
Problem Dependent View on Structured Thresholding Bandit Problems

James Cheshire¹ Pierre Ménard¹ Alexandra Carpentier¹

Abstract

We investigate the *problem dependent regime* in the stochastic *Thresholding Bandit problem (TBP)* under several *shape constraints*. In the *TBP* the objective of the learner is to output, at the end of a sequential game, the set of arms whose means are above a given threshold. The vanilla, unstructured, case is already well studied in the literature. Taking K as the number of arms, we consider the case where (i) the sequence of arm's means $(\mu_k)_{k=1}^K$ is monotonically increasing (*MTBP*) and (ii) the case where $(\mu_k)_{k=1}^K$ is concave (*CTBP*). We consider both cases in the *problem dependent* regime and study the probability of error - i.e. the probability to mis-classify at least one arm. In the fixed budget setting, we provide upper and lower bounds for the probability of error in both the concave and monotone settings, as well as associated algorithms. In both settings the bounds match in the *problem dependent* regime up to universal constants in the exponential.

1. Introduction

Stochastic multi-armed bandit problems model situations in which a learner faces multiple unknown probability distributions, or “arms”, and has to sequentially sample these arms.

In this paper, we focus on the Thresholding Bandit Problem (*TBP*), a *Combinatorial Pure Exploration (CPE)* bandit setting introduced by [Chen et al. \(2014\)](#). The learner is presented with $[K] = \{1, \dots, K\}$ arms, each following an unknown distribution ν_k with unknown mean μ_k . We focus on the *fixed budget* variant of this problem. Given a budget $T > 0$, the learner samples the arms sequentially for a total of T times and then aims at predicting the set of arms whose mean is above a known threshold $\tau \in \mathbb{R}$. We will measure the learner's performance by the *probability*

¹Otto von Guericke University Magdeburg. Correspondence to: James Cheshire <james.cheshire@ovgu.de>.

of error - i.e. the probability that the learner mis-classifies at least one arm - and consider therefore the *problem dependent regime*.

The focus of this paper is on *structured, shape constrained TBP*. More precisely, we study the influence of some classical *structures, in the form of a shape constraint* on the *sequence of means of the arms*, on the *TBP* problem. That is, we study how classical shape constraints influence the probability of error. A related study was performed by [Cheshire et al. \(2020\)](#) for the problem independent (overall worst-case) regime, and we aim at extending this study to the *problem dependent regime*. We will aim at finding the problem dependent quantities that have an impact on the optimal probability of error, and at providing matching upper and lower bounds.

We will discuss three structured *TBPs* in this paper; among those, we recall existing results of one, and provide results for two. Here is a short overview.

Vanilla, unstructured case *TBP* The vanilla, unstructured case is the simplest *TBP* where we only assume that the distributions of the arms are sub-Gaussian - also related to the TOP-M¹ setting. The *TBP* is already well studied in the literature - both in a fixed budget and in a fixed confidence context - and we only introduce it here to provide a benchmark for later structured problems. We recall here results in the problem dependent, fixed budget, setting, which is most relevant for this paper. [Locatelli et al. \(2016\)](#) prove that up to multiplicative constants, and additives $\log(TK)$ terms, in the exponential, the optimal probability of regret in this problem is $\exp(-\frac{T}{\sum_{i:\Delta_i>0} \Delta_i^{-2}})$, where $\Delta_i = |\tau - \mu_i|$. We present their results for completeness and comparison to the bounds under additional shape constraints in Table 5.3 - see also Subsection 3.1. The *TBP* in the problem dependent regime is also studied by [Mukherjee et al. \(2017\)](#) and [Zhong et al. \(2017\)](#), however they consider a problem complexity based also upon variance making their results not so relevant to our setting. The *problem independent* regime for the *TBP* is studied by [Cheshire et al. \(2020\)](#), we also present their results in Ta-

¹In the TOP-M setting, the objective of the learner is to output the M arms with highest means. A popular version of it is the TOP-1 or “best arm identification” problem where the aim is to find the arm that realises the maximum.

ble 5.3 for comparison across the different regimes.

Monotone constraint, *MTBP*. We then consider the problem where on top of assuming that the distributions are sub-Gaussian, we assume that the sequence of means $(\mu_k)_{k \in [K]}$ is monotone - this is problem *MTBP*. This specific instance of the *TBP* is introduced within the context of drug dosing by Garivier et al. (2017). In this paper, the authors provide an algorithm for the fixed confidence setting that is optimal asymptotically, in the fixed confidence regime. However the definition of the algorithms, as well as the provided optimal error bound, are defined in an implicit way and not so easy to relate in a simple way to the gaps Δ_i moreover it is not clear how to translate a result from the fixed confidence setting to the fixed budget one. On the other hand, the shape constraint on the means of the arms implies that the *MTBP* is related to *noisy binary search*, i.e. inserting an element into its correct place within an ordered list when only noisy labels of the elements are observed, see Feige et al. (1994). They describe an algorithm structurally similar to ours, using a binary tree with infinite extension however they consider a simpler setting where the probability of correct labeling is fixed as some $\delta > \frac{1}{2}$ and go on to show that there exists an algorithm that will correctly insert an element with probability at least $1 - \delta$ in $\mathcal{O}(\log(\frac{K}{\delta}))$ steps. For further literature on the related yet different problem of noisy binary search, see Feige et al. (1994), Ben-Or & Hassidim (2008), Emamjomeh-Zadeh et al. (2016), Nowak (2011). Again, these papers consider settings with more structural assumptions than our own and are focused on the problem independent, fixed confidence regime. The *problem independent* regime for the *MTBP* is studied by Cheshire et al. (2020), we also present their results in Table 5.3 for comparison across the different regimes.

In this work, we prove that, up to universal multiplicative constants and additive $\log(K)$ terms in the exponential, the optimal error probability is $\exp(-T \min_k \Delta_k^2)$, which highlights the somewhat surprising fact that this structured monotone *TBP* problem is akin to a one armed *TBP*- see Subsection 3.2. We provide the Problem Dependent Monotone *TBP* (**PD-MTB**) algorithm that matches this bound, see Section 4.

Concave constraint, *CTBP*. We next consider the problem where on top of assuming that the distributions are sub-Gaussian, we assume that the sequence of means $(\mu_k)_{k \in [K]}$ is concave - this is problem *CTBP*. Again, in the problem independent regime the *CTBP* has been studied by Cheshire et al. (2020). In the problem dependent regime however, to the best of our knowledge, the *CTBP* has not been studied in the literature. However the related problems of estimating a concave function and optimising a concave function are well studied in the literature. Both problems are considered primarily in the continuous regime

which makes comparison to the K -armed bandit setting difficult. The problem of estimating a concave function has been thoroughly studied in the noiseless setting, and also in the noisy setting, see e.g. Simchowitz et al. (2018), where a continuous set of arms is considered, under Hölder smoothness assumptions. The problem of optimising a convex function in noise without access to its derivative - namely zeroth order noisy optimisation - has also been extensively studied. See e.g. Nemirovski & Yudin. (1983)[Chapter 9], and Wang et al. (2018); Agarwal et al. (2011); Liang et al. (2014) to name a few, all of them in a continuous setting with dimension d . The focus of this literature is however very different to ours and Cheshire et al. (2020), as the main difficulty under their assumption is to obtain a good dependence in the dimension d , and with this in mind logarithmic factors are not very relevant.

In this work, we prove that, up to universal multiplicative constants and additive $\log(K)$ terms in the exponential, the optimal error probability is $\exp(-T \min_k \Delta_k^2)$, which highlights the somewhat surprising fact that this structured concave *TBP* problem is also akin to a one armed *TBP*- see Subsection 3.3. We provide the Problem Dependent Concave *TBP* (**PD-CTB**) algorithm that matches this bound, see Section 4.

Organisation of the paper This paper is structured as follows. In Section 2 we formally introduce the *TBP* setting along with the monotone and concave shape constraints. We also describe the performance criterion - probability of error, we will be primarily using for the duration of the paper. Following this, upper and lower bounds on probability of error for all shape constraints are presented in Section 3. Descriptions of algorithms achieving said upper bounds can be found in Section 4. The results are discussed and compared to related work in Section 5. In Appendix ?? we conduct some preliminary experiments to explore how our theoretical results translate in practice. All proofs are found in the Appendix.

2. Setting

Problem formulation The learner is presented with a K -armed bandit problem $\nu = \{\nu_1, \dots, \nu_K\}$, with $K \geq 3$, where ν_k is the unknown distribution of arm k .

Let $\sigma^2 \geq 0$. We remind the learner that distribution ν of mean μ is said to be σ^2 -sub-Gaussian if for all $t \in \mathbb{R}$ we have,

$$\mathbb{E}_{X \sim \nu} [e^{t(X-\mu)}] \leq \exp\left(\frac{\sigma^2 t^2}{2}\right).$$

In particular the Gaussian distributions with variance smaller than σ^2 and the distributions with absolute values bounded by σ are σ^2 -sub-Gaussian.

Let $\mathcal{B} := \mathcal{B}(K, \sigma^2)$ be the set of all bandit problems as presented above, i.e. where the distributions ν_k of the arms

are all σ^2 sub-Gaussian.

In what follows, we assume that all $\nu \in \mathcal{B}$, and we write μ_k for the mean of arm k . Let $\tau \in \mathbb{R}$ be a fixed threshold known to the learner. We aim to devise an algorithm which classifies arms as above or below threshold τ based on their means. That is, the learner aims at finding the vector $Q \in \{-1, 1\}^K$ that encodes the true classification, i.e. $Q_k = 2\mathbb{1}_{\{\mu_k \geq \tau\}} - 1$ with the convention $Q_k = 1$ if arm k is above the threshold and $Q_k = -1$ otherwise. The *fixed budget* bandit sequential learning setting goes as follows: the learner has a budget $T > 0$ and at each round $t \leq T$, the learner pulls an arm $k_t \in [K]$ and observes a sample $Y_t \sim \nu_{k_t}$, conditionally independent from the past. After interacting with the bandit problem and expending their budget, the learner outputs a vector $\hat{Q} \in \{-1, 1\}^K$ and the aim is that it matches the unknown vector Q as well as possible.

Unstructured case TBP In the *problem dependent* regime, for $\bar{\Delta} \in \mathbb{R}_+^K$, we consider the following class of problems

$$\mathcal{B}^{\bar{\Delta}} = \{\nu \in \mathcal{B} : \forall k \in [K], |\mu_k - \tau| = \bar{\Delta}_k\}.$$

Monotone case MTBP We denote by \mathcal{B}_m the set of bandit problems,

$$\mathcal{B}_m := \{\nu \in \mathcal{B} : \mu_1 \leq \mu_2 \leq \dots \leq \mu_K\},$$

where the learner is given the additional information that the sequence of means $(\mu_k)_{k \in [K]}$ is a monotonically increasing sequence. We denote by $\Delta\mathcal{B}_m = \{\bar{\Delta} \in \mathbb{R}_+^K : \exists \nu \in \mathcal{B}_m, \forall k \in [K], |\mu_k - \tau| = \bar{\Delta}_k\}$ the set of possible vectors of gaps in \mathcal{B}_m - i.e. the set of sequences $\bar{\Delta}$ that would correspond to at least one problem in \mathcal{B}_m . In the *problem dependent* regime, for $\bar{\Delta} \in \Delta\mathcal{B}_m$, we consider the following class of problems

$$\mathcal{B}_m^{\bar{\Delta}} = \{\nu \in \mathcal{B}_m : \forall k \in [K], |\mu_k - \tau| = \bar{\Delta}_k\}.$$

Concave case CTBP We will denote by \mathcal{B}_c the set of bandit problems,

$$\mathcal{B}_c := \left\{ \nu \in \mathcal{B} : \forall 1 < k < K - 1, \frac{1}{2}\mu_{k-1} + \frac{1}{2}\mu_{k+1} \leq \mu_k \right\},$$

where the learner is given the additional information that the sequence of means $(\mu_k)_{k \in [K]}$ is concave. We denote by $\Delta\mathcal{B}_c = \{\bar{\Delta} \in \mathbb{R}_+^K : \exists \nu \in \mathcal{B}_c, \forall k \in [K], |\mu_k - \tau| = \bar{\Delta}_k, \exists l : \mu_l \geq \tau\}$ the set of possible vectors of gaps in \mathcal{B}_c where at least one arm is above threshold - i.e. the set of sequences $\bar{\Delta}$ that would correspond to at least one problem in \mathcal{B}_c where at least one arm is above threshold. In the *problem independent* regime, for $\bar{\Delta} \in \Delta\mathcal{B}_c$, we consider the following class of problems

$$\mathcal{B}_c^{\bar{\Delta}} := \left\{ \nu \in \mathcal{B}_c : \forall k < K, |\mu_k - \tau| \in \left[\frac{\bar{\Delta}_k}{2}, 3\frac{\bar{\Delta}_k}{2} \right] \right\}.$$

Remark 1. The classes of problems $\mathcal{B}^{\bar{\Delta}}, \mathcal{B}_m^{\bar{\Delta}}, \mathcal{B}_c^{\bar{\Delta}}$ contain bandit problems in resp. $\mathcal{B}, \mathcal{B}_m, \mathcal{B}_c$ that are 'local' around $\bar{\Delta}$ in the sense that while the sign of $\mu_k - \tau$ is arbitrary - although severely restricted by the shape constraint when it comes to $\mathcal{B}_m^{\bar{\Delta}}, \mathcal{B}_c^{\bar{\Delta}}$ - the gap of arm k is fixed to being - approximately, for the concave case set $\mathcal{B}_c^{\bar{\Delta}} - \bar{\Delta}_k$. This implies that in each case and on top of the respective shape constraint, we restrict ourselves to a small class of problems whose complexity is entirely characterised by $\bar{\Delta}$, in a problem dependent sense.

Strategy A strategy is a sequence of functions that maps the information gathered in the past to an arm and finally to a classification. Precisely, if we denote by I_t the information available to the player at time t , that is $I_t = \{Y_1, Y_2, \dots, Y_t\}$, with the convention $I_0 = \emptyset$. Then a strategy $\pi = ((\pi_t)_{t \in [T]}, \hat{Q}^\pi)$ is given by a sampling rule $\pi_t(I_{t-1}) = k_t \in [K]$ and a classification rule $\hat{Q}^\pi(I_T) = \hat{Q} \in \{-1, 1\}^K$.

Minimax expected regret The *problem independent*, *fixed budget* objective of the learner following the strategy π is then to minimize the expected simple regret of this classification for $\hat{Q} := \hat{Q}^\pi$:

$$r_T^{\nu, \pi} = \mathbb{E}_\nu \left[\max_{\{k \in [K] : \hat{Q}_k^\pi \neq Q_k\}} \Delta_k \right],$$

where $\Delta_k := |\tau - \mu_k|$ is the gap of arm k , and where \mathbb{E}_ν is defined as the expectation on problem ν and \mathbb{P}_ν the probability. However, the focus of this paper is on the *problem dependent* regime where, as usual, we consider as a performance criterion rather the related *probability of error*

$$e_T^{\nu, \pi} = \mathbb{P}_\nu \left(\exists k \in [K] : \hat{Q}_k^\pi \neq Q_k \right).$$

When it is clear from the context we will remove the dependence on the bandit problem ν and/or the strategy π . Note that if we denote by $\bar{\Delta}_{\min} = \min_{k \in [K]} \bar{\Delta}_k$ the minimum of the gaps then

$$r_T^{\nu, \pi} \geq \bar{\Delta}_{\min} e_T^{\nu, \pi}.$$

Consider a set of bandit problems $\tilde{\mathcal{B}} \subset \mathcal{B}$. The minimax optimal probability of error on $\tilde{\mathcal{B}}$ is then

$$e_T^*(\tilde{\mathcal{B}}) := \inf_{\pi} \sup_{\nu \in \tilde{\mathcal{B}}} e_T^{\nu, \pi}.$$

We will study this quantity over the local classes $\mathcal{B}^{\bar{\Delta}}, \mathcal{B}_m^{\bar{\Delta}}, \mathcal{B}_c^{\bar{\Delta}}$.

Remark 2. As argued above, the classes $\mathcal{B}^{\bar{\Delta}}, \mathcal{B}_m^{\bar{\Delta}}, \mathcal{B}_c^{\bar{\Delta}}$ contain only bandit problems that satisfy their respective shape constraint and whose complexity is entirely characterised by $\bar{\Delta}$, in a problem dependent sense. Studying the minimax probability of error over these very restricted classes

is therefore a very meaningful way of studying the problem dependent regime of structured TBP problems - and we expect this probability of error to heavily depend on $\bar{\Delta}$. The focus of this paper is to characterise this dependence in a tight manner.

3. Minimax rates

In this section we present upper and lower bounds on probability of error for all three shape constraints. Given a vector $\bar{\Delta} \in \mathbb{R}_+^K$ we denote $\bar{\Delta}_{\min} = \min_{k \in [K]} \bar{\Delta}_k$.

3.1. Problem dependent unstructured setting TBP

The unstructured thresholding bandit in the problem dependent regime has already been considered in the literature. We remind results from [Locatelli et al. \(2016\)](#), where they provide tight upper and lower bounds over $e_T^*(\mathcal{B}^{\bar{\Delta}})$, for any $\bar{\Delta} \in \mathbb{R}_+^K$. In our context they prove that

$$\begin{aligned} \exp\left(-\frac{3}{\sigma^2} \frac{T}{H} - 4\sigma^{-2} \log(12(\log T + 1)K)\right) &\leq e_T^*(\mathcal{B}^{\bar{\Delta}}) \\ &\leq \exp\left(-\frac{1}{64\sigma^2} \frac{T}{H} + 2 \log((\log T + 1)K)\right), \end{aligned}$$

where $H = \sum_{i: \bar{\Delta}_i > 0} 1/\bar{\Delta}_i^2$ - see Theorems 1 and 2 by [Locatelli et al. \(2016\)](#). This implies that up to multiplicative universal constants and whenever $T \geq H\sigma^2 \log(\log(T) + K)$, it holds that

$$-\log\left(e_T^*(\mathcal{B}^{\bar{\Delta}})\right) \asymp \frac{1}{\sigma^2} \frac{T}{H},$$

and upper and lower bound match up to universal multiplicative constants in the exponential of the error probability. The quantity H is therefore the problem dependent quantity that characterises the difficulty of the problem. Note that of course, the APT algorithm by [Locatelli et al. \(2016\)](#) does not take any information on the class - $\bar{\Delta}$, but also σ^2 - as parameters, and is essentially parameter free.

In this paper, we won't therefore discuss further this unstructured setting - the reminder provided here is only to be taken as a benchmark for the rest of the paper. We will on the other hand focus on the structured problems - monotone and concave and study how the minimax error probability evolves, in particular depending on the problem-dependent quantities $\bar{\Delta}$.

3.2. Problem dependent monotone setting

Given a class of problems $\mathcal{B}_m^{\bar{\Delta}}$ for some $\bar{\Delta} \in \Delta\mathcal{B}_m$, the following theorem provides a lower bound on the probability of error for any strategy π . The proof of [Theorem 3](#) can be found in [Appendix ??](#).

Theorem 3. *Let $\bar{\Delta} \in \Delta\mathcal{B}_m$. For any strategy π there exists a monotone bandit problem $\nu \in \mathcal{B}_m^{\bar{\Delta}}$ such that*

$$e_T^{\nu, \pi} \geq \frac{1}{4} \exp\left(-\frac{T\bar{\Delta}_{\min}^2}{\sigma^2}\right).$$

Now the following theorem gives an upper bound on the probability of error for the [PD-MTB](#) algorithm. The proof of [Theorem 4](#) can be found in [Appendix ??](#).

Theorem 4. *Let $\nu \in \mathcal{B}_m$ associated with arm gaps Δ , and assume that $T > 36 \log(K)$. The algorithm [PD-MTB](#) satisfies the following bound on error probability:*

$$e_T^{\nu, \text{PD-MTB}} \leq \exp\left(-c_{\text{mon}} \frac{T\Delta_{\min}^2}{\sigma^2} + c'_{\text{mon}} \log(K)\right)$$

where $c_{\text{mon}} = 1/48$ and $c'_{\text{mon}} = 12$.

The parameter free algorithm [PD-MTB](#) is described in [Sections 4](#) - see also [Appendix ??](#).

The assumption on T is reasonable as in the monotone setting it is clear no algorithm can gain enough information in less than $\log(K)$ pulls, see [Cheshire et al. \(2020\)](#). Note that combining both bounds yields that whenever $T > 36 \log(K)/\bar{\Delta}_{\min}^2$:

$$-\log\left(e_T^*(\mathcal{B}_m^{\bar{\Delta}})\right) \asymp \frac{1}{\sigma^2} T \bar{\Delta}_{\min}^2,$$

and upper and lower bound match up to universal multiplicative constants in the exponential of the error probability. Perhaps surprisingly, the number of arms plays no role in this rate - as long as we assume that $T > 36 \log(K)/\bar{\Delta}_{\min}^2$. Only the minimal arm gap appears, and this amounts to saying that when $T > 36 \log(K)/\bar{\Delta}_{\min}^2$, this problem is not more difficult - in order, up to universal multiplicative constants in the exponential - than a one-armed TBP with gap $\min_k \Delta_k$! And that in a sense, even if we knew in our monotone problem the position of all means but one - the arm with minimal gap - with respect to the threshold, the problem would not be significantly easier.

3.3. Problem dependent concave setting

Given a class of problems $\mathcal{B}_c^{\bar{\Delta}}$ for some $\bar{\Delta} \in \Delta\mathcal{B}_c$ the following theorem provides a lower bound on the probability of error for any strategy π . The proof of [Theorem 5](#) can be found in [Appendix ??](#).

Theorem 5. *Let $\bar{\Delta} \in \Delta\mathcal{B}_c$. For any strategy π there exists a problem $\nu \in \mathcal{B}_c^{\bar{\Delta}}$ such that*

$$e_T^{\nu, \pi} \geq \frac{1}{4} \exp\left(-9 \frac{T\bar{\Delta}_{\min}^2}{\sigma^2}\right).$$

Now the following theorem gives an upper bound on the probability of error for the [PD-CTB](#) algorithm. The proof of [Theorem 6](#) can be found in [Appendix ??](#).

Theorem 6. *Let $\nu \in \mathcal{B}_c$ with associated gaps Δ and assume $T > 108 \log(K)$. The algorithm [PD-CTB](#) has the following bound on error;*

$$e_T^{\nu, \text{PD-CTB}} \leq 3 \exp\left(-c_{\text{con}} \frac{T\Delta_{\min}^2}{\sigma^2} + c'_{\text{con}} \log(K)\right)$$

where $c_{\text{con}} = 1/576$ and $c'_{\text{con}} = 12$.

The parameter free algorithm **PD-CTB** is described in Sections 4 - see also Appendix ??.

The assumption on T is reasonable as in the monotone setting it is clear no algorithm can gain enough information in less than $\log(K)$ pulls, see [Cheshire et al. \(2020\)](#). Note that combining both bounds yields that whenever $T > 108 \frac{\log(K)}{\Delta_{\min}^2}$:

$$-\log\left(e_T^*(\mathcal{B}_m^{\Delta})\right) \asymp \frac{1}{\sigma^2} T \Delta_{\min}^2,$$

and upper and lower bound match up to universal multiplicative constants in the exponential of the error probability. Similar comments can be made here as in the case of the monotone *TBP* in Section 3.2: the convex *TBP* is also as difficult as a one-armed *TBP* with gap $\min_k \Delta_k$.

4. Optimal algorithms in the problem dependent regime

4.1. Monotone case *MTBP*

We assume in this section, without loss of generality, instead of considering K arms, we consider for technical reasons $K+2$ arms adding two deterministic arms 0 and $K+1$ with respective means $\mu_0 = -\infty$ and $\mu_{K+1} = +\infty$. While we assume that the distributions of the original K arms are σ^2 -sub-Gaussian the addition of two such arms will not invalidate our proofs, see Appendix ?. We do this to ensure that, after re-indexing of the arms and adapting the number of arms, $\tau \in [\mu_1, \mu_K]$.

To match a minimax rate as described in Section 3 we will utilise a modified version of the *MTB* algorithm described by [Cheshire et al. \(2020\)](#). The algorithm **PD-MTB** performs a random walk on the set of arms $[K]$ as a binary tree. We consider the binary tree as [Cheshire et al. \(2020\)](#) with an specific extension akin to that by [Feige et al. \(1994\)](#).

Binary Tree We associate to each problem $\nu \in \mathcal{B}_m$ a binary tree. Precisely we consider a binary tree with nodes of the form $v = \{L, M, R\}$ where $\{L, M, R\}$ are indexes of arms and we note respectively $v(l) = L, v(r) = R, v(m) = M$. The tree is built recursively as follows: the root is $\text{root} = \{1, \lfloor (1+K)/2 \rfloor, K\}$, and for a node $v = \{L, M, R\}$ with $L, M, R \in \{1, \dots, K\}$ the left child of v is $L(v) = \{L, M_l, M\}$ and the right child is $R(v) = \{M, M_r, R\}$ with $M_l = \lfloor (L+M)/2 \rfloor$ and $M_r = \lfloor (M+R)/2 \rfloor$ as the middle index between. The leaves of the tree will be the nodes $\{v = \{L, M, R\} : R = L+1\}$. If a node v is a leaf we set $R(v) = L(v) = \emptyset$. We consider the tree up to maximum depth $H = \lfloor \log_2(K) \rfloor + 1$. We note $P(l(v)) = P(r(v))$ the parent of the two children and let $|v|$ denote the depth of node v in the tree, with $|\text{root}| = 0$. We adopt the convention $P(\text{root}) = \text{root}$.

Extended Binary Tree We extend the above Binary tree in the following manner. For a leaf v we replace the condition $R(v) = L(v) = \emptyset$ with the following: for any leaf $v = \{L, M, R\}$ we set $R(v) = \tilde{v}$ where $\tilde{v} = \{L, M, R\}$ and set $L(v) = \emptyset$. Note that \tilde{v} is also a leaf therefore iterative application this relation will lead to an infinite extension. The result being that each leaf in our original binary tree is now the root of an infinite chain of identical nodes, see Figure 1. For practical purposes we need only consider such an extension up to depth T and can simply cut the tree at this depth.

Remark 7. We set $L(v) = \emptyset$ for some leaf v during the extension of the binary tree as by construction all leaves of the original binary tree are of the form $\{v = \{L, M, R\} : R = L+1 \text{ and } M = L\}$.

In order to predict the right classification we want to find the arm whose mean is the one just above the threshold τ . Finding this arm is equivalent to inserting the threshold into the (sorted) list of means, which can be done with a binary search in the aforementioned binary tree. But in our setting we only have access to estimates of the means which can be very unreliable if the mean is close to the threshold. Because of this there is a high chance we will make a mistake on some step of the binary search. For this reason we must allow **PD-MTB** to backtrack and this is why **PD-MTB** performs a binary search *with corrections*.

PD-MTB algorithm First, define the following integers

$$T_1 := \lceil 6 \log(K) \rceil \quad T_2 := \left\lfloor \frac{T}{3T_1} \right\rfloor. \quad (1)$$

The algorithm **PD-MTB** is then essentially a random walk on said binary tree moving one step per iteration for a total of T_1 steps. Let $v_1 = \text{root}$ and for $t < T_1$ let v_t denote the current node, the algorithm samples arms $\{v_t(j) : j \in \{l, m, r\}\}$ each T_2 times. Let the sample mean of arm $v_t(j)$ be denoted $\hat{\mu}_{j,t}$. **PD-MTB** will use these estimates to decide which node to explore next. If an error is detected - i.e. the interval between left and right-most sample mean does not contain the threshold, then the algorithm backtracks to the parent of the current node, otherwise **PD-MTB** acts as the deterministic binary search for inserting the threshold τ in the sorted list of means. More specifically, if there is an anomaly, $\tau \notin [\hat{\mu}_{l,t}, \hat{\mu}_{r,t}]$, then the next node is the parent $v_{t+1} = P(v_t)$, otherwise if $\tau \in [\hat{\mu}_{l,t}, \hat{\mu}_{m,t}]$ the the next node is the left child $v_{t+1} = L(v_t)$ and if $\tau \in [\hat{\mu}_{m,t}, \hat{\mu}_{r,t}]$ the next node is the right child $v_{t+1} = R(v_t)$. If at time t , $\tau \in [\hat{\mu}_{l,t}, \hat{\mu}_{r,t}]$ and the node v_t is a leaf, that is $v(r) = v(l) + 1$, then due to the extension of our binary tree $R(v_t) = L(v_t) = \tilde{v}_t$ where \tilde{v} is a duplicate of v_t . Hence $v_{t+1} = \tilde{v}_t$. Via this mechanism the **PD-MTB** algorithm essentially gives additional preference the the node v_t . See **PD-MTB** for details. We now formally

state the parameter free **PD-MTB** algorithm (Problem Dependent Monotone Thresholding Bandit Algorithm). We rely on the assumption $T > 36 \log(K)$, see Theorem 4 to ensure $T_2 \geq 1$.

Algorithm 1 **PD-MTB**

Initialization: $v_1 = \text{root}$
for $t = 1 : T_1$ **do**
 sample T_2 times each arm in v_t
if $\tau \notin [\hat{\mu}_{l,t}, \hat{\mu}_{r,t}]$ **then**
 $v_{t+1} = P(v_t)$
else if $\hat{\mu}_{m,t} \leq \tau \leq \hat{\mu}_{r,t}$ **then**
 $v_{t+1} = R(v_t)$
else if $\hat{\mu}_{l,t} \leq \tau \leq \hat{\mu}_{m,t}$ **then**
 $v_{t+1} = L(v_t)$
end if
end for
 Set $\hat{k} = v_{T_1+1}(r)$
return $(\hat{k}, \hat{Q}) : \hat{Q}_k = 2\mathbb{1}_{\{k \geq \hat{k}\}} - 1$

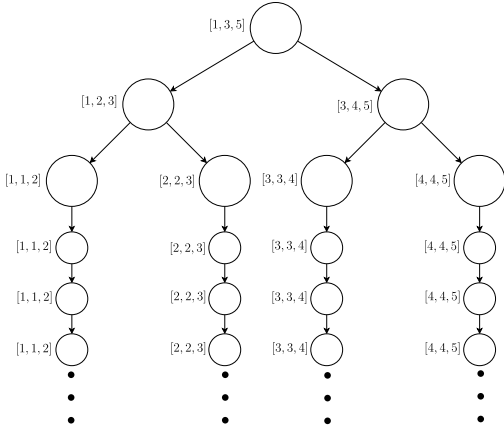


Figure 1. Extended binary tree for $K = 5$

Remark 8 (Adaptation of **PD-MTB** to a non-increasing sequence, **PD-DEC-MTB**). **PD-MTB** is applied for a monotone non-decreasing sequence $(\mu_k)_{k \in [K]}$, and it is easy to adapt it to a monotone non-increasing sequence $(\mu_k)_{k \in [K]}$. In this case, we transform the label of arm k into $K - k$, and apply **PD-MTB** to the newly labeled problem - where the mean sequence is now non-decreasing. We refer to this modification as **PD-DEC-MTB**.

Remark 9 (Relaxing the monotone assumption). By inspecting the proof of Theorem 4 in Appendix ?? we can obtain the same guarantee for a larger class of problem than one with increasing means. Indeed we only need that there exists an arm for which all the arms before it have a mean below the threshold and all arm after have a mean above the threshold. Precisely the bound of Theorem 4 holds also

for problems that belongs to

$$\mathcal{B}_{r_m} := \{\nu \in \mathcal{B} : \exists k \in [1, K], \forall j \leq k \mu_j \leq \tau, \forall j \geq k + 1 \mu_j \geq \tau\}.$$

Note the same remark also applies for problems with monotone non-increasing sequence.

4.2. Concave case CTBP

We assume in this section, without loss of generality, instead of considering K arms, we consider for technical reasons $K + 2$ arms adding two deterministic arms 0 and $K + 1$ with respective means $\mu_0 = \mu_{K+1} = -\infty$. While we assume that the distributions of the original K arms are σ^2 -sub-Gaussian the addition of two such arms will not invalidate our proofs, see Appendix ?. We do this to ensure that after re-indexing $\tau > \mu_1, \mu_K$.

As in the monotone case we construct a binary tree to span the arms of the bandit problem. The construction of this tree is identical to that described in Section 4.1 but without the infinite extension. We will use a variant off the **PD-MTB** Algorithm, **Grad-Explore** to move around the tree. The difference is that **Grad-Explore** bases its movement off the estimated gradients of the arms as opposed to their sample means. The objective of **Grad-Explore** is to find an arm with corresponding mean above threshold. Once such an arm has been identified we split our problem into two “relaxed monotone” bandit problems - see Remark 9, one increasing and one decreasing. We then run **PD-MTB** and **PD-DEC-MTB** respectively. We split our budget evenly across the three algorithms: **Grad-Explore**, **PD-MTB** and **PD-DEC-MTB**.

Grad-Explore algorithm As with **PD-MTB** the algorithm **Grad-Explore** is essentially a random walk on the said binary tree moving one step per iteration for a total of T_1 steps. Let $v_1 = \text{root}$ and for $t < T_1$ let v_t denote the current node, the algorithm samples arms $\{v_t(l), v_t(l) + 1, v_t(m), v_t(m) + 1, v_t(r), v_t(r) + 1\}$ each T_2 times. As in Section 4.1, we adopt the convention that the arm $K + 1$ is a Dirac distribution at $-\infty$. Let the sample mean of arm $v_t(j)$ be denoted $\hat{\mu}_{j,t}$ and the sample mean of arm $v_t(j) + 1$ be denoted $\hat{\mu}_{j+1,t}$. Let the estimated local gradient at arm j , that is $\hat{\mu}_{j,t} - \hat{\mu}_{j+1,t}$ denote $\hat{\nabla}_{j,t}$. **Grad-Explore** will use these estimates to decide which node to explore next. If an error is detected - i.e. the left most or right most gradient is negative or positive respectively, then the algorithm backtracks to the parent of the current node, otherwise **Grad-Explore** acts as the deterministic binary search for the maximum mean, $\max_{i \in [K]} \mu_i$. More specifically, if there is an anomaly, $(\hat{\nabla}_{l,t}, \hat{\nabla}_{r,t}) \notin (\mathbb{R}_+, \mathbb{R}_-)$, then the next node is the parent $v_{t+1} = P(v_t)$, otherwise if $\hat{\nabla}_{m,t} < 0$ the next node is the

left child $v_{t+1} = L(v_t)$ and if $\hat{\nabla}_{m,t} \geq 0$ the next node is the right child $v_{t+1} = R(v_t)$. See Algorithm 2 for details.

Algorithm 2 Grad-Explore

Initialization: $v_1 = \text{root}$
for $t = 1 : T_1$ **do**
 $S_{t+1} = S_t$
 for each $k \in v_t$ **sample** $\frac{T_2}{12}$ **times the arms** $k, k + 1$
 if $\exists k \in \{l, m, r\} : \hat{\mu}_k > \tau$ **then**
 Append arm k to the list S_{t+1}
 $v_{t+1} = v_t$
 else if $(\hat{\nabla}_{l,t}, \hat{\nabla}_{r,t}) \notin (\mathbb{R}_+, \mathbb{R}_-)$ **then**
 $v_{t+1} = P(v_t)$
 else if $\hat{\nabla}_{m,t} \geq 0$ **then**
 $v_{t+1} = R(v_t)$
 else if $\hat{\nabla}_{m,t} < 0$ **then**
 $v_{t+1} = L(v_t)$
 end if
end for

Algorithm 3 PD-CTB

run Grad-Explore
output list S_{T_1}
if $|S_{T_1}| \leq \frac{T_1}{4}$ **then**
 return $\hat{Q} = \{-1\}^K$
else
 $\hat{k} = \text{Median}(S_{T_1})$
 $l = \text{output of PD-DEC-MTB on set of arms } [1, \hat{k}] \text{ budget: } \frac{T}{3}$
 $r = \text{output of PD-MTB on set of arms } [\hat{k}, K] \text{ budget: } \frac{T}{3}$
 return $\hat{Q} : \hat{Q}_k = 1 - 2\mathbb{1}_{k < l} - 2\mathbb{1}_{k > r}$
end if

For the arms whose means are below threshold, due to the concave property gradients are essentially greater than $\bar{\Delta}_{\min}$ and can easily be estimated. Above threshold however gradients are less than $\bar{\Delta}_{\min}$ and are relatively hard to estimate. Therefore, although on the face Grad-Explore is in part a binary search for the arm with maximum mean, in reality this is not feasible. The true utility of Grad-Explore to the learner is to act as a binary search for the "set" of arms above threshold. If we refer to nodes containing an arm $k : \mu_k > \tau$ as "good nodes" the idea behind Grad-Explore is to spend a sufficient amount of time in exploring this set of nodes and adding "good arms" - i.e ones with a corresponding mean above threshold, to the list S . We can then output such an arm with high probability when outputting the median of S_{T_1} .

Once we have identified our arm above threshold we split our problem into two bandit problems where the classification can be done by binary search, see Remark 9 and 8. We can thus then apply PD-MTB and PD-DEC-MTB. Precisely, the complete procedure, namely PD-CTB (Problem

Dependent- Concave Threshold Bandits), is detailed in Algorithm 3.

5. Discussion

5.1. Algorithms PD-MTB and PD-CTB

Both the PD-MTB and PD-CTB are based upon a binary search with corrections, this allows them to exploit the structure of the shape constraints reducing the problems to sets of arms with cardinality of order $\log(K)$, something in sharp contrast to existing algorithms for the vanilla setting. The difference between PD-MTB and PD-CTB is that while PD-MTB works exclusively on a binary tree based upon the classification of an arms mean above or below threshold, the sub algorithm Grad-Explore of PD-CTB bases a binary tree on positive or negative gradient. Therefore PD-MTB acts as a search for the point the arms cross threshold while Grad-Explore acts as a search for the arm $k^* = \arg \max_k (\bar{\Delta}_k)$. Another more subtle difference is that on a "good decision" at time t - i.e when the sample means are well concentrated up to $\bar{\Delta}_{\min}$, PD-MTB will make a step in the right direction. The same cannot be said for Grad-Explore as we can only guarantee that the increments between arms are greater than $\bar{\Delta}_{\min}$ for arms below threshold, this is a direct result of the concave property. Therefore the true utility of Grad-Explore is not to find k^* but to find any arm $k : \mu_k > \tau$.

It is worth noting that both algorithms described in this paper are parameter free, being adaptive not only to the hardness of the problem characterised by the gaps $\bar{\Delta}$, but also to the underlying sub-Gaussian assumption parameter σ^2 .

5.2. Problem classes and optimality

In the monotone and concave settings we consider a very narrow class of problems and argue our classes are relevant for characterising the problem dependent regime - i.e. are narrow enough.

- In the monotone setting this is obvious as the class of problems is defined by a specific vector $\bar{\Delta} \in \mathbb{R}_+^K$, so that all problems in this class have a similar complexity, bear in mind that our algorithms do not need to know $\bar{\Delta}_{\min}$ or any aspect of $\bar{\Delta}$. In fact, when constructing our lower bound, we just need a class with two problems where, given a first problem, we simply switch the arm with minimal gap $\bar{\Delta}_{\min}$ from below to above threshold in order to obtain the second problem - see the proof of Theorem 3.
- In the concave setting this approach is unfeasible as under the concave constraints the class of problems defined by a specific vector of gaps $\bar{\Delta} \in \mathbb{R}_+^K$ has very often cardinality 1 which is nonsensical for a lower bound. Instead, given a specific vector $\bar{\Delta} \in \mathbb{R}_+^K$ we consider a class of problems with gaps within a proportional tolerance of $\bar{\Delta}$. This class is designed to be

as narrow as possible while still containing multiple problems which disagree on the placement of certain arms above or below threshold. In fact, when constructing our lower bound, we just need a class with two problems where, starting from a first problem, we simply flip the arm with minimal gap and translate other means vertically in such a way to preserve concavity - see the proof of Theorem 3.

In both cases, we prove that for T large enough, the problem dependent optimal probability of error is of order

$$\exp(-T\bar{\Delta}_{\min}^2/\sigma^2),$$

up to universal multiplicative constants inside and outside the exponential. This implies that from a problem dependent perspective, both problems are as difficult as a one armed bandit problem where we just want to decide whether the arm with minimal gap $\bar{\Delta}_{\min}$ is up or down the threshold, which is quite surprising - as the number of arms plays therefore no role asymptotically. While the lower bounds are relatively simple, the upper bounds are more interesting and challenging.

5.3. Comparison of rates between settings

Table 5.3 presents a comparison of results across the problem independent and dependent regimes. Although the results are not immediately comparable between the regimes, of particular interest is the difference in rates across the monotone and concave settings in the problem independent regime compared to the lack of difference between said rates in the problem dependent regime.

problem:	independent	dependent
Unconstrained	$\sqrt{\frac{K \log K}{T}}$	$\exp\left(-\frac{T}{H}\right)$
Monotone	$\sqrt{\frac{\log K \vee 1}{T}}$	$\exp\left(-T\bar{\Delta}_{\min}^2\right)$
Concave	$\sqrt{\frac{\log \log K \vee 1}{T}}$	$\exp\left(-T\bar{\Delta}_{\min}^2\right)$

Table 1. Order of the optimal problem dependent probability of error, and of the problem independent expected simple regret for the three structured *TBP*, in the case of all four structural assumptions on the means of the arms considered in this paper. All results are given up to universal multiplicative constants both in and outside the exponential. The first line concerns the problem independent setting and the simple regret, see Cheshire et al. (2020). The second line concerns the problem dependent setting and the probability of error, the main focus of this paper. The results for the monotone and concave are novel and can be found in this paper, see Section 3. The results for the unstructured setting are by Locatelli et al. (2016), where they take $H = \sum_{i=1}^K \bar{\Delta}_i^{-2}$

In both the monotone and concave setting an initial lower bound is one which does not depend upon K - imagine the setting in which a learner places their entire budget on

the two arms either side of the threshold. We show that in the problem dependent regime a binary search with corrections can match this bound, up to a $\log(K)$ term which disappears for large T . The intuition behind this is that as the depth of the tree is only $\log(K)$ the binary search can quickly find the point of interest and spend the majority of its time there. As both the concave and monotone problems can be solved with a binary search they therefore have the same rate.

In the problem independent regime the situation is slightly more nuanced. In terms of lower bounds one is no longer restricted to a narrow class of problems and can consider a number of different problems, all close in terms of distributional distance but nevertheless disagreeing on the classification of certain arms above or below threshold. The cardinality of these sets differs between the monotone and concave setting - being $\log(K)$ and $\log \log(K)$ respectively. This then leads to a difference in the lower bound. Upper bounds naturally must follow suit, while an adaptation of the standard binary search is still optimal in the monotone case in the concave case an algorithm using a binary search on a log scale is required. The above is by no means a rigorous explanation but hopefully gives the reader some intuition behind the differences in rates between the problem dependent and independent regimes, for more detail refer to Cheshire et al. (2020).

Acknowledgements

The work of J. Cheshire is supported by the Deutsche Forschungsgemeinschaft (DFG) DFG - 314838170, GRK 2297 MathCoRe. The work of P. Ménard is supported by the SFI Sachsen-Anhalt for the project RE-BCI ZS/2019/10/102024 by the Investitionsbank Sachsen-Anhalt. The work of A. Carpentier is partially supported by the Deutsche Forschungsgemeinschaft (DFG) Emmy Noether grant MuSyAD (CA 1488/1-1), by the DFG - 314838170, GRK 2297 MathCoRe, by the DFG GRK 2433 DAEDALUS (384950143/GRK2433), by the DFG CRC 1294 'Data Assimilation', Project A03, and by the UFA-DFH through the French-German Doktorandenkolleg CDFA 01-18 and by the UFA-DFH through the French-German Doktorandenkolleg CDFA 01-18 and by the SFI Sachsen-Anhalt for the project RE-BCI.

References

- Agarwal, A., Foster, D. P., Hsu, D. J., Kakade, S. M., and Rakhlin, A. Stochastic convex optimization with bandit feedback. In *Advances in Neural Information Processing Systems*, pp. 1035–1043, 2011.
- Ben-Or, M. and Hassidim, A. The bayesian learner is optimal for noisy binary search (and pretty good for quantum as well). In *2008 49th Annual IEEE Symposium on*

- Foundations of Computer Science*, pp. 221–230. IEEE, 2008.
- Chen, S., Lin, T., King, I., Lyu, M. R., and Chen, W. Combinatorial pure exploration of multi-armed bandits. In *Advances in Neural Information Processing Systems*, pp. 379–387, 2014.
- Cheshire, J., Menard, P., and Carpentier, A. The influence of shape constraints on the thresholding bandit problem. In Abernethy, J. and Agarwal, S. (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 1228–1275. PMLR, 09–12 Jul 2020.
- Emamjomeh-Zadeh, E., Kempe, D., and Singhal, V. Deterministic and probabilistic binary search in graphs. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pp. 519–532. ACM, 2016.
- Feige, U., Raghavan, P., Peleg, D., and Upfal, E. Computing with noisy information. *SIAM Journal on Computing*, 23(5):1001–1018, 1994.
- Garivier, A., Ménard, P., Rossi, L., and Menard, P. Thresholding bandit for dose-ranging: The impact of monotonicity. *arXiv preprint arXiv:1711.04454*, 2017.
- Liang, T., Narayanan, H., and Rakhlin, A. On zeroth-order stochastic convex optimization via random walks. *arXiv preprint arXiv:1402.2667*, 2014.
- Locatelli, A., Gutzeit, M., and Carpentier, A. An optimal algorithm for the thresholding bandit problem. In *International Conference on Machine Learning*, pp. 1690–1698. PMLR, 2016.
- Mukherjee, S., Purushothama, N. K., Sudarsanam, N., and Ravindran, B. Thresholding bandits with augmented ucb. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 2515–2521. AAAI Press, 2017.
- Nemirovski, A. and Yudin., D. Problem complexity and method efficiency in optimization. *Wiley, New York*, 1983.
- Nowak, R. D. The geometry of generalized binary search. *IEEE Transactions on Information Theory*, 57(12):7893–7906, 2011.
- Simchowitz, M., Jamieson, K., Suchow, J. W., and Griffiths, T. L. Adaptive sampling for convex regression. *arXiv preprint arXiv:1808.04523*, 2018.
- Wang, Y., Du, S., Balakrishnan, S., and Singh, A. Stochastic zeroth-order optimization in high dimensions. In Storkey, A. and Perez-Cruz, F. (eds.), *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 1356–1365. PMLR, 09–11 Apr 2018.
- Zhong, J., Huang, Y., and Liu, J. Asynchronous parallel empirical variance guided algorithms for the thresholding bandit problem. *arXiv preprint arXiv:1704.04567*, 2017.