence: ob2105@columbia.edu, ygur14@gsb.columbia.edu, assaf@gsb.columbia.edu.

routes; as well as many other instances. In all the above examples decisions can be adjusted on a weekly, daily or hourly basis (if not more frequently), and the history of observations may be used to optimize current and future performance. To do so effectively, the decision maker tries to balance between the acquisition cost of new information (*exploration*) that may be used to improve future decisions and rewards, and the generation of instantaneous rewards based on the existing information (*exploitation*).

A widely studied paradigm that captures the tension between exploration and exploitation is that of multi armed bandits (MAB), originally proposed in the context of drug testing by Thompson (1933), and placed in a general setting by Robbins (1952). The original setting has a gambler choosing among K slot machines at each round of play, and upon that selection observing a reward realization. In this classical formulation the rewards are assumed to be independent and identically distributed according to an unknown distribution that characterizes each machine. The objective is to maximize the expected sum of (possibly discounted) rewards received over a given (possibly infinite) time horizon. Since their inception, MAB problems with various modifications have been studied extensively in Statistics, Economics, Operations Research, and Computer Science, and are used to model a plethora of dynamic optimization problems under uncertainty; examples include clinical trials (Zelen 1969), strategic pricing (Bergemann and Valimaki 1996), investment in innovation (Bergemann and Hege 2005), packet routing (Awerbuch and Kleinberg 2004), on-line auctions (Kleinberg and Leighton 2003), assortment selection (Caro and Gallien 2007), and on-line advertising (Pandey et al. 2007), to name but a few. For overviews and further references cf. the monographs by Berry and Fristedt (1985), Gittins (1989) for Bayesian / dynamic programming formulations, and Cesa-Bianchi and Lugosi (2006) that covers the machine learning literature and the so-called adversarial setting.

Since the set of MAB instances in which one can identify the optimal policy is extremely limited, a typical yardstick to measure performance of a candidate policy is to compare it to a benchmark: an *oracle* that at each time instant selects the arm that maximizes expected reward. The difference between the performance of the policy and that of the oracle is called the *regret*. When the growth of the regret as a function of the horizon T is *sub-linear*, the policy is *long-run average optimal*: its long run average performance converges to that of the oracle. Hence the first order objective is to develop policies with this characteristic. The precise rate of growth of the regret as a function of T provides a refined measure of policy performance. Lai and Robbins (1985) is the first paper that provides a sharp characterization of the regret growth rate in the context of the traditional (stationary random rewards) setting, often referred to as the *stochastic* MAB problem. Most of the literature has followed this path with the objective of designing policies that exhibit the “slowest possible” rate of growth in the regret (often referred to as *rate optimal* policies).

In many application domains, several of which were noted above, temporal changes in the structure of

the reward distribution are an intrinsic characteristic of the problem. These are ignored in the traditional stochastic MAB formulation, but there have been several attempts to extend that framework. The origin of this line of work can be traced back to Gittins and Jones (1974) who considered a case where only the state of the chosen arm can change, giving rise to a rich line of work (see, e.g., Gittins 1979, and Whittle 1981). In particular, Whittle (1988) introduced the term *restless bandits*; a model in which the states (associated with the reward distributions) of the arms change in each step according to an arbitrary, yet known, stochastic process. Considered a notoriously hard class of problems (cf. Papadimitriou and Tsitsiklis 1994), this line of work has led to various approximation approaches, see, e.g., Bertsimas and Nino-Mora (2000), and relaxations, see, e.g., Guha and Munagala (2007) and references therein.

Departure from the stationarity assumption that has dominated much of the MAB literature raises fundamental questions as to how one should model temporal uncertainty in rewards, and how to benchmark performance of candidate policies. One extreme view, is to allow the reward realizations of arms to be selected at any point in time by an *adversary*. These ideas have their origins in game theory with the work of Blackwell (1956) and Hannan (1957), and have since seen significant development; Foster and Vohra (1999) and Cesa-Bianchi and Lugosi (2006) provide reviews of this line of research. Within this so called *adversarial* formulation, the efficacy of a policy over a given time horizon T is often measured relative to a benchmark which is defined by the single best action one could have taken in hindsight (after seeing all reward realizations). The single best action benchmark represents a *static* oracle, as it is constrained to a single (static) action. For obvious reasons, this static oracle can perform quite poorly relative to a “dynamic oracle” that follows the optimal *dynamic* sequence of actions, as the latter optimizes the (expected) reward at each time instant over all possible actions.¹ Thus, a potential limitation of the adversarial framework is that even if a policy has a “small” regret relative to a static oracle, there is no guarantee with regard to its performance relative to the dynamic oracle.

Main contributions. At a high level, the main contribution of this paper lies in fully characterizing the (regret) complexity of a broad class of MAB problems with non-stationary reward structure by establishing a direct link between the extent of reward “variation” and the minimal achievable worst-case regret. More specifically, the paper’s contributions are along four dimensions. On the modeling side we formulate a class of non-stationary reward structures that is quite general, and hence can be used to realistically capture a variety of real-world type phenomena, yet remain mathematically tractable. The main constraint that we impose on the evolution of the mean rewards is that their variation over the relevant time horizon is bounded by a *variation budget* V_T ; a concept that was recently introduced in Besbes et al. (2013) in the context of non-stationary stochastic approximation. This limits the power of nature compared to the adversarial setup discussed above where rewards can be picked to maximally

¹Under non-stationary reward structure it is immediate that the single best action may be sub-optimal in a large number of decision epochs, and the gap between the performance of the static and the dynamic oracles can grow linearly with T .

damage the policy at each instance within $\{1, \dots, T\}$. Nevertheless, this constraint still allows for a very rich class of temporal changes. In particular, this class extends most of the treatment in the non-stationary stochastic MAB literature which mainly focuses on a finite (known) number of changes in the mean reward values, see, e.g., Garivier and Moulines (2011) and references therein (see also Auer et al. (2002) in the adversarial context). It is also consistent with more extreme settings, such as the one treated in Slivkins and Upfal (2008) where reward distributions evolve according to a Brownian motion and hence the regret is linear in T (we explain these connections in more detail in §5).

The second dimension of contribution lies in the analysis domain. For the class of non-stationary reward distributions described above, we establish lower bounds on the performance of *any* non-anticipating policy relative to the *dynamic* oracle, and show that these bounds can be achieved, uniformly over the class of admissible reward distributions, by a suitable policy construction. The term “achieved” is meant in the sense of the order of the regret as a function of the time horizon T , the variation budget V_T , and the number of arms K . More precisely, our policies are shown to be minimax optimal up to a term that is logarithmic in the number of arms, and the regret is sublinear and is of the order of $(KV_T)^{1/3} T^{2/3}$. Auer et al. (2002), in the adversarial setting, and Garivier and Moulines (2011) in the stochastic setting, considered non-stationary rewards where the identity of the best arm can change a *finite* number of times; the regret in these instances (relative to a dynamic oracle) is shown to be of order \sqrt{T} . Our analysis complements these results by treating a broader and more flexible class of temporal changes in the reward distributions, yet still establishing optimality results and showing that sublinear regret is achievable. When V_T increases with the time horizon T , our results provide a spectrum of orders of the minimax regret ranging between order $T^{2/3}$ (when V_T is a constant independent of T) and order T (when V_T grows linearly with T), mapping allowed variation to best achievable performance.

With the analysis described above we shed light on the exploration-exploitation trade off that is a characteristic of the non-stationary reward setting, and the change in this trade off compared to the stationary setting. In particular, our results highlight the tension that exists between the need to “remember” and “forget.” This is characteristic of several algorithms that have been developed in the adversarial MAB literature, e.g., the family of exponential weight methods such as EXP3, EXP3.S and the like; see, e.g., Auer et al. (2002), and Cesa-Bianchi and Lugosi (2006). In a nutshell, the fewer past observations one retains, the larger the stochastic error associated with one’s estimates of the mean rewards, while at the same time using more past observations increases the risk of these being biased.

One interesting observation drawn in this paper is a connection between the adversarial MAB setting, and the non-stationary environment studied here. In particular, as in Besbes et al. (2013), it is seen that optimal policy in the adversarial setting may be suitably calibrated to perform near-optimally in the non-stationary stochastic setting. This will be further discussed after the main results are established.

Structure of the paper. §2 introduces the basic formulation of the stochastic non-stationary MAB problem. In §3 we provide a lower bound on the regret that any admissible policy must incur relative to a dynamic oracle. §4 introduces a policy that achieves that lower bound. §5 contains a brief discussion. Proofs can be found in the Appendix.

2 Problem Formulation

Let $\mathcal{K} = \{1, \dots, K\}$ be a set of arms. Let $\mathcal{T} = \{1, 2, \dots, T\}$ denote the sequence of decision epochs faced by the decision maker. At any epoch $t \in \mathcal{T}$, a decision-maker pulls one of the K arms. When pulling arm $k \in \mathcal{K}$ at epoch $t \in \mathcal{T}$, a reward $X_t^k \in [0, 1]$ is obtained, where X_t^k is a random variable with expectation

$$\mu_t^k = \mathbb{E}[X_t^k].$$

We denote the best possible expected reward at decision epoch t by μ_t^* , i.e.,

$$\mu_t^* = \max_{k \in \mathcal{K}} \{\mu_t^k\}.$$

Changes in the expected rewards of the arms. We assume the expected reward of each arm μ_t^k may change at any decision point. We denote by μ^k the sequence of expected rewards of arm k : $\mu^k = \{\mu_t^k\}_{t=1}^T$. In addition, we denote by μ the sequence of vectors of all K expected rewards: $\mu = \{\mu^k\}_{k=1}^K$. We assume that the expected reward of each arm can change an arbitrary number of times, but bound the total variation of the expected rewards:

$$\sum_{t=1}^{T-1} \sup_{k \in \mathcal{K}} |\mu_t^k - \mu_{t+1}^k|. \quad (1)$$

Let $\{V_t : t = 1, 2, \dots\}$ be a non-decreasing sequence of positive real numbers such that $V_1 = 0$, $KV_t \leq t$ for all t , and for normalization purposes set $V_2 = 2 \cdot K^{-1}$. We refer to V_T as the *variation budget* over \mathcal{T} . We define the corresponding *temporal uncertainty set*, as the set of reward vector sequences that are subject to the variation budget V_T over the set of decision epochs $\{1, \dots, T\}$:

$$\mathcal{V} = \left\{ \mu \in [0, 1]^{K \times T} : \sum_{t=1}^{T-1} \sup_{k \in \mathcal{K}} |\mu_t^k - \mu_{t+1}^k| \leq V_T \right\}.$$

The variation budget captures the constraint imposed on the non-stationary environment faced by the decision-maker. While limiting the possible evolution in the environment, it allows for many different forms in which the expected rewards may change: continuously, in discrete shocks, and of a changing rate (for illustration, Figure 1 depicts two different variation patterns that correspond to the same variation budget). In general, the variation budget V_T is designed to depend on the number of pulls T .

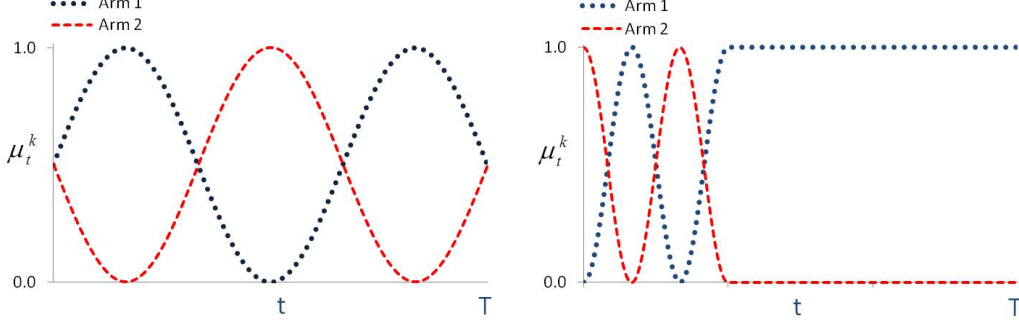


Figure 1: Two instances of variation in the expected rewards of two arms: (*Left*) Continuous variation in which a fixed variation budget (that equals 3) is spread over the whole horizon. (*Right*) “Compressed” instance in which the same variation budget is “spent” in the first third of the horizon.

Admissible policies, performance, and regret. Let U be a random variable defined over a probability space $(\mathbb{U}, \mathcal{U}, \mathbf{P}_u)$. Let $\pi_1 : \mathbb{U} \rightarrow \mathcal{K}$ and $\pi_t : [0, 1]^{t-1} \times \mathbb{U} \rightarrow \mathcal{K}$ for $t = 2, 3, \dots$ be measurable functions. With some abuse of notation we denote by $\pi_t \in \mathcal{K}$ the action at time t , that is given by

$$\pi_t = \begin{cases} \pi_1(U) & t = 1, \\ \pi_t(X_{t-1}^\pi, \dots, X_1^\pi, U) & t = 2, 3, \dots, \end{cases}$$

The mappings $\{\pi_t : t = 1, \dots, T\}$ together with the distribution \mathbf{P}_u define the class of admissible policies. We denote this class by \mathcal{P} . We further denote by $\{\mathcal{H}_t, t = 1, \dots, T\}$ the filtration associated with a policy $\pi \in \mathcal{P}$, such that $\mathcal{H}_1 = \sigma(U)$ and $\mathcal{H}_t = \sigma\left(\left\{X_j^\pi\right\}_{j=1}^{t-1}, U\right)$ for all $t \in \{2, 3, \dots\}$. Note that policies in \mathcal{P} are non-anticipating, i.e., depend only on the past history of actions and observations, and allow for randomized strategies via their dependence on U .

We define the *regret* under policy $\pi \in \mathcal{P}$ compared to a *dynamic* oracle as the worst-case difference between the expected performance of pulling at each epoch t the arm which has the highest expected reward at epoch t (the dynamic oracle performance) and the expected performance under policy π :

$$\mathcal{R}^\pi(\mathcal{V}, T) = \sup_{\mu \in \mathcal{V}} \left\{ \sum_{t=1}^T \mu_t^* - \mathbb{E}^\pi \left[\sum_{t=1}^T \mu_t^\pi \right] \right\},$$

where the expectation $\mathbb{E}^\pi[\cdot]$ is taken with respect to the noisy rewards, as well as to the policy’s actions. In addition, we denote by $\mathcal{R}^*(\mathcal{V}, T)$ the minimal worst-case regret that can be guaranteed by an admissible policy $\pi \in \mathcal{P}$:

$$\mathcal{R}^*(\mathcal{V}, T) = \inf_{\pi \in \mathcal{P}} \mathcal{R}^\pi(\mathcal{V}, T).$$

$\mathcal{R}^*(\mathcal{V}, T)$ is the best achievable performance. In the following sections we study the magnitude of $\mathcal{R}^*(\mathcal{V}, T)$. We analyze the magnitude of this quantity by establishing upper and lower bounds; in these bounds we refer to a constant C as *absolute* if it is independent of K , V_T , and T .

3 Lower bound on the best achievable performance

We next provide a lower bound on the the best achievable performance.

Theorem 1 *Assume that rewards have a Bernoulli distribution. Then, there is some absolute constant $C > 0$ such that for any policy $\pi \in \mathcal{P}$ and for any $T \geq 1$, $K \geq 2$ and $V_T \in [K^{-1}, K^{-1}T]$,*

$$\mathcal{R}^\pi(\mathcal{V}, T) \geq C (KV_T)^{1/3} T^{2/3}.$$

We note that when reward distributions are stationary, there are known policies such as UCB1 and ε -greedy (Auer et al. 2002) that achieve regret of order \sqrt{T} in the stochastic setup. When the environment is non-stationary and the reward structure is defined by the class \mathcal{V} , then no policy may achieve such a performance and the best performance must incur a regret of at least order $T^{2/3}$. This additional complexity embedded in the stochastic non-stationary MAB problem compared to the stationary one will be further discussed in §5.

Remark 1 (Growing variation budget) Theorem 1 holds when V_T is increasing with T . In particular, when the variation budget is linear in T , the regret grows linearly and long run average optimality is not achievable. This also implies the observation of Slivkins and Upfal (2008) about linear regret in an instance in which expected rewards evolve according to a Brownian motion.

The driver of the change in the best achievable performance (relative to the one established in a stationary environment) is the optimal exploration-exploitation balance. Beyond the tension between exploring different arms and capitalizing on the information already collected, captured by the “classical” exploration-exploitation trade-off, a second tradeoff is introduced by the non-stationary environment, between “remembering” and “forgetting”: estimating the expected rewards is done based on past observations of rewards. While keeping track of more observations may decrease the variance of mean rewards estimates, the non-stationary environment implies that “old” information is potentially less relevant and creates a bias that stems from possible changes in the underlying rewards. The changing rewards give incentive to dismiss old information, which in turn encourages enhanced exploration. The proof of Theorem 1 emphasizes these two tradeoffs and their impact on achievable performance. At a high level the proof of Theorem 1 builds on ideas of identifying a worst-case “strategy” of nature (e.g., Auer et al. 2002, proof of Theorem 5.1) adapting them to our setting. While the proof is deferred to the appendix, we next describe the key ideas.

Selecting a subset of feasible reward paths. We define a subset of vector sequences $\mathcal{V}' \subset \mathcal{V}$ and show that when μ is drawn randomly from \mathcal{V}' , any admissible policy must incur regret of order $(KV_T)^{1/3} T^{2/3}$. We define a partition of the decision horizon \mathcal{T} into batches $\mathcal{T}_1, \dots, \mathcal{T}_m$ of size $\tilde{\Delta}_T$ each (except, possibly the last batch):

$$\mathcal{T}_j = \left\{ t : (j-1)\tilde{\Delta}_T + 1 \leq t \leq \min \{ j\tilde{\Delta}_T, T \} \right\}, \quad \text{for all } j = 1, \dots, m, \quad (2)$$

where $m = \lceil T/\tilde{\Delta}_T \rceil$ is the number of batches. In \mathcal{V}' , in every batch there is exactly one “good” arm with expected reward $1/2 + \varepsilon$ for some $0 < \varepsilon \leq 1/4$, and all the other arms have expected reward $1/2$. The “good” arm is drawn independently in the beginning of each batch according to a discrete uniform distribution over $\{1, \dots, K\}$. Thus, the identity of the “good” arm can change only between batches. See Figure 2 for a description and a numeric example of possible realizations of a sequence μ that is randomly drawn from \mathcal{V}' . Since there are m batches we obtain a set \mathcal{V}' of K^m possible, equally probable

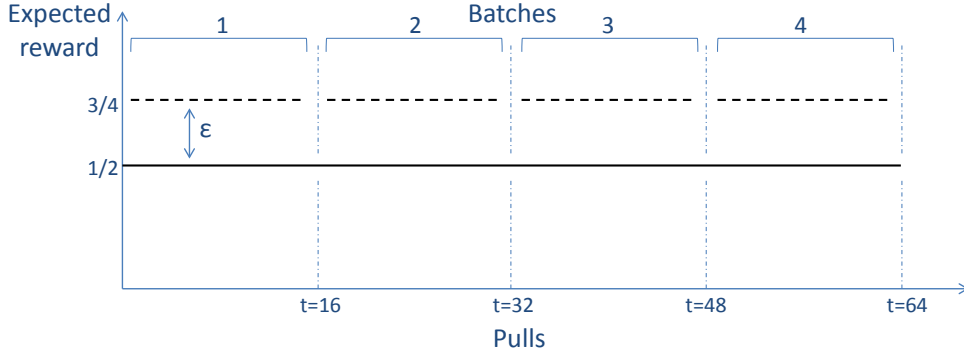


Figure 2: (Drawing a sequence from \mathcal{V}' .) A numerical example of possible realizations of expected rewards. Here $T = 64$, and we have 4 decision batches, each contains 16 pulls. We have K^4 possible realizations of reward sequences. In every batch one arm is randomly and independently drawn to have an expected reward of $1/2 + \varepsilon$, where in this example $\varepsilon = 1/4$. This example corresponds to a variation budget of $V_T = \varepsilon\tilde{\Delta}_T = 1$.

realizations of μ . By selecting ε such that $\varepsilon T/\tilde{\Delta}_T \leq V_T$, any $\mu \in \mathcal{V}'$ is composed of expected reward sequences with a variation of at most V_T , and therefore $\mathcal{V}' \subset \mathcal{V}$. Given the draws under which expected reward sequences are generated, nature prevents any accumulation of information from one batch to another, since at the beginning of each batch a new “good” arm is drawn independently of the history.

Countering possible policies. For the sake of simplicity, the discussion in this paragraph assumes a variation budget that is fixed and independent of T (the proof of the theorem details the more general treatment for a variation budget that depends on T). The proof of Theorem 1 establishes that under the setting described above, if $\varepsilon \approx 1/\sqrt{\tilde{\Delta}_T}$ no admissible policy can identify the “good” arm with high probability within a batch. Since there are $\tilde{\Delta}_T$ epochs in each batch, the regret that any policy must incur along a batch is of order $\tilde{\Delta}_T \cdot \varepsilon \approx \sqrt{\tilde{\Delta}_T}$, which yields a regret of order $\sqrt{\tilde{\Delta}_T} \cdot T/\tilde{\Delta}_T \approx T/\sqrt{\tilde{\Delta}_T}$ throughout the whole horizon. Selecting the smallest feasible $\tilde{\Delta}_T$ such that the variation budget constraint is satisfied leads to $\tilde{\Delta}_T \approx T^{2/3}$, yielding a regret of order $T^{2/3}$ throughout the horizon.

4 A near-optimal policy

In this section we apply the ideas underlying the lower bound in Theorem 1 to develop a rate optimal policy for the non-stationary MAB problem with a variation budget. Consider the following policy:

Rexp3. Inputs: a positive number γ , and a batch size Δ_T .

1. Set batch index $j = 1$
 2. Repeat while $j \leq \lceil T/\Delta_T \rceil$:
 - (a) Set $\tau = (j - 1) \Delta_T$
 - (b) Initialization: for any $k \in \mathcal{K}$ set $w_t^k = 1$
 - (c) Repeat for $t = \tau + 1, \dots, \min\{T, \tau + \Delta_T\}$:
 - For each $k \in \mathcal{K}$, set

$$p_t^k = (1 - \gamma) \frac{w_t^k}{\sum_{k'=1}^K w_t^{k'}} + \frac{\gamma}{K}$$
 - Draw an arm k' from \mathcal{K} according to the distribution $\{p_t^k\}_{k=1}^K$
 - Receive a reward $X_t^{k'}$
 - For arm k' set $\hat{X}_t^{k'} = X_t^{k'}/p_t^{k'}$, and for any $k \neq k'$ set $\hat{X}_t^k = 0$. For all $k \in \mathcal{K}$ update:

$$w_{t+1}^k = w_t^k \exp \left\{ \frac{\gamma \hat{X}_t^k}{K} \right\}$$
 - (d) Set $j = j + 1$, and return to the beginning of step 2
-

Clearly $\pi \in \mathcal{P}$. The Rexp3 policy uses Exp3, a policy introduced by Freund and Schapire (1997) for solving a worst-case sequential allocation problem, as a subroutine, restarting it every Δ_T epochs.

Theorem 2 *Let π be the Rexp3 policy with a batch size $\Delta_T = \left\lceil (K \log K)^{1/3} (T/V_T)^{2/3} \right\rceil$ and with $\gamma = \min \left\{ 1, \sqrt{\frac{K \log K}{(e-1)\Delta_T}} \right\}$. Then, there is some absolute constant \bar{C} such that for every $T \geq 1$, $K \geq 2$, and $V_T \in [K^{-1}, K^{-1}T]$:*

$$\mathcal{R}^\pi(\mathcal{V}, T) \leq \bar{C} (K \log K \cdot V_T)^{1/3} T^{2/3}.$$

Theorem 2 is obtained by establishing a connection between the regret relative to the single best action in the adversarial setting, and the regret with respect to the dynamic oracle in non-stationary stochastic setting with variation budget. Several classes of policies, such as exponential-weight policies (including Exp3) and polynomial-weight policies, have been shown to achieve regret of order \sqrt{T} with respect to the single best action in the adversarial setting (see Auer et al. (2002) and chapter 6 of Cesa-Bianchi and Lugosi (2006) for a review). While in general these policies tend to perform well numerically, there is no guarantee for its performance with respect to the dynamic oracle studied in this paper (see also Hartland et al. (2006) for a study of the empirical performance of one class of

algorithms), since the single best action itself may incur linear (with respect to T) regret relative to the dynamic oracle. The proof of Theorem 2 shows that *any* policy that achieves regret of order \sqrt{T} with respect to the single best action in the adversarial setting, can be used as a subroutine to obtain near-optimal performance with respect to the dynamic oracle in our setting.

Rexp3 emphasizes the two tradeoffs discussed in the previous section. The first tradeoff, information acquisition versus capitalizing on existing information, is captured by the subroutine policy Exp3. In fact, any policy that achieves a good performance compared to a single best action benchmark in the adversarial setting must balance exploration and exploitation, and therefore the loss incurred by experimenting on sub-optimal arms is indeed balanced with the gain of better estimation of expected rewards. The second tradeoff, “remembering” versus “forgetting,” is captured by restarting Exp3 and forgetting any acquired information every Δ_T pulls. Thus, old information that may slow down the adaptation to the changing environment is being discarded.

Taking Theorem 1 and Theorem 2 together, we have characterized the minimax regret (up to a multiplicative factor, logarithmic in the number of arms) in a full spectrum of variations V_T :

$$\mathcal{R}^*(\mathcal{V}, T) \asymp (KV_T)^{1/3} T^{2/3}.$$

Hence, we have quantified the impact of the extent of change in the environment on the best achievable performance in this broad class of problems. For example, for the case in which $V_T = C \cdot T^\beta$, for some absolute constant C and $0 \leq \beta < 1$ the best achievable regret is of order $T^{(2+\beta)/3}$.

4.1 Numerical Results

We illustrate the upper bound on the regret by a numerical experiment that measures the average regret that is incurred by Rexp3, in the presence of changing environments.

Setup. We consider instances where two arms are available: $\mathcal{K} = \{1, 2\}$. The reward X_t^k associated with arm k at epoch t has a Bernoulli distribution with a changing expectation μ_t^k :

$$X_t^k = \begin{cases} 1 & \text{w.p. } \mu_t^k \\ 0 & \text{w.p. } 1 - \mu_t^k \end{cases}$$

for all $t = 1, \dots, T$, and for any pulled arm $k \in \mathcal{K}$. The evolution patterns of μ_t^k , $k \in \mathcal{K}$ will be specified below. At each epoch $t \in \mathcal{T}$ the policy selects an arm $k \in \mathcal{K}$. Then, the binary rewards are generated, and X_t^k is observed. The pointwise regret that is incurred at epoch t is $X_t^k - X_t^{k_t^*}$, where $k_t^* = \arg \max_{k \in \mathcal{K}} \mu_t^k$. We note that while the pointwise regret at epoch t is not necessarily positive, its expectation is. Summing over the whole horizon and replicating 20,000 times for each instance of changing rewards, the average regret approximates the expected regret compared to the dynamic oracle.

First stage (Fixed variation, different time horizons). The objective of the first part of the simulation is to measure the growth rate of the average regret incurred by the policy, as a function of the horizon length, under a fixed variation budget. We use two basic instances. In the first instance (displayed on the left side of Figure 1) the expected rewards are sinusoidal:

$$\mu_t^1 = \frac{1}{2} + \frac{1}{2} \sin\left(\frac{V_T \pi t}{T}\right), \quad \mu_t^2 = \frac{1}{2} + \frac{1}{2} \sin\left(\frac{V_T \pi t}{T} + \pi\right)$$

for all $t = 1, \dots, T$. In the second instance (depicted on the right side of Figure 1) similar sinusoidal evolution of the expected reward is “compressed” into the first third of the horizon, where in the rest of the horizon the expected rewards remain fixed:

$$\mu_t^1 = \begin{cases} \frac{1}{2} + \frac{1}{2} \sin\left(\frac{3V_T \pi t}{T} + \frac{\pi}{2}\right) & \text{if } t < \frac{T}{3} \\ 0 & \text{otherwise} \end{cases} \quad \mu_t^2 = \begin{cases} \frac{1}{2} + \frac{1}{2} \sin\left(\frac{3V_T \pi t}{T} - \frac{\pi}{2}\right) & \text{if } t < \frac{T}{3} \\ 1 & \text{otherwise} \end{cases}$$

for all $t = 1, \dots, T$. Both instances describe different changing environments under the same (fixed) variation budget $V_T = 3$. While in the first instance the variation budget is spent throughout the whole horizon, in the second one the same variation budget is spent only over the first third of the horizon. For different values of T (between 3000 and 40000) and for both variation instances we estimated the regret through 20,000 replications (the average performance trajectory of Rexp3 for $T = 5000$ is depicted in the upper-left and upper-right plots of Figure 3).

Discussion of the first stage. The first part of the simulation illustrates the decision process of the policy, as well as the order $T^{2/3}$ growth rate of the regret. The upper parts of Figure 3 describe the performance trajectory of the policy. One may observe that the policy identifies the arm with the higher expected rewards, and selects it with higher probability. The Rexp3 policy adjusts to changes in the expected rewards and updates the probabilities of selecting each arm according to the received rewards. While the policy adapts quickly to the changes in the expected rewards (and in the identity of the “better” arm), it keeps experimenting with the sub-optimal arm (the policy’s trajectory doesn’t reach the one of the dynamic oracle). The Rexp3 policy balances the remembering-forgetting tradeoff using the restarting points, occurring every Δ_T epochs. The exploration-exploitation tradeoff is balanced throughout each batch by the subroutine policy Exp3. While Exp3 explores at an order of $\sqrt{\Delta_T}$ epochs in each batch, restarting it every Δ_T (V_T is fixed, therefore one has an order of $T^{1/3}$ batches, each batch with an order of $T^{2/3}$ epochs) yields an exploration rate of order $T^{2/3}$.

The lower-left and lower-right parts in Figure 3 show plots of the natural logarithm of the averaged regret as a function of the natural logarithm of the horizon length. All the standard errors of the data points in these log-log plots are lower than 0.004. These plots detail the linear dependence between the natural logarithm of the averaged regret, and the natural logarithm of T . In both cases the slope of the linear fit for increasing values of T supports the $T^{2/3}$ dependence of the minimax regret.

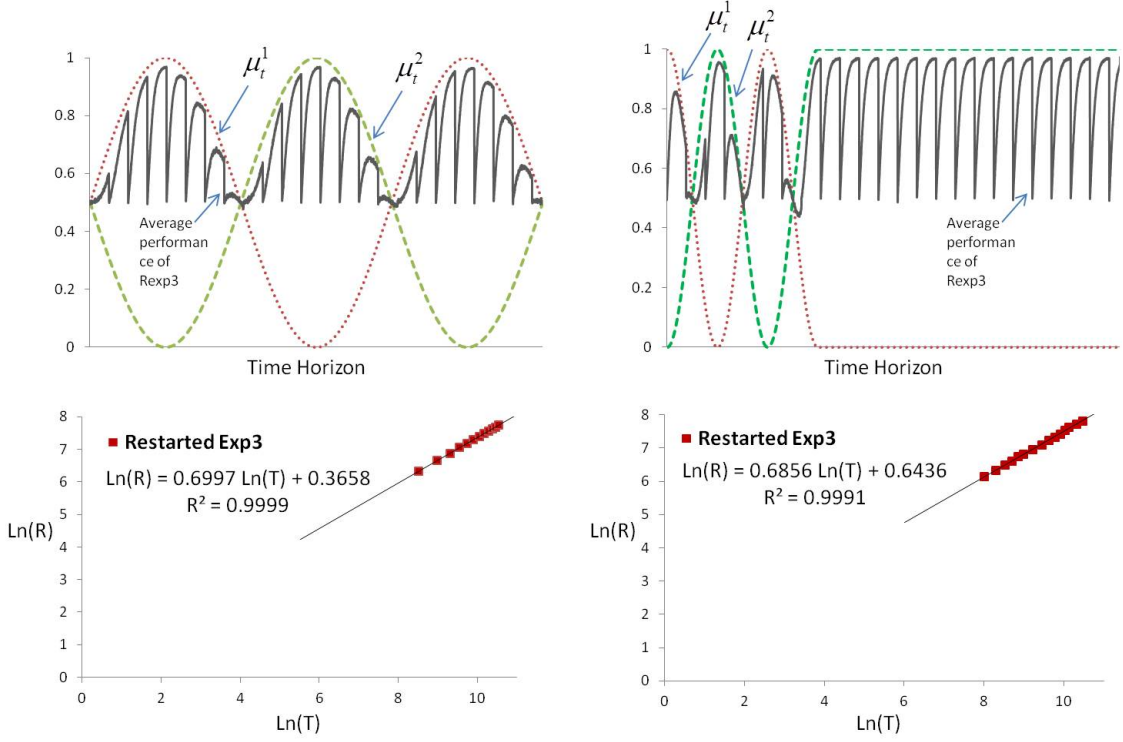


Figure 3: Numerical simulation of the performance of Rexp3, in two complementary instances: (*Upper left*) The average performance trajectory in the presence of sinusoidal expected rewards, with a fixed variation budget $V_T = 3$. (*Upper right*) The average performance trajectory under an instance in which the same variation budget is “spent” only over the first third of the horizon. In both of the instances the average performance trajectory of the policy is generated along $T = 5,000$ epochs. (*Bottom*) Log-log plots of the averaged regret as a function of the horizon length T .

Second stage (Increasing the variation). The objective of the second part of the simulation is to measure how the growth rate of the averaged regret (as a function of T) established in the first part changes when the variation increases. For this purpose we used a variation budget of the form $V_T = 3T^\beta$. Using first instance of sinusoidal variation, we repeated the first step for different values of β between 0 (implying a constant variation, that was simulated at the first stage) and 1 (implying linear variation). The upper plots of Figure 4 depicts the average performance trajectories of the Rexp3 policy under different variation budgets. The different slopes, representing different growth rate of the regret for different values of β appear in the table and the plot, at the bottom of Figure 4.

Discussion of the second stage. The second part of the simulation illustrates the way variation affects the policy decision process and the minimax regret. Since Δ_T is of order $(T/V_T)^{2/3}$, holding T fixed and increasing V_T affects the decision process and in particular the batch size of the policy. This is illustrated at the top plots of Figure 4. The slopes that were estimated for each β value (in the variation structure $V_T = 3T^\beta$) ranging from 0.1 to 0.9 describing the linear log-log dependencies (the case of $\beta = 0.0$ is already depicted at the bottom-left plot in Figure 3) are summarized in Table 1. The

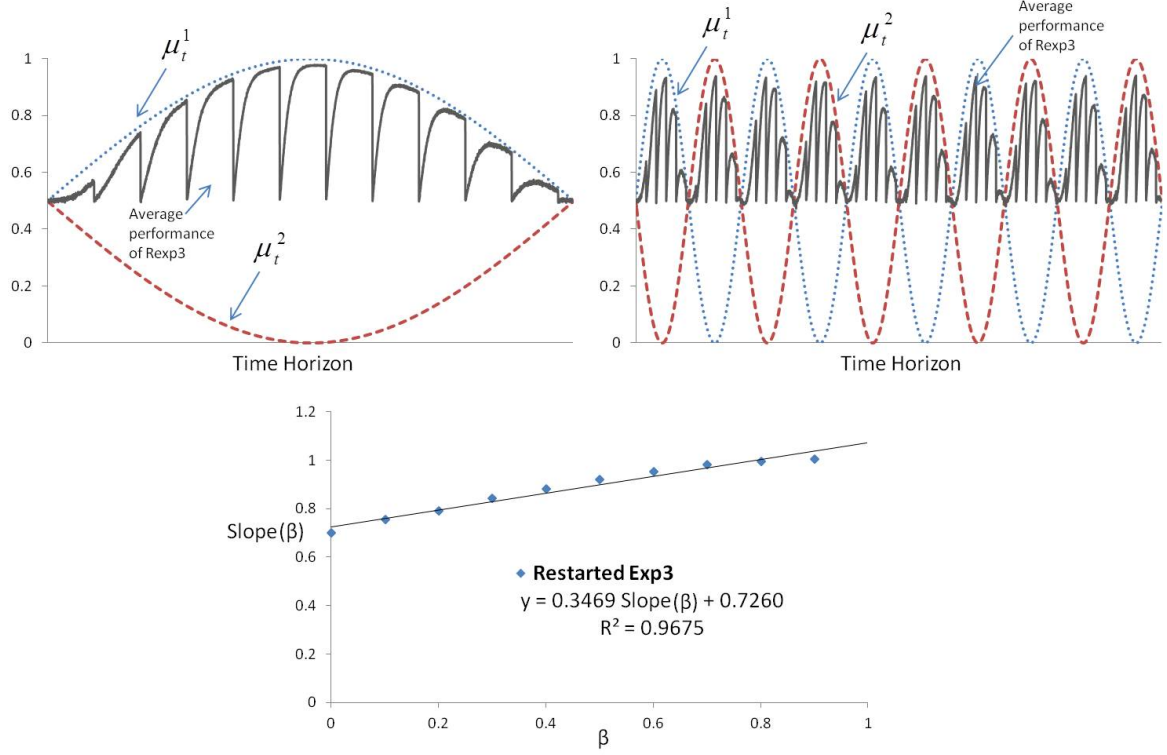


Figure 4: Variation and performance: (*Upper left*) The averaged performance trajectory for $V_T = 1$, and $T = 5000$. (*Upper right*) The averaged performance trajectory for $V_T = 10$, and $T = 5000$. (*Bottom*) The slope of the linear fit between the data points of Table 1 imply the growth rate $V_T^{1/3}$.

β value	Estimated slope
0.0	0.6997
0.1	0.7558
0.2	0.7915
0.3	0.8421
0.4	0.8801
0.5	0.9210
0.6	0.9519
0.7	0.9813
0.8	0.9942
0.9	1.0036

Table 1: **Estimated slopes for growing variation budgets.** The estimated log-log slopes obtained for different β values in the variation structure $V_T = 3T^\beta$.

bottom part of Figure 4 show the slope of the linear fit between the data points of Table 1, illustrates the growth rate of the regret when the variation (as a function of T) increases, supports the $V_T^{1/3}$ dependence of the minimax regret, and emphasizes the full spectrum of minimax regret rates (of order $V_T^{1/3}T^{2/3}$) that are obtained for different variation levels.

5 Discussion

Contrasting with traditional (stationary) MAB problems. The tight bounds that were established on the minimax regret in our stochastic non-stationary MAB problem allows one to quantify the “price of non-stationarity,” which mathematically captures the added complexity embedded in changing rewards versus stationary ones. While Theorem 1 and Theorem 2 together characterize minimax regret of order $V_T^{1/3}T^{2/3}$, the characterized minimax regret in the stationary stochastic setting is of order $\log T$ in the case where rewards are guaranteed to be “well separated” one from the other, and of order \sqrt{T} when expected rewards can be arbitrarily close to each other (see Lai and Robbins (1985) and Auer et al. (2002) for more details). Contrasting the different regret growth rates quantifies the “price,” in terms of best achievable performance, of non-stationary rewards compared to stationary ones, as a function of the variation that is allowed in the non-stationary case. Clearly, this comparison shows that additional complexity is introduced even when the allowed variation is fixed and independent of the horizon length.

Contrasting with other non-stationary MAB instances. The class of MAB problems with non-stationary rewards that is formulated in the current chapter extends other MAB formulations that allow rewards to change in a more structured manner. We already discussed in Remark 1 the consistency of our results (in the case where the variation budget grows linearly with the time horizon) with the setting treated in Slivkins and Upfal (2008) where reward evolve according to a Brownian motion and hence the regret is linear in T . Two other representative studies are those of Garivier and Moulines (2011), that study a stochastic MAB problems in which expected rewards may change a finite number of times, and Auer et al. (2002) that formulate an adversarial MAB problem in which the identity of the best arm may change a finite number of times. Both studies suggest policies that, utilizing the prior knowledge that the number of changes must be finite, achieve regret of order \sqrt{T} relative to the best sequence of actions. However, the performance of these policies can deteriorate to regret that is linear in T when the number of changes is allowed to depend on T . When there is a finite variation (V_T is fixed and independent of T) but not necessarily a finite number of changes, we establish that the best achievable performance deteriorate to regret of order $T^{2/3}$. In that respect, it is not surprising that the “hard case” used to establish the lower bound in Theorem 1 describes a nature’s strategy that allocates the allowed variation over a large (as a function of T) number of changes in the expected rewards.

A Proofs

Proof of Theorem 1. At a high level the proof adapts a general approach of identifying a worst-case nature “strategy” (see proof of Theorem 5.1 in Auer et al. (2002) which analyze the worst-case regret

relative to a single best action benchmark in a fully adversarial environment), extending these ideas appropriately to our setting. Fix $T \geq 1$, $K \geq 2$, and $V_T \in [K^{-1}, K^{-1}T]$. In what follows we restrict nature to the class $\mathcal{V}' \subseteq \mathcal{V}$ that was described in §3, and show that when μ is drawn randomly from \mathcal{V}' , any policy in \mathcal{P} must incur regret of order $(KV_T)^{1/3} T^{2/3}$.

Step 1 (Preliminaries). Define a partition of the decision horizon \mathcal{T} to $m = \left\lceil \frac{T}{\tilde{\Delta}_T} \right\rceil$ batches $\mathcal{T}_1, \dots, \mathcal{T}_m$ of size $\tilde{\Delta}_T$ each (except perhaps \mathcal{T}_m) according to (2). For some $\varepsilon > 0$ that will be specified shortly, define \mathcal{V}' to be the set of reward vectors sequences μ such that:

- $\mu_t^k \in \{1/2, 1/2 + \varepsilon\}$ for all $k \in \mathcal{K}$, $t \in \mathcal{T}$
- $\sum_{k \in \mathcal{K}} \mu_t^k = K/2 + \varepsilon$ for all $t \in \mathcal{T}$
- $\mu_t^k = \mu_{t+1}^k$ for any $(j-1)\tilde{\Delta}_T + 1 \leq t \leq \min\{j\tilde{\Delta}_T, T\} - 1$, $j = 1, \dots, m$, for all $k \in \mathcal{K}$

For each sequence in \mathcal{V}' in any epoch there is exactly one arm with expected reward $1/2 + \varepsilon$ where the rest of the arms have expected reward $1/2$, and expected rewards cannot change within a batch. Let $\varepsilon = \min\{1/4, V_T \tilde{\Delta}_T / T\}$. Then, for any $\mu \in \mathcal{V}'$ one has:

$$\sum_{t=1}^{T-1} \sup_{k \in \mathcal{K}} |\mu_t^k - \mu_{t+1}^k| \leq \sum_{j=1}^{m-1} \varepsilon = \left(\left\lceil \frac{T}{\tilde{\Delta}_T} \right\rceil - 1 \right) \cdot \varepsilon \leq \frac{T\varepsilon}{\tilde{\Delta}_T} \leq V_T,$$

where the first inequality follows from the structure of \mathcal{V}' . Therefore, $\mathcal{V}' \subset \mathcal{V}$.

Step 2 (Single batch analysis). Fix some policy $\pi \in \mathcal{P}$, and fix a batch $j \in \{1, \dots, m\}$. Let k_j denote the “good” arm of batch j . We denote by $\mathbb{P}_{k_j}^j$ the probability distribution conditioned on arm k_j being the “good” arm in batch j , and by \mathbb{P}_0 the probability distribution with respect to random rewards (i.e. expected reward $1/2$) for each arm. We further denote by $\mathbb{E}_{k_j}^j[\cdot]$ and $\mathbb{E}_0[\cdot]$ the respective expectations. Assuming binary rewards, we let X denote a vector of $|\mathcal{T}_j|$ rewards, i.e. $X \in \{0, 1\}^{|\mathcal{T}_j|}$. We denote by N_k^j the number of times arm k was selected in batch j . In the proof we use Lemma A.1 from Auer et al. (2002) that characterizes the difference between the two different expectations of some function of the observed rewards vector:

Lemma 1 *Let $f : \{0, 1\}^{|\mathcal{T}_j|} \rightarrow [0, M]$ be a bounded real function. Then, for any $k \in \mathcal{K}$:*

$$\mathbb{E}_k^j[f(X)] - \mathbb{E}_0[f(X)] \leq \frac{M}{2} \sqrt{-\mathbb{E}_0[N_k^j] \log(1 - 4\varepsilon^2)}.$$

Recalling that k_j denotes the “good” arm of batch j , one has

$$\mathbb{E}_{k_j}^j[\mu_t^\pi] = \left(\frac{1}{2} + \varepsilon \right) \mathbb{P}_{k_j}^j\{\pi_t = k_j\} + \frac{1}{2} \mathbb{P}_{k_j}^j\{\pi_t \neq k_j\} = \frac{1}{2} + \varepsilon \mathbb{P}_{k_j}^j\{\pi_t = k_j\},$$

and therefore,

$$\mathbb{E}_{k_j}^j \left[\sum_{t \in \mathcal{T}_j} \mu_t^\pi \right] = \frac{|\mathcal{T}_j|}{2} + \sum_{t \in \mathcal{T}_j} \varepsilon \mathbb{P}_{k_j}^j\{\pi_t = k_j\} = \frac{|\mathcal{T}_j|}{2} + \varepsilon \mathbb{E}_{k_j}^j[N_{k_j}^j]. \quad (3)$$

In addition, applying Lemma 1 with $f(X) = N_{k_j}^j$ (clearly $N_{k_j}^j \in \{0, \dots, |\mathcal{T}_j|\}$) we have:

$$\mathbb{E}_{k_j}^j [N_{k_j}^j] \leq \mathbb{E}_0 [N_{k_j}^j] + \frac{|\mathcal{T}_j|}{2} \sqrt{-\mathbb{E}_0 [N_{k_j}^j] \log(1 - 4\varepsilon^2)}.$$

Summing over arms, one has:

$$\begin{aligned} \sum_{k_j=1}^K \mathbb{E}_{k_j}^j [N_{k_j}^j] &\leq \sum_{k_j=1}^K \mathbb{E}_0 [N_{k_j}^j] + \sum_{k_j=1}^K \frac{|\mathcal{T}_j|}{2} \sqrt{-\mathbb{E}_0 [N_{k_j}^j] \log(1 - 4\varepsilon^2)} \\ &\stackrel{(a)}{\leq} |\mathcal{T}_j| + \frac{|\mathcal{T}_j|}{2} \sqrt{-\log(1 - 4\varepsilon^2) |\mathcal{T}_j| K} \\ &\stackrel{(b)}{\leq} \tilde{\Delta}_T + \frac{\tilde{\Delta}_T}{2} \sqrt{-\log(1 - 4\varepsilon^2) \tilde{\Delta}_T K}, \end{aligned} \quad (4)$$

for any $j \in \{1, \dots, m\}$, where: (a) holds since $\sum_{k_j=1}^K \mathbb{E}_0 [N_{k_j}^j] = |\mathcal{T}_j|$, and thus by Cauchy-Schwarz inequality $\sum_{k_j=1}^K \sqrt{\mathbb{E}_0 [N_{k_j}^j]} \leq \sqrt{|\mathcal{T}_j| K}$; and (b) holds since $|\mathcal{T}_j| \leq \tilde{\Delta}_T$ for all $j \in \{1, \dots, m\}$.

Step 3 (Regret along the horizon). Let $\tilde{\mu}$ be a random sequence of expected rewards vectors, in which in every batch the “good” arm is drawn according to an independent uniform distribution over the set \mathcal{K} . Clearly, every realization of $\tilde{\mu}$ is in \mathcal{V}' . In particular, taking expectation over $\tilde{\mu}$, one has:

$$\begin{aligned} \mathcal{R}^\pi(\mathcal{V}', T) &= \sup_{\mu \in \mathcal{V}'} \left\{ \sum_{t=1}^T \mu_t^* - \mathbb{E}^\pi \left[\sum_{t=1}^T \mu_t^\pi \right] \right\} \geq \mathbb{E}^{\tilde{\mu}} \left[\sum_{t=1}^T \tilde{\mu}_t^* - \mathbb{E}^\pi \left[\sum_{t=1}^T \tilde{\mu}_t^\pi \right] \right] \\ &\geq \sum_{j=1}^m \left(\sum_{t \in \mathcal{T}_j} \left(\frac{1}{2} + \varepsilon \right) - \frac{1}{K} \sum_{k_j=1}^K \mathbb{E}^\pi \mathbb{E}_{k_j}^j \left[\sum_{t \in \mathcal{T}_j} \tilde{\mu}_t^\pi \right] \right) \\ &\stackrel{(a)}{\geq} \sum_{j=1}^m \left(\sum_{t \in \mathcal{T}_j} \left(\frac{1}{2} + \varepsilon \right) - \frac{1}{K} \sum_{k_j=1}^K \left(\frac{|\mathcal{T}_j|}{2} + \varepsilon \mathbb{E}^\pi \mathbb{E}_{k_j}^j [N_{k_j}^j] \right) \right) \\ &\geq \sum_{j=1}^m \left(\sum_{t \in \mathcal{T}_j} \left(\frac{1}{2} + \varepsilon \right) - \frac{|\mathcal{T}_j|}{2} - \frac{\varepsilon}{K} \mathbb{E}^\pi \sum_{k_j=1}^K \mathbb{E}_{k_j}^j [N_{k_j}^j] \right) \\ &\stackrel{(b)}{\geq} \sum_{j=1}^m \left(|\mathcal{T}_j| \varepsilon - \frac{\varepsilon}{K} \left(\tilde{\Delta}_T + \frac{\tilde{\Delta}_T}{2} \sqrt{-\log(1 - 4\varepsilon^2) \tilde{\Delta}_T K} \right) \right) \\ &\stackrel{(c)}{\geq} T\varepsilon - \frac{T\varepsilon}{K} - \frac{T\varepsilon}{2K} \sqrt{-\log(1 - 4\varepsilon^2) \tilde{\Delta}_T K} \\ &\stackrel{(d)}{\geq} \frac{T\varepsilon}{2} - \frac{T\varepsilon^2}{K} \sqrt{\log(4/3) \tilde{\Delta}_T K}, \end{aligned}$$

where: (a) holds by (3); (b) holds by (4); (c) holds since $\sum_{j=1}^m |\mathcal{T}_j| = T$ and since $m \geq T/\tilde{\Delta}_T$; and (d) holds since $4\varepsilon^2 \leq 1/4$, and since $-\log(1 - x) \leq 4\log(4/3)x$ for all $x \in [0, 1/4]$, and because $K \geq 2$. Set $\tilde{\Delta}_T = \left\lceil K^{1/3} \left(\frac{T}{V_T} \right)^{2/3} \right\rceil$. Recall that $\varepsilon = \min \left\{ 1/4, V_T \tilde{\Delta}_T / T \right\}$. Suppose first that $V_T \tilde{\Delta}_T / T \leq 1/4$. Then, $\varepsilon = V_T \tilde{\Delta}_T / T \geq (KV_T / T)^{1/3}$, and one has

$$\mathcal{R}^\pi(\mathcal{V}', T) \geq \frac{1}{2} \cdot (KV_T)^{1/3} T^{2/3} - \sqrt{\log(4/3)} \cdot (KV_T)^{1/3} T^{2/3} \geq \frac{1}{8} \cdot (KV_T)^{1/3} T^{2/3}.$$

On the other hand, if $V_T \tilde{\Delta}_T / T \geq 1/4$, one has $\varepsilon = 1/4$, and therefore

$$\mathcal{R}^\pi(\mathcal{V}', T) \geq \frac{\frac{T(KV_T)^{1/3}}{4} - \frac{T^{4/3} \sqrt{\log(4/3)}}{16}}{(KV_T)^{1/3}} \geq \frac{T^{4/3}}{8(KV_T)^{1/3}} \geq \frac{1}{8} \cdot (KV_T)^{1/3} T^{2/3},$$

where the last two inequalities hold by $T \geq KV_T$. Thus, since $\mathcal{V}' \subset \mathcal{V}$, we have established that:

$$\mathcal{R}^\pi(\mathcal{V}, T) \geq \mathcal{R}^\pi(\mathcal{V}', T) \geq \frac{1}{8} \cdot (KV_T)^{1/3} T^{2/3}.$$

This concludes the proof. ■

Proof of Theorem 2 The structure of the proof is as follows. First, breaking the decision horizon to a sequence of batches of size Δ_T each, we analyze the difference in performance between the the single best action and the performance of the dynamic oracle in a single batch. Then, we plug in a known performance guarantee for Exp3 relative to the single best action in the adversarial setting, and sum over batches to establish the regret of Rexp3 with respect to the dynamic oracle.

Step 1 (Preliminaries). Fix $T \geq 1$, $K \geq 2$, and $V_T \in [K^{-1}, K^{-1}T]$. Let π be the Rexp3 policy described in §4, tuned by $\gamma = \min \left\{ 1, \sqrt{\frac{K \log K}{(e-1)\Delta_T}} \right\}$ and a batch size $\Delta_T \in \{1, \dots, T\}$ (to be specified later on). We break the horizon \mathcal{T} into a sequence of batches $\mathcal{T}_1, \dots, \mathcal{T}_m$ of size Δ_T each (except, possibly \mathcal{T}_m) according to (2). Let $\mu \in \mathcal{V}$, and fix $j \in \{1, \dots, m\}$. We decompose the regret in batch j :

$$\mathbb{E}^\pi \left[\sum_{t \in \mathcal{T}_j} (\mu_t^* - \mu_t^\pi) \right] = \underbrace{\sum_{t \in \mathcal{T}_j} \mu_t^* - \mathbb{E} \left[\max_{k \in \mathcal{K}} \left\{ \sum_{t \in \mathcal{T}_j} X_t^k \right\} \right]}_{J_{1,j}} + \underbrace{\mathbb{E} \left[\max_{k \in \mathcal{K}} \left\{ \sum_{t \in \mathcal{T}_j} X_t^k \right\} \right] - \mathbb{E}^\pi \left[\sum_{t \in \mathcal{T}_j} \mu_t^\pi \right]}_{J_{2,j}}. \quad (5)$$

The first component, $J_{1,j}$, corresponds to the expected loss associated with using a single action over the batch. The second component, $J_{2,j}$, corresponds to the expected regret with respect to the best static action in batch j .

Step 2 (Analysis of $J_{1,j}$ and $J_{2,j}$). Defining $\mu_{T+1}^k = \mu_T^k$ for all $k \in \mathcal{K}$, we denote by $V_j = \sum_{t \in \mathcal{T}_j} \max_{k \in \mathcal{K}} |\mu_{t+1}^k - \mu_t^k|$ the variation in expected rewards along batch j . We note that

$$\sum_{j=1}^m V_j = \sum_{j=1}^m \sum_{t \in \mathcal{T}_j} \max_{k \in \mathcal{K}} |\mu_{t+1}^k - \mu_t^k| \leq V_T. \quad (6)$$

Let k_0 by an arm with the best expected performance (the best static strategy) over batch \mathcal{T}_j , i.e., $k_0 \in \arg \max_{k \in \mathcal{K}} \left\{ \sum_{t \in \mathcal{T}_j} \mu_t^k \right\}$. Then,

$$\max_{k \in \mathcal{K}} \left\{ \sum_{t \in \mathcal{T}_j} \mu_t^k \right\} = \sum_{t \in \mathcal{T}_j} \mu_t^{k_0} = \mathbb{E} \left[\sum_{t \in \mathcal{T}_j} X_t^{k_0} \right] \leq \mathbb{E} \left[\max_{k \in \mathcal{K}} \left\{ \sum_{t \in \mathcal{T}_j} X_t^k \right\} \right], \quad (7)$$

and therefore, one has:

$$\begin{aligned}
J_{1,j} &= \sum_{t \in \mathcal{T}_j} \mu_t^* - \mathbb{E} \left[\max_{k \in \mathcal{K}} \left\{ \sum_{t \in \mathcal{T}_j} X_t^k \right\} \right] \stackrel{(a)}{\leq} \sum_{t \in \mathcal{T}_j} (\mu_t^* - \mu_t^{k_0}) \\
&\leq \Delta_T \max_{t \in \mathcal{T}_j} \left\{ \mu_t^* - \mu_t^{k_0} \right\} \stackrel{(b)}{\leq} 2V_j \Delta_T,
\end{aligned} \tag{8}$$

for any $\mu \in \mathcal{V}$ and $j \in \{1, \dots, m\}$, where (a) holds by (7) and (b) holds by the following argument: otherwise there is an epoch $t_0 \in \mathcal{T}_j$ for which $\mu_{t_0}^* - \mu_{t_0}^{k_0} > 2V_j$. Indeed, let $k_1 = \arg \max_{k \in \mathcal{K}} \mu_{t_0}^k$. In such case, for all $t \in \mathcal{T}_j$ one has $\mu_t^{k_1} \geq \mu_{t_0}^{k_1} - V_j > \mu_{t_0}^{k_0} + V_j \geq \mu_t^{k_0}$, since V_j is the maximal variation in batch \mathcal{T}_j . This however, implies that the expected reward of k_0 is dominated by an expected reward of another arm throughout the whole period, and contradicts the optimality of k_0 .

In addition, Corollary 3.2 in Auer et al. (2002) points out that the regret with respect to the single best action of the batch, that is incurred by Exp3 with the tuning parameter $\gamma = \min \left\{ 1, \sqrt{\frac{K \log K}{(e-1)\Delta_T}} \right\}$, is bounded by $2\sqrt{e-1}\sqrt{\Delta_T K \log K}$. Therefore, for each $j \in \{1, \dots, m\}$ one has

$$J_{2,j} = \mathbb{E} \left[\max_{k \in \mathcal{K}} \left\{ \sum_{t \in \mathcal{T}_j} X_t^k \right\} - \mathbb{E}^\pi \left[\sum_{t \in \mathcal{T}_j} \mu_t^\pi \right] \right] \stackrel{(a)}{\leq} 2\sqrt{e-1}\sqrt{\Delta_T K \log K}, \tag{9}$$

for any $\mu \in \mathcal{V}$, where (a) holds since within each batch arms are pulled according to Exp3(γ).

Step 3 (Regret throughout the horizon). Summing over $m = \lceil T/\Delta_T \rceil$ batches we have:

$$\begin{aligned}
\mathcal{R}^\pi(\mathcal{V}, T) &= \sup_{\mu \in \mathcal{V}} \left\{ \sum_{t=1}^T \mu_t^* - \mathbb{E}^\pi \left[\sum_{t=1}^T \mu_t^\pi \right] \right\} \stackrel{(a)}{\leq} \sum_{j=1}^m \left(2\sqrt{e-1}\sqrt{\Delta_T K \log K} + 2V_j \Delta_T \right) \\
&\stackrel{(b)}{\leq} \left(\frac{T}{\Delta_T} + 1 \right) \cdot 2\sqrt{e-1}\sqrt{\Delta_T K \log K} + 2\Delta_T V_T. \\
&= \frac{2\sqrt{e-1}\sqrt{K \log K} \cdot T}{\sqrt{\Delta_T}} + 2\sqrt{e-1}\sqrt{\Delta_T K \log K} + 2\Delta_T V_T.
\end{aligned}$$

where: (a) holds by (5), (8), and (9); and (b) follows from (6). Selecting $\Delta_T = \left\lceil (K \log K)^{1/3} (T/V_T)^{2/3} \right\rceil$, we establish:

$$\begin{aligned}
\mathcal{R}^\pi(\mathcal{V}, T) &\leq 2\sqrt{e-1} (K \log K \cdot V_T)^{1/3} T^{2/3} + 2\sqrt{e-1} \sqrt{\left((K \log K)^{1/3} (T/V_T)^{2/3} + 1 \right) K \log K} \\
&\quad + 2 \left((K \log K)^{1/3} (T/V_T)^{2/3} + 1 \right) V_T \\
&\stackrel{(a)}{\leq} (6\sqrt{e-1} + 4) (K \log K \cdot V_T)^{1/3} T^{2/3},
\end{aligned}$$

where (a) follows from $K \geq 2$ and $V_T \in [K^{-1}, K^{-1}T]$. This concludes the proof. \blacksquare

References

- Auer, P., N. Cesa-Bianchi, and P. Fischer (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning* 47, 235–246.
- Auer, P., N. Cesa-Bianchi, Y. Freund, and R. E. Schapire (2002). The non-stochastic multi-armed bandit problem. *SIAM journal of computing* 32, 48–77.
- Awerbuch, B. and R. D. Kleinberg (2004). Addaptive routing with end-to-end feedback: distributed learning and geometric approaches. In *Proceedings of the 36th ACM Symposium on Theory of Computing (STOC)*, 45–53.
- Bergemann, D. and U. Hege (2005). The financing of innovation: Learning and stopping. *RAND Journal of Economics* 36 (4), 719–752.
- Bergemann, D. and J. Valimaki (1996). Learning and strategic pricing. *Econometrica* 64, 1125–1149.
- Berry, D. A. and B. Fristedt (1985). *Bandit problems: sequential allocation of experiments*. Chapman and Hall.
- Bertsimas, D. and J. Nino-Mora (2000). Restless bandits, linear programming relaxations, and primal dual index heuristic. *Operations Research* 48(1), 80–90.
- Besbes, O., Y. Gur, and A. Zeevi (2013). Non-stationary stochastic optimization. *Working paper*.
- Blackwell, D. (1956). An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics* 6, 1–8.
- Caro, F. and G. Gallien (2007). Dynamic assortment with demand learning for seasonal consumer goods. *Management Science* 53, 276–292.
- Cesa-Bianchi, N. and G. Lugosi (2006). *Prediction, Learning, and Games*. Cambridge University Press, Cambridge, UK.
- Foster, D. P. and R. Vohra (1999). Regret in the on-line decision problem. *Games and Economic Behaviour* 29, 7–35.
- Freund, Y. and R. E. Schapire (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. System Sci.* 55, 119–139.
- Garivier, A. and E. Moulines (2011). On upper-confidence bound policies for switching bandit problems. In *Algorithmic Learning Theory*, pp. 174–188. Springer Berlin Heidelberg.
- Gittins, J. C. (1979). Bandit processes and dynamic allocation indices (with discussion). *Journal of the Royal Statistical Society, Series B* 41, 148–177.
- Gittins, J. C. (1989). *Multi-Armed Bandit Allocation Indices*. John Wiley and Sons.
- Gittins, J. C. and D. M. Jones (1974). *A dynamic allocation index for the sequential design of experiments*. North-Holland.
- Guha, S. and K. Munagala (2007). Approximation algorithms for partial-information based stochastic control with markovian rewards. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 483–493.
- Hannan, J. (1957). *Approximation to bayes risk in repeated plays, Contributions to the Theory of Games, Volume 3*. Princeton University Press, Cambridge, UK.
- Hartland, C., S. Gelly, N. Baskiotis, O. Teytaud, and M. Sebag (2006). Multi-armed bandit, dynamic environments and meta-bandits. *NIPS-2006 workshop, Online trading between exploration and exploitation, Whistler, Canada*.

- Kleinberg, R. D. and T. Leighton (2003). The value of knowing a demand curve: Bounds on regret for online posted-price auctions. *In Proceedings of the 44th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 594–605.
- Lai, T. L. and H. Robbins (1985). Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics* 6, 4–22.
- Pandey, S., D. Agarwal, D. Chakrabarti, and V. Josifovski (2007). Bandits for taxonomies: A model-based approach. *In SIAM International Conference on Data Mining*.
- Papadimitriou, C. H. and J. N. Tsitsiklis (1994). The complexity of optimal queueing network control. *In Structure in Complexity Theory Conference*, 318–322.
- Robbins, H. (1952). Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society* 55, 527–535.
- Slivkins, A. and E. Upfal (2008). Adapting to a changing environment: The brownian restless bandits. *In Proceedings of the 21st Annual Conference on Learning Theory (COLT)*, 343–354.
- Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 25, 285–294.
- Whittle, P. (1981). Arm acquiring bandits. *The Annals of Probability* 9, 284–292.
- Whittle, P. (1988). Restless bandits: Activity allocation in a changing world. *Journal of Applied Probability* 25A, 287–298.
- Zelen, M. (1969). Play the winner rule and the controlled clinical trials. *Journal of the American Statistical Association* 64, 131–146.