



HAL
open science

Online Sign Identification: Minimization of the Number of Errors in Thresholding Bandits

Reda Ouhamma, Rémy Degenne, Pierre Gaillard, Vianney Perchet

► **To cite this version:**

Reda Ouhamma, Rémy Degenne, Pierre Gaillard, Vianney Perchet. Online Sign Identification: Minimization of the Number of Errors in Thresholding Bandits. NeurIPS 2021 - 35th International Conference on Neural Information Processing Systems, Dec 2021, Virtual, Canada. pp.1-25, 10.5555/3540261.3541681 . hal-03363014

HAL Id: hal-03363014

<https://inria.hal.science/hal-03363014v1>

Submitted on 4 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Online Sign Identification: Minimization of the Number of Errors in Thresholding Bandits

Reda Ouhamma

Univ. Lille, Inria, CNRS,
Centrale Lille, UMR 9189 CRIStAL,
F-59000 Lille, France
reda.ouhamma@gmail.com

Rémy Degenne

Univ. Lille, Inria, CNRS,
Centrale Lille, UMR 9189 CRIStAL,
F-59000 Lille, France
remy.degenne@inria.fr

Pierre Gaillard

Univ. Grenoble Alpes, Inria, CNRS,
Grenoble INP, LJK, 38000 Grenoble, France
pierre.gaillard@inria.fr

Vianney Perchet

Crest, Ensaë & Criteo AI Lab
vianney.perchet@normalesup.org

Abstract

In the fixed budget thresholding bandit problem, an algorithm sequentially allocates a budgeted number of samples to different distributions. It then predicts whether the mean of each distribution is larger or lower than a given threshold. We introduce a large family of algorithms (containing most existing relevant ones), inspired by the Frank-Wolfe algorithm, and provide a thorough yet generic analysis of their performance. This allowed us to construct new explicit algorithms, for a broad class of problems, whose losses are within a small constant factor of the non-adaptive oracle ones. Quite interestingly, we observed that adaptive methods empirically greatly out-perform non-adaptive oracles, an uncommon behavior in standard online learning settings, such as regret minimization. We explain this surprising phenomenon on an insightful toy problem.

1 Introduction and related work

In a stochastic multi-armed bandit problem, a decision maker sequentially samples from different distributions in order to optimize a loss that depends on the unknown parameters of those distributions. As a consequence, a tradeoff arises between gathering more samples from any possible distribution (to enhance the estimation of relevant parameters) and optimizing the allocation to minimize the final loss. We can distinguish two main categories of losses, focusing on “exploitation” vs “exploration”. The former directly depends on the whole allocation of samples and the typical example is regret minimization (we refer to the recent monographs [27, 6, 34] that cover this setting almost exhaustively). The later is a bit different; after the budget of samples is exhausted, the algorithms must answer one or several “questions” (on the different distribution) and its loss is related to the number of mistakes made; the typical application being best-arm identification and variants [3, 24].

We investigate a class of pure exploration problems, called “thresholding bandit” [28, 35]. The key property of this class is that a question is asked about each distribution, and the probability of making a mistake decreases with the total information gathered on that distribution solely. The typical question the algorithm must answer is “is the mean of the distribution above or below some threshold?” (say, 0, for simplicity); giving the wrong answer can either incur a unit cost - independently from the distribution -, or a data-dependent cost (say, the distance to the threshold that represents the “risk” of that distribution). A typical application of thresholding bandits is crowdsourcing [7] where the

objective is to distinguish workers with positive (vs. negative) efficiency; another one is bandit binary classification [18].

Some care must be taken when designing a performance criterion for a thresholding bandit problem, since any non-stupid algorithm will eventually answer all questions correctly (hence have a 0 loss) if it has enough samples. Furthermore, if distributions are sub-Gaussian (a rather mild assumption that we are going to make), the probability of making a single mistake decreases exponentially fast with the number of samples. As a consequence, the focus must be on controlling the exponential decay constant. We illustrate that issue on the unit cost problem described as follows. There are K different σ -sub-Gaussian distributions; the mean of distribution k is denoted by μ_k and the (variance-normalized) gap of distribution k to the threshold 0 is denoted by $\Delta_k := |\mu_k|/\sqrt{2\sigma^2}$. The algorithm has a budget of T samples to (sequentially) allocate to those distributions and, based on the $N_{k,T}$ samples of distribution k , it must decide the sign of μ_k ; any mistake has a cost of one. We denote by $E_k \in \{0, 1\}$ an indicator of a wrong sign prediction of μ_k after exhausting the budget of T samples. The loss is then $L_T^1 := \sum_k E_k$. It is not difficult to see that the expected number of mistakes could be of order $\sum_{k=1}^K \exp(-N_{k,T} \Delta_k^2)$.

In particular, sampling evenly across distributions ($N_{k,T} = T/K$) gives $\mathbb{E}[L_T^1] \approx \sum_k \exp(-\frac{T}{K} \Delta_k^2)$, which has an exponential decay in T . However, this uniform allocation is far from being optimal in term of the exponential decay constant. Computing an (approximate) optimal fixed allocation in hindsight is not difficult: just optimize the upper-bound of $\mathbb{E}[L_T^1]$. Since even the uniform allocation has a loss decaying exponentially, the performance of an algorithm should be measured not with respect to $\mathbb{E}[L_T^1]$ (see [24]) but rather in terms of $-\log(\mathbb{E}[L_T^1])/T$. The oracle that uses knowledge of the gaps Δ_k to optimize its fixed allocation verifies

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \log(\mathbb{E}[L_T^1]) \leq -\frac{1}{\sum_k 1/\Delta_k^2}.$$

This unit cost framework has been investigated recently [35] with a simple yet effective algorithm called LSA (Logarithmic-Sample Algorithm) designed exclusively for this problem; it samples the distribution with the smallest current index defined as $\alpha N_{k,t} \hat{\Delta}_{k,t}^2 + \log N_{k,t}$, where $\hat{\Delta}_{k,t}$ is the empirical estimate of Δ_k and α is some parameter to be chosen. LSA is "optimal up to a constant", but the constant is unfortunately in the exponential decay, as it was proved that¹

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \log(\mathbb{E}[L_T^1]) \leq -\frac{1}{16020} \frac{1}{\sum_k 1/\Delta_k^2} \quad \text{for LSA.}$$

As we shall see, this result can be drastically improved with our more refined and general analysis (that implies choosing a totally different input parameter $\alpha = 1$ instead of $1/10$ as suggested originally).

1.1 Contributions

We investigate the thresholding bandit problem with a weighted number of errors loss. Our contributions are twofold: 1) a generic method to design algorithms, with a generic proof, showing good performance on the weighted number of errors loss. 2) new lower-bounds and counter-intuitive results for the unit cost problem.

A generic algorithm with performance guarantees We propose a Frank-Wolfe inspired method to design bandit algorithms. We develop a proof technique to obtain loss bounds for the type of algorithms that our method produces, which we apply to the thresholding bandit with losses

$$L_T = \sum_{k=1}^K a_k E_k \quad \text{or} \quad L_T^\Delta = \sum_{k=1}^K \Delta_k E_k, \quad (1)$$

where $(a_k)_{k \in [K]}$ are known costs. The class of algorithms we analyze includes both LSA and APT (Anytime Parameter-free Thresholding) [35, 28]. We obtain precise non-asymptotic loss bounds for $\mathbb{E}[L_T]$; for instance, we improve the original bound of LSA by a factor 4005 (and APT by a factor 8). More importantly, we get a new algorithm whose expected error for the unit cost problem is within a

¹See Remark 1 [35]. This bound implies that LSA - with the specified choice of $\alpha = 0.1$ needs 16000 times more samples than the oracle to achieve the same performances.

factor 4 of the oracle. We emphasize again that those “constant” factors are in the exponential (and are not mere multiplicative constants).

Interestingly, this class of algorithms are *not* driven either by the “optimism under uncertainty” principle, a standard technique in multi-armed bandit [4] nor “Explore-then-commit / Successive Elimination” [31, 12].

New insights on the thresholding bandit problem First, the optimal allocation provided by the oracle of [35] in the unit cost problem has a M-shape (see Figure 1) because of two concurrent phenomena. On the one hand, the arms close to the threshold should not be pulled too much because their sign is difficult (if not impossible) to identify and it is a waste of budget. On the other hand, the signs of the arms far from the threshold are quickly well estimated and therefore should not be chosen too often either. The middle arms are the ones that need to be pulled the most frequently. As T gets larger, more and more budget is allocated to difficult arms. In section 2.2, we provide a lower-bound that shows that this M shape is actually impossible to achieve for a sequential algorithm. Typically, the hollow inside of the M shape corresponds to arms whose sign cannot be well-estimated. In particular, it is not possible to distinguish arms that are very close to the threshold from the arms that are at the top of the M and should be pulled the most frequently according to the oracle.

Our second insight is corroborated by numerical simulations in Section 4. We show empirically that our algorithms not only match but also surpass the optimal non-adaptive sampling of the oracle. We conjecture that our algorithms take advantage of the chance due to noise that can move its estimate of the arm away from the threshold. In particular, when all the gaps Δ_k are equal, the non-adaptive optimal allocation should be uniform, which is significantly outperformed by adaptive algorithms. This suggests that adaptivity is crucial for this problem and may inspire future research directions to the multi-armed bandit community in order to prove theoretical guarantees for such phenomena.

1.2 Additional related work

Zero-one loss Most of the literature on thresholding bandits [28, 30, 8] aims at minimizing the probability of making any sign error, i.e., minimizing the loss

$$L_T^* = \mathbb{I}\{\exists k \in [K], E_k = 1\} = \max_k E_k. \quad (2)$$

We already mentioned the algorithm APT [28], that gets an exponential decay of that loss (variants include variance estimation [36] and/or delayed feedbacks). Other algorithms exist, but based on the optimism principle [23, 30]. Unfortunately they suffer from a degraded exponential decay constant (by a factor bigger than 1000).

Another part of the literature focuses on the fixed confidence framework, where the objective is to answer some questions with some fixed probability of mistake (and obviously with a minimal sample budget). For instance, an objective could be to return any arm above some threshold as soon as possible [21, 9], or the one closest to the threshold [16], or just identifying that one arm is above that threshold [25], or even to control false discovery rates and variants [19, 18].

Global loss, dynamic allocation and outliers detection The loss considered in thresholding bandits can be seen as a variant of a “global loss” (i.e., essentially non-linear) that has been extensively studied in the bandit literature [1, 2, 29]. However, the major difference is, again, that the optimal allocation is time dependent and that the loss converges exponentially fast to zero (no matter the algorithm). Similarly, Frank Wolfe algorithms have been introduced in this setting [5, 13]; even though our algorithms share some similarities, they are intrinsically different for the same reasons.

Similarly, the problem investigated could be seen as a special case of bandit resource allocations [26, 7, 32, 11, 14] but where the loss is always decreasing with respect to the budget allocated per resource (hence again leading to a zero loss exponentially fast).

Finally the global objective of thresholding bandits is to obtain a synthetic view of how the means of distributions are spread on the real line (which ones are above/below some threshold). In that aspect, this problem sheds some similarities with outlier detection in multi-armed bandits [22, 38, 37].

2 Preliminaries

We describe here the weighted number of errors setting, in which an error on arm k has a known cost $a_k > 0$. The sum-of-gaps setting will be briefly investigated in section 3.3. The environment is composed of $K > 1$ arms and an algorithm sequentially pulls them. After pulling arm $k \in [K]$, it observes a sample from a distribution ν_k with mean μ_k , and that sample is independent of past observations. The distribution ν_k is supposed σ -sub-Gaussian, that is

$$\forall \lambda \in \mathbb{R} : \mathbb{E}_{X \sim \nu_k} [\exp(\lambda(X - \mu_k))] \leq \exp(\sigma^2 \lambda^2 / 2).$$

The total number of rounds (and samples) T is known in advance and called the horizon. After pulling T arms, the task of the algorithm is to classify the arms depending on whether $\mu_k > \theta$ or not, where θ is a known threshold that we conveniently set to 0 (although it could be any other value, even different from arm to arm, without significant change to the analysis). Let $s_k \in \{-1, 1\}$ be the sign of $\mu_k - \theta$, equal to 1 iff $\mu_k - \theta > 0$. The algorithm returns for all arms an estimated sign $\hat{s}_k \in \{-1, 1\}$. The objective is to minimize the expected weighted number of missclassified arms, where a mistake on arm k has a known cost $a_k > 0$,

$$L_T = \sum_{k=1}^K a_k \mathbb{I}\{\hat{s}_k \neq s_k\} = \sum_{k=1}^K a_k E_k. \quad (3)$$

Note that the linear form of the loss is quite general: since $E_k \in \{0, 1\}$, any separable loss $\sum_k f_k(E_k)$ is the sum of a constant and $\sum_k a_k E_k$ for some costs a_k .

We conclude this description of the problem with notations used in the design of algorithms. Let $N_{k,t}$ and $\hat{\mu}_{k,t} = \frac{1}{N_{k,t}} \sum_{s=1}^t \mathbb{I}\{i_s = k\} X_s$ be the number of times the learner has pulled arm k up to round t (included) and the subsequent empirical mean of arm k respectively. Define further $\hat{\Delta}_{k,t} = |\hat{\mu}_{k,t} - \theta| / \sqrt{2\sigma^2}$ and $\Delta_k = |\mu_k - \theta| / \sqrt{2\sigma^2}$, respectively the empirical and the true (variance-normalized) gap of arm k to the threshold after t rounds.

2.1 The benchmarks: a lower bound and a non-adaptive oracle

Following the proof of [35] in a slightly more generic fashion (using exponential families with one parameter instead of Bernoulli distribution), we obtain a lower bound on the performance of any algorithm (see appendix A) from which we get Theorem 1.

Theorem 1. (Similar to Theorem 20 in [35]) *Let $(\Delta_1, \dots, \Delta_K)$ be a sequence of gaps. Then for any algorithm and time horizon $T \geq K$, there exists an instance in which all arms $k \in [K]$ have Gaussian distributions with variance σ^2 and mean in $\{\Delta_k, -\Delta_k\}$ such that*

$$\mathbb{E}[L_T] \geq \frac{1}{4} \min_{\sum_k N_k = T} \sum_{k=1}^K a_k e^{-4N_k \Delta_k^2}.$$

We now deriving an optimal but unrealistic oracle, which requires prior knowledge of the gaps as input. Consider the algorithm that pulls each arm $N_{k,T}$ times, a number fixed in advance, then returns the sign of the empirical mean $\hat{\mu}_{k,T}$. Using Hoeffding's inequality, the expected loss verifies:

$$\mathbb{E}[L_T] = \sum_{k=1}^K a_k \mathbb{P}((\hat{\mu}_{k,T} - \theta)(\mu_k - \theta) < 0) \leq \sum_{k=1}^K a_k e^{-N_{k,T} \Delta_k^2} \quad (4)$$

We define the *non-adaptive oracle* as the allocation N_T which minimizes that upper bound. Its error probability has the same form as the lower bound of Theorem 1, but has a different constant in the exponential (1 instead of 4). We can solve that minimization problem and make the error bound more explicit. To that end, suppose that the arms are ordered such that $a_1 \Delta_1^2 \leq \dots \leq a_K \Delta_K^2$. There is a set $S = \{k_0, k_0 + 1, \dots, K\}$ and a constant C_S such that the oracle non-adaptive algorithm has $N_{k,T} = 0$ for $k \notin S$ and $N_{k,T} = (C_S + \log(a_k \Delta_k^2)) / \Delta_k^2$ for $k \in S$ (see appendix B for details). The expected loss of that non-adaptive oracle is

$$\mathbb{E}[L_T] \leq \sum_{k \notin S} a_k + \sum_{k \in S} a_k \exp\left(-\frac{T + \sum_{j \in S} \frac{1}{\Delta_j^2} \log\left(\frac{a_k \Delta_k^2}{a_j \Delta_j^2}\right)}{\sum_{j \in S} \frac{1}{\Delta_j^2}}\right). \quad (5)$$

The oracle is not pulling arms $1, \dots, k_0 - 1$. These are the arms which are too close to the threshold (in a distance weighted by a_k) and thus too hard to classify to be worth trying. Giving up on those arms is not something that a non-oracle algorithm can do. Figure 1 illustrates on an example ($\mu_k = (-1)^k (k/K)^2$, $k = 1, \dots, 50$, and $T = 500$) the shape of the optimal allocation (arms near the threshold should not be drawn) as well as the empirical sampling distributions of several algorithms that pull all arms. In Appendix G, we illustrate how this optimal allocation evolves with the horizon T .

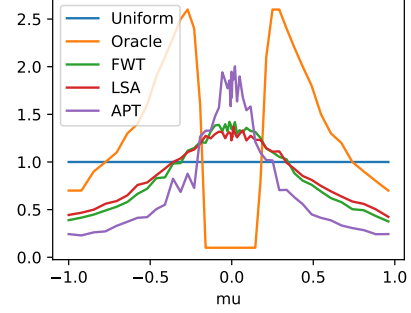


Figure 1: Optimal and empirical sampling distributions with respect to μ .

2.2 A good algorithm must pull all arms

We provide a new lower bound for the thresholding bandit with unit-cost problem, to support the claim that it is not possible to avoid pulling the arms which are close to the threshold. Consider the following 4 Gaussian bandit models (with variances 1) with means

$$\begin{aligned} \mu_{+\varepsilon} &= (\varepsilon, \dots, \varepsilon, \mu_{K_0+1}, \dots, \mu_K), & \mu'_+ &= (\mu_{K_0+1}, \dots, \mu_{K_0+1}, \mu_{K_0+1}, \dots, \mu_K), \\ \mu_{-\varepsilon} &= (-\varepsilon, \dots, -\varepsilon, \mu_{K_0+1}, \dots, \mu_K), & \mu'_- &= (-\mu_{K_0+1}, \dots, -\mu_{K_0+1}, \mu_{K_0+1}, \dots, \mu_K). \end{aligned}$$

where $0 < \varepsilon < \mu_{K_0+1} \leq \dots \leq \mu_K$, the value μ_{K_0+1} is large enough for the oracle to pull all arms on μ'_+ and $\varepsilon \leq \sqrt{\log(2)/(2T)}$.

Lemma 1. If $\mathbb{E}_{\tilde{\mu}}[L_T] \leq c_1 \min_k \sum_{N_k=T} e^{-c_0 N_k \Delta_k^2}$ for constants c_0, c_1 on $\tilde{\mu} \in \{\mu'_+, \mu'_-\}$, then

$$\max_{\mu \in \{\mu_{+\varepsilon}, \mu_{-\varepsilon}\}} \mathbb{E}_{\mu} \left[\sum_{k=1}^{K_0} N_{k,T} \right] \geq \frac{1}{2(\mu_{K_0+1} - \varepsilon)^2} \left(c_0 \frac{T + H^{\log}}{H} + \log \frac{K_0}{32c_1 H} \right).$$

where $H = \frac{K_0}{\Delta_{K_0+1}^2} + \sum_{k=K_0+1}^K \frac{1}{\Delta_k^2}$ and $H^{\log} = \frac{K_0}{\Delta_{K_0+1}^2} \log \frac{1}{\Delta_{K_0+1}^2} + \sum_{k=K_0+1}^K \frac{1}{\Delta_k^2} \log \frac{1}{\Delta_k^2}$.

The proof is postponed to Appendix A. In a few words, if an algorithm has an expected loss close to the loss of the non-adaptive oracle, then it must pull linearly the arms which are close to the threshold.

3 Algorithm and upper-bound

We introduce and analyse a new class of algorithms for the thresholding bandit problem that we call *index-based* algorithms. That class unifies several existing algorithms, including APT [28] and LSA [35]. An index-based algorithm pulls the minimum of K quantities, one for each arm, that each depends only on the rewards and pull counts of the respective arm (it does not change when pulling other arms). In particular, we consider algorithms for which the sampled arm is $i_{t+1} \in \arg \min_{k \in [K]} F(N_{k,t}, N_{k,t} \hat{\Delta}_{k,t}^2; a_k)$ for a function $F: \mathbb{N} \times \mathbb{R}_+ \times \mathbb{R}_+^* \rightarrow \mathbb{R}$ that depends on the pull counts, the information about the sign and the weight of the arm.

After T rounds, the algorithm recommends the sign of the arms at the round $t_{\max} \in [T]$ at which $\min_{k \in [K]} F(N_{k,t}, N_{k,t} \hat{\Delta}_{k,t}^2; a_k)$ was maximal. This rule is used as opposed to returning the sign of all arms at time T to facilitate the analysis, which is based on the observation that there is a small probability of error when all arms have high index. The time t_{\max} should be close to T : in particular, only one arm is sampled (possibly several times) between t_{\max} and T (see Appendix C). In Sec. 3.2, we provide a generic analysis for index-based algorithms satisfying the assumption below.

Assumption 1. The index function $F(n, x; a): \mathbb{N} \times \mathbb{R}_+ \times \mathbb{R}_+^* \rightarrow \mathbb{R}$ is non-decreasing in n and x and $\lim_{n \rightarrow +\infty} F(n, ny; a) = +\infty$ for all $y > 0, a > 0$.

Intuitively, algorithms that verify Assumption 1 prefer pulling arms that were pulled the least (smallest n) and whose quantity of information about the sign ($n \hat{\Delta}_{k,n}^2$) is small. This class includes several algorithms from the thresholding bandits literature: APT [28] for $F(n, x; a_k) = x$ and LSA [35] for $F(n, x; a_k) = x + \log(n)$ (these algorithms are only defined for $a_k = 1$). We now propose a generic method for designing an index-based algorithm.

Algorithm 1 Index-based algorithm for thresholding bandit

Inputs: an index function $F : \mathbb{N} \times \mathbb{R}_+ \times \mathbb{R}_+^* \rightarrow \mathbb{R}$; $a_1, \dots, a_K \in \mathbb{R}_+^*$; $\sigma > 0$; and $\theta \in \mathbb{R}$

For $t = 1, \dots, T$ do

- for all $k \in [K]$ define

$$N_{k,t-1} = \sum_{s=1}^{t-1} \mathbb{I}\{k = i_s\}, \hat{\mu}_{k,t-1} = \frac{1}{N_{k,t-1}} \sum_{s=1}^{t-1} \mathbb{I}\{k = i_s\} X_s, \text{ and } \hat{\Delta}_{k,t-1}^2 = \frac{1}{2\sigma^2} (\hat{\mu}_{k,t-1} - \theta)^2$$

- pull $i_t \in \arg \min_{k \in [K]} F(N_{k,t-1}, N_{k,t-1} \hat{\Delta}_{k,t-1}^2; a_k)$.
- observe $X_t \sim \nu_{i_t}$

Define $t_{\max} = \max_{t \in [T]} \min_{k \in [K]} F(N_{k,t}, N_{k,t} \hat{\Delta}_{k,t}^2; a_k)$
Return for each $k \in [K]$ the sign $\hat{s}_k = \text{sign}(\hat{\mu}_{k,t_{\max}} - \theta)$

3.1 Frank-Wolfe for Thresholding bandits

Our strategy to minimize the expected loss is inspired by the Frank-Wolfe algorithm [15] and aims at controlling an upper-bound on the loss, such as the right hand side of Inequality (4). Let's write that function as $B(N_T) = \sum_{k=1}^K a_k e^{-N_{k,T} \Delta_k^2}$. The high-level idea is to sequentially estimate its gradient and move to the minimizer of its linear approximation. If the gaps were known, we could compute at time $t + 1$ the gradient of the bound with respect to N_t , $\nabla B(N_t) = (-a_k \Delta_k^2 e^{-N_{k,t} \Delta_k^2})_k$ and use the Frank-Wolfe algorithm. The algorithm would pull $i_{t+1} \in \arg \min_u u^\top \nabla B(N_t)$ for u in the simplex, which is simply $\arg \min_{k \in [K]} (-a_k \Delta_k^2 e^{-N_{k,t} \Delta_k^2})$. The gaps are however unknown. We therefore compute an estimate of the gaps $\hat{\Delta}_{k,t}$, with which we form the estimated gradient

$$\hat{\nabla} B(N_t)_k = -a_k \hat{\Delta}_k^2 e^{-N_{k,t} \hat{\Delta}_k^2} = -\exp\left(-\left(N_{k,t} \hat{\Delta}_k^2 - \log(N_{k,t} \hat{\Delta}_k^2) + \log\left(\frac{N_{k,t}}{a_k}\right)\right)\right).$$

This gives a natural choice for the index function of our algorithm $F(n, x; a_k) = x - \log x + \log(n/a_k)$. However, the latter is decreasing in x for $x \in (0, 1)$, which in addition to violating Assumption 1, may lead to instability in the initial phase when the gaps Δ_k are poorly estimated by $\hat{\Delta}_{k,n}$. We therefore propose a slight modification that preserves the asymptotic behavior of F and we call the resulting algorithm FWT (Frank-Wolfe for Thresholding bandits):

$$F(n, x; a_k) = \max\{x, 1\} - \log(\max\{x, 1\}) + \log(n/a_k). \quad (\text{FWT})$$

Recovering APT Using different upper-bounds B on the expected loss may lead to different algorithms. In particular, we highlight a link between our Frank-Wolfe inspired method and the APT algorithm of [28], which was designed to minimize the loss

$$L_T = \sum_{k=1}^K a_k \mathbb{I}\{\hat{s}_k \neq s_k\} = \sum_{k=1}^K a_k E_k.$$

Following our method with the choice $B(N_t) = \max_{k \in [K]} e^{-N_{k,t} \Delta_k^2}$ results in exactly the same sampling rule as the one of the APT algorithm (the recommendation rule differs slightly since we recommend the sign at t_{\max} and not at T). Indeed, the derivative of B with respect to $N_{k,t}$ is nonzero (and negative) if and only if $N_{k,t} \Delta_k^2 = \arg \min_j N_{j,t} \Delta_j^2$ (ignoring the case in which there are several argmins, for which the tie breaking can be arbitrary). This leads to the choice $F(n, x; a_k) = x$ in Algorithm 1, which then pulls $i_{t+1} = \arg \min_{k \in [K]} N_{k,t} \Delta_k^2$. This is the sampling rule of APT.

3.2 Loss upper bound

We provide a loss upper bound that is valid for all index-based algorithms that verify Assumption 1. We then give a compact summary of the analysis outline and the resulting loss bounds.

Theorem 2. Let $K \geq 1$, $a_1, \dots, a_K > 0$, $T \geq 1$, and $\sigma > 0$. Let $F : \mathbb{N} \times \mathbb{R} \times \mathbb{R}_+^* \rightarrow \mathbb{R}$ that satisfies Assumption 1. Let $C_1, \dots, C_K > \max_k F(0, 0; a_k)$. For all $j, k \in [K]$, define

- $t_j(C_k)$ a solution of the equation $F(t, t\Delta_j^2; a_j) = C_k$,
- $S_k \subseteq [K]$ and $t_{j,0}(C_k) \in \mathbb{R}_+$, a set and values such that for $i \notin S_k$, $\mathbb{P}(\exists n \leq t_{i,0}(C_k), F(n, n\hat{\Delta}_{n,i}^2; a_i) \geq C) = 1$.

Then the expected loss of Algorithm 1 is upper-bounded as

$$\mathbb{E}[L_T^A] \leq \sum_{k=1}^K a_k \left(e \cdot \exp \left(-\frac{\frac{1}{2} \left(T - \sum_{j \notin S_k} t_{j,0}(C_k) \right) - \sum_{j \in S_k} t_j(C_k)}{\sum_{j \in S_k} 1/\Delta_j^2} \right) + T \cdot e^{-t_k(C_k)\Delta_k^2} \right).$$

Refer to Appendix D for the proof. It is composed of two parts:

1. First we establish that for any arm $j \in [K]$, with large probability, there is a time $\tau_j(C_k)$ such that $F(\tau_j(C_k), \tau_j(C_k)\hat{\Delta}_{\tau_j(C_k),j}^2; a_j) \geq C_k$. We prove that for all $j, k \in [K]$, $\tau_j(C_k)$ has an exponential tail then use the fact that the algorithm pulls the minimal index to control the probability that the minimum never reaches C_k .
2. We show that if an arm's index is large, then the probability of mistake on it is small.

The times $t_j(C_k)$ of Theorem 2 are the smallest numbers of samples such that $t_j(C_k) \geq \tau_j(C_k)$ with high enough probability. By determining those times, we derive explicit bounds for algorithms that verify Assumption 1. In particular we derive a bound for the variant of APT which returns the sign at the time t_{\max} when the minimal index was maximal.

Corollary 1. Suppose that for all $k \in [K]$, $a_k = 1$. For all $T \in \mathbb{N}^*$,

$$\mathbb{E}[L_T^{\text{APT}}] \leq 2K\sqrt{e \cdot T} \cdot \exp \left(-\frac{1}{4} \frac{T}{\sum_{j=1}^K 1/\Delta_j^2} \right).$$

Refer to Appendix D.3 for the proof. Since $\max_k E_k \leq \sum_k E_k$, the bound of Corollary 1 is also a bound on the zero-one loss, which we can compare to the result of [28]. Our result shows a 1/4 factor in the exponential instead of the worse 1/32 constant of the original paper.

LSA and FWT Theorem 2 applies to LSA and FWT with the following times:

- LSA: $t_j(C_k) = W(e^{C_k} \Delta_j^2) / \Delta_j^2$ and $t_{j,0}(C_k) = e^{C_k}$,
- FWT: $t_j(C_k) = \log(e^{C_k} a_j \Delta_j^2) / \Delta_j^2$ and $t_{j,0}(C_k) = a_j e^{C_k - 1}$,

where W is the Lambert W function, which verifies $|W(x) - (\log x - \log \log x)| \leq \log(1 + 1/e)$ for $x \geq e$. Therefore, for the two algorithms, the times $t_j(C_k)$ are close (equal up to the $\log \log$ terms in W), thus their bounds are close as well. Note that LSA is only defined for $a_j = 1$ for all j . In contrast to LSA, our bound for FWT has the notable property that, in the regime where $T \geq 2 \sum_{j=1}^K \frac{1}{\Delta_j^2} (2 + \log \frac{a_j \Delta_j^2 \max_i a_i \Delta_i^2}{(\min_k a_k \Delta_k^2)^2} - \log \frac{T}{e})$, we recover the same exponent as in the non-adaptive oracle loss bound (5) (up to a factor 1/4). Indeed we show that for such T

$$\mathbb{E}[L_T^{\text{FWT}}] \leq 2\sqrt{eT} \sum_{k=1}^K a_k \exp \left(-\frac{1}{4} \frac{T + 2 \sum_{j=1}^K \frac{1}{\Delta_j^2} \log \frac{a_k \Delta_k^2}{a_j \Delta_j^2}}{\sum_{j=1}^K 1/\Delta_j^2} \right). \quad (6)$$

In the same regime of large T , the bound that we obtain for LSA is of the same order, but less explicit due to the function W . The latter is still impressive since the original theorem of [35] for LSA exhibits an exponent significantly looser, of order $\exp \left(-\frac{1}{16020} \frac{T}{\sum_{j=1}^K 1/\Delta_j^2} \right)$, i.e. 4005-times worse than our bound. We finally derive a bound for our newly introduced algorithm.

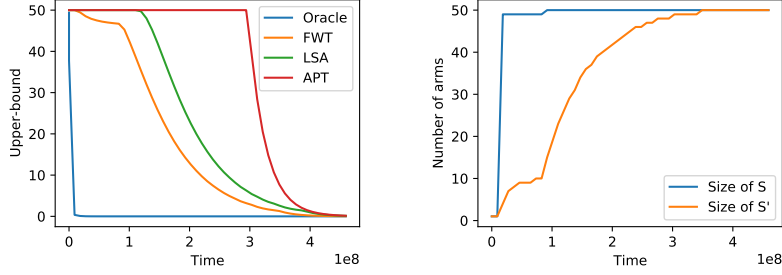


Figure 2: [left] Comparison of the upper-bounds of Corollaries 1, 2 and that of the optimal non-adaptive oracle of Eq. (5) (blue) when the gaps are of the form $\Delta_i = (i/K)^2$. [right] Evolution over time of the size of the optimal sets S (blue) and S' (orange) that minimize the bound of Corollary 2.

Corollary 2. Let S, S' be two sets with $S' \subseteq S \subseteq [K]$ and let $C \in \mathbb{R}$ be such that $C \geq 1 + \max_{k \in S} \log \frac{1}{a_k \Delta_k^2}$. Then, for all $T \geq 1$

$$\begin{aligned} \mathbb{E}[L_T^{\text{FWT}}] &\leq \sum_{k \notin S'} a_k + e \sum_{k \in S'} a_k \exp \left(- \frac{\frac{1}{2}(T - \sum_{j \notin S} a_j e^{C-1}) + \sum_{j \in S} \frac{1}{\Delta_j^2} \log \frac{1}{a_j \Delta_j^2}}{\sum_{j \in S} 1/\Delta_j^2} + C \right) \\ &\quad + T \sum_{k \in S'} a_k \exp(-C + \log(1/(a_k \Delta_k^2))) . \end{aligned}$$

Figure 2 compares the upper-bounds of Corollary 1 (APT), Theorem 2 (see also Equation (9) in the Appendix) (LSA), and Corollary 2 (FWT) for the particular case $\Delta_i = (i/K)^2$ and $a_i = 1$, for $i = 1, \dots, K = 50$. See also Figure 4 in the supplementary material for $\Delta_i = i/K$. We can see that while the bounds of LSA, APT, and FWT are asymptotically similar, that of FWT starts to be significant for much smaller T . On the right, we can see the importance of the set S' in Corollary 2: the bound first ignores all the arms, and suffers a loss of 1 and then adds them one by one as soon as they can be classified. The bound derived in [35] for LSA is not represented on the figures, since it is still bigger than K for the considered range of T .

3.3 The sum-of-gaps objective

We show that our method applies for the sum-of-gaps objective $\sum_{k=1}^K \Delta_k E_k$. This is not a particular case of the setting discussed previously since a_k was known to the algorithm, while Δ_k is unknown. It serves as a proof of concept for the extensibility of our method. The index given by FWT in this setting is $F(n, x) = x' - \frac{3}{2} \log(x') + \frac{3}{2} \log(n)$, where $x' = \max(x, \frac{3}{2})$. We can then bound the sum-of-gaps loss using our generic analysis by proceeding similarly to Theorem (2).

Corollary 3. (FWT for the sum-of-gaps objective) In the regime where $T \geq 2 \sum_{j=1}^k \frac{1}{\Delta_j^2} \left(3 + 3 \log \frac{\Delta_j \max_i \Delta_i}{(\min_i \Delta_i)^2} - \log \frac{T}{e} \right)$, we show that

$$\mathbb{E} \left[\sum_{k=1}^K \Delta_k E_k \right] \leq 2\sqrt{eT} \sum_k \Delta_k \exp \left(- \frac{1}{2} \frac{\frac{T}{2} + \sum_j \frac{3}{2} \frac{1}{\Delta_j^2} \log \frac{\Delta_k^2}{\Delta_j^2}}{\sum_j 1/\Delta_j^2} \right) .$$

See Appendix E for the proof and for a different bound that is valid for all times T . This can be useful for applications in which errors are more tolerated for arms that are close to the threshold.

4 Beating the oracle? The benefits of adaptivity.

We argue that in some situations adaptive algorithms can greatly outperform the non-adaptive oracle of Section 2.1, i.e., the cost of non-adaptivity can be much higher than the cost of learning. The algorithms in the family we considered are all adaptive in the sense that they adapt their drawing strategy as more information is observed, at the cost of learning the parameter μ_k . We illustrate the benefits of adaptivity in the following toy example.

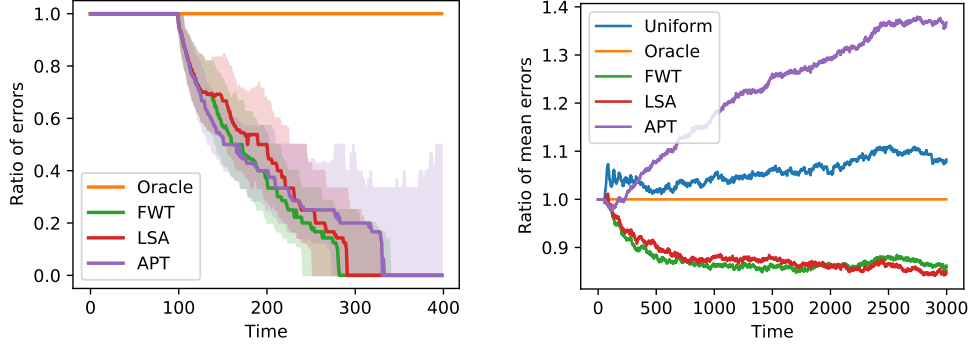


Figure 3: [left] Median (and 0.25, 0.75 empirical quantiles obtained on 500 runs) of the ratio between the error suffered by each algorithm and that of the optimal non-adaptive oracle ($\mu_k = (-1)^k, k = 1, \dots, 100$). [right] Ratio of the averaged errors (over 500 runs) of each algorithm with that of the oracle ($\mu_k = (-1)^k (k/K)^2, k = 1, \dots, 50$).

The “optimal” non-adaptive algorithm may be worse than adaptive algorithms. Consider the following parametric problem. An arm distribution is parametrized by $x \in \mathbb{R}$ and is supported on $\{0, x\}$; a sample of that distribution is equal to 0 or x , each with probability $1/2$. We assume that all arms have non-zero parameter and we will compute the optimal non-adaptive allocation.

We make the convention that if an algorithm sees only zeros for one arm, it returns any sign with probability $1/2$. The error probability of a non-adaptive allocation N_T^k for arm k is half of the probability of seeing only zeros (since if anything else is observed, the arm can be classified with perfect accuracy). Hence the total error is

$$\mathbb{E}[L_T] = \frac{1}{2} \sum_{k=1}^K \frac{1}{2^{N_{k,T}}} \geq \frac{K}{2^{(T/K)+1}},$$

which is minimized with the uniform allocation: $N_{k,T} = \frac{T}{K}$ for all $k \in [K]$.

Consider now an adaptive procedure that sample each arm in turn, but stops sampling an arm as soon as it sees a non-zero value. We crudely prove an upper bound for its number of errors, by remarking that it is zero if the algorithm classifies all arms correctly and smaller than K otherwise. The number of samples required to perfectly classify an arm follows a geometric distribution with parameter $1/2$. As a consequence, the number of required samples to classify all arms correctly follows a negative binomial $\text{NB}(K, 1/2)$. Let Z be such a negative binomial random variable. The expected number of errors of the adaptive procedure is up to $K\mathbb{P}(Z > T)$. It then verifies

$$\mathbb{E}[L_T] \leq K\mathbb{P}(Z > T) \leq K e^{-(\log(2)/2)T} \mathbb{E} e^{(\log(2)/2)Z} = \frac{K}{2^{T/2}} \left(1 + \frac{1}{\sqrt{2}}\right)^K,$$

where the value $\log(2)/2$ is chosen for simplicity (in $[0, \log 2]$). In the regime where T is large, this is of order $1/2^{T/2}$, which for $K > 2$ is much smaller than $1/2^{T/K}$ for the uniform allocation.

This toy example differs drastically from more realistic situations, as one non-zero sample for an arm is sufficient to know the sign of the expectation perfectly. We therefore consider empirically more reasonable frameworks, closer to those analyzed in the paper: the distributions of K arms are either $\mathcal{N}(1, 1)$ or $\mathcal{N}(-1, 1)$. Since all gaps Δ_i are equal, the optimal non-adaptive oracle is also the uniform sampling. The results are illustrated on the left part of Figure 3 and highlight the fact that all the adaptive algorithms considered (APT, LSA or FWT) drastically outperform the oracle. The right part of the figure shows the same phenomenon on another example in which the gaps are not constant. In particular, we can see that FWT and LSA have similar performance while APT (not designed for this purpose) generally suffers from a larger error. This result was corroborated by most of our experiments. We refer to Appendix G for more details.

Discussion

An interesting research direction is to consider objective functions more general than (1). In particular, we believe that our approach can be generalized to losses of the form $L_T = \sum_{k=1}^K f(\Delta_k, E_k)$ under certain regularity assumptions on f . Moreover, we focused on separable losses (hence linear wlog) and the index based algorithms we analyze reflect that separability. An obvious and intriguing direction for further work is to replace that assumption. One might for example want to design an algorithm that minimizes the probability of making more than a given number of mistakes.

The fact that adaptive algorithms can beat non-adaptive oracles has already been observed empirically for fixed confidence identification [33, 10], although only in cases where the non-adaptive oracle was worse only for small times and was still asymptotically optimal. The phenomenon we observe for fixed budget thresholding is much more significant and remains to be explained by theoretical arguments. Currently, the best theoretical bound for adaptive algorithms is still a factor $1/4$ away in the exponent from the non-adaptive oracle bound.

Acknowledgments and Disclosure of Funding

V. Perchet acknowledges support from the French National Research Agency (ANR) under grant number #ANR-19-CE23-0026 as well as the support grant, as well as from the grant "Investissements d'Avenir" (LabEx Ecodec/ANR-11-LABX-0047)". R. Ouhamma also acknowledges support from Ecole polytechnique under the AMX funding. P. Gaillard and R. Degenne were supported by the French government under management of Agence Nationale de la Recherche as part of the "Investissements d'avenir" program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute).

References

- [1] Agarwal, A., Foster, D. P., Hsu, D. J., Kakade, S. M., and Rakhlin, A. Stochastic convex optimization with bandit feedback. In Shawe-Taylor, J., Zemel, R. S., Bartlett, P. L., Pereira, F., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 24*, pp. 1035–1043. Curran Associates, Inc., 2011.
- [2] Agrawal, S. and Devanur, N. R. Bandits with concave rewards and convex knapsacks. In *Proceedings of the Fifteenth ACM Conference on Economics and Computation, EC '14*, pp. 989–1006, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2565-3.
- [3] Audibert, J.-Y., Bubeck, S., and Munos, R. Best arm identification in multi-armed bandits. In *COLT*, pp. 41–53, 2010.
- [4] Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- [5] Berthet, Q. and Perchet, V. Fast rates for bandit optimization with upper-confidence frank-wolfe. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 2225–2234. 2017.
- [6] Bubeck, S., Cesa-Bianchi, N., et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- [7] Chen, X., Lin, Q., and Zhou, D. Statistical decision making for optimal budget allocation in crowd labeling. *The Journal of Machine Learning Research*, 16(1):1–46, 2015.
- [8] Cheshire, J., Menard, P., and Carpentier, A. The influence of shape constraints on the thresholding bandit problem. In *Conference on Learning Theory*, pp. 1228–1275. PMLR, 2020.
- [9] Degenne, R. and Koolen, W. M. Pure exploration with multiple correct answers. *arXiv preprint arXiv:1902.03475*, 2019.
- [10] Degenne, R., Koolen, W. M., and Ménard, P. Non-asymptotic pure exploration by solving games. *arXiv preprint arXiv:1906.10431*, 2019.

- [11] Devanur, N. R., Jain, K., Sivan, B., and Wilkens, C. A. Near optimal online algorithms and fast approximation algorithms for resource allocation problems. *Journal of the ACM (JACM)*, 66(1): 7, 2019.
- [12] Even-Dar, E., Mannor, S., Mansour, Y., and Mahadevan, S. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7(6), 2006.
- [13] Fontaine, X., Berthet, Q., and Perchet, V. Regularized contextual bandits. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2144–2153. PMLR, 2019.
- [14] Fontaine, X., Mannor, S., and Perchet, V. An adaptive stochastic optimization algorithm for resource allocation. In *Algorithmic Learning Theory*, pp. 319–363. PMLR, 2020.
- [15] Frank, M. and Wolfe, P. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.
- [16] Garivier, A., Ménard, P., Rossi, L., and Menard, P. Thresholding bandit for dose-ranging: The impact of monotonicity. *arXiv preprint arXiv:1711.04454*, 2017.
- [17] Hoorfar, A. and Hassani, M. Inequalities on the lambert w function and hyperpower function. *J. Inequal. Pure and Appl. Math*, 9(2):5–9, 2008.
- [18] Jain, L. and Jamieson, K. G. A new perspective on pool-based active classification and false-discovery control. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [19] Jamieson, K. and Jain, L. A bandit approach to multiple testing with false discovery control. *arXiv preprint arXiv:1809.02235*, 2018.
- [20] Janson, S. Tail bounds for sums of geometric and exponential variables. *Statistics & Probability Letters*, 135:1–6, 2018.
- [21] Kano, H., Honda, J., Sakamaki, K., Matsuura, K., Nakamura, A., and Sugiyama, M. Good arm identification via bandit feedback. *Machine Learning*, 108(5):721–745, 2019.
- [22] Katariya, S., Tripathy, A., and Nowak, R. Maxgap bandit: Adaptive algorithms for approximate ranking. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [23] Katz-Samuels, J. and Scott, C. Feasible arm identification. In *International Conference on Machine Learning*, pp. 2535–2543. PMLR, 2018.
- [24] Kaufmann, E., Cappé, O., and Garivier, A. On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.
- [25] Kaufmann, E., Koolen, W., and Garivier, A. Sequential test for the lowest mean: From thompson to murphy sampling. *arXiv preprint arXiv:1806.00973*, 2018.
- [26] Koopman, B. O. The optimum distribution of effort. *Journal of the Operations Research Society of America*, 1(2):52–63, 1953.
- [27] Lattimore, T. and Szepesvári, C. *Bandit algorithms*. Cambridge University Press, 2020.
- [28] 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.
- [29] Mannor, S., Perchet, V., and Stoltz, G. Approachability in unknown games: Online learning meets multi-objective optimization. In *Conference on Learning Theory*, pp. 339–355. PMLR, 2014.
- [30] Mukherjee, S., Naveen, K. P., Sudarsanam, N., and Ravindran, B. Thresholding bandits with augmented ucb. *arXiv preprint arXiv:1704.02281*, 2017.

- [31] Perchet, V., Rigollet, P., Chassang, S., Snowberg, E., et al. Batched bandit problems. *Annals of Statistics*, 44(2):660–681, 2016.
- [32] Salehi, M. A., Smith, J., Maciejewski, A. A., Siegel, H. J., Chong, E. K., Apodaca, J., Briceno, L. D., Renner, T., Shestak, V., Ladd, J., et al. Stochastic-based robust dynamic resource allocation for independent tasks in a heterogeneous computing system. *Journal of Parallel and Distributed Computing*, 97:96–111, 2016.
- [33] Simchowitz, M., Jamieson, K., and Recht, B. The simulator: Understanding adaptive sampling in the moderate-confidence regime. In *Conference on Learning Theory*, pp. 1794–1834. PMLR, 2017.
- [34] Slivkins, A. et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286, 2019.
- [35] Tao, C., Blanco, S., Peng, J., and Zhou, Y. Thresholding bandit with optimal aggregate regret. In *Advances in Neural Information Processing Systems*, pp. 11664–11673, 2019.
- [36] Zhong, J., Huang, Y., and Liu, J. Asynchronous parallel empirical variance guided algorithms for the thresholding bandit problem. *arXiv preprint arXiv:1704.04567*, 2017.
- [37] Zhu, Y., Katariya, S., and Nowak, R. Robust outlier arm identification. In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 11566–11575. PMLR, 13–18 Jul 2020.
- [38] Zhuang, H., Wang, C., and Wang, Y. Identifying outlier arms in multi-armed bandit. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 5210–5219, 2017.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A] It is a purely theoretical paper
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See Appendix G.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the description of each figure.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix G.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]

- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Lower Bounds

We follow the method of [35]. Let \hat{s}_k be the estimated sign of μ . The expected loss on problem μ is

$$\mathbb{E}[L_T(\mu)] = \sum_{k=1}^K a_k \mathbb{P}_\mu \{\hat{s}_k \neq s_k\}.$$

For each arm $k \in [K]$, we define two values $\mu_k, \tilde{\mu}_k \in \mathbb{R}$, with $\mu_k < \theta < \tilde{\mu}_k$. Let $\mu = (\mu_k)_{k \in [K]}$. For some fixed one-parameter exponential family, we denote by $\text{KL}(a, b)$ the Kullback-Leibler divergence between distributions with mean a and b . We recall the for Gaussians with variance σ^2 , $\text{KL}(a, b) = \frac{(a-b)^2}{2\sigma^2}$.

Theorem 3. *For any algorithm,*

$$\sup_{S \in \mathcal{P}([K])} \mathbb{E}[L_T(\mu_S)] \geq \frac{1}{4} \min_{N: \sum_k N_k = T} \sum_{k=1}^K a_k \exp(-N_k \max\{\text{KL}(\mu_k, \tilde{\mu}_k), \text{KL}(\tilde{\mu}_k, \mu_k)\})$$

In particular, for Gaussians with variance σ^2 and $\tilde{\mu}^k = \theta + (\theta - \mu^k)$,

$$\sup_{S \in \mathcal{P}([K])} \mathbb{E}[L_T(\mu_S)] \geq \frac{1}{4} \min_{N: \sum_k N_k = T} \sum_{k=1}^K a_k \exp(-4N_k \Delta_k^2).$$

Proof. Given a vector $\lambda \in \mathbb{R}^K$ with $\lambda_k \in \{\mu_k, \tilde{\mu}_k\}$ for all $k \in [K]$ and $S \subseteq [K]$, let λ_S be such that $\lambda_{k,S} \in \{\mu_k, \tilde{\mu}_k\}$ and $\lambda_{k,S} \neq \lambda_k$ for $k \in S$ and $\lambda_{j,S} = \lambda_j$ for $j \notin S$.

For $S \in \mathcal{P}([K])$, let $S \pm i$ be equal to $S \cup \{i\}$ if $i \notin S$ and to $S \setminus \{i\}$ otherwise. Also, we denote by $(s_k(\lambda))$ be the signs of (λ_k) . Then the following holds

$$\begin{aligned} \sup_{S \in \mathcal{P}([K])} \mathbb{E}[L_T(\mu_S)] &\geq \frac{1}{2^K} \sum_{S \in \mathcal{P}([K])} \mathbb{E}[L_T(\mu_S)] \\ &= \frac{1}{2^K} \sum_{S \in \mathcal{P}([K])} \sum_{k=1}^K a_k \mathbb{P}_{\mu_S} \{\hat{s}_k \neq s_k(\mu_S)\} \\ &= \frac{1}{2^{K+1}} \sum_{S \in \mathcal{P}([K])} \sum_{k=1}^K a_k \mathbb{P}_{\mu_S} \{\hat{s}_k \neq s_k(\mu_S)\} + a_k \mathbb{P}_{\mu_{S \pm k}} \{\hat{s}_k \neq s_k(\mu_{S \pm k})\} \\ &= \frac{1}{2^{K+1}} \sum_{S \in \mathcal{P}([K])} \sum_{k=1}^K a_k \mathbb{P}_{\mu_S} \{\hat{s}_k \neq s_k(\mu_S)\} + a_k \mathbb{P}_{\mu_{S \pm k}} \{\hat{s}_k = s_k(\mu_S)\}. \end{aligned}$$

For each arm k , we can bound the sum of the two probabilities from below. Let $\mathcal{E}_{k,S} = \{\hat{s}_k \neq s_k(\mu_S)\}$.

$$\mathbb{P}_{\mu_S}(\mathcal{E}_{k,S}) + \mathbb{P}_{\mu_{S \pm k}}(\overline{\mathcal{E}_{k,S}}) \geq \frac{1}{2} \exp(-\mathbb{E}_{\mu_S}[N_{k,T}] \text{KL}(\mu_{k,S}, \mu_{k,S \pm k})),$$

so that, when plugged back in the previous equation, we get

$$\begin{aligned} \sup_{S \in \mathcal{P}([K])} \mathbb{E}[L_T(\mu_S)] &\geq \frac{1}{2^{K+1}} \sum_S \sum_{k=1}^K \frac{1}{2} a_k \exp(-\mathbb{E}_{\mu_S}[N_{k,T}] \text{KL}(\mu_{k,S}, \mu_{k,S \pm k})) \\ &\geq \frac{1}{4} \frac{1}{2^K} \sum_S \min_{N: \sum_k N_k = T} \sum_{k=1}^K a_k \exp(-N_k \text{KL}(\mu_{k,S}, \mu_{k,S \pm k})) \\ &\geq \frac{1}{4} \min_{N: \sum_k N_k = T} \sum_{k=1}^K a_k \exp(-N_k \max\{\text{KL}(\mu_k, \tilde{\mu}_k), \text{KL}(\tilde{\mu}_k, \mu_k)\}) \end{aligned}$$

□

A.1 Lower bound on the number of pulls of arms close to zero

Consider a Gaussian bandit model (with variances $\sigma = 1$) with vector of means $\mu_{+\varepsilon} = (\varepsilon, \dots, \varepsilon, \mu_{K_0+1}, \dots, \mu_K)$, in which arms $1, \dots, K_0$ have mean $\varepsilon > 0$ and arms $K - K_0 + 1, \dots, K$ have mean greater than ε . Let μ'_+ be equal to $\mu_{+\varepsilon}$ except that $\mu'_{+,j} = \mu_{K_0+1}$ for $j \in [K_0]$. We suppose that μ_{K_0+1} is large enough for the non-adaptive oracle to pull all arms on μ'_+ . We also define $\mu_{-\varepsilon} = (-\varepsilon, \dots, -\varepsilon, \mu_{K_0+1}, \dots, \mu_K)$ and $\mu'_- = (-\mu_{K_0+1}, \dots, -\mu_{K_0+1}, \mu_{K_0+1}, \dots, \mu_K)$.

Lemma 2. If an algorithm verifies $\mathbb{E}_{\tilde{\mu}}[L_T] \leq c_1 \min_{N: \sum_k N_k = T} \sum_k e^{-c_0 N_k \tilde{\Delta}_k^2}$ for constants c_0, c_1 on all Gaussian problems with variance 1, for all mean vectors $\tilde{\mu}$ with gaps $\tilde{\Delta}$, then for $\varepsilon \leq \sqrt{\log(2)/(2T)}$,

$$\max_{\mu \in \{\mu_{+\varepsilon}, \mu_{-\varepsilon}\}} \mathbb{E}_{\mu} \left[\sum_{k=1}^{K_0} N_{k,T} \right] \geq \frac{1}{2(\mu_{K_0+1} - \varepsilon)^2} \left(c_0 \frac{T + H^{\log}}{H} + \log \frac{K_0}{32c_1 H} \right).$$

where $H = \frac{K_0}{\Delta_{K_0+1}^2} + \sum_{k=K_0+1}^K \frac{1}{\Delta_k^2}$ and $H^{\log} = \frac{K_0}{\Delta_{K_0+1}^2} \log \frac{1}{\Delta_{K_0+1}^2} + \sum_{k=K_0+1}^K \frac{1}{\Delta_k^2} \log \frac{1}{\Delta_k^2}$.

Proof. We will prove that the number of pulls of arms $1, \dots, K_0$ cannot be too small. Formally, let $n_{\varepsilon} = \sum_{k=1}^{K_0} \mathbb{E}_{\mu_{+\varepsilon}} N_{k,T}$ be the expected number of pulls under $\mu_{+\varepsilon}$ of the arms with mean ε . We aim at showing that that number cannot be zero. We first prove that

$$\mathbb{P}_{\mu_{+\varepsilon}}(L_T > K_0/2) \geq \frac{1}{4}. \quad (7)$$

This follows from the basic inequalities,

$$\mathbb{P}_{\mu_{+\varepsilon}}(L_T(\mu_{+\varepsilon}) > K_0/2) + \mathbb{P}_{\mu_{-\varepsilon}}(L_T(\mu_{+\varepsilon}) \leq K_0/2) \geq \frac{1}{2} e^{-n_{\varepsilon} \text{KL}(\varepsilon, -\varepsilon)} \geq \frac{1}{2} e^{-2\varepsilon^2 T}$$

In particular, for $\varepsilon \leq \sqrt{\frac{\log 2}{2T}}$, $\max\{\mathbb{P}_{\mu_{+\varepsilon}}(L_T(\mu_{+\varepsilon}) > K_0/2), \mathbb{P}_{\mu_{-\varepsilon}}(L_T(\mu_{-\varepsilon}) > K_0/2)\} \geq \frac{1}{4}$, and either Inequality 7 either holds for $\mu_{+\varepsilon}$, or we just need to switch the role of ε and $-\varepsilon$ in this proof. Suppose now that we inequality 7 holds for $\mu_{+\varepsilon}$.

The Kullback-Leibler divergence between $\mu_{+\varepsilon}$ and μ'_+ is

$$\begin{aligned} \sum_k \mathbb{E}_{\mu} [N_{k,t}] \text{KL}(\mu_{k,+\varepsilon}, \mu'_{k,+}) &= \sum_{k=1}^{K_0} \mathbb{E}_{\mu} [N_{k,t}] \text{KL}(\varepsilon, \mu_{k_0+1}) \geq \text{kl}(\mathbb{P}_{\mu_{+\varepsilon}}(L_T > K_0/2), \mathbb{P}_{\mu'_+}(L_T > K_0/2)) \\ &\geq \text{kl}\left(\frac{1}{4}, \frac{2}{K_0} \mathbb{E}_{\mu'_+}[L_T]\right) \geq \frac{1}{4} \log \frac{K_0}{2\mathbb{E}_{\mu'_+}[L_T]} - \log 2, \end{aligned}$$

We have proved that $n_{\varepsilon} \geq \frac{1}{\text{KL}(\varepsilon, \mu_{K_0+1})} \left(\frac{1}{4} \log \frac{K_0}{2\mathbb{E}_{\mu'_+}[L_T]} - \log 2 \right)$ and the final result is obtained by using the explicit form for the bound on $\mathbb{E}_{\mu'_+}[L_T]$. \square

B Non-adaptive oracle

The objective of this section is to explicit the solution of

$$\min_{\sum_k N_k = T} \sum_k a_k e^{-N_k \Delta_k^2}.$$

Introducing the Lagrange multiplier $\gamma \in \mathbb{R}$, it is straightforward that the solution is such that all N_k which are nonzero verify $\frac{\partial}{\partial N_k} (\sum_j a_j e^{-N_j \Delta_j^2}) = \gamma$. Then there exists a set S and a constant $\gamma_S(T) > 0$ for which $k \notin S \implies N_k = 0$ and for $k \in S$, $N_k \neq 0$ and

$$a_k \Delta_k^2 e^{-N_k \Delta_k^2} = \gamma_S(T).$$

That is, $N_k = \frac{1}{\Delta_k^2} (\gamma_S(T) + \log(a_k \Delta_k^2))$.

We remark that $k \in S$ iff $\frac{1}{\Delta_k^2}(\gamma_S(T) + \log(a_k \Delta_k^2)) > 0$, which then implies that if $a_1 \Delta_1^2 \leq \dots \leq a_K \Delta_K^2$, then $S = \{k_0, k_0 + 1, \dots, K\}$ for some $k_0 \in [K]$.

Using the condition $\sum_k N_k = T$ to determine $\gamma_S(T)$, we get

$$\sum_{k=k_0}^K \frac{1}{\Delta_k^2} (\gamma_S(T) + \log(a_k \Delta_k^2)) = T \quad \implies \quad \gamma_S(T) = \frac{T + \sum_{k=k_0}^K \frac{1}{\Delta_k^2} \log \frac{1}{a_k \Delta_k^2}}{\sum_{k=k_0}^K \frac{1}{\Delta_k^2}}.$$

Finally, we can characterize k_0 . Notice that $k \in S$ iff $\frac{1}{\Delta_k^2}(\gamma_S(T) + \log(a_k \Delta_k^2)) > 0$, i.e. iff

$$\frac{1}{\Delta_k^2} \left(\frac{T + \sum_{j=k_0}^K \frac{1}{\Delta_j^2} \log \frac{1}{a_j \Delta_j^2}}{\sum_{j=k_0}^K \frac{1}{\Delta_j^2}} + \log(a_k \Delta_k^2) \right) > 0 \quad \Leftrightarrow \quad T > \sum_{j=k_0}^K \frac{1}{\Delta_j^2} \log \frac{a_j \Delta_j^2}{a_k \Delta_k^2}.$$

Finally, let $H_k = \sum_{j=k+1}^K \frac{1}{\Delta_j^2} \log \frac{a_j \Delta_j^2}{a_k \Delta_k^2}$, with $H_0 = +\infty$ and $H_K = 0$. Then k_0 is the unique element of $[K]$ such that $H_{k_0} < T \leq H_{k_0-1}$.

C Properties of index-based algorithms

An algorithm is index-based if, at any round t , it pulls $k_t = \arg \min_k I_{N_{k,t-1}}^k$ where the index $I_{N_{k,t}}^k$ depends only on the number of pulls and on rewards of arm k . That index does not change when other arms are pulled.

For $C \geq 0$, let $\mathcal{F}_C \triangleq \{\exists T' \leq T, \forall k \in [K], I_{N_{k,T'}}^k \geq C\}$ be the event that at some time before T , all arm indices are above a value C . And let $\tau_k(C) = \min\{n | I_n^k \geq C\}$ be the minimal number of pulls of arm k such that its index becomes greater than C .

We start with two immediate remarks about index-based algorithms.

Lemma 3. If $I_{N_{j,t}}^j \geq C$, then at the next time t' when an index-based algorithm pulls arm j , it necessarily holds that $\min_k I_{N_{k,t'}-1}^k \geq C$.

Lemma 4. If $\min_k I_{N_{k,t}}^k \geq C$ then for all k , by definition of $\tau_k(C)$, $N_{k,t} \geq \tau_k(C)$.

This next lemma explicits \mathcal{F}_C using $(\tau_k(C))_{k \in [K]}$.

Lemma 5. An index-based algorithm verifies $\mathcal{F}_C = \{\sum_k \tau_k(C) \leq T\}$.

Proof. We first prove the inclusion $\mathcal{F}_C \subseteq \{\sum_k \tau_k(C) \leq T\}$. At the time T' defined in \mathcal{F}_C , it holds $\min_k I_{N_{k,T'}}^k \geq C$. The results then follows from lemma 4: $\sum_k \tau_k(C) \leq \sum_k N_{k,T'} = T' \leq T$.

We now prove $\{\sum_k \tau_k(C) \leq T\} \subseteq \mathcal{F}_C$. If there is no j with $N_{j,T} > \tau_j(C)$, we have $T = \sum_k N_{k,T} \leq \sum_k \tau_k(C) \leq T$. Hence there is equality and we have $N_{k,T} = \tau_k(C)$ for all k and \mathcal{F}_C is true for $T' = T$.

If there is some j such that $N_{j,T} > \tau_j(C)$, then after the time at which arm j was pulled $\tau_j(C)$ times it verified $I_{N_{j,t}}^j \geq C$. Arm j is again pulled at least once at some time t' , and at that time we have by lemmas 3 that for all k , $I_{N_{k,t'}-1}^k \geq I_{N_{j,t'}-1}^j = I_{N_{j,t}}^j \geq C$. Stated otherwise, the event \mathcal{F}_C happens. \square

Lemma 6. Let $t_{\max} = \arg \max_{t \in [T]} \min_{k \in [K]} I_{N_{k,t}}^k$. Then for all arms except at most one, $N_{k,t_{\max}} = N_{k,T}$.

Proof. The algorithm switches arm only if the index of the pulled arm becomes strictly greater than the minimal index of the others. As a consequence, the value of the minimal index at times of arm changes is increasing. If two or more arms are pulled since t_{\max} , there is an arm change later than t_{\max} and the minimal index value at that time is higher than at t_{\max} . This is a contradiction. \square

D Loss upper bound

Notation: we analyze index functions slightly more general than $F(n, x; a_k)$. Each arm has a potentially different index function $F_k(n, x)$.

Consider then algorithms which obey -a slightly more generic- Assumption 1, *i.e.* whose index can be written as $I_n^k = F_k(n, n\hat{\Delta}_{n,k}^2)$, where each F_k is non-decreasing in both variables, and $\lim_{n \rightarrow +\infty} F_k(n, ny) = +\infty$ for $y > 0$. And recall the theorem we aim at proving.

Theorem. Let $K \geq 1$, $a_1, \dots, a_K > 0$, $T \geq 1$. For $k \in [K]$, let $F_k : \mathbb{N} \times \mathbb{R} \rightarrow \mathbb{R}$ satisfying Assumption 1. Let $C_1, \dots, C_K > \max_k F_k(0, 0)$. For all $j, k \in [K]$, define

- $t_j(C_k)$ a solution of the equation $F_j(t, t\Delta_j^2) = C_k$,
- $S_k \subseteq [K]$ and $t_{j,0}(C_k) \in \mathbb{R}_+$, a set and values such that for $i \notin S_k$, $\mathbb{P}(\exists n \leq t_{i,0}(C_k), F_i(n, n\hat{\Delta}_{n,i}^2) \geq C) = 1$.

Then the expected loss of Algorithm 1 is upper-bounded as

$$\mathbb{E}[L_T^A] \leq \sum_{k=1}^K a_k \left(e \cdot \exp \left(-\frac{\frac{1}{2} \left(T - \sum_{j \notin S_k} t_{j,0}(C_k) \right) - \sum_{j \in S_k} t_j(C_k)}{\sum_{j \in S_k} 1/\Delta_j^2} \right) + T \cdot e^{-t_k(C_k)\Delta_k^2} \right).$$

Proof. Set $(C_k)_{1 \leq k \leq K} \in \mathbb{R}^+$, so that it immediately follows

$$\begin{aligned} \mathbb{E}[L_T] &= \sum_k a_k \mathbb{P}(\hat{s}_k \neq s_k) \\ &\leq \sum_{k=1}^K a_k \mathbb{P}(\hat{s}_k \neq s_k \wedge \exists T' \in [T], I_{N_{T'}}^k \geq C_k) + a_k \mathbb{P}(\forall T' \in [T], I_{N_{T'}}^k < C_k) \\ &\leq \sum_{k=1}^K a_k \mathbb{P}(\hat{s}_k \neq s_k \wedge \exists T' \in [T], I_{N_{T'}}^k \geq C_k) + a_k \mathbb{P}(\overline{\mathcal{F}_{C_k}}), \end{aligned}$$

where $\overline{\mathcal{E}}$ stands for the complement of an event \mathcal{E} . The proof then proceeds in two steps, which are proved in subsections D.1 and D.2, in order to control both probabilities introduced above:

1. Lemma 9: with large probability, there is some time t for which the index $I_{N_{k,t}}^k$ is large for all $k \in [K]$ (and all arms are well explored)
2. Theorem 5: if $I_{N_{k,t}}^k$ is large, then there is a small probability of mistake for arm k .

□

D.1 With large probability, all arms are well explored

Technical tools used in this section First, we state a simple but useful lemma.

Lemma 7. Let $\Delta_{n,k}^\delta = \Delta_k - \sqrt{\frac{\log(1/\delta)}{n}}$. For all $\delta \in (0, 1)$ such that $\Delta_{n,k}^\delta \geq 0$, Hoeffding's inequality implies that $\mathbb{P}(\hat{\Delta}_{n,k} < \Delta_{n,k}^\delta) \leq \delta$.

The next lemma will be used to bound the sum of exponentially tailed distributions.

Lemma 8 ([20]). Let Z_1, \dots, Z_K be independent random variables and $a_1, \dots, a_K \in \mathbb{R}^+$ be such that for all $k \in [K]$ and $x \in \mathbb{R}^+$, $\mathbb{P}(Z_k \geq x) \leq e^{-a_k x}$. Then for all $\lambda \geq 0$,

$$\mathbb{P}\left(\sum_k Z_k \geq \lambda \sum_k \frac{1}{a_k}\right) \leq e^{1-\lambda}.$$

Corollary 4. Let Y_1, \dots, Y_K be independent random variables and $y_1, \dots, y_K \in \mathbb{R}$, $a_1, \dots, a_K \in \mathbb{R}^+$ be such that for all $k \in [K]$ and $x \in \mathbb{R}^+$, $\mathbb{P}(Y_k \geq y_k + x) \leq e^{-a_k x}$. Then for all $x \geq \sum_k y_k$,

$$\mathbb{P}\left(\sum_k Y_k \geq x\right) \leq e \times \exp\left(-\frac{x - \sum_k y_k}{\sum_k 1/a_k}\right).$$

The corollary is a direct application of Lemma 8 to $Z_k = Y_k - y_k$.

Main proof We know from Lemma 5 that any index based algorithm verifies $\mathcal{F}_C = \{\sum_k \tau_k(C) \leq T\}$. Hence, to prove that \mathcal{F}_C happens with great probability, it suffices to show that $\sum_k \tau_k(C)$ has an exponential tail.

We derive a bound on $\mathbb{P}(\sum_k \tau_k(C) > T)$. To that end, we bound individually for each arm $\mathbb{P}(\tau_k(C) \geq t_k + x)$ for some t_k to be defined and $x \geq 0$, and conclude by Corollary 4.

Theorem 4. *Under Assumption 1, the algorithm verifies, for all $x \geq 0$,*

$$\mathbb{P}(\sqrt{\tau_k(C)} > \sqrt{t_k(C)} + x) \leq \exp(-\Delta_k^2 x^2),$$

with $t_k(C)$ solution to $F_k(t, t\Delta_k^2) = C$, if such a solution exists. Otherwise, if $C < F_k(0, 0)$ and no solution exists, $\tau_k(C) = 0$ with probability 1.

Proof. We first bound $\mathbb{P}(\tau_k(C) > t_k + x)$, where t_k is chosen later, and $x \geq 0$.

$$\begin{aligned} \mathbb{P}(\tau_k(C) > t_k + x) &= \mathbb{P}(\forall n \leq t_k + x, I_n^k < C) \leq \mathbb{P}(I_{t_k+x}^k < C) \\ &= \mathbb{P}(F_k(t_k + x, (t_k + x)\hat{\Delta}_{t_k+x,k}^2) < C) \end{aligned}$$

First, by monotonicity of F_k , this probability equals zero if $C \leq F_k(t_k + x, 0)$. If $C > F_k(t_k + x, 0)$, we define $\delta_{t_k+x,k,C}$ such that $\Delta_{t_k+x,k}^{\delta_{t_k+x,k,C}} = \inf\{\Delta \geq 0 \mid C < F_k(t_k + x, (t_k + x)\Delta^2)\}$. In the following, we write $F_k(n, \cdot)^{-1}(C) := \inf\{x \mid C \leq F_k(n, x)\}$. If $F_k(n, \cdot)$ is increasing, this is its inverse, but we only suppose that $F_k(n, \cdot)$ is non-decreasing. Note that $x < F_k(n, \cdot)^{-1}(C)$ implies that $F_k(n, x) < C$. With that definition, $(t_k + x)(\Delta_{t_k+x,k}^{\delta_{t_k+x,k,C}})^2 = F_k(t_k + x, \cdot)^{-1}(C)$. As a consequence, we get

$$\begin{aligned} \mathbb{P}(\tau_k(C) > t_k + x) &\leq \mathbb{P}(I_{t_k+x}^k < C) = \mathbb{P}(F_k(t_k + x, (t_k + x)\hat{\Delta}_{t_k+x,k}^2) < C) \\ &\leq \mathbb{P}((t_k + x)\hat{\Delta}_{t_k+x,k}^2 < F_k(t_k + x, \cdot)^{-1}(C)) \\ &= \mathbb{P}(\hat{\Delta}_{t_k+x,k}^2 < (\Delta_{t_k+x,k}^{\delta_{t_k+x,k,C}})^2) \\ &\leq \delta_{t_k+x,k,C}, \end{aligned}$$

where by definition,

$$\begin{aligned} \delta_{t_k+x,k,C} &= \exp\left(- (t_k + x)(\Delta_k - \Delta_{t_k+x,k}^{\delta_{t_k+x,k,C}})^2\right) \\ &= \exp\left(- \left(\sqrt{\Delta_k^2(t_k + x)} - \sqrt{F_k(t_k + x, \cdot)^{-1}(C)}\right)^2\right). \end{aligned}$$

We intend to prove an exponential decrease with x . In order to have it, we will set t_k such that the exponential is equal to 1 for $x = 0$, and then decreases as x grows. Let then t_k be such that $C \leq F_k(t_k, t_k\Delta_k^2)$. It exists as soon as $C \geq F_k(0, 0)$ (where the later is non-positive for specific algorithms we will consider). For all $t \geq t_k$, $F_k(t, t\Delta_k^2) \geq F_k(t_k, t_k\Delta_k^2) \geq C$, which leads to $\sqrt{\Delta_k^2 t} - \sqrt{F_k(t, \cdot)^{-1}(C)} \geq 0$. Note that since F_k is non-decreasing in the first variable we have $F_k(t + x, \cdot)^{-1}(C) \leq F_k(t, \cdot)^{-1}(C)$ for all $t, x \geq 0$, and

$$\begin{aligned} \delta_{t_k+x,k,C} &= \exp\left(- \left(\sqrt{\Delta_k^2(t_k + x)} - \sqrt{F_k(t_k + x, \cdot)^{-1}(C)}\right)^2\right) \\ &\leq \exp\left(- \left(\sqrt{\Delta_k^2(t_k + x)} - \sqrt{F_k(t_k, \cdot)^{-1}(C)}\right)^2\right) \\ &\leq \exp\left(-\Delta_k^2 (\sqrt{t_k + x} - \sqrt{t_k})^2\right). \end{aligned}$$

Let $Y_k = \max(\sqrt{\tau_k(C)}, \sqrt{t_k})$; we have proved that for all $x \geq 0$,

$$\mathbb{P}(Y_k > \sqrt{t_k + x}) \leq \exp\left(-\Delta_k^2 (\sqrt{t_k + x} - \sqrt{t_k})^2\right).$$

By setting $x = 2\sqrt{\lambda t_k} + \lambda$ for $\lambda \geq 0$, we get $\mathbb{P}(Y_k > \sqrt{t_k} + \lambda) \leq \exp(-\Delta_k^2 \lambda^2)$. \square

Lemma 9. For all $C > 0$, t_k such that $F(t_k, t_k \Delta_k^2) \geq C$.

$$\mathbb{P}\left(\sum_k \tau_k(C) \geq T\right) \leq e \times \exp\left(-\frac{T - \sum_k t_k}{\sum_k 1/\Delta_k^2}\right).$$

Proof. Rewrite the event $\{\sum_k \tau_k(C) > T\}$ using $Y_k = \max(\sqrt{\tau_k(C)}, \sqrt{t_k})$, so that:

$$\begin{aligned} \mathbb{P}\left(\sum_k \tau_k(C) > T\right) &= \mathbb{P}\left(\sum_k ((Y_k - \sqrt{t_k}) + \sqrt{t_k})^2 > T\right) \\ &\leq \mathbb{P}\left(2 \sum_k (Y_k - \sqrt{t_k})^2 + 2 \sum_k t_k > T\right) \\ &\leq \mathbb{P}\left(\sum_k (Y_k - \sqrt{t_k})^2 > T/2 - \sum_k t_k\right). \end{aligned} \quad (8)$$

We now apply Lemma 8 to $(Y_k - \sqrt{t_k})^2$, which verifies $\mathbb{P}((Y_k - \sqrt{t_k})^2 \geq x) \leq \exp(-\Delta_k^2 x)$. From Equation (8), we obtain

$$\mathbb{P}\left(\sum_k \tau_k(C) > T\right) \leq \mathbb{P}\left(\sum_k (Y_k - \sqrt{t_k})^2 > T/2 - \sum_k t_k\right) \leq e \times \exp\left(-\frac{T/2 - \sum_k t_k}{\sum_k 1/\Delta_k^2}\right).$$

□

Remark We can actually derive a tighter bound than (8)

$$\begin{aligned} \mathbb{P}\left(\sum_k \tau_k(C) > T\right) &= \mathbb{P}\left(\sum_k ((Y_k - \sqrt{t_k}) + \sqrt{t_k})^2 > T\right) \\ &\leq \mathbb{P}\left(\left(\sqrt{\sum_k (Y_k - \sqrt{t_k})^2} + \sqrt{\sum_k t_k}\right)^2 > T\right) \\ &= \mathbb{P}\left(\sum_k (Y_k - \sqrt{t_k})^2 > (\sqrt{T} - \sqrt{\sum_k t_k})^2\right). \end{aligned}$$

Roughly speaking, to get it, just write $\|(Y - \sqrt{t}) + \sqrt{t}\|^2 \leq (\|Y - \sqrt{t}\| + \|\sqrt{t}\|)^2$. To get Equation (8), we further use $(\|Y - \sqrt{t}\| + \|\sqrt{t}\|)^2 \leq 2\|Y - \sqrt{t}\|^2 + 2\|\sqrt{t}\|^2$. In the case of APT (at least) it leads to the *same* final bound on the algorithm because when we optimize further down, we set $\sum_k t_k = T/4$ no matter which of these inequalities we use, value for which resulting exponents are equal.

Corollary 5. Suppose now that there is a set S_C such that for $k \notin S_C$, $\mathbb{P}(\tau_k(C) > t_k) = 0$. Then we can refine Lemma 9 to

$$\mathbb{P}\left(\sum_k \tau_k(C) > T\right) \leq \mathbb{P}\left(\sum_{k \in S_C} \tau_k(C) > T - \sum_{k \notin S_C} t_k\right) \leq e \times \exp\left(-\frac{(T - \sum_{k \notin S_C} t_k)/2 - \sum_{k \in S_C} t_k}{\sum_{k \in S_C} 1/\Delta_k^2}\right).$$

D.2 When an arm index is large, the probability of mistake is small

The goal of this section is to bound $\mathbb{P}(\hat{s}_k \neq s_k \wedge \exists T' \in [T], I_{N_{T'}}^k \geq C)$. We define the random variable $E_{k,t} = \mathbb{I}\{(\hat{\mu}_{k,t} - \theta)(\mu_k - \theta) < 0\}$; it is equal to 1 iff there is an error on the sign at time t . The algorithm makes a mistake on arm k is $E_{k,t_{\max}} = 1$ since it returns the sign at that time.

Theorem 5. *The algorithm using F_k for its index definition verifies*

$$\begin{aligned} \mathbb{P}(\hat{s}_k \neq s_k \wedge \exists T' \in [T], I_{N_{T'}}^k \geq C) &\leq \mathbb{P}(\exists T' \in [T], E_{k,T'} = 1 \wedge I_{N_{T'}}^k \geq C) \\ &\leq T \cdot \inf\{e^{-n_k \Delta_k^2} \mid F_k(n_k, n_k \Delta_k^2) < C\}. \end{aligned}$$

Proof. We use Lemma 12 (below): find n_k as large as possible such that $F_k(n_k, n_k \Delta_k^2) < C$. Then, since the algorithm returns the sign of the arm at the time at which its index was maximal, we get $\mathbb{P}(\hat{s}_k \neq s_k \wedge \exists T' \in [T], I_{N_{T'}}^k \geq C) \leq \mathbb{P}(\exists T' \in [T], E_{k,T'} = 1 \wedge I_{N_{T'}}^k \geq C) \leq T \cdot e^{-n_k \Delta_k^2}$. \square

Lemma 10. For any $\delta_k \in (0, 1)$, with probability at least $1 - \delta_k$, and for all $n \in [T]$, it holds

$$\sqrt{N_{k,t}}(\hat{\mu}_t^k - \mu^k) \leq \sqrt{\log\left(\frac{T}{\delta_k}\right)}.$$

Proof. This is a direct implication of Hoeffding's inequality with a union bound for time-uniformity. \square

Define $n_k = \frac{1}{\Delta_k^2} \log\left(\frac{T}{\delta_k}\right)$. Consider the following three facts (their definition will be useful for the following proofs):

1. If the concentration holds, then $\sqrt{N_{k,t}}(\hat{\mu}_t^k - \mu^k) \leq \sqrt{n_k} \Delta_k$.
2. If there is a mistake at time t , then we have
 - (a) $\hat{\mu}_t^k - \mu^k \geq \Delta_k$.
 - (b) $(\hat{\mu}_t^k - \mu^k)^2 \geq \hat{\Delta}_{N_{k,t},k}^2 + \Delta_k^2$.
3. If $I_{N_{k,t}}^k > C$ then $F_k(N_{k,t}, N_{k,t} \hat{\Delta}_{t,k}^2) > C$ and $N_{k,t} \hat{\Delta}_{t,k}^2 \geq F_k(N_{k,t}, \cdot)^{-1}(C)$.

Lemma 11. If at time t , concentration holds and there is a mistake (1 and 2 are true), then $N_{k,t} \leq n_k$ and $F_k(N_{k,t}, (n_k - N_{k,t}) \Delta_k^2) \geq I_{N_{k,t}}^k$.

Proof. First point: combine 1 and 2(a). Second point: use 1, then 2(b), then the definition of $I_{N_{k,t}}^k$:

$$n_k \Delta_k^2 \geq N_{k,t} (\hat{\mu}_t^k - \mu^k)^2 \geq N_{k,t} \hat{\Delta}_{t,k}^2 + N_{k,t} \Delta_k^2 \geq F_k(N_{k,t}, \cdot)^{-1}(I_{N_{k,t}}^k) + N_{k,t} \Delta_k^2.$$

\square

Lemma 12. If at time t , all three ‘‘if’’ are true, then $F_k(n_k, n_k \Delta_k^2) \geq C$.

Proof. Use the monotonicity of F_k in the inequality of Lemma 11. We have $N_{k,t} \leq n_k$ and $n_k - N_{k,t} \leq n_k$. \square

D.3 Examples

In this section, we explicit Theorem 2 for certain algorithms from the literature and for our algorithm, FWT.

D.3.1 APT ([28])

This algorithm (in its variant that stops at t_{\max}) corresponds to $F(n, x) = x$. To apply Theorem 2 we find that $t_j(C_k) = \frac{C_k}{\Delta_j^2}$ is solution, then:

$$\mathbb{E}[L_T] \leq \sum_{k=1}^K e a_k \exp\left(-\left(\frac{T/2}{\sum_{j=1}^K 1/\Delta_j^2} - C_k\right)\right) + T \cdot a_k \exp(-C_k)$$

An optimal C_k is such that

$$e \cdot \exp\left(-\frac{T/2}{\sum_{j=1}^K 1/\Delta_j^2} + C_k\right) = T \cdot \exp(-C_k)$$

Then the bound becomes:

$$\mathbb{E}[L_T] \leq 2\sqrt{e \cdot T} \sum_{k=1}^K a_k \exp\left(-\frac{1}{4} \frac{T}{\sum_{k=1}^K 1/\Delta_k^2}\right)$$

D.3.2 LSA ([35])

This algorithm corresponds to $F(n, x) = x + \log n$, the stopping time $t_j(C_k)$ is solution to $t\Delta_j^2 + \log t = C_k$. This equation has a closed form solution: $t_j(C_k) = \frac{1}{\Delta_j^2} W(\Delta_j^2 \exp(C_k))$, the loss bound becomes

$$\mathbb{E}[L_T] \leq \sum_{k=1}^K e \cdot a_k \exp \left(- \left(\frac{\frac{T}{2} - \sum_{j=1}^K \frac{1}{\Delta_j^2} W(\Delta_j^2 \exp(C_k))}{\sum_{j=1}^K 1/\Delta_j^2} \right) \right) + T a_k \exp(-W(\Delta_k^2 \exp(C_k)))$$

We can use an inequality on the lambert function (*cf.* Corollary 2.4 in [17]), For all $x \geq e$ we have

$$\log x - \log \log x \leq W(x) \leq \log x - \log \log x + \log(1 + e^{-1}),$$

this entails that if $\forall k, j \quad C_k + \log \Delta_j^2 \geq 1$ we obtain the more accurate bound

$$\begin{aligned} \mathbb{E}[L_T] &\leq \sum_{k=1}^K (1 + e) \cdot a_k \exp \left(- \left(\frac{\frac{T}{2} + \sum_{j=1}^K \frac{1}{\Delta_j^2} (\log \frac{1}{\Delta_j^2} + \log(C_k + \log \Delta_j^2))}{\sum_{j=1}^K 1/\Delta_j^2} \right) + C_k \right) \\ &\quad + \sum_{k=1}^K \frac{T(C_k + \log \Delta_k^2)}{\Delta_k^2} a_k \exp(-C_k). \end{aligned} \quad (9)$$

D.3.3 FWT (our algorithm)

Our algorithm corresponds to $F_k(n, x) = \max(1, x) - \log \max(1, x) + \log n - \log a_k$. In order to find the times $t_j(C_k)$ of Theorem 2 we solve the equation:

$$\max(1, n\Delta_k^2) - \log \max(1, n\Delta_k^2) + \log(n\Delta_k^2) = C + \log a_k + \log \Delta_k^2.$$

Let $D_k = C + \log a_k + \log \Delta_k^2$. We want a solution to $\max(1, x) - \log \max(1, x) + \log x = D_k$. The function on the left, which we now denote by \mathcal{I} , is increasing and bijective from \mathbb{R}^+ to \mathbb{R} .

- If $D_k \geq 1$, $\mathcal{I}(D_k) = D_k$.
- If $D_k \leq 1$, $\mathcal{I}(e^{D_k-1}) = D_k$.
- For all $D_k > 0$, $\mathcal{I}(e^{D_k-1}) \geq D_k$.

Moreover, we have $F_k(a_k e^{C-1}) \geq C$, it comes $\mathbb{P}(\tau_k(C) > a_k e^{C-1}) = 0$.

Let $C_k > 0$ and $S_k \subseteq \{j \in [K] \mid C_k + \log a_j \Delta_j^2 \geq 1\}$. Let S' be a set such that for all $k \in S'$, $k \in S_k$.

$$\begin{aligned} \mathbb{E}[L_T] &\leq \sum_{k \notin S'} a_k + \sum_{k \in S'} a_k \mathbb{P}(\overline{\mathcal{F}_{C_k}}) + \sum_{k \in S'} a_k \mathbb{P}(\overline{\mathcal{E}_k(T)} \cap \mathcal{F}_{C_k}) \\ &\leq \sum_{k \notin S'} a_k + e \sum_{k \in S'} a_k \exp \left(- \frac{\frac{1}{2}(T - \sum_{j \notin S_k} a_j e^{C_k-1}) - \sum_{j \in S_k} \frac{1}{\Delta_j^2} \log a_j \Delta_j^2}{\sum_{j \in S_k} 1/\Delta_j^2} + C_k \right) \\ &\quad + T \sum_{k \in S'} a_k \exp(-C_k - \log a_k \Delta_k^2). \end{aligned}$$

Large T The value of C_k which equalizes the two terms indexed by k verifies

$$C_k = \frac{1}{2} \frac{\frac{1}{2}(T - \sum_{j \notin S_k} a_j e^{C_k-1}) - \sum_{j \in S_k} \frac{1}{\Delta_j^2} \log a_j \Delta_j^2}{\sum_{j \in S_k} 1/\Delta_j^2} - \frac{1}{2} \log a_k \Delta_k^2 + \frac{1}{2} \log \frac{T}{e}.$$

The latter can be chosen if T is big enough such that $S_k = [K]$ for all arms, this is the case if $T \geq 2 \sum_{j=1}^K \frac{1}{\Delta_j^2} (2 + \log \frac{a_j \Delta_j^2 \max_i a_i \Delta_i^2}{(\min_k a_k \Delta_k^2)^2} - \log \frac{T}{e^3})$, we get the bound

$$\mathbb{E}[L_T] \leq 2\sqrt{eT} \sum_k a_k \exp \left(- \frac{1}{2} \frac{T/2 - \sum_j \frac{1}{\Delta_j^2} \log \frac{a_j \Delta_j^2}{a_k \Delta_k^2}}{\sum_j 1/\Delta_j^2} \right)$$

Up to a 1/4 factor, this is the exponent of the optimal non-adaptive oracle (*cf.* Eq. 4).

General T We choose a set $S \subseteq [K]$ and set $C_k = C$, a common value still to be determined, for all $k \in S'$. Then for all $k \in S'$, we set $S_k = S$. We impose $C \geq 1 + \max_{j \in S} \log \frac{1}{a_j \Delta_j^2}$, such that the condition $S_k \subseteq \{j \in [K] \mid C_k + \log a_j \Delta_j^2 \geq 1\}$ is verified. We get that for all $S' \subseteq S \subseteq [K]$ and $C \geq 1 + \max_{j \in S} \log \frac{1}{a_j \Delta_j^2}$,

$$\begin{aligned} \mathbb{E}[L_T] \leq & \sum_{k \notin S'} a_k + e \sum_{k \in S'} a_k \exp \left(-\frac{\frac{1}{2}(T - \sum_{j \notin S} a_j e^{C-1}) + \sum_{j \in S} \frac{1}{\Delta_j^2} \log \frac{1}{a_j \Delta_j^2}}{\sum_{j \in S} 1/\Delta_j^2} + C \right) \\ & + T \sum_{k \in S'} a_k \exp(-C - \log a_k \Delta_k^2). \end{aligned}$$

E Extention to sum-of-gaps

The global loss we investigate in this section is

$$L_T = \sum_k \Delta_k E_k.$$

First we write the Frank-Wolfe index: $\arg \min_k N_{k,t} \hat{\Delta}_{k,t}^2 - \frac{3}{2} \log(N_{k,t} \hat{\Delta}_{k,t}^2) + \frac{3}{2} \log(N_{k,t})$, then we slightly modify it to comply with Assumption 1 (see explanation above Eq. FWT):

$$I_{N_{k,t}}^k = \max\left(\frac{3}{2}, N_{k,t} \hat{\Delta}_{k,t}^2\right) - \frac{3}{2} \log \max\left(\frac{3}{2}, N_{k,t} \hat{\Delta}_{k,t}^2\right) + \frac{3}{2} \log N_{k,t}.$$

This corresponds to the function $F(n, x) = \frac{3}{2} \log n + \max(\frac{3}{2}, x) - \frac{3}{2} \log \max(\frac{3}{2}, x)$, used for all arms.

Solving $F(n, n\Delta_k^2) = C$ gives rise to two cases:

- $n = \frac{1}{\Delta_k^2}(C + \frac{3}{2} \log \Delta_k^2)$ if $C + \frac{3}{2} \log \Delta_k^2 \geq \frac{3}{2}$,
- $n = \frac{3}{2} \exp(\frac{2}{3}C - 1)$ otherwise. In that case, $n \leq \frac{3}{2} \frac{1}{\Delta_k^2}$.

Also, $n = \frac{3}{2} \exp(\frac{2}{3}C - 1)$ is solution to $F(n, 0) = C$.

Consider $C_k > 0$, let $S_k = \{j \in [K] \mid C_k + \frac{3}{2} \log \Delta_j^2 \geq \frac{3}{2}\}$ and S' be a set such that for all $k \in S'$, $k \in S_k$, then

$$\begin{aligned} \mathbb{E}[L_T] \leq & \sum_{k \notin S'} \Delta_k + \sum_{k \in S'} \Delta_k \mathbb{P}(\overline{\mathcal{F}_{C_k}}) + \sum_{k \in S'} \Delta_k \mathbb{P}(\overline{\mathcal{E}_k(T)} \cap \mathcal{F}_{C_k}) \\ \leq & \sum_{k \notin S'} \Delta_k + e \sum_{k \in S'} \Delta_k \exp \left(-\frac{\frac{1}{2}(T - \sum_{j \notin S_k} \frac{3}{2} \Delta_j e^{\frac{2}{3}C_k - 1}) + \sum_{j \in S_k} \frac{3}{2} \frac{1}{\Delta_j^2} \log \frac{1}{\Delta_j^2}}{\sum_{j \in S_k} 1/\Delta_j^2} + C_k \right) \\ & + \sum_{k \in S'} T \Delta_k \exp \left(-C_k - \frac{3}{2} \log \Delta_k^2 \right) \end{aligned}$$

Large T The values of C_k that optimize the r.h.s of the previous inequality verify:

$$C_k = \frac{1}{2} \frac{\frac{1}{2}(T - \sum_{j \notin S_k} \frac{3}{2} \Delta_j e^{\frac{2}{3}C_k - 1}) + \sum_{j \in S_k} \frac{3}{2} \frac{1}{\Delta_j^2} \log \frac{1}{\Delta_j^2}}{\sum_{j \in S_k} 1/\Delta_j^2} - \frac{1}{2} \log \Delta_k^2 + \frac{1}{2} \log \frac{T}{e}$$

If T is big enough such that $S_k = [K]$ for all arms, which happens once $T \geq 2 \sum_{j=1}^k \frac{1}{\Delta_j^2} \left(3 + 3 \log \frac{\Delta_j \max_i \Delta_i}{(\min_k \Delta_i)^2} - \log \frac{T}{e} \right)$, we get the bound:

$$\mathbb{E}[L_T] \leq 2\sqrt{eT} \sum_k \Delta_k \exp \left(-\frac{\frac{T}{2} + \sum_j \frac{3}{2} \frac{1}{\Delta_j^2} \log \frac{\Delta_k^2}{\Delta_j^2}}{\sum_j 1/\Delta_j^2} \right).$$

F Upper-bounds comparison

Figure 4 compares the upper-bounds of Corollary 1 (APT), Equation (9) (LSA), and Corollary 2 (FWT) for the particular case $\Delta_i = i/K$ and $a_i = 1$, for all $i = 1, \dots, K$ and $K = 50$. We observe a behavior similar to that of Figure 2.

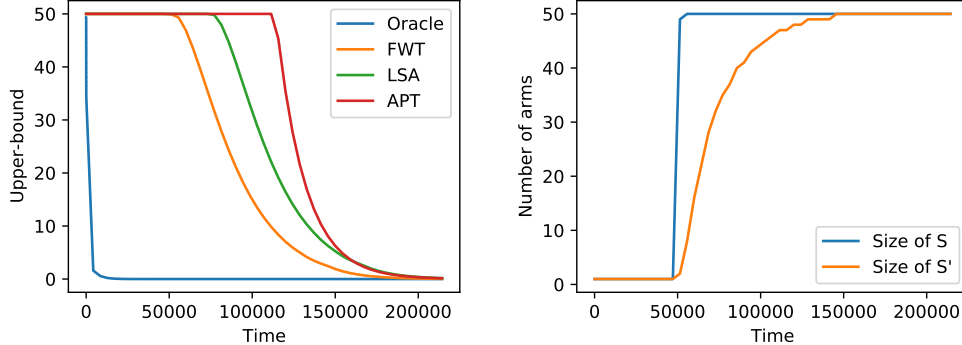


Figure 4: [left] Comparison of the upper-bounds of Corollaries 1, 2 and that of the optimal non-adaptive oracle of Eq. (5) (blue) when the gaps are of the form $\Delta_i = i/K$. [right] Evolution over time of the size of the optimal sets S (blue) and S' (orange) that minimize the bound of Corollary 2.

G Additional Experiments

In this section, we illustrate on synthetic data the performance of APT, LSA, and FWT. The implemented algorithms respectively correspond to Algorithm 1 with the following choices:

$$F(n, x; 1) = x \quad (\text{APT})$$

$$F(n, x; 1) = x + \log(n) \quad (\text{LSA})$$

$$F(n, x; 1) = (1 + \sqrt{x})^2 - \log(1 + \sqrt{x})^2 + \log(n) \quad (\text{FWT})$$

Note that we used a slightly different version for APT than the one proposed in the analysis $F(n, x; 1) = \max\{1, x\} - \log \max\{1, x\} + \log(n)$. The analysis and experiments work similarly for both versions. But the $(1 + \sqrt{x})^2$ version performs slightly better empirically while the $\max\{1, x\}$ version provides cleaner theoretical results. The experiments are averaged over 500 runs and consider arm distributions of the form $\nu_k = \mathcal{N}(\mu_k, 1)$. The gaps are thus $\Delta_k = |\mu_k|$, for $k = 1, \dots, K$. The performance criterion is the sum of errors defined in Equation 4 with $a_k = 1$. The experiments were run on a personal laptop with Intel Core i5, Dual Core, 3.1 GHz.

Since most of the tested experiments obtained similar performance, we only provide the results of a few experiments. Although our theoretical upper bounds are slightly better for FWT, LSA and FWT generally have similar performance, while APT underperforms. This last point is not surprising since, although we provide in Corollary 1 an upper bound for APT that appears asymptotically similar to those of LSA and FWT, APT was not designed to minimize the sum of errors. APT was made to minimize the probability of making at least one error and thus focuses too much on arms with very small gaps that are very difficult to classify.

Figure 5 shows the performance of the algorithms together with the non-adaptive oracle of Section 4. Figure 6 plots the ratio of error with respect to the non-adaptive oracle. Interestingly, in all of our experiments, APT and FWT perform better than it.

Figures 7 and 8 represent the optimal non-adaptive sampling distribution if the means were known and the empirical sampling distribution of the algorithms for different numbers of iterations. As we can see, for the initial phase, the arms that are closest to the threshold should ideally not be drawn. Yet, as our lower bound in Section A.1 illustrates, this is not possible for sequential algorithms. All arms must be sampled. We can see that this is indeed the case for all algorithms: the closer the arms are to the threshold, the more likely they are to be sampled.

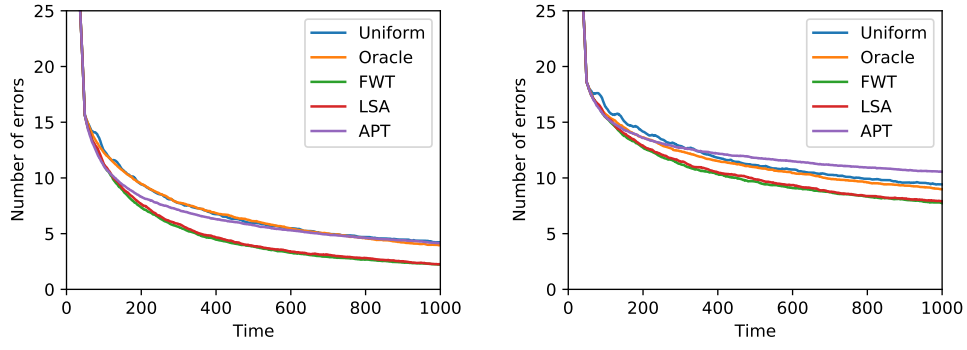


Figure 5: Sum of errors over time of the different algorithms when $\mu_k = \frac{(-1)^k k}{K}$ [left] and $\mu_k = \frac{(-1)^k k^2}{K^2}$ [right] for $K = 50$ arms.

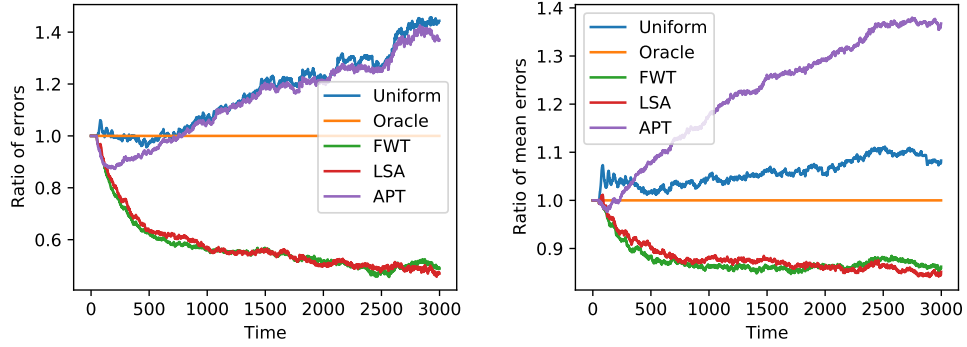


Figure 6: Ratio of improvement with respect to the non adaptive oracle sampling when $\mu_k = \frac{(-1)^k k}{K}$ [left] and $\mu_k = \frac{(-1)^k k^2}{K^2}$ [right] for $K = 50$ arms.

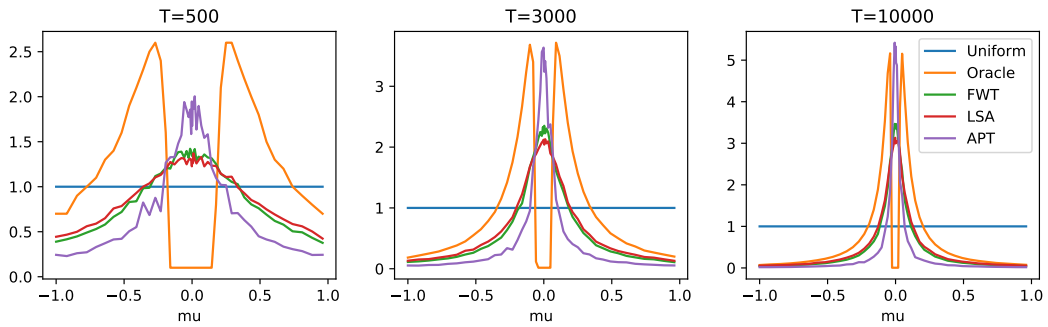


Figure 7: Optimal sampling distribution and empirical sampling distribution with respect to μ when $\mu_k = \frac{(-1)^k k^2}{K^2}$ for $K = 50$ arms.

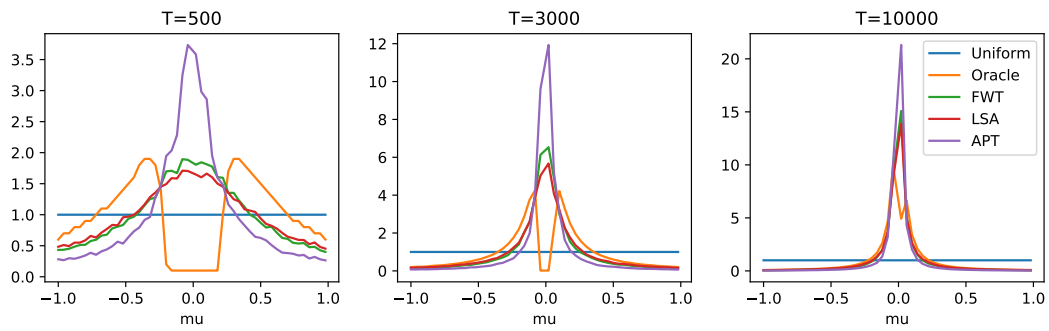


Figure 8: Optimal sampling distribution and empirical sampling distribution with respect to μ when $\mu_k = \frac{(-1)^k k}{K}$ for $K = 50$ arms.