
Asymptotically Optimal Bandits under Weighted Information

Matias I. Müller

Division of Decision and Control Systems
KTH Royal Institute of Technology
SE-100 44 Stockholm, Sweden
mimr2@kth.se

Cristian R. Rojas

Division of Decision and Control Systems
KTH Royal Institute of Technology
SE-100 44 Stockholm, Sweden
crrro@kth.se

Abstract

We study the problem of regret minimization in a multi-armed bandit setup where the agent is allowed to play multiple arms at each round by spreading the resources usually allocated to only one arm. At each iteration the agent selects a normalized power profile and receives a Gaussian vector as outcome, where the unknown variance of each sample is inversely proportional to the power allocated to that arm. The reward corresponds to a linear combination of the power profile and the outcomes, resembling a linear bandit. By spreading the power, the agent can choose to collect information much faster than in a traditional multi-armed bandit at the price of reducing the accuracy of the samples. This setup is fundamentally different from that of a linear bandit—the regret is known to scale as $\Theta(\sqrt{T})$ for linear bandits, while in this setup the agent receives a much more detailed feedback, for which we derive a tight $\log(T)$ problem-dependent lower-bound. We propose a Thompson-Sampling-based strategy, called Weighted Thompson Sampling (WTS), that designs the power profile as its posterior belief of each arm being the best arm, and show that its upper bound matches the derived logarithmic lower bound. Finally, we apply this strategy to a problem of control and system identification, where the goal is to estimate the maximum gain (also called \mathcal{H}_∞ -norm) of a linear dynamical system based on batches of input-output samples.

1 Introduction

A stochastic multi-armed bandits (MAB, or simply bandit) is a fundamental instance of reinforcement learning which originated in the early 1900's [1], but it only recognized as an experiment design problem in the 1950's [2]. MABs capture an essential sequential decision problem subject to an exploration-exploitation trade-off. In a general MAB setup, an agent is presented with a finite set of actions from which it can take only one at every round, revealing a random outcome/pay-off/reward whose mean is fixed and unknown to the agent. The goal of the agent is to sample arms sequentially so as to maximize the expected cumulative reward or, equivalently, to minimize the regret of not choosing the arm with the largest mean. Numerous extensions and variants of the classical stochastic MAB problem [3] have been recently investigated. These extensions are motivated by problems arising in various fields, including online recommendation systems (search engines, display ads, etc.), financial portfolio design [4, 5], adaptive network routing [6], web crawling [7], clinical trials [8, 1], and, recently, in control theory [9] and system identification [10, 11].

In this paper we consider a bandit problem in which, at every round, the agent can spread the available resources to many arms instead of allocating all the resources to only one arm. This setup allows the agent to collect more information from the environment, since it perceives a pay-off from every arm where it has allocated resources, but at the cost of reducing the accuracy of the samples. More

specifically, the outcomes from each arm have a fixed mean for every experiment, but the variance of the samples is inversely proportional to the amount of resources allocated to each arm. This poses an extra difficulty to the agent, which not only has to balance the exploration-exploitation trade-off, but now it also has to balance the precision invested in every sample. The sequence of outcomes received by the agent at every round resembles the data collected in a traditional MAB along many experiments, where sampling a given arm many times emulates one noisy sample from its mean but with smaller variance. By not allocating any resource to an arm, the variance of its sample diverges, which can be interpreted as missing or non-informative data, recovering the traditional MAB setup where the pay-offs are only revealed to the agent when an arm is selected. The class of distributions addressed in this problem are Gaussian, two-dimensional, and where the minimum arm's variance (*i.e.*, the one generating the data when all the resources are allocated to only one arm) is possibly unknown, leading to a 1-parameter or 2-parameter Gaussian bandit. Gaussian distributions are particularly suitable for this class of problems, since their means and variances are decoupled, *i.e.*, they can be fixed independently.

The bandit problem presented in this work is motivated by a fundamental problem in model-based control, where it is necessary to estimate the largest ℓ_2 -gain of the difference between a linear time-invariant dynamical system (the system we want to control) and its model, also known as the modeling error. Estimating this quantity efficiently in a model-free fashion involves collecting input-output data from the modeling error when the input applied to the system produces an output of large gain (in Euclidean norm). It is well known that the input signal achieving the largest gain is a sinusoid whose frequency is equal to the one at which the frequency response of the modeling error is the largest, also known as the peak frequency. The problem of data collection can be addressed as a regret-minimization problem of choosing as input to the system a sinusoid of frequency different from the peak frequency.

Our contributions

The main contributions of our work are:

1. the introduction of a new MAB problem called a bandit under weighted information, which allows resource-spreading policies;
2. a lower bound on the asymptotic regret incurred by a wide class of uniformly efficient policies;
3. a bandit algorithm, WTS (Weighted Thompson Sampling), that can efficiently spread the resources at every round;
4. a thorough theoretical derivation of an upper bound on the asymptotic regret incurred by WTS for the 2-parameter Gaussian bandit under weighted information, which matches the lower bound on the regret, thus establishing the asymptotic optimality of WTS for these classes of problems; and
5. an application of WTS to the gain estimation problem in control and system identification, where the goal is to estimate the \mathcal{H}_∞ -norm of an unknown system from input-output data.

The remainder of this paper is organized as follows: Section 2 formalizes the MAB problem and its relation to linear bandits. Section 3 introduces the notion of uniformly efficient policies and derives a lower bound for the regret of policies in this class. Section 4 presents WTS together with its proof of optimality, while Section 5 presents the application to the gain estimation problem in control and system identification. Conclusions are presented in Section 6. For brevity, all the proofs appear in the supplementary material.

1.1 Related work

Fundamental limitations on the frequentist regret depend on the precise structure of the problem, and on the prior knowledge available to the agent, which was originally formalized by [12] for a general class of problems. The framework introduced in this work is very similar to the bandit problem with Gaussian rewards [13], which is one of the most fundamental stochastic MABs. The first finite-time analysis for 1-parameter Gaussian bandits appeared in [14], where UCB-Normal was introduced and proved to attain logarithmic regret. Later, k1-UCB was introduced in [15], where it was proved to match the lower bound of Lai and Robbins [3]. Regarding the 2-parameter Gaussian bandit, [16]

introduced UCB-V, an algorithm that uses variance estimates to construct confidence intervals, which is shown to attain logarithmic but suboptimal regret. Using a similar idea, [17] proved the optimality of their improved version named ISM.

Thompson Sampling (TS) [1] is a famous heuristic algorithm that was re-discovered in the last decade due to its excellent empirical performance [18]. The first thorough theoretical analysis of TS was carried out in [19], where a logarithmic (yet suboptimal) upper bound was derived for the regret incurred by TS for bounded distributions. A matching upper bound for TS was derived in [20] for Bernoulli bandits, which was later extended in [21] to 1-dimensional Gaussian bandits under a Jeffreys prior. [22] has developed a careful analysis for TS under a 2-parameter (mean and variance) Gaussian bandit, concluding that its optimality crucially depends on the choice of the prior. In fact, the algorithm does not achieve optimality when, for example, a Jeffreys prior is employed, nor even a logarithmic asymptotic regret. Problem-independent upper bounds for TS under bounded rewards can be found in [23].

The idea of playing more than one arm per round given a certain budget is not new in machine learning. Examples of problems where the budget is a stopping time appear in [24, 25], however none of them considers weighted information being fed back to the agent: the feedback is the same as in a stochastic MAB that samples many arms per round. The first 1-parameter Gaussian lower bound for a bandit with weighted information was presented in [10], together with WTS. Later, TS was proven to be asymptotically optimal for the 2-parameter Gaussian bandit, but when the outcomes are 2-dimensional and the regret is a nonlinear function of these outcomes [11].

With respect to the application to gain estimation in control, some iterative approaches based on the power-iterations method of numerical linear algebra have been already proposed in the control community [26, 27], where the input signal is allowed to be designed in a sequential manner based on data from previous rounds. The specific problem of \mathcal{H}_∞ -norm estimation has gained some attention in computer science; for example, [28] has derived sharp asymptotic bounds on the error incurred by a method that firstly fits an FIR (finite impulse response) filter of L coefficients to N -length data, in terms of N .

2 Problem formulation

Consider a stochastic multi-armed bandit problem (MAB) in which an agent can sample from K different independent distributions $(\nu_k^t)_{k=1}^K$ at every round $t \in \mathbb{N}$. The means of these distributions are unknown but fixed and denoted as $\boldsymbol{\mu} := (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_K)$, where $\boldsymbol{\mu}_k \in \mathbb{R}^2$, for each $k = 1, \dots, K$; however, unlike in a standard stochastic MAB setup, the variance of these distributions will depend on the action taken by the agent. We define $k^* := \arg \min_{k=1, \dots, K} \|\boldsymbol{\mu}_k\|$, $\boldsymbol{\mu}^* := \boldsymbol{\mu}_{k^*}$, and the gaps $\Delta_k := \|\boldsymbol{\mu}^*\| - \|\boldsymbol{\mu}_k\|$, where $\|\cdot\|$ is the Euclidean norm. At every round $t \in \mathbb{N}$, the agent chooses a *power profile*¹ $p^{t, \boldsymbol{\pi}} := (p_1^{t, \boldsymbol{\pi}}, p_2^{t, \boldsymbol{\pi}}, \dots, p_K^{t, \boldsymbol{\pi}})$ according to policy $\boldsymbol{\pi}$, and it receives a random outcome tuple $X^t := (X_1^t, X_2^t, \dots, X_K^t)$, where $X_k^t \in \mathbb{R}^2$ satisfies

$$X_k^t \sim \mathcal{N} \left(\boldsymbol{\mu}_k, \frac{\sigma_k^2}{2p_k^{t, \boldsymbol{\pi}}} \mathbf{I}_2 \right) \quad (1)$$

for a fixed and unknown vector of variances $\boldsymbol{\sigma} := (\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2)$, where \mathbf{I}_2 denotes the identity matrix in $\mathbb{R}^{2 \times 2}$. In other words, the Gaussian outcome distribution ν_k^t is parametrized by the amount of power $p_k^{t, \boldsymbol{\pi}}$ allocated to that particular arm, which only affects the variance of the distribution. The power profiles are restricted to fulfil $p^{t, \boldsymbol{\pi}} \in \Lambda$ at every round, where Λ is the simplex $\Lambda := \{\xi \in [0, 1]^K : \sum_{k=1}^K \xi_k = 1\}$. The selection of $\boldsymbol{\pi}$ may be adaptive and depend on the power levels and observed outcomes in previous experiments. Let $\boldsymbol{\Pi}$ denote the set of policies $\boldsymbol{\pi}$ such that the power profile $p^{t, \boldsymbol{\pi}}$ assigned under policy $\boldsymbol{\pi}$ at experiment t is \mathcal{F}_t -measurable, where \mathcal{F}_t is the sigma-algebra generated by the previous observations $(p^{1, \boldsymbol{\pi}}, X^1, p^{2, \boldsymbol{\pi}}, X^2, \dots, p^{t-1, \boldsymbol{\pi}}, X^{t-1})$.

We are interested on finding a policy $\boldsymbol{\pi} \in \boldsymbol{\Pi}$ minimizing the frequentist regret of not allocating all the power at arm k^* (having the largest mean in magnitude). To this end, we introduce the following

¹The dependence of the power profiles on $\boldsymbol{\pi}$ is made explicit in the notation, however, we might drop the superscript when this dependence is clear.

notion of expected cumulative regret:

$$\mathbb{E} \{R_{\mu, \sigma}^{\pi}(T)\} = \sum_{t=1}^T \sum_{k=1}^K \|\mu_k\| (p_k^{k, \star} - \mathbb{E} \{p_k^{t, \pi}\}), \quad (2)$$

where $p^{t, \star}$ is the power profile assigned by an oracle policy that knows k^* in hindsight, *i.e.*, assigning $p_{k^*}^{t, \star} = 1$ and $p_k^{t, \star} = 0$, for every suboptimal arm $k \neq k^*$, at every round $t \in \mathbb{N}$. Observe that the dependence of the regret on μ and σ is made explicit. We formalize this problem as the regret minimization (RM) problem:

$$(RM) \quad \min_{\pi \in \Pi} \mathbb{E} \{R_{\mu, \sigma}(T)\}.$$

For the sake of comparison, we introduce a sub-class of policies in Π that sample only one arm k_t at each round $t \in \mathbb{N}$, by applying the power profile $p_{k_t}^{t, \pi} = 1$ and $p_k^{t, \pi} = 0$ for every other arm $k \neq k_t$. This class is denoted as $\Pi_{NS} \subset \Pi$, where NS stands for non-spreading, in the sense that our problem considers bandits that can “spread” the available budget to many arms per round. Certainly, algorithms in Π_{NS} correspond to traditional MAB policies where only one arm per round is sampled.

2.1 Relation to linear bandits

The expected cumulative regret in (2) resembles the one measuring the performance of a linear bandit, *i.e.*, where the reward is a noisy linear function of the distribution means [13]. To see this, define the reward perceived by the agent at every round $t \in \mathbb{N}$ as the linear combination

$$\sum_{k=1}^K p_k^{t, \pi} (\|X_k^t\|^2 - \sigma_k^2), \quad (3)$$

where we subtract the variance of the outcomes because we are not interested in that part of the (expected) reward generated by the variance of the outcomes. The agent then aims to maximize the expected cumulative reward

$$\mathbb{E} \left\{ \sum_{t=1}^T \sum_{k=1}^K p_k^{t, \pi} (\|X_k^t\|^2 - \sigma_k^2) \right\} = \sum_{t=1}^T \sum_{k=1}^K \mathbb{E} \{p_k^{t, \pi}\} \|\mu_k\|^2, \quad (4)$$

which is maximized by the oracle policy. Since such an optimal policy is independent of the exponent in $\|\mu_k\|^2$, one can instead consider $\max_{\pi \in \Pi} \sum_{t=1}^T \mathbb{E} \{p_k^{t, \pi}\} \|\mu_k\|$, which is equivalent to solving (RM).

The difference between solving (RM) and a traditional linear bandit lies in the more detailed feedback we receive in (RM) at every round. More precisely, the agent measures each of the outcomes in the linear combination (3) which provides information about all distributions $(\nu_k^t)_k$ simultaneously. As shown in Section 3, this crucial difference allows us to dramatically reduce the rate of growth of the expected cumulative regret $\mathbb{E} \{R_{\mu, \sigma}^{\pi}(T)\}$ in (2) from $\Theta(\sqrt{T})$ [29, 13] to $\Theta(\log T)$. The problem addressed in this work thus presents a new instance of a stochastic MAB that lies between linear bandits and traditional MABs, since the cost function is linear in the actions but the information received by the agent resembles the one collected in a traditional MAB after many rounds.

3 Regret lower bound

In this section we introduce a class of algorithms that formally address the exploration-exploitation dilemma successfully enough in order to quickly discriminate the optimal arm. Following the notion originally introduced by [3] for traditional MABs, we call this family the set of *uniformly efficient policies*. In the context of bandits with weighted information, this notion has already been introduced in [10], and it corresponds to the spreading version of the one in [3].

Definition 1 (Uniform efficiency). *A policy π is said to be uniformly efficient if the cumulative power $p_k^{t, \pi}$ allocated at every suboptimal arm $k \neq k^*$ satisfies $\mathbb{E} \{\sum_{t=1}^T p_k^{t, \pi}\} = o(T^\alpha)$, for every $\alpha > 0$. The set of uniformly efficient strategies is denoted as Π^* . The subset of non-spreading but uniformly efficient policies in Π_{NS} is denoted as Π_{NS}^* . \triangle*

Table 1: Regret lower bounds.

	$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}\{R_{\mu, \sigma}^{\pi}(T)\}}{\log T} \geq$	
	Known $\sigma = (\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2)$	Unknown $\sigma = (\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2)$
$\pi \in \Pi_{\text{NS}}^*$	$\sum_{k \neq k^*} \frac{\sigma_k^2}{\Delta_k}$ [30]	$\sum_{k \neq k^*} \frac{\Delta_k}{\log\left(1 + \frac{\Delta_k^2}{\sigma_k^2}\right)}$ [11]
$\pi \in \Pi^*$	$\sum_{k \neq k^*} \frac{\sigma_k^2}{\Delta_k}$ [10]	$\sum_{k \neq k^*} \frac{\sigma_k^2}{\Delta_k}$

In their seminal paper [3], Lai and Robbins derived a tight asymptotic (in T) regret lower bound satisfied by every non-spreading uniformly efficient policy. Following this idea, [11] derived a lower bound for the regret in (2), for policies in Π_{NS}^* , showing that $\liminf_{T \rightarrow \infty} \mathbb{E}\{R_{\mu, \sigma}^{\pi}(T)\} / \log T = \sum_{k \neq k^*} \sigma_k^2 / \Delta_k$ when σ is known, and that $\liminf_{T \rightarrow \infty} \mathbb{E}\{R_{\mu, \sigma}^{\pi}(T)\} / \log T = \sum_{k \neq k^*} \Delta_k / \log(1 + \Delta_k^2 / \sigma_k^2)$ when σ is unknown. This result suggests that it becomes easier for the agent to identify suboptimal arms when the variance of the noise in the outcomes is known in hindsight. For bandits under weighted information, [10] had already reported that $\liminf_{T \rightarrow \infty} \mathbb{E}\{R_{\mu, \sigma}^{\pi}(T)\} / \log T = \sum_{k \neq k^*} \sigma_k^2 / \Delta_k$ for every $\pi \in \Pi^*$ when σ is known, suggesting that there may not be any performance improvement (asymptotically) when spreading strategies are preferred. The following result describes the missing lower bound: when σ is unknown, one can recover the same lower bound when resources can be spread among arms, as if σ were revealed to the agent in hindsight.

Theorem 1 (Regret lower bound). *Consider (RM) with unknown parameters μ and $\sigma = (\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2)$. Under every uniformly efficient policy $\pi \in \Pi^*$, the expected cumulative regret, as defined in (2), satisfies*

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}\{R_{\mu, \sigma}^{\pi}(T)\}}{\log T} \geq \sum_{k \neq k^*} \frac{\sigma_k^2}{\Delta_k}. \quad (5)$$

△

The result in Theorem 1 completes the set of lower bounds for policies in $(\Pi^*, \Pi_{\text{NS}}^*)$ when σ is known or not to the agent. It also formalizes the idea that sampling many arms simultaneously at every round can lead to better performance, in terms of cumulative expected regret, when the variances in σ are unknown. Table 1 puts this result in perspective by showing the different regret lower bounds under known/unknown variance for uniformly efficient policies in the set of spreading and non-spreading strategies Π_{NS}^* and Π^* , respectively. The lower-right entry, originally unknown, is filled by Theorem 1.

4 The Weighted Thompson Sampling (WTS) algorithm

Weighted Thompson Sampling (WTS) is a spreading strategy in Π originally introduced in [10]. This algorithm is based on the algorithm by Thompson [1], called Thompson Sampling (TS), and it corresponds to a Bayesian policy that condenses gathered information in the form of posterior distributions of each arm being optimal but, instead of randomizing the action (as TS does), WTS uses this posterior distribution as a power profile p^t . In this section we develop a thorough analysis of WTS which overlaps with some of the ideas used to study TS for the solution of (RM) in [11].

Just as with TS, WTS starts with a prior distribution $\rho^1 = (\rho_k^1)_{k=1}^K$ encoding our confidence on each arm being the optimal arm, *i.e.*, $\rho_k^1 = \text{Prob}\{\tilde{k}_1^* = k \mid \mathcal{F}_1\}$, where \tilde{k}_t^* is a random variable whose distribution (conditioned to \mathcal{F}_t) captures our confidence of k being optimal given the history, and $\mathcal{F}_1 = \{\phi, \Omega\}$ (*i.e.*, no data). Recall that \mathcal{F}_t denotes the sigma algebra generated by $(p^{1, \pi}, X^1, p^{2, \pi}, X^2, \dots, p^{t-1, \pi}, X^{t-1})$ under policy $\pi \in \Pi$. At every round $t \in \mathbb{N}$, WTS applies the power profile $p^{t, \text{WTS}} = \rho^t$ and then updates its belief of each arm's optimality ρ^{t+1} based on the history up to round t . Obtaining ρ^t in closed form is, in general, not possible. This can be overcome by introducing the posterior means

$\tilde{\boldsymbol{\mu}}_k(t) \sim f_{\boldsymbol{\mu}_k | \mathcal{F}_t}$, where $f_{\boldsymbol{\mu}_k | \mathcal{F}_t}$ is the posterior distribution on the mean $\boldsymbol{\mu}_k$ given the data up to round $t - 1$. As discussed in [31], if $\|\tilde{\boldsymbol{\mu}}^*(t)\| := \max_k \|\tilde{\boldsymbol{\mu}}_k(t)\|$, then

$$\rho_k^t := \mathbb{P}\{\tilde{k}_t^* = k | \mathcal{F}_t\} = \mathbb{E}\{\mathbb{1}\{\tilde{k}_t^* = k\} | \mathcal{F}_t\} = \mathbb{P}\{\|\tilde{\boldsymbol{\mu}}^*(t)\| = \|\tilde{\boldsymbol{\mu}}_k(t)\| | \mathcal{F}_t\}, \quad (6)$$

which means that ρ^t is completely determined by $(f_{\boldsymbol{\mu}_1 | \mathcal{F}_t}, f_{\boldsymbol{\mu}_2 | \mathcal{F}_t}, \dots, f_{\boldsymbol{\mu}_K | \mathcal{F}_t})$. It is important to mention that finding a closed form for the mapping from $(f_{\boldsymbol{\mu}_1 | \mathcal{F}_t}, f_{\boldsymbol{\mu}_2 | \mathcal{F}_t}, \dots, f_{\boldsymbol{\mu}_K | \mathcal{F}_t})$ to ρ^t involves the calculation of very complicated K -dimensional integrals, that can however be approximated by Monte Carlo calculations [47], as we explain in Appendix E.

It is clear that a key component in WTS is to be able to obtain $f_{\boldsymbol{\mu}_k | \mathcal{F}_t}$ at every round $t \in \mathbb{N}$, and a method to achieve this task involves computing sufficient statistics for the distribution of the outcomes. To obtain $f_{\boldsymbol{\mu}_k | \mathcal{F}_t}$ we select a prior distribution representing our total ignorance on what the real value of $(\boldsymbol{\mu}, \boldsymbol{\sigma})$ is by considering a uniform improper distribution $f_{\boldsymbol{\mu}_k, \sigma_k^2 | \mathcal{F}_1}(\mathbf{m}, \zeta^2) \propto 1$, for every $(\mathbf{m}, \zeta^2) \in \mathbb{R}^2 \times (0, \infty)$ and for every $k = 1, \dots, K$. The posterior $f_{\boldsymbol{\mu}_k | \mathcal{F}_t}(\mathbf{m})$ is then obtained as $\int_0^\infty f_{\boldsymbol{\mu}_k, \sigma_k^2}(\mathbf{m}, \zeta^2) d\zeta^2$. This paper shows, in Section 4.2, that WTS is optimal under such a prior, even though other choices of priors may also attain optimality [22].

4.1 Sufficient Statistics

Here we show that the sample mean, the sample variance and the trajectory of the power profiles $(p^{1, \pi}, p^{2, \pi}, \dots, p^{t, \pi})$ are sufficient statistics for the data collected up to round t . We additionally characterize the distribution of the sample mean, the sample variance, and the posterior mean distribution $f_{\boldsymbol{\mu}_k | \mathcal{F}_t}$. These results are presented in the following lemma.

Lemma 1. *For each arm $k \in \{1, \dots, K\}$,*

$$\bar{\mathbf{x}}_k(t) := \frac{\sum_{\ell=1}^t p_k^\ell X_k^\ell}{\sum_{\ell=1}^t p_k^\ell}, \quad S_k(t) := \sum_{\ell=1}^t p_k^\ell \|X_k^\ell - \bar{\mathbf{x}}_k(t)\|^2 \quad (7)$$

are sufficient statistics for $(X_k^\ell)_{\ell=1}^t$ conditioned on $p_k^1, p_k^2, \dots, p_k^t$. Furthermore, conditioned on the trajectory of the power profiles, $\bar{\mathbf{x}}_k(t) \sim \mathcal{N}(\boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I}_2 / (2 \sum_{\ell=1}^t p_k^\ell))$ and $2S_k(t) / \sigma_k^2 \sim \chi_{2(t-1)}^2$ are statistically independent. \triangle

The following result states that the posterior means follow a bivariate t-student distribution [32] that is symmetrically distributed around the empirical mean $\bar{\mathbf{x}}_k(t)$ with variance proportional to $S_k(t)/t$, at every round $t \in \mathbb{N}$ and for every arm $k = 1, \dots, K$.

Lemma 2. *Consider the improper uniform prior distribution $f_{\boldsymbol{\mu}_k, \sigma_k^2 | \mathcal{F}_1} \propto 1$. Then, the posterior distribution $f_{\boldsymbol{\mu}_k | \mathcal{F}_t}$ at round $t \geq 4$ is*

$$f_{\boldsymbol{\mu}_k | \mathcal{F}_t}(\tilde{\boldsymbol{\mu}}) = \frac{\sum_{\ell=1}^{t-1} p_k^\ell (t-3)}{\pi S_k(t-1)} \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell \|\tilde{\boldsymbol{\mu}} - \bar{\mathbf{x}}_k(t-1)\|^2}{S_k(t-1)} \right)^{-t+2}, \quad (8)$$

for every $k = 1, \dots, K$. \triangle

We can now formalize the WTS algorithm introduced at the beginning of this section by using Lemmas 1 and 2 to update ρ^t in a tractable fashion. In fact, having a closed-form expression for the posterior means in (8) allows us to approximate ρ^t by means of (6) and Monte Carlo simulations, as explained in Appendix E. We summarize WTS in Algorithm 1, where we remark that, by Lemma 2, each arm needs to be sampled at least 3 times.

4.2 Optimality of WTS

Here we provide a theoretical analysis of WTS, showing that its upper bound matches the regret lower bound predicted by Theorem 1.

Theorem 2 (Upper bound). *Consider (RM) under unknown $(\boldsymbol{\mu}, \boldsymbol{\sigma})$, and the improper uniform prior distribution $f_{\boldsymbol{\mu}_k, \sigma_k^2 | \mathcal{F}_1} \propto 1$. Then, the regret, as defined in (2), incurred by WTS satisfies*

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}\{R_{\boldsymbol{\mu}, \boldsymbol{\sigma}}^{\text{WTS}}(T)\}}{\log T} \leq \sum_{k \neq k^*} \frac{\sigma_k^2}{\Delta_k}. \quad (9)$$

\triangle

Algorithm 1 WTS: proposed implementation via sufficient statistics

- 1: Input: $T, \rho^1 = (1/K, 1/K, \dots, 1/K)$ (prior distribution for \tilde{k}_1^*)
 - 2: **for** $t = 1$ to 3 **do**
 - 3: Apply the power profile $p^{t, \text{WTS}} = \rho^1$ and collect the outcome X^t
 - 4: Update the sufficient statistics $\bar{x}_k(t)$ and $S_k(t)$ according to Lemma 1
 - 5: **end for**
 - 6: Obtain $(f_{\mu_1 | \mathcal{F}_4}, f_{\mu_2 | \mathcal{F}_4}, \dots, f_{\mu_K | \mathcal{F}_4})$ via Lemma 2 and compute ρ^3
 - 7: **for** $t = 4$ to T **do**
 - 8: Apply the power profile $p^{t, \text{WTS}} = \rho^1$ and collect the outcome X^t
 - 9: Update the sufficient statistics $\bar{x}_k(t)$ and $S_k(t)$ according to Lemma 1
 - 10: Obtain $(f_{\mu_1 | \mathcal{F}_{t+1}}, f_{\mu_2 | \mathcal{F}_{t+1}}, \dots, f_{\mu_K | \mathcal{F}_{t+1}})$ via Lemma 2 and compute ρ^{t+1}
 - 11: **end for**
-

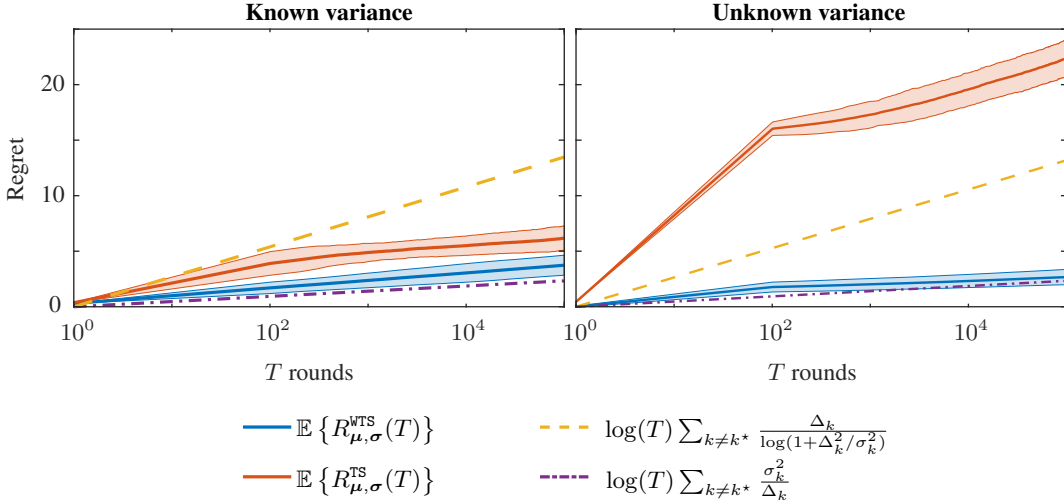


Figure 1: Summary of the performances attained by WTS and TS, in the light of the lower bounds of Table 1, under known and unknown σ . The shaded areas represent 2 standard deviations. The implementation of WTS and TS under known σ is described in Appendix F.

The reason why WTS is able to attain the same performance regardless its knowledge of σ follows from the fact that, for every $k = 1, \dots, K$, the outcomes from arm k provide the same information about σ_k^2 despite how small p_k^t is at round t . In fact, the sample variance $2S_k(t)/t$ concentrates exponentially around σ_k^2 (much faster than the sample means $\bar{x}_k(t)$), meaning that WTS can quickly make decisions as if it knew the value of σ in hindsight.

4.3 Simulation study 1: regret of WTS and TS under known/unknown variance

We consider an arbitrary system with prefixed values for μ, σ, K , on which we test TS and WTS for known and unknown σ . The values of μ, σ, K are chosen so as to highlight the different lower bounds in Table 1, and their values are detailed in Appendix F. We run 300 Monte Carlo simulations to approximate the expected cumulative regret for each of them during $T = 10^5$ rounds. In line with the results derived in [10, 11], the simulations, depicted in Figure 1, show that TS and WTS are matching algorithms for each of the two setups (known and unknown σ), *i.e.*, their asymptotic expected cumulative regret matches each of the lower bounds summarized in Table 1. When σ is known (left), the asymptotic expected regret is the same for TS (red) and WTS (blue), meaning that there is no improvement, asymptotically, in spreading the resources at every round. Both algorithms match the same lower bound (purple) prescribed for both Π^* and Π_{NS}^* under known σ [10]. When the noise is unknown (right), such lower bounds are different, where we can observe that TS (red) and WTS (blue) match their respective lower bounds (yellow and purple, respectively).

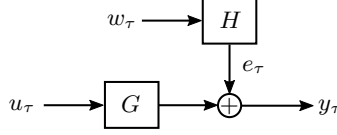


Figure 2: Mathematical model of a linear system with additive non-white noise.

5 Application to gain estimation of dynamical systems

In this section we study the extension of WTS to an important problem in control and system identification [33, 34]. One of the fundamental steps in designing a model-based controller for a dynamical linear and time-invariant (LTI) system is to estimate the mismatch between the real system and its LTI model in terms of the largest gain it can induce, also known as the induced ℓ_2 -gain or \mathcal{H}_∞ -norm of the system [35]. The system depicted in Figure 2 is suitable for studying the problem of estimating such a quantity from input-output data. To see this, let G_0 and \hat{G} denote a system and its model, respectively, and define $G = G_0 - \hat{G}$. Then measurements from G can be collected by exciting the real system G_0 and its model \hat{G} with the same input and then subtracting their respective outputs.

Following the setup described in [10, 36], let $G(e^{j\omega})$ and $H(e^{j\omega})$ denote the transfer functions of the unknown, causal, stable LTI systems G and H in Fig. 2, where w is Gaussian $\mathcal{N}(0, 1)$ white noise, and $j := \sqrt{-1}$. We are interested on estimating the \mathcal{H}_∞ -norm of G , defined as $\beta := \|G\|_\infty = \max_{\omega \in [0, \pi]} |G(e^{j\omega})|$. It is well known [37] that estimating β from input-output data crucially depends on the input signal. In fact, applying a sinusoidal of frequency equal to $\arg \max_{\omega \in [0, \pi]} |G(e^{j\omega})|$ provides the best possible data for this task because it maximizes the ratio $\|Gu\|/\|u\|$ [38, 39]. In line with this argument, we study the gain estimation problem as two sub-problems: the one of generating the input signal and that of using the collected data in a point-estimator of β . The data is collected in a sequence of independent experiments², where at each round $t \in \mathbb{N}$, the agent is allowed to design the input signal $u^t = (u_\tau^t)_{\tau=0}^{N-1}$ and collect the output $y^t = (y_\tau^t)_{\tau=0}^{N-1}$, for some prefixed N .

We delegate the problem of optimally generating the input signal u^t in an adaptive fashion to WTS. Since G is LTI, we restrict the class of input signals to a sum of sinusoids (*i.e.*, a multisine) of K predetermined frequencies $(2\pi k/(2K+1))_{k=1}^K$ which discretize the frequency range $[0, \pi]$ into K equispaced frequencies, where the agent can choose the amplitude of each sinusoid. This can be achieved by choosing $N = 2K + 1$ and then letting the agent design u^t in the frequency domain. More specifically, if $U^t(j\omega) := \sum_{\tau=0}^{N-1} u_\tau^t e^{-j\omega\tau}$ denotes the discrete Fourier transform [40] of u^t , then the agent designs $|U^t(j\omega_k)|^2 = p_k^t$ at the frequencies of interest $\omega_k = 2\pi k/(2K+1)$. By applying an input u^t satisfying³ $|U^t(j\omega_k)|^2 = p_k^t$, the input output data (u^t, y^t) satisfies

$$\frac{Y^t(j\omega_k)}{U^t(j\omega_k)} = G(e^{j\omega_k}) + \frac{E^t(j\omega_k)}{U^t(j\omega_k)}, \quad (10)$$

where $E^t(j\omega_k)$ is a circularly symmetric complex zero-mean and white sequence whose real and imaginary parts are independent for every $k = 1, \dots, K$ and $t \in \mathbb{N}$ [41]. Moreover, its real and imaginary parts are Gaussian with variance $|H(e^{j\omega_k})|^2/2$. The $(2K+1)$ -length noisy output is completely revealed to the agent after every experiment, who can then define the sequence of outcomes $X_k^t := [\text{Re}\{Y^t(j\omega_k)/U^t(j\omega_k)\} \quad \text{Im}\{Y^t(j\omega_k)/U^t(j\omega_k)\}]^\top$, $k = 1, \dots, K$. In consequence, $X_k^t \sim \mathcal{N}(\mu_k, \sigma_k^2/(2p_k^t)\mathbf{I}_2)$, where $\mu_k = [\text{Re}G(e^{j\omega_k}) \quad \text{Im}G(e^{j\omega_k})]^\top$ and $\sigma_k^2 := |H(e^{j\omega_k})|^2$, which recovers the bandit problem under weighted information defined in Section 2. The goal of the agent collecting the data is then to find the power profile p^t that maximizes $\|Gu\|$ in the class of multisine inputs which, for K large enough, reasonably approximates $\|G\|_\infty \approx \max_{k=1, \dots, K} |G(e^{j\omega_k})|$.

The point-estimator for $\beta = \|G\|_\infty$ is independent of the underlying data-collection algorithm and, for the sake of simplicity, it is chosen as $\hat{\beta}_t := \|\hat{x}_{\hat{k}_t}(t)\|$, with $\hat{k}_t = \arg \max_{k=1, \dots, K} \sum_{\ell=1}^t p_k^{\ell, \text{WTS}}$.

²Statistically independent experiments can be achieved by either waiting long enough between experiments (so that the natural response of G due to initial conditions decays exponentially to zero), by using a controller to bring the state of the system to zero, or, if possible, by manually resetting G .

³This can be achieved by, for example, taking the inverse Fourier transform of the sequence $(0, \sqrt{p_1^t}, \sqrt{p_2^t}, \dots, \sqrt{p_K^t}, \sqrt{p_K^t}, \dots, \sqrt{p_2^t}, \sqrt{p_1^t})$, or by setting $u^t = (u_\tau^t)_{\tau=1}^N = \sum_{k=1}^K \sqrt{p_k^t} \sin(\omega_k \tau)$.

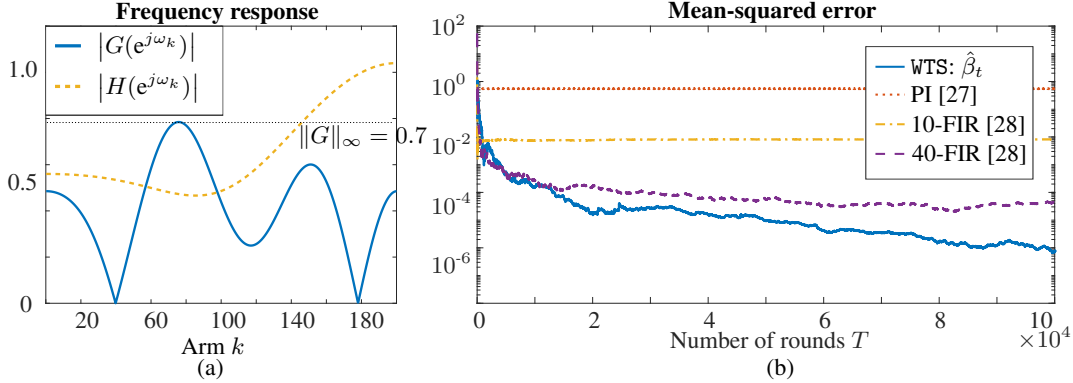


Figure 3: Simulation study in Section 5.1. The left-hand figure (a) depicts the frequency response of filters G and H in Fig. 2, while the right-hand figure (b) shows the performance attained by different state-of-the-art estimators and our proposed method.

5.1 Simulation study 2: WTS in the gain estimation problem

Consider the configuration shown in Fig. 2, where G and H have the frequency responses depicted in Fig. 3a. This simulation study is particularly confusing for any \mathcal{H}_∞ -norm estimation algorithm since the frequency response of G is almost completely covered by the one of H . We implement WTS with $K = 200$ arms (to avoid large discretization errors) with an estimator $\hat{\beta}_t$, together with the power iterations (PI) method [26, 27], and a model-based approach that first derives an FIR model of lengths 10 (10-FIR) and 40 (40-FIR), and then approximates β by the \mathcal{H}_∞ -norm of the model. To average the squared error, we use 10 Monte Carlo simulations, each of them running 10^5 rounds. The results are illustrated in Fig. 3b, where we observe that the proposed method (WTS + $\hat{\beta}_t$) attains the best performance among all algorithms. We observe that, as predicted by [27], PI leads to asymptotically biased estimations under the presence of noise. On the other hand, model-based approaches lead inherently to biased estimations since their model structure does not (in general) include the one generating the data. The main source of confusion for the other three algorithms is the large variance at high frequency, which misleads these methods to believe that large outcomes are a consequence of large gains at high frequency.

6 Conclusions

We have introduced a new instance of the stochastic multi-armed bandit problem in which the agent can play many arms per round by spreading a budget usually allocated to only one arm. This setup is called a bandit problem under weighted information. By spreading the available resources, the agent receives a Gaussian sequence as outcome, where the variance of each entry is inversely-proportional to the amount of resources allocated to that arm. The goal of the agent is to minimize the regret of not allocating all the available resources to the arm that has the largest mean in magnitude. We have derived a tight lower bound for the regret incurred by a wide class of bandit algorithms and introduced WTS (Weighted Thompson Sampling), a Thompson-Sampling-based strategy whose regret upper bound has been shown, in this work, to match the predicted lower bound. We have illustrated the importance of this new bandit problem by means of a fundamental problem in automatic control and system identification, also known as the *gain estimation problem*, where the goal is to find the induced ℓ_2 -gain (or \mathcal{H}_∞ -norm) of an unknown system in a model-free fashion. By arguing that such a problem can be decomposed into a data-collection and point-estimation problems, we have shown that collecting data via WTS and then using a point-estimator on this data outperforms state-of-the-art methods for solving this problem. Future work involves studying the effect of weighted information from non-Gaussian distributions, as well as deriving problem-independent and probability bounds for the regret.

References

- [1] W. R. Thompson, "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples," *Biometrika*, vol. 25, no. 3, pp. 285–294, 1933.
- [2] H. Robbins, "Some aspects of the sequential design of experiments," *Bulletin of the American Mathematical Society*, vol. 58, no. 5, pp. 527–535, 1952.
- [3] T. L. Lai and H. Robbins, "Asymptotically efficient adaptive allocation rules," *Advances in applied mathematics*, vol. 6, pp. 4–22, 1985.
- [4] W. Shen, J. Wang, Y. Gang Jiang, and H. Zha, "Portfolio choices with orthogonal bandit learning," in *Twenty-Fourth International Joint Conference on Artificial Intelligence (IJCAI 2015)*, 2015, pp. 974–980.
- [5] X. Huo and F. Fu, "Risk-aware multi-armed bandit problem with application to portfolio selection," *Royal society open science*, no. 171377, 2017.
- [6] M. S. Talebi, Z. Zou, R. Combes, A. Proutiere, and M. Johansson, "Stochastic online shortest path routing: The value of feedback," *IEEE Transactions on Automatic Control*, vol. 63, no. 4, pp. 915–930, 2018.
- [7] A. Kolobov, S. Bubeck, and J. Zimmert, "Online learning for active cache synchronization," in *37th International Conference on Machine Learning (ICML)*, 2020, pp. 974–980.
- [8] W. H. Press, "Bandit solutions provide unified ethical models for randomized clinical trials and comparative effectiveness research," *Proceedings of the National Academy of Sciences*, vol. 106, no. 52, pp. 22 387–22 392, 2009.
- [9] T. Verstraeten, E. Bargiacchi, P. J. K. Libin, J. Helsen, D. M. Roijers, and A. Nowé, "Multi-agent Thompson Sampling for bandit applications with sparse neighbourhood structures," *Nature, Scientific Reports*, p. 6728, 2020.
- [10] M. I. Müller, P. E. Valenzuela, A. Proutiere, and C. R. Rojas, "A stochastic multi-armed bandit approach to nonparametric \mathcal{H}_∞ -norm estimation," in *Proceedings of the 56th IEEE Conference on Decision and Control (CDC)*, 2017.
- [11] M. I. Müller and C. R. Rojas, "Gain estimation of dynamical linear systems using Thompson Sampling," in *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- [12] A. N. Burnetas and M. N. Katehakis, "Optimal adaptive policies for sequential allocation problems," *Advances in Applied Mathematics*, vol. 17, no. 7, pp. 122–142, 1996.
- [13] T. Lattimore and C. Szepesvári, *Bandit Algorithms*. Cambridge University Press, 2020.
- [14] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, "The nonstochastic multiarmed bandit problem," *SIAM J. Comput.*, vol. 32, no. 1, pp. 48–77, 2002.
- [15] O. Cappé, A. Garivier, O.-A. Maillard, R. Munos, and G. Stoltz, "Kullback-Leibler upper confidence bounds for optimal sequential allocation," *The Annals of Statistics*, vol. 41, no. 3, pp. 1516–1541, 2013.
- [16] J.-Y. Audibert, R. Munos, and C. Szepesvári, "Exploration-exploitation tradeoff using variance estimates in multi-armed bandits," *Theoretical Computer Science*, vol. 410, no. 19, pp. 1876–1902, 2009.
- [17] W. Cowan, J. Honda, and M. N. Katehakis, "Normal bandits of unknown means and variances," *Journal of Machine Learning Research*, vol. 18, no. 154, pp. 1–28, 2018.
- [18] O. Chapelle and L. Li, "An empirical evaluation of Thompson Sampling," in *Advances in Neural Information Processing Systems (NIPS) 24*, 2011.
- [19] S. Agrawal and N. Goyal, "Analysis of thompson sampling for the multi-armed bandit problem," in *Proceedings of the 25th Annual Conference on Learning Theory (COLT)*, 2012, pp. 39.1–39.26.

- [20] E. Kaufmann, N. Korda, and R. Munos, “Thompson Sampling: An asymptotically optimal finite-time analysis,” in *Proceedings of the 23rd International Conference on Algorithmic Learning Theory (ALT)*. Springer, 2012, pp. 199–213.
- [21] N. Korda, E. Kaufmann, and R. Munos, “Thompson Sampling for 1-dimensional exponential family bandits,” in *Advances in Neural Information Processing Systems (NIPS) 26*, 2013.
- [22] J. Honda and A. Takemura, “Optimality of Thompson Sampling for Gaussian bandits depends on priors,” in *Proceedings of the 17th international conference on Artificial Intelligence and Statistics (AISTATS)*, 2014.
- [23] S. Agrawal and N. Goyal, “Further optimal regret bounds for thompson sampling,” *Journal of the ACM*, vol. 64, 09 2012.
- [24] D. P. Zhou and C. J. Tomlin, “Budget-constrained multi-armed bandits with multiple plays,” in *Proceedings of The Thirty-Second Association for the Advancement of Artificial Intelligence Conference (AAAI)*, 2018, pp. 4572–4579.
- [25] L. Tran-Thanh, A. Chapman, E. Munoz de Cote, A. Rogers, and N. Jennings, “Epsilon–first policies for budget–limited multi-armed bandits,” in *Proceedings of The Twenty-Fourth Association for the Advancement of Artificial Intelligence Conference (AAAI)*, 2010, pp. 1211–1216.
- [26] B. Wahlberg, M. B. Syberg, and H. Hjalmarsson, “Non-parametric methods for \mathcal{L}_2 -gain estimation using iterative experiments,” *Automatica*, vol. 46, no. 8, pp. 1376 – 1381, 2010.
- [27] C. R. Rojas, T. Oomen, H. Hjalmarsson, and B. Wahlberg, “Analyzing iterations in identification with application to nonparametric \mathcal{H}_∞ -norm estimation,” *Automatica*, vol. 48, no. 11, pp. 2776–2790, 2012.
- [28] S. Tu, R. Boczar, and B. Recht, “On the approximation of Toeplitz operators for nonparametric \mathcal{H}_∞ -norm estimation,” in *Proceedings of the American Control Conference (ACC)*, 2018.
- [29] V. Dani, T. P. Hayes, and S. M. Kakade, “The price of bandit information for online optimization,” in *Proceedings of the 21st Annual Conference on Learning Theory (COLT)*, 2008, pp. 355–366.
- [30] S. Bubeck and N. Cesa-Bianchi, “Regret analysis of stochastic and nonstochastic multi-armed bandit problems,” *Foundations and Trends in Machine Learning*, vol. 5, no. 1, pp. 1–122, 2012.
- [31] E. Kaufmann, “Analyse de stratégies bayésiennes et fréquentistes pour l’allocation séquentielle de ressources,” Ph.D. dissertation, Paris Institute of Technology, Paris, France, 2014.
- [32] A. Papoulis and U. Pillai, *Probability, Random Processes, and Stochastic Processes*, 4th ed. McGraw-Hill, 2002.
- [33] K. Zhou, J. C. Doyle, and K. Glover, *Robust and Optimal Control*. Prentice-Hall, 1996.
- [34] L. Ljung, *System Identification: Theory for the User*, 2nd ed. Prentice Hall, 1999.
- [35] A. van der Schaft, *L_2 Gain and Passivity Techniques in Nonlinear Control*, 3rd ed. Springer, 2017.
- [36] M. I. Müller, J. Milošević, H. Sandberg, and C. R. Rojas, “A risk-theoretical approach to \mathcal{H}_2 -optimal control under covert attacks,” in *Proceedings of the 57th IEEE Conference on Decision and Control (CDC)*, 2018, pp. 4553–4558.
- [37] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice-Hall, 1993.
- [38] F. W. Fairman, *Linear Control Theory: The State Space Approach*. John Wiley & Sons, 1998.
- [39] T. Kailath, *Linear Systems*. Prentice Hall, 1980.
- [40] A. V. Oppenheim and R. W. Schaffer, *Discrete-time Signal Processing*. Prentice-Hall, 2009.

- [41] J. C. Agüero, J. I. Yuz, G. C. Goodwin, and R. Delgado, “On the equivalence of time and frequency domain maximum likelihood estimation,” *Automatica*, vol. 46, no. 2, pp. 260–270, 2010.
- [42] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Wiley-Interscience, 2006.
- [43] G. H. Hardy, J. E. Littlewood, and G. Pólya, *Inequalities*, 2nd ed. Cambridge University Press, 1988.
- [44] R. Durrett, *Probability: Theory and Examples*, 4th ed. Cambridge University Press, 2010.
- [45] R. A. Horn and C. R. Johnson, *Matrix Analysis*, 2nd ed. Cambridge University Press, 2012.
- [46] T. Söderström, *Discrete-time Stochastic Systems*. Springer, 2002.
- [47] G. Fishman, *Monte Carlo: Concepts, Algorithms, and Applications*. Springer, 1996.

Appendix

Table of Contents

A Proof of Theorem 1: regret lower bound	13
B Proof of Lemma 1: sufficient statistics	15
B.1 Sufficiency of the statistics	16
B.2 Distributions and statistical independence	16
C Proof of Lemma 2: posterior mean distribution	18
D Proof of Theorem 2: regret upper bound	18
D.1 Complementary upper bounds	20
D.2 Technical Lemmas	28
E Approximating ρ^t via Monte Carlo simulations	32
F Details of simulation study 1 in Section 4.3	33

A Proof of Theorem 1: regret lower bound

Before proving this result, we need to introduce the following technical lemma [10], which we re-state here for completeness, together with its respective proof.

Lemma 3. *Let ν and ν' denote two different distributions for a random variable, and let \mathbb{P} and \mathbb{P}' denote the respective probability measures. Then, for every measurable event B , the KL-divergence between ν and ν' satisfies*

$$\mathbb{D}\{\nu||\nu'\} \geq d(\mathbb{P}\{B\}, \mathbb{P}'\{B\}), \quad (11)$$

where $d(q, s) = q \log(q/s) + (1 - q) \log((1 - q)/(1 - s))$ is the binary entropy function [42], corresponding to the KL-divergence between two Bernoulli variables of means q and s . \triangle

Proof. Let $\mathbb{E}\{\cdot\}$ and $\mathbb{E}'\{\cdot\}$ denote the expectation operator under \mathbb{P} and \mathbb{P}' , respectively. Introducing a change of measure, and by Jensen's inequality [43], we have that

$$\begin{aligned} \mathbb{P}'\{B\} &= \mathbb{E}'\{\mathbb{1}\{B\}\} = \int \left(\frac{d\mathbb{P}'}{d\mathbb{P}} \right) \mathbb{1}\{B\} d\mathbb{P} = \int e^{\ell(\nu') - \ell(\nu)} \mathbb{1}\{B\} d\mathbb{P} \\ &= \mathbb{E}\left\{ e^{\ell(\nu') - \ell(\nu)} \mathbb{1}\{B\} \right\} \\ &= \mathbb{E}\left\{ e^{\ell(\nu') - \ell(\nu)} \mid B \right\} \mathbb{P}\{B\} \\ &\geq e^{\mathbb{E}\{\ell(\nu') - \ell(\nu) \mid B\}} \mathbb{P}\{B\}, \end{aligned} \quad (12)$$

from which we conclude that

$$\mathbb{E}\{\ell(\nu) - \ell(\nu') \mid B\} \geq \log \frac{\mathbb{P}\{B\}}{\mathbb{P}'\{B\}}. \quad (13)$$

It now follows that

$$\begin{aligned}
\mathbb{D}\{\nu|\nu'\} &= \mathbb{E}\{\ell(\nu) - \ell(\nu')\} \\
&= \mathbb{E}\{\ell(\nu) - \ell(\nu')|B\}\mathbb{P}\{B\} + \mathbb{E}\{\ell(\nu) - \ell(\nu')|B^c\}\mathbb{P}\{B^c\} \\
&\geq \mathbb{P}\{B\} \log \frac{\mathbb{P}\{B\}}{\mathbb{P}'\{B\}} + \mathbb{P}\{B^c\} \log \frac{\mathbb{P}\{B^c\}}{\mathbb{P}'\{B^c\}} \\
&= \mathbb{P}\{B\} \log \frac{\mathbb{P}\{B\}}{\mathbb{P}'\{B\}} + (1 - \mathbb{P}\{B\}) \log \frac{1 - \mathbb{P}\{B\}}{1 - \mathbb{P}'\{B\}} \\
&= d(\mathbb{P}\{B\}, \mathbb{P}'\{B\}).
\end{aligned} \tag{14}$$

To prove the regret lower bound, we follow the standard approach by [3, 12].

Without loss of generality, let $\boldsymbol{\mu} := (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)$ satisfy $\|\boldsymbol{\mu}_1\| > \|\boldsymbol{\mu}_k\|$, $k = 2, \dots, K$, meaning that arm 1 is assumed to be the (unique) optimal arm. The log-likelihood function of $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ given the outcomes (X^1, X^2, \dots, X^T) is

$$\begin{aligned}
\ell(\boldsymbol{\mu}, \boldsymbol{\sigma}) &= \log \nu(X^1, \dots, X^T) \\
&= \sum_{t=1}^T \sum_{k=1}^K \log \left(\frac{p_k^{t,\pi}}{\pi \sigma_k^2} e^{-\frac{p_k^{t,\pi}}{\sigma_k^2} \|X_k^t - \boldsymbol{\mu}_k\|^2} \right) \\
&= \sum_{t=1}^T \sum_{k=1}^K \left(-\log \pi + \log p_k^{t,\pi} - \log \sigma_k^2 - \frac{p_k^{t,\pi} \|X_k^t - \boldsymbol{\mu}_k\|^2}{\sigma_k^2} \right).
\end{aligned} \tag{15}$$

Now consider confusion parameters $\boldsymbol{\mu}' := (\boldsymbol{\mu}'_1, \dots, \boldsymbol{\mu}'_K)$ and $\boldsymbol{\sigma}' := (\sigma_1'^2, \dots, \sigma_K'^2)$, and let \mathcal{P}_{t+1} denote the σ -algebra generated by (p^1, p^2, \dots, p^t) . Then, the expected difference of the likelihood function at $(\boldsymbol{\mu}, \boldsymbol{\sigma})$ and at $(\boldsymbol{\mu}', \boldsymbol{\sigma}')$ satisfies

$$\begin{aligned}
&\mathbb{E}_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \{ \ell(\boldsymbol{\mu}, \boldsymbol{\sigma}) - \ell(\boldsymbol{\mu}', \boldsymbol{\sigma}') | \mathcal{P}_{t+1} \} \\
&= \mathbb{E}_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \left\{ \sum_{t=1}^T \sum_{k=1}^K \left(\frac{p_k^{t,\pi} \|X_k^t - \boldsymbol{\mu}'_k\|^2}{\sigma_k'^2} - \frac{p_k^{t,\pi} \|X_k^t - \boldsymbol{\mu}_k\|^2}{\sigma_k^2} + \log \frac{\sigma_k'^2}{\sigma_k^2} \right) \middle| \mathcal{P}_{t+1} \right\} \\
&= \sum_{t=1}^T \sum_{k=1}^K \left(\frac{p_k^{t,\pi} \left(\frac{\sigma_k^2}{p_k^{t,\pi}} + \|\boldsymbol{\mu}_k - \boldsymbol{\mu}'_k\|^2 \right)}{\sigma_k'^2} - 1 + \log \frac{\sigma_k'^2}{\sigma_k^2} \right) \\
&= \sum_{t=1}^T \sum_{k=1}^K \left(\frac{p_k^{t,\pi} \|\boldsymbol{\mu}_k - \boldsymbol{\mu}'_k\|^2}{\sigma_k'^2} - 1 + \frac{\sigma_k^2}{\sigma_k'^2} + \log \frac{\sigma_k'^2}{\sigma_k^2} \right),
\end{aligned} \tag{16}$$

where $\mathbb{E}_{\boldsymbol{\mu}, \boldsymbol{\sigma}}$ denotes the expectation under the distribution of the outcomes determined by $(\boldsymbol{\mu}, \boldsymbol{\sigma})$. Now, let parameter $\boldsymbol{\mu}'$ satisfy $\boldsymbol{\mu}'_k = \boldsymbol{\mu}_k$, for every $k \neq a$, and $\|\boldsymbol{\mu}'_a\| > \|\boldsymbol{\mu}_1\|$, for a fixed $a \neq 1$. This means that arm 1 is optimal in $(\boldsymbol{\mu}, \boldsymbol{\sigma})$ but suboptimal in $(\boldsymbol{\mu}', \boldsymbol{\sigma}')$. It then follows that

$$\begin{aligned}
\mathbb{D}\{\boldsymbol{\mu}, \boldsymbol{\sigma} | \boldsymbol{\mu}', \boldsymbol{\sigma}'\} &= \mathbb{E}\{\ell(\boldsymbol{\mu}, \boldsymbol{\sigma}) - \ell(\boldsymbol{\mu}, \boldsymbol{\sigma}')\} \\
&= \mathbb{E}\{\mathbb{E}_{\boldsymbol{\mu}, \boldsymbol{\sigma}}\{\ell(\boldsymbol{\mu}, \boldsymbol{\sigma}) - \ell(\boldsymbol{\mu}, \boldsymbol{\sigma}') | \mathcal{P}_{t+1}\}\} \\
&= \frac{\mathbb{E}\{p_a^{t,\pi}\} \|\boldsymbol{\mu}_a - \boldsymbol{\mu}'_a\|^2}{\sigma_a'^2} - 1 + \frac{\sigma_a^2}{\sigma_a'^2} + \log \frac{\sigma_a'^2}{\sigma_a^2}.
\end{aligned} \tag{17}$$

On the other hand, let \mathbb{P} and \mathbb{P}' (and $\mathbb{E}\{\cdot\}$ with $\mathbb{E}'\{\cdot\}$) denote probability measures for $(X^t)_t$ (and expectations) under $(\boldsymbol{\mu}, \boldsymbol{\sigma})$ and $(\boldsymbol{\mu}', \boldsymbol{\sigma}')$, respectively, and recall that $\boldsymbol{\pi} \in \boldsymbol{\Pi}^*$. Invoking Lemma 3 with $B = B_T := \{\sum_{t=1}^T p_1^{t,\pi} < T - T^\gamma\}$, for some $0 < \gamma < 1$, yields

$$\begin{aligned}
\mathbb{P}\{B_T\} &= \mathbb{P}\left\{T - \sum_{t=1}^T p_1^{t,\pi} > T^\gamma\right\} \stackrel{(a)}{=} \mathbb{P}\left\{\sum_{k=2}^K \sum_{t=1}^T p_k^{t,\pi} > T^\gamma\right\} \\
&\stackrel{(b)}{\leq} \frac{1}{T^\gamma} \sum_{k=2}^K \sum_{t=1}^T \mathbb{E}\{p_k^{t,\pi}\} \\
&\stackrel{(c)}{=} o(T^{\alpha-\gamma}),
\end{aligned} \tag{18}$$

for every $\alpha > 0$, where (a) follows from $\sum_{k=1}^K p_k^{t,\pi} = 1$ for every $t \in \{1, \dots, T\}$, (b) is Markov's inequality [44], and where (c) is a consequence of π being uniformly efficient according to Definition 1, *i.e.*, the expected-cumulative power allocated at arms $2, \dots, K$ by $\pi \in \Pi^*$ is $\mathfrak{o}(T^\alpha)$. Similarly, under $(\boldsymbol{\mu}, \boldsymbol{\sigma}')$ (where arm 1 is suboptimal) we have that

$$\begin{aligned} \mathbb{P}'\{B_T^c\} &= \mathbb{P}'\left\{T - \sum_{t=1}^T p_1^{t,\pi} < T^\gamma\right\} \leq \frac{1}{T - T^\gamma} \sum_{t=1}^T \mathbb{E}'\{p_1^{t,\pi}\} \\ &= \frac{\mathfrak{o}(T^\alpha)}{T - T^{\alpha-1}} \\ &= \mathfrak{o}(T^{\alpha-1}), \end{aligned} \quad (19)$$

for every $\alpha > 0$. Now for $\alpha < \gamma$, (18) and (19) lead us to

$$\begin{aligned} d(\mathbb{P}\{B_T\}, \mathbb{P}'\{B_T\}) &\geq \mathbb{P}\{B_T\} \log \frac{\mathbb{P}\{B_T\}}{\mathbb{P}'\{B_T\}} + (1 - \mathbb{P}\{B_T\}) \log \frac{1 - \mathbb{P}\{B_T\}}{1 - \mathbb{P}'\{B_T\}} \\ &= \mathfrak{o}(T^{\alpha-\gamma}) \log \frac{\mathfrak{o}(T^{\alpha-\gamma})}{1 - \mathfrak{o}(T^{\alpha-1})} + (1 - \mathfrak{o}(T^{\alpha-\gamma})) \log \frac{1 - \mathfrak{o}(T^{\alpha-\gamma})}{\mathfrak{o}(T^{\alpha-1})} \\ &\geq \mathfrak{o}(1) + (1 - \alpha) \log T. \end{aligned} \quad (20)$$

Thus, combining this result with (17) and Lemma 3 yields

$$\sum_{t=1}^T \frac{\mathbb{E}\{p_a^{t,\pi}\} \|\boldsymbol{\mu}_a - \boldsymbol{\mu}'_a\|^2}{\sigma_a'^2} - 1 + \frac{\sigma_a^2}{\sigma_a'^2} + \log \frac{\sigma_a'^2}{\sigma_a^2} \geq \mathfrak{o}(1) + (1 - \alpha) \log T, \quad (21)$$

implying hat

$$\sum_{t=1}^T \mathbb{E}\{p_a^{t,\pi}\} \geq \frac{\sigma_a'^2}{\|\boldsymbol{\mu}_a - \boldsymbol{\mu}'_a\|^2} \left[\left(1 - \frac{\sigma_a^2}{\sigma_a'^2} - \log \frac{\sigma_a'^2}{\sigma_a^2}\right) T + \mathfrak{o}(1) + (1 - \alpha) \log T \right]. \quad (22)$$

It now follows that the confusion parameters $(\boldsymbol{\mu}', \boldsymbol{\sigma}')$ maximizing the lower bound in (22) are obtained by solving

$$\sup_{\boldsymbol{\mu}'_a, \sigma_a'^2: \|\boldsymbol{\mu}'_a\| > \|\boldsymbol{\mu}_1\|} \frac{\sigma_a'^2}{\|\boldsymbol{\mu}_a - \boldsymbol{\mu}'_a\|^2} \left[\left(1 - \frac{\sigma_a^2}{\sigma_a'^2} - \log \frac{\sigma_a'^2}{\sigma_a^2}\right) T + \mathfrak{o}(1) + (1 - \alpha) \log T \right],$$

where a necessary condition for optimality is that the term multiplying T must be zero since

$$1 - \frac{\sigma_a^2}{\sigma_a'^2} - \log \frac{\sigma_a'^2}{\sigma_a^2} \leq 0,$$

which is achieved iff $\sigma_a'^2 = \sigma_a^2$. The supremum is attained by choosing $\boldsymbol{\mu}'_a = \|\boldsymbol{\mu}_1\| \boldsymbol{\mu}_a / \|\boldsymbol{\mu}_a\|$, that is, a vector in the direction of $\boldsymbol{\mu}_a$ with the magnitude of $\boldsymbol{\mu}_1$, allowing us to conclude

$$\liminf_{T \rightarrow \infty} \frac{\sum_{t=1}^T \mathbb{E}\{p_a^{t,\pi}\}}{\log T} \geq \frac{\sigma_a^2(1 - \alpha)}{(\|\boldsymbol{\mu}_1\| - \|\boldsymbol{\mu}_a\|)^2} = \frac{\sigma_a^2(1 - \alpha)}{\Delta_a^2}. \quad (23)$$

The above result is valid for every $a = 2, 3, \dots, K$, so the proof is finished by taking $\alpha \rightarrow 0^+$ and observing that $\mathbb{E}\{R_{\boldsymbol{\mu}, \boldsymbol{\sigma}^\pi}^\pi(T)\} = \sum_{a=2}^K \sum_{t=1}^T \mathbb{E}\{p_a^{t,\pi}\} \Delta_a$.

B Proof of Lemma 1: sufficient statistics

To prove this result, we first show that $\bar{\mathbf{x}}_k(t)$, $S_k(t)$ and (p^1, p^2, \dots, p^t) are sufficient statistics for the distribution of the outcomes. This is followed by the characterization of their distributions and statistical independence.

B.1 Sufficiency of the statistics

The density of the sequence $(X_k^t)_t$, for every arm $k \in \{1, \dots, K\}$, is given by

$$\begin{aligned}
\nu(X_k^1, \dots, X_k^t | p_k^1, \dots, p_k^t) &= \frac{(\prod_{\ell=1}^t p_k^\ell)}{(\pi \sigma_k^2)^n} e^{-\frac{1}{\sigma_k^2} \sum_{\ell=1}^t p_k^\ell \|X_k^\ell - \boldsymbol{\mu}_k\|^2} \\
&= \frac{(\prod_{\ell=1}^t p_k^\ell)}{(\pi \sigma_k^2)^n} e^{-\frac{1}{\sigma_k^2} \sum_{\ell=1}^t p_k^\ell \|X_k^\ell - \boldsymbol{\mu}_k + \bar{\mathbf{x}}_k(t) - \bar{\mathbf{x}}_k(t)\|^2} \\
&= \frac{(\prod_{\ell=1}^t p_k^\ell)}{(\pi \sigma_k^2)^n} e^{-\frac{1}{\sigma_k^2} \sum_{\ell=1}^t p_k^\ell (\|X_k^\ell - \bar{\mathbf{x}}_k(t)\|^2 + \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\|^2 - 2[X_k^\ell - \bar{\mathbf{x}}_k(t)]^\top [\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k(t)])},
\end{aligned}$$

where

$$\begin{aligned}
&\sum_{\ell=1}^t p_k^\ell [X_k^\ell - \bar{\mathbf{x}}_k(t)]^\top [\boldsymbol{\mu}_k - \bar{\mathbf{x}}_k(t)] \\
&= \sum_{\ell=1}^t \left(X_k^\ell - \frac{\sum_{i=1}^t p_k^i X_k^i}{\sum_{j=1}^t p_k^j} \right)^\top \left(\boldsymbol{\mu}_k - \frac{\sum_{i=1}^t p_k^i X_k^i}{\sum_{j=1}^t p_k^j} \right) \\
&= \frac{1}{\left(\sum_{j=1}^t p_k^j \right)^2} \sum_{\ell=1}^t p_k^\ell \left(\sum_{j=1}^t p_k^j X_k^\ell - \sum_{i=1}^t p_k^i X_k^i \right)^\top \left(\sum_{j=1}^t p_k^j \boldsymbol{\mu}_k - \sum_{i=1}^t p_k^i X_k^i \right) \\
&= \frac{1}{\left(\sum_{j=1}^t p_k^j \right)^2} \sum_{\ell=1}^t \sum_{i=1}^t \sum_{j=1}^t p_k^\ell p_k^i p_k^j (X_k^\ell - X_k^i)^\top (\boldsymbol{\mu}_k - X_k^j) \\
&= \frac{1}{\left(\sum_{j=1}^t p_k^j \right)^2} \sum_{j=1}^t p_k^j (\boldsymbol{\mu}_k - X_k^j)^\top \underbrace{\sum_{\ell=1}^t \sum_{i=1}^t p_k^\ell p_k^i (X_k^\ell - X_k^i)}_{=0} \\
&= 0.
\end{aligned} \tag{24}$$

Therefore,

$$\begin{aligned}
\nu(X_k^1, \dots, X_k^t | p_k^1, \dots, p_k^t) &= \frac{(\prod_{\ell=1}^t p_k^\ell)}{(\pi \sigma_k^2)^n} e^{-\frac{1}{\sigma_k^2} \sum_{\ell=1}^t p_k^\ell (\|X_k^\ell - \bar{\mathbf{x}}_k(t)\|^2 + \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\|^2)} \\
&= \frac{(\prod_{\ell=1}^t p_k^\ell)}{(\pi \sigma_k^2)^n} e^{-\frac{1}{\sigma_k^2} (S_k(t) + \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\|^2 \sum_{\ell=1}^t p_k^\ell)},
\end{aligned} \tag{25}$$

which means that it is possible to keep track of the distribution of the sequence (X_k^1, \dots, X_k^t) (given the trajectory of allocated power) by just considering $S_k(t)$ and $\bar{\mathbf{x}}_k(t)$, for every $k \in \{1, \dots, K\}$.

B.2 Distributions and statistical independence

To simplify the notation, we write p_k^t instead of $p_k^{t, \text{WTS}}$. Additionally, in this section, all the expectations (and covariance matrices) are conditioned on $\mathcal{P}_{t+1} = \sigma(p_k^1, p_k^2, \dots, p_k^t)$. Notice that $\bar{\mathbf{x}}_k(t)$ is Gaussian distributed, since (X_k^1, \dots, X_k^t) is an independent Gaussian sequence, with mean

$\mathbb{E}\{\bar{\mathbf{x}}_k(t)\} = \boldsymbol{\mu}_k$ and whose covariance matrix is

$$\begin{aligned}
\text{cov}\{\bar{\mathbf{x}}_k(t)\} &= \mathbb{E}\left\{(\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k)(\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k)^\top\right\} \\
&= \frac{1}{(\sum_{\ell=1}^t p_k^\ell)^2} \mathbb{E}\left\{\left(\sum_{i=1}^t p_k^i (X_k^i - \boldsymbol{\mu}_k)\right) \left(\sum_{i=1}^t p_k^i (X_k^i - \boldsymbol{\mu}_k)\right)^\top\right\} \\
&= \frac{1}{(\sum_{\ell=1}^t p_k^\ell)^2} \mathbb{E}\left\{\sum_{i=1}^t (p_k^i)^2 (X_k^i - \boldsymbol{\mu}_k)(X_k^i - \boldsymbol{\mu}_k)^\top\right\} \\
&= \frac{1}{(\sum_{\ell=1}^t p_k^\ell)^2} \sum_{i=1}^t (p_k^i)^2 \frac{\sigma_k^2}{2p_k^i} \mathbf{I}_2 \\
&= \frac{\sigma_k^2}{2\sum_{\ell=1}^t p_k^\ell} \mathbf{I}_2. \tag{26}
\end{aligned}$$

To characterize the distribution of $S_k(t) = \sum_{\ell=1}^t p_k^\ell \|X_k^\ell - \bar{\mathbf{x}}_k(t)\|^2$ we exploit the fact that the data is statistically independent across arms and rounds. To this end, let us rewrite $S_k(t) = \sum_{\ell=1}^t [p_k^\ell (X_k^\ell(1) - \bar{\mathbf{x}}_k(t,1))^2 + p_k^\ell (X_k^\ell(2) - \bar{\mathbf{x}}_k(t,2))^2]$, where $X_k^t(i)$ and $\bar{\mathbf{x}}_k(t,i)$, $i \in \{1,2\}$ denote the i -th component of the \mathbb{R}^2 -vectors X_k^t and $\bar{\mathbf{x}}_k(t)$, respectively. Since the terms $(X_k^t(i) - \bar{\mathbf{x}}_k(t,i))^2$ in $S_k(t)$ are iid, we only need to characterize one of them. Let us introduce $s_k(t) := \sum_{\ell=1}^t p_k^\ell (X_k^\ell(1) - \bar{\mathbf{x}}_k(t,1))^2$ and, for k and t fixed, let $\mathbf{X} := [X_k^1(1) \ \dots \ X_k^t(1)]^\top$, $\mathbf{p} := [p_k^1 \ \dots \ p_k^t]^\top$. Additionally, let $\boldsymbol{\mu}_k(1)$ denote the first component of $\boldsymbol{\mu}_k \in \mathbb{R}^2$. Notice that $(\mathbf{X} - \boldsymbol{\mu}_k(1)\mathbf{1}) \sim \mathcal{N}(\mathbf{0}_t, (\sigma_k^2/2)\text{diag}\{\mathbf{p}\})$, for every $i = 1, \dots, t$, since the outcomes $(X_k^1, X_k^2, \dots, X_k^t)$ are Gaussian, statistically independent, and with individual distribution $X_k^t(1) \sim \mathcal{N}(\boldsymbol{\mu}_k(1), \sigma_k^2/2)$. Under this definition, $\bar{\mathbf{x}}_k(t,1) = \mathbf{1}^\top \mathbf{X} / (\mathbf{1}^\top \mathbf{p})$ and, hence, the first component of $\boldsymbol{\mu}_k \in \mathbb{R}^2$, $s_k(t)$ can be written as

$$\begin{aligned}
s_k(t) &= \mathbf{X}^\top \left(\mathbf{I}_t - \frac{\mathbf{1}\mathbf{p}^\top}{\mathbf{1}^\top \mathbf{p}} \right)^\top \text{diag}\{\mathbf{p}\} \left(\mathbf{I}_t - \frac{\mathbf{1}\mathbf{p}^\top}{\mathbf{1}^\top \mathbf{p}} \right) \mathbf{X} \\
&= \mathbf{X}^\top \left(\text{diag}\{\mathbf{p}\} - \frac{\mathbf{p}\mathbf{p}^\top}{\mathbf{1}^\top \mathbf{p}} \right) \mathbf{X} \\
&= \frac{\sigma_k^2}{2} \mathbf{E}^\top (\text{diag}\{\mathbf{p}\})^{-1/2} \left(\text{diag}\{\mathbf{p}\} - \frac{\mathbf{p}\mathbf{p}^\top}{\mathbf{1}^\top \mathbf{p}} \right) (\text{diag}\{\mathbf{p}\})^{-1/2} \mathbf{E} \\
&= \frac{\sigma_k^2}{2} \mathbf{E}^\top \underbrace{\left(\mathbf{I}_t - \frac{\mathbf{p}^{1/2} \mathbf{p}^{\top/2}}{\mathbf{1}^\top \mathbf{p}} \right)}_{:=\mathbf{A}} \mathbf{E}, \tag{27}
\end{aligned}$$

where $\mathbf{E} \sim \mathcal{N}(\mathbf{0}_t, \mathbf{I}_t)$, and where $\mathbf{p}^{1/2} := [\sqrt{p_k^1} \ \dots \ \sqrt{p_k^t}]^\top$. Notice that $\text{Rank } \mathbf{A} = t - 1$ since $\mathbf{p}^{1/2} \mathbf{p}^{\top/2}$ is a rank-1 perturbation [45], and that \mathbf{A} corresponds to an idempotent matrix since

$$\left(\mathbf{I}_t - \frac{\mathbf{p}^{1/2} \mathbf{p}^{\top/2}}{\mathbf{1}^\top \mathbf{p}} \right)^2 = \mathbf{I}_t - 2 \frac{\mathbf{p}^{1/2} \mathbf{p}^{\top/2}}{\mathbf{1}^\top \mathbf{p}} + \frac{\mathbf{p}^{1/2} \mathbf{p}^{\top/2} \mathbf{p}^{1/2} \mathbf{p}^{\top/2}}{(\mathbf{1}^\top \mathbf{p})^2} = \mathbf{I}_t - \frac{\mathbf{p}^{1/2} \mathbf{p}^{\top/2}}{\mathbf{1}^\top \mathbf{p}}. \tag{28}$$

Then, by [46, Lemma B.2], we have that $s_k(t)/(\sigma_k^2/2) \sim \chi_{t-1}^2$ and therefore, because both terms in $S_k(t)$ are iid, $S_k(t)/(\sigma_k^2/2) \sim \chi_{2(t-1)}^2$.

It now remains to prove conditional independence of $\bar{\mathbf{x}}_k(t)$ and $S_k(t)$, for every $k \in \{1, \dots, K\}$, and $t \in \{1, \dots, T\}$, given the trajectory of the power profiles. To show this, observe that

$$\begin{aligned}
\text{cov}(X_k^i - \bar{\mathbf{x}}_k(t), \bar{\mathbf{x}}_k(t)) &= \mathbb{E}\left\{(X_k^i - \bar{\mathbf{x}}_k(t))(\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k)^\top\right\} \\
&= (\sum_{\ell=1}^t p_k^\ell)^{-2} \mathbb{E}\left\{\sum_{j=1}^t \sum_{\ell=1}^t p_j p_\ell (X_k^i - X_k^j)(X_\ell - \boldsymbol{\mu}_k)\right\} \\
&= 0, \tag{29}
\end{aligned}$$

and that X_k^i and $\bar{\mathbf{x}}_k(t)$ are jointly Gaussian distributed. This implies that $\{X_k^i - \bar{\mathbf{x}}_k(t)\}_{i=1}^t$ are statistically independent of $\bar{\mathbf{x}}_k(t)$ and, therefore, the result follows because $S_k(t)$ is a function of the former.

C Proof of Lemma 2: posterior mean distribution

Recall that, from Lemma 1, $\bar{\mathbf{x}}_k(t)$ and $S_k(t)$ are conditionally independent (given the trajectory of allocated powers). In consequence, for every $k \in \{1, \dots, K\}$ and $t = 3, 4, \dots$, the posterior distribution of $(\boldsymbol{\mu}_k, \sigma_k^2)$ given $(\bar{\mathbf{x}}_k(t-1), S_k(t-1), p_k^1, p_k^2, \dots, p_k^{t-1})$ is

$$\begin{aligned} & f_{\boldsymbol{\mu}_k, \sigma_k^2} | \bar{\mathbf{x}}_k(t-1)=\mathbf{x}, S_k(t-1)=s, p_k^1, p_k^2, \dots, p_k^{t-1}(\boldsymbol{\mu}_k, \sigma_k^2) \\ &= \frac{f_{\bar{\mathbf{x}}_k(t-1), S_k(t-1)} | \boldsymbol{\mu}_k, \sigma_k^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s) \cdot f_{\boldsymbol{\mu}_k, \sigma_k} | p_k^1, p_k^2, \dots, p_k^{t-1}(\boldsymbol{\mu}_k, \sigma_k^2)}{\int_0^\infty \int_{\mathbb{R}^2} f_{\bar{\mathbf{x}}_k(t-1), S_k(t-1)} | \boldsymbol{\mu}_k, \sigma_k^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s) \cdot \underbrace{f_{\boldsymbol{\mu}_k, \sigma_k} | p_k^1, p_k^2, \dots, p_k^{t-1}(\boldsymbol{\mu}_k, \sigma_k^2)}_{\propto 1} d\boldsymbol{\mu} d(\sigma^2)} \\ &= \frac{f_{\bar{\mathbf{x}}_k(t-1), S_k(t-1)} | \boldsymbol{\mu}_k, \sigma_k^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s) \cdot f_{\boldsymbol{\mu}_k, \sigma_k} | \mathcal{F}_1(\boldsymbol{\mu}_k, \sigma_k^2)}{\int_0^\infty \int_{\mathbb{R}^2} f_{\bar{\mathbf{x}}_k(t-1), S_k(t-1)} | \boldsymbol{\mu}_k, \sigma_k^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s) \cdot f_{\boldsymbol{\mu}_k, \sigma_k} | \mathcal{F}_1(\boldsymbol{\mu}_k, \sigma_k^2) d\boldsymbol{\mu} d(\sigma^2)} \\ &= \frac{f_{\bar{\mathbf{x}}_k(t-1), S_k(t-1)} | \boldsymbol{\mu}_k, \sigma_k^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s)}{\int_0^\infty \int_{\mathbb{R}^2} f_{\bar{\mathbf{x}}_k(t-1), S_k(t-1)} | \boldsymbol{\mu}_k, \sigma_k^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s) d\boldsymbol{\mu} d(\sigma^2)} \end{aligned} \quad (30)$$

where, from Lemma 1,

$$f_{\bar{\mathbf{x}}_k(t-1), S_k(t-1)} | \boldsymbol{\mu}_k, \sigma_k^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s) = \frac{(\sum_{\ell=1}^t p_k^\ell) s^{t-2} e^{-\frac{1}{\sigma_k^2} \left(s + \sum_{\ell=1}^t p_k^\ell \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right)}}{\pi \Gamma(t-1) (\sigma_k^2)^t}. \quad (31)$$

Thus, the integral term in (30) is given by

$$\begin{aligned} & \int_0^\infty \int_{\mathbb{R}^2} f_{\bar{\mathbf{x}}_k(t), S_k(t)} | \boldsymbol{\mu}_k = \boldsymbol{\mu}, \sigma_k^2 = \sigma^2, p_k^1, p_k^2, \dots, p_k^{t-1}(\mathbf{x}, s) d\boldsymbol{\mu} d(\sigma^2) \\ &= \frac{(\sum_{\ell=1}^t p_k^\ell) s^{t-2}}{\pi \Gamma(t-1)} \int_0^\infty \frac{e^{-s/\sigma^2}}{(\sigma^2)^t} \int_{\mathbb{R}^2} e^{-\frac{1}{\sigma^2} \sum_{\ell=1}^t p_k^\ell \|\mathbf{x} - \boldsymbol{\mu}\|^2} d\boldsymbol{\mu} d(\sigma^2) \\ &= \frac{s^{t-2}}{\Gamma(t-1)} \int_0^\infty \frac{e^{-s/\sigma^2}}{(\sigma^2)^{t-1}} d(\sigma^2) \\ &= \frac{\Gamma(t-2)}{\Gamma(t-1)}, \end{aligned} \quad (32)$$

implying that

$$f_{\boldsymbol{\mu}_k, \sigma_k^2} | \bar{\mathbf{x}}_k(t-1)=\mathbf{x}, S_k(t-1)=s, p_k^1, p_k^2, \dots, p_k^{t-1}(\boldsymbol{\mu}_k, \sigma_k^2) = \frac{\sum_{\ell=1}^t p_k^\ell s^{t-2} e^{-\frac{1}{\sigma_k^2} \left(s + \sum_{\ell=1}^t p_k^\ell \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right)}}{\pi \Gamma(t-2) (\sigma_k^2)^t}.$$

Therefore, the posterior distribution of $\boldsymbol{\mu}_k$ given $(\bar{\mathbf{x}}_k(t), S_k(t), p_k^1, p_k^2, \dots, p_k^{t-1})$ is

$$\begin{aligned} f_{\boldsymbol{\mu}_k} | \bar{\mathbf{x}}_k(t-1)=\mathbf{x}, S_k(t-1)=s, p_k^1, p_k^2, \dots, p_k^{t-1}(\boldsymbol{\mu}_k) &= \frac{\sum_{\ell=1}^t p_k^\ell s^{t-2}}{\pi \Gamma(t-2)} \int_0^\infty \frac{e^{-\frac{1}{\sigma^2} \left(s + \sum_{\ell=1}^t p_k^\ell \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right)}}{(\sigma^2)^t} d(\sigma^2) \\ &= \frac{\sum_{\ell=1}^t p_k^\ell (t-2)}{\pi s} \left(1 + \frac{\sum_{\ell=1}^t p_k^\ell \|\mathbf{x} - \boldsymbol{\mu}_k\|^2}{s} \right)^{-t+1}. \end{aligned}$$

D Proof of Theorem 2: regret upper bound

We proceed by decomposing the regret $R_{\boldsymbol{\mu}, \sigma}^{\text{WTS}}(T)$ into three different terms which we upper bound separately in three different lemmas presented in Section D.1. It is worth to mention that splitting

the regret into these 3 terms is a standard technique in finite-time regret analysis, and our choice is inspired directly by [22, 11]. Nevertheless, upper-bounding each of these terms involves completely different techniques than the ones employed in [22, 11] since the direct application of the techniques in those references to bandits under weighted information does not hold.

The technical results supporting the three upper bounds, such as concentration inequalities, can be found in Section D.2.

Without loss of generality, let $k^* = 1$. Also, let $0 < \epsilon < \min_k(\|\boldsymbol{\mu}_1\| - \|\boldsymbol{\mu}_k\|)/2$ and define the following events:

$$\begin{aligned}\mathcal{A}_t &:= \{\|\tilde{\boldsymbol{\mu}}^*(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon\}, \\ \mathcal{B}_{k,t} &:= \{\|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \epsilon/2 \quad \text{and} \quad S_k(t) \leq t(\sigma_k^2 + \epsilon)\},\end{aligned}$$

where $\tilde{\boldsymbol{\mu}}_k(t) \sim f_{\boldsymbol{\mu}_k} | \mathcal{F}_t$ as in Lemma 2, and where $\|\tilde{\boldsymbol{\mu}}^*(t)\| := \max_k \|\tilde{\boldsymbol{\mu}}_k(t)\|$. Observe that the events $\mathcal{B}_{k,t}$ are defined in terms of $\bar{\mathbf{x}}_k(t)$ and $S_k(t)$, unlike the ones in the proof of optimality of TS [11] where their counterparts $\mathcal{B}_t(k)$ are defined in terms of $\bar{\mathbf{x}}_k(t-1)$ and $S_k(t-1)$. This choice of events $\mathcal{B}_{k,t}$ is essential in proving this theorem.

Since all prior mean distributions are equal, we have that $\rho_k^1 = \rho_k^2 = K^{-1}$, for every $k \in \{1, \dots, K\}$, according to Algorithm 1. By the definition of $p_k^{t,\text{WTS}}$, it holds that

$$\mathbb{E}\{p_k^{t,\text{WTS}}\} = \mathbb{E}\{\rho_k^t\} = \mathbb{E}\left\{\mathbb{P}\left\{\tilde{k}_t^* = k \mid \mathcal{F}_t\right\}\right\} = \mathbb{P}\left\{\tilde{k}_t^* = k\right\}, \quad (33)$$

where $\mathbb{P}\left\{\tilde{k}_t^* = k\right\}$ is the unconditional probability of arm k being optimal. Then, by (6), the expected cumulative regret under WTS satisfies

$$\begin{aligned}\mathbb{E}\{R_{\boldsymbol{\mu},\boldsymbol{\sigma}}^{\text{WTS}}(T)\} &= \mathbb{E}\left\{\sum_{t=1}^T \sum_{k=2}^K \Delta_k p_k^{t,\text{WTS}}\right\} \\ &= \sum_{k=2}^K \Delta_k (\rho_k^1 + \rho_k^2) + \mathbb{E}\left\{\sum_{t=3}^T \sum_{k=2}^K \Delta_k p_k^{t,\text{WTS}}\right\} \\ &= 2K^{-1} \sum_{k=2}^K \Delta_k + \sum_{t=3}^T \sum_{k=2}^K \Delta_k \mathbb{P}\left\{\tilde{k}_t^* = k\right\} \\ &\leq 2\Delta_{\max} + \sum_{t=3}^T \sum_{k=2}^K \Delta_k \mathbb{P}\left\{\tilde{k}_t^* = k\right\} \\ &= 2\Delta_{\max} + \sum_{k=2}^K \Delta_k \left(\sum_{t=3}^T \mathbb{P}\left\{\tilde{k}_t^* = k, \mathcal{A}_t\right\} + \mathbb{P}\left\{\tilde{k}_t^* = k, \mathcal{A}_t^c\right\}\right) \\ &\leq 2\Delta_{\max} + \sum_{k=2}^K \Delta_k \left(\sum_{t=3}^T \mathbb{P}\left\{\tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t}\right\} + \mathbb{P}\left\{\tilde{k}_t^* = k, \mathcal{B}_{k,t}^c\right\}\right. \\ &\quad \left. + \mathbb{P}\left\{\tilde{k}_t^* = k, \mathcal{A}_t^c\right\}\right) \\ &\leq 2\Delta_{\max} + \sum_{k=2}^K \Delta_k \sum_{t=3}^T \left(\mathbb{P}\left\{\tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t}\right\} + \mathbb{P}\left\{\tilde{k}_t^* = k, \mathcal{B}_{k,t}^c\right\}\right) \\ &\quad + \Delta_{\max} \sum_{t=3}^T \mathbb{P}\left\{\tilde{k}_t^* \neq 1, \mathcal{A}_t^c\right\},\end{aligned} \quad (34)$$

where \mathcal{A}_t^c and $\mathcal{B}_{k,t}^c$ denote the complements of \mathcal{A}_t and $\mathcal{B}_{k,t}$, respectively. Now, by means of Lemma 4, Lemma 5 and Lemma 6 (defined in the following section), the expected cumulative regret under WTS

satisfies

$$\mathbb{E} \{R_{\mu, \sigma}^{\text{WTS}}(T)\} \leq \sum_{k=2}^K \Delta_k \frac{\sigma_k^2 + \epsilon}{(\Delta_k - 2\epsilon)^2} \log T + \text{O}(1) + \text{O}(\epsilon^{-2}) + \text{O}(\epsilon^{-2}), \quad (35)$$

where the O terms are independent of T , implying that

$$\frac{\mathbb{E} \{R_{\mu, \sigma}^{\text{WTS}}(T)\}}{\log T} \leq \sum_{k=2}^K \Delta_k \frac{\sigma_k^2 + \epsilon}{(\Delta_k - 2\epsilon)^2} + \frac{\text{O}(1) + \text{O}(\epsilon^{-2}) + \text{O}(\epsilon^{-2})}{\log T}. \quad (36)$$

The numerator of the second term is independent of T , so the result now follows from choosing $\epsilon = \epsilon(T) < \log^{-a} T$, for some $0 < a < 1/2$ and taking the limit $T \rightarrow \infty$:

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E} \{R_{\mu, \sigma}^{\text{WTS}}(T)\}}{\log T} \leq \sum_{k=2}^K \frac{\sigma_k^2}{\Delta_k}. \quad (37)$$

D.1 Complementary upper bounds

Here we provide the complementary lemmas, and their respective proofs, that complete the proof of Theorem 2.

Lemma 4. *Consider the conditions of Theorem 2 and the events*

$$\begin{aligned} \mathcal{A}_t &= \{\|\tilde{\mu}^*(t)\| \geq \|\mu_1\| - \epsilon\}, \\ \mathcal{B}_{k,t} &= \{\|\bar{\mathbf{x}}_k(t) - \mu_k\| \leq \epsilon/2 \quad \text{and} \quad S_k(t) \leq t(\sigma_k^2 + \epsilon)\}, \end{aligned}$$

for every $k = 1, \dots, K$ and $t \in \mathbb{N}$, defined for some fixed ϵ such that $0 < \epsilon < \min_{k \neq k^*} \Delta_k/2$. Then,

$$\sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t} \right\} \leq \frac{\sigma_k^2 + \epsilon}{(\Delta_k - 2\epsilon)^2} \log T + \text{O}(\epsilon^{-2}), \quad (38)$$

for every $k \neq k^*$, where the $\text{O}(\epsilon^{-2})$ term does not depend on T . \triangle

Proof. Fix k and define the random variable $z_t := \sum_{\ell=1}^t p_k^\ell$, representing the cumulative power allocated at arm k up to round t . Define also $c(\epsilon) := \frac{\sigma_k^2 + \epsilon}{(\Delta_k - 2\epsilon)^2} > 0$, for every $\epsilon > 0$ (notice that $c(\epsilon)$ is $\text{O}(1)$ as $\epsilon \rightarrow 0$, and independent of T). Then,

$$\begin{aligned} \sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t} \right\} &= \sum_{t=3}^T \left(\mathbb{P} \left\{ \tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t}, z_{t-1} < c(\epsilon) \log T \right\} \right. \\ &\quad \left. + \mathbb{P} \left\{ \tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t}, z_{t-1} \geq c(\epsilon) \log T \right\} \right) \quad (39) \end{aligned}$$

We proceed to upper bound each sum in (39) separately. Since, for every event \tilde{E} , $\mathbb{P} \left\{ \tilde{E} \right\} = \mathbb{E} \left\{ \mathbb{1} \left\{ \tilde{E} \right\} \right\} = \mathbb{E} \left\{ \mathbb{E} \left\{ \mathbb{1} \left\{ \tilde{E} \right\} \mid \mathcal{F}_t \right\} \right\} = \mathbb{E} \left\{ \mathbb{P} \left\{ \tilde{E} \mid \mathcal{F}_t \right\} \right\}$, the first sum satisfies

$$\begin{aligned} \sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* = k, z_{t-1} < c(\epsilon) \log T \right\} &= \mathbb{E} \left\{ \sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* = k, z_{t-1} < c(\epsilon) \log T \mid \mathcal{F}_t \right\} \right\} \\ &\stackrel{(a)}{=} \mathbb{E} \left\{ \sum_{t=3}^T \mathbb{1} \left\{ z_{t-1} < c(\epsilon) \log T \right\} \mathbb{P} \left\{ \tilde{k}_t^* = k \mid \mathcal{F}_t \right\} \right\} \\ &= \mathbb{E} \left\{ \sum_{t=3}^T p_k^t \mathbb{1} \left\{ z_{t-1} < c(\epsilon) \log T \right\} \right\} \\ &\leq \mathbb{E} \left\{ \sum_{t=3}^T p_k^t \mathbb{1} \left\{ z_t < 1 + c(\epsilon) \log T \right\} \right\} \end{aligned}$$

$$\begin{aligned}
&\leq 1 + c(\epsilon) \log T \\
&= \frac{\sigma_k^2 + \epsilon}{(\Delta_k - 2\epsilon)^2} \log T + \mathcal{O}(1), \tag{40}
\end{aligned}$$

where (a) follows from z_{t-1} being \mathcal{F}_t -measurable and because of Lemma 8, which guarantees that $\sum_{\ell=1}^t p_k^\ell$ eventually reaches $1 + c(\epsilon) \log T$. This first bound provides the desired rate of growth in the right-hand side of (38), so it remains to prove that the second sum in (39) grows sub-logarithmically, at most, with T large, and polynomially in ϵ^{-2} as $\epsilon \rightarrow 0$.

Deriving an upper bound for the second term in (39) involves bounding the deviation of $\tilde{\boldsymbol{\mu}}_k(t)$ around $\bar{\boldsymbol{x}}_k(t-1)$. However, the definition of $\mathcal{B}_{k,t}$ only provides information about $\bar{\boldsymbol{x}}_k(t)$ and $S_k(t)$. This issue is tackled by formalizing the intuitive property in which $\bar{\boldsymbol{x}}_k(t)$ is *close* to $\bar{\boldsymbol{x}}_k(t-1)$ with high probability. More precisely, and since the events $\{\tilde{k}_t^* = k\}$ and \mathcal{A}_t imply event $\{\|\tilde{\boldsymbol{\mu}}_k(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon\}$ (which follows from (6)), the second term in (39) satisfies

$$\begin{aligned}
&\mathbb{P} \left\{ \tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t}, z_{t-1} \geq c(\epsilon) \log T \right\} \\
&\leq \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon, \|\bar{\boldsymbol{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \epsilon/2, S_k(t) \leq t(\sigma_k^2 + \epsilon), z_{t-1} \geq c(\epsilon) \log T \right\} \\
&= \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon, \underbrace{\|\bar{\boldsymbol{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \epsilon/2, \|\bar{\boldsymbol{x}}_k(t) - \bar{\boldsymbol{x}}_k(t-1)\| \leq \epsilon/2}_{\implies \|\bar{\boldsymbol{x}}_k(t-1) - \boldsymbol{\mu}_k\| \leq \epsilon}, \right. \\
&\quad \left. \underbrace{S_k(t) \leq t(\sigma_k^2 + \epsilon)}_{\implies S_k(t-1) \leq t(\sigma_k^2 + \epsilon)}, z_{t-1} \geq c(\epsilon) \log T \right\} \\
&\quad + \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon, \underbrace{\|\bar{\boldsymbol{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \epsilon/2, \|\bar{\boldsymbol{x}}_k(t) - \bar{\boldsymbol{x}}_k(t-1)\| \geq \epsilon/2}_{\implies \|X_k^t - \boldsymbol{\mu}_k\| \geq z_{t-1}\epsilon/2p_k^t}, \right. \\
&\quad \left. S_k(t) \leq t(\sigma_k^2 + \epsilon), z_{t-1} \geq c(\epsilon) \log T \right\} \\
&\leq \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon, \|\bar{\boldsymbol{x}}_k(t-1) - \boldsymbol{\mu}_k\| \leq \epsilon, S_k(t-1) \leq t(\sigma_k^2 + \epsilon), z_{t-1} \geq c(\epsilon) \log T \right\} \\
&\quad + \mathbb{P} \left\{ \|X_k^t - \boldsymbol{\mu}_k\| \geq z_{t-1}\epsilon/2p_k^t \right\}, \tag{41}
\end{aligned}$$

which follows from the fact that

$$\epsilon/2 \leq \|\bar{\boldsymbol{x}}_k(t) - \bar{\boldsymbol{x}}_k(t-1)\| = \left\| \bar{\boldsymbol{x}}_k(t) - \frac{z_t \bar{\boldsymbol{x}}_k(t) - p_k^t X_k^t}{z_{t-1}} \right\| = \frac{p_k^t}{z_{t-1}} \|X_k^t - \bar{\boldsymbol{x}}_k(t)\|$$

and that $\{S_k(t)\}_t$ is non-decreasing in t . Now, because $\epsilon < \min_{i \neq k^*} \Delta_i/2 \leq \Delta_k/2$, meaning that $\mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon, \|\bar{\boldsymbol{x}}_k(t-1) - \boldsymbol{\mu}_k\| \leq \epsilon \right\} \leq \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t) - \bar{\boldsymbol{x}}_k(t-1)\| \geq \Delta_k - 2\epsilon \right\}$, the first

term in (41) satisfies

$$\begin{aligned}
& \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon, \|\bar{\mathbf{x}}_k(t-1) - \boldsymbol{\mu}_k\| \leq \epsilon, S_k(t-1) \leq t(\sigma_k^2 + \epsilon), z_{t-1} \geq c(\epsilon) \log T \right\} \\
& \leq \mathbb{E} \left\{ \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t) - \bar{\mathbf{x}}_k(t-1)\| \geq \Delta_k - 2\epsilon \mid \mathcal{F}_t \right\} \right. \\
& \quad \left. \times \mathbb{1} \left\{ \|\bar{\mathbf{x}}_k(t-1) - \boldsymbol{\mu}_k\| \leq \epsilon, S_k(t-1) \leq t(\sigma_k^2 + \epsilon), z_{t-1} \geq c(\epsilon) \log T \right\} \right\} \\
& = \mathbb{E} \left\{ \left(1 + \frac{z_{t-1}(\Delta_k - 2\epsilon)^2}{S_k(t-1)} \right)^{-t+3} \right. \\
& \quad \left. \times \mathbb{1} \left\{ \|\bar{\mathbf{x}}_k(t-1) - \boldsymbol{\mu}_k\| \leq \epsilon, S_k(t-1) \leq t(\sigma_k^2 + \epsilon), z_{t-1} \geq c(\epsilon) \log T \right\} \right\} \\
& \leq \mathbb{E} \left\{ \left(1 + c^{-1}(\epsilon) \frac{z_{t-1}}{t} \right)^{-t+3} \mathbb{1} \left\{ z_{t-1} \geq c(\epsilon) \log T \right\} \right\}, \tag{42}
\end{aligned}$$

while, because p_k^t is \mathcal{F}_t -measurable, the second term satisfies

$$\begin{aligned}
\sum_{t=3}^T \mathbb{P} \left\{ \|X_k^t - \boldsymbol{\mu}_k\| \geq \frac{z_{t-1}\epsilon}{2p_k^t} \right\} &= \sum_{t=3}^T \mathbb{E} \left\{ \mathbb{P} \left\{ \|X_k^t - \boldsymbol{\mu}_k\| \geq \frac{z_{t-1}\epsilon}{2p_k^t} \mid \mathcal{F}_t \right\} \right\} \\
&\stackrel{(a)}{\leq} \sum_{t=3}^T \mathbb{E} \left\{ e^{-z_{t-1}^2 \epsilon^2 / 4\sigma_k^2 p_k^t} \right\} \\
&\leq \mathbb{E} \left\{ \sum_{t=3}^T p_k^t e^{-z_{t-1}^2 \epsilon^2 / 4\sigma_k^2} \right\} \\
&= \mathbb{E} \left\{ \sum_{i=0}^T \sum_{t=3}^T p_k^t e^{-z_{t-1}^2 \epsilon^2 / 4\sigma_k^2} \mathbb{1} \{z_t \in (i, i+1)\} \right\} \\
&\leq \sum_{i=1}^T e^{-(i-1)^2 \epsilon^2 / 4\sigma_k^2} \mathbb{E} \left\{ \sum_{t=3}^T p_k^t \mathbb{1} \{z_t \in (i, i+1)\} \right\} \\
&\leq e^{\epsilon^2 / 4\sigma_k^2} \sum_{i=1}^T e^{-i\epsilon^2 / 4\sigma_k^2} 2 \\
&\leq 2e^{\frac{\min_{i \neq k} \Delta_i^2}{16\sigma_k^2}} \sum_{i=1}^{\infty} e^{-i\epsilon^2 / 4\sigma_k^2} \\
&= 2e^{\frac{\min_{i \neq k} \Delta_i^2}{16\sigma_k^2}} \frac{e^{-\epsilon^2 / 4\sigma_k^2}}{1 - e^{-\epsilon^2 / 4\sigma_k^2}} \\
&= \mathcal{O}(\epsilon^{-2}), \tag{43}
\end{aligned}$$

where (a) follows from the fact that $1 + x \leq e^x$, $x \in \mathbb{R}$, holds, in particular, for $x = -1 + 1/p_k^t$. In consequence, the second sum in (39) can be upper bounded as

$$\begin{aligned}
& \sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t}, z_t \geq c(\epsilon) \log T \right\} \\
& \leq \mathcal{O}(\epsilon^{-2}) + \mathbb{E} \left\{ \sum_{t=3}^T \left(1 + c^{-1}(\epsilon) \frac{z_{t-1}}{t} \right)^{-t+3} \mathbb{1} \{z_t \geq c(\epsilon) \log T\} \right\} \\
& \leq \mathcal{O}(\epsilon^{-2}) + (1 + c^{-1}(\epsilon))^3 \mathbb{E} \left\{ \sum_{t=3}^T \left(1 + \frac{\log T}{t} \right)^{-t} \mathbb{1} \{z_t \geq c(\epsilon) \log T\} \right\}
\end{aligned}$$

$$\begin{aligned}
&\leq \mathbf{O}(\epsilon^{-2}) + (1 + c^{-1}(\epsilon))^3 \mathbb{E} \left\{ \sum_{t: z_{t-1} \geq c(\epsilon) \log T}^T \left(1 + \frac{\log T}{t-1}\right)^{-t} \right\} \\
&\leq \mathbf{O}(\epsilon^{-2}) + (1 + c^{-1}(\epsilon))^3 \sum_{t=\lfloor c(\epsilon) \log T \rfloor}^T \left(1 + \frac{\log T}{t}\right)^{-t} \\
&\stackrel{(b)}{\leq} \mathbf{O}(\epsilon^{-2}) + (1 + c^{-1}(\epsilon))^3 \sum_{t=\lfloor c(\epsilon) \log T \rfloor}^T e^{\frac{-t \log T}{t+\log T}} \\
&\leq \mathbf{O}(\epsilon^{-2}) + (1 + c^{-1}(\epsilon))^3 \int_{\lfloor c(\epsilon) \log T \rfloor - 1}^T e^{\frac{-t \log T}{t+\log T}} dt \\
&\leq \mathbf{O}(\epsilon^{-2}) + (1 + c^{-1}(\epsilon))^3 \int_{c(\epsilon) \log(T) - 2}^T e^{\frac{-t \log T}{t+\log T}} dt, \tag{44}
\end{aligned}$$

where (a) follows from the fact that $(1 + x/y)^y = e^{y \log(1+x/y)} \geq e^{y \frac{x/y}{1+x/y}}$. Deriving an upper bound for the right-hand side of the equation above involves solving a very delicate integral. We do this by finding an upper and lower bound of an anti-derivative for $e^{\frac{-x \log T}{x+\log T}}$. Such an anti-derivative can be obtained using integration by parts:

$$\begin{aligned}
\int e^{\frac{-x \log T}{x+\log T}} dx &= (x + \log T) e^{\frac{-x \log T}{x+\log T}} + \int (x + \log T) e^{\frac{-x \log T}{x+\log T}} \frac{\log^2 T}{(x + \log T)^2} dx \\
&\stackrel{(a)}{=} (x + \log T) e^{\frac{-x \log T}{x+\log T}} + \int \frac{\log^2 T}{\log(T) - u} e^{-u} du \\
&\stackrel{(b)}{=} (x + \log T) e^{\frac{-x \log T}{x+\log T}} - T^{-1} \log^2(T) \int \frac{e^v}{v} dv,
\end{aligned}$$

where (a) follows from $u = \log(T)x/(x + \log T)$ and (b) from $v = \log(T) - u$. The function $f(v) := \log v + \sum_{n=1}^{\infty} \frac{1}{n \cdot n!} v^n$ is an anti-derivative for the function $v \mapsto e^v/v$, since

$$\frac{df}{dv} = v^{-1} + \sum_{n=1}^{\infty} \frac{v^{n-1}}{n!} = \frac{1}{v} \left(1 + \sum_{n=1}^{\infty} \frac{v^n}{n!} \right) = \frac{e^v}{v}. \tag{45}$$

Additionally, for every $v > 0$, it holds that

$$\log v \leq f(v) \leq 2e^v,$$

hence, by reversing the change of variables $v = v(u(x))$, we have that

$$f(v(u(x))) = f\left(\frac{\log^2 T}{x + \log T}\right)$$

is an anti-derivative for $e^{\frac{-x \log T}{x+\log T}}$ that satisfies

$$\log \frac{\log^2 T}{x + \log T} \leq f\left(\frac{\log^2 T}{x + \log T}\right) \leq 2e^{\frac{\log^2 T}{x+\log T}}. \tag{46}$$

Therefore, the definite integral in the right-hand of (44) can be upper bounded, by means of (46), as

$$\begin{aligned}
& \int_{c(\epsilon) \log T - 2}^T e^{\frac{-t \log T}{t + \log T}} dt \\
&= (t + \log T) e^{\frac{-t \log T}{t + \log T}} \Bigg|_{t=c(\epsilon) \log(T) - 2}^T - T^{-1} \log^2(T) f\left(\frac{\log^2 T}{t + \log T}\right) \Bigg|_{t=c(\epsilon) \log(T) - 2}^T \\
&= (T + \log T) e^{\frac{-T \log T}{T + \log T}} - (c(\epsilon) \log(T) - 2 + \log T) e^{\frac{-(c(\epsilon) \log(T) - 2) \log T}{c(\epsilon) \log(T) - 2 + \log T}} \\
&\quad - T^{-1} \log^2(T) f\left(\frac{\log^2 T}{T + \log T}\right) + T^{-1} \log^2 T f\left(\frac{\log^2 T}{c(\epsilon) \log(T) - 2 + \log T}\right) \\
&\stackrel{(a)}{\leq} (T + \log T) e^{\frac{-T \log T}{T + \log T}} - T^{-1} \log^2(T) f\left(\frac{\log^2 T}{T + \log T}\right) \\
&\quad + T^{-1} \log^2 T f\left(\frac{\log^2 T}{c(\epsilon) \log(T) - 2 + \log T}\right) \\
&\leq 2T^{1 - \frac{T}{T + \log T}} + T^{-1} \log^3(2T) + T^{\frac{\log T}{c(\epsilon) \log(T) - 2 + \log T} - 1} \log^2 T \\
&\leq 4 + T^{-1} \log^3(2T) + T^{-1} \log^2 T \\
&\leq 4 + 2 \cdot 3^3/e^3 + (2/e)^2 \\
&= O(\epsilon^0), \tag{47}
\end{aligned}$$

where (a) follows because the second term in the line above is always negative, since T is assumed to satisfy $c(\epsilon) \log T - 2 > 0$. The constants hidden in $O(\epsilon^0)$ are universal (*i.e.*, independent of T), therefore, putting together (40), (44) and (47) yields

$$\begin{aligned}
\sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, \mathcal{A}_t, \mathcal{B}_{k,t} \right\} &\leq \frac{\sigma_k^2 + \epsilon}{(\Delta_k - 2\epsilon)^2} \log T + O(1) + O(\epsilon^{-2}) + O(1) \\
&= \frac{\sigma_k^2 + \epsilon}{(\Delta_k - 2\epsilon)^2} \log T + O(\epsilon^{-2}),
\end{aligned}$$

concluding the proof. ■

Lemma 5. Consider the conditions of Theorem 2. Let $\epsilon > 0$, and define the events

$$\mathcal{B}_{k,t} = \{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \epsilon/2 \quad \text{and} \quad S_k(t) \leq t(\sigma_k^2 + \epsilon) \},$$

for every $k = 1, \dots, K$, and $t \in \mathbb{N}$. Then,

$$\sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, \mathcal{B}_{k,t}^c \right\} = O(\epsilon^{-2}), \tag{48}$$

for every $k \neq k^*$, where the $O(\epsilon^{-2})$ term is independent of T . △

Proof. Recall that $h(x) = x - \log(1 + x)$. Fix k and let $z_t := \sum_{\ell=1}^t p_k^\ell$ denote the (random) cumulative allocated power at frequency k up to round t . Define \mathcal{P}_{t+1} as the σ -algebra generated by the power's trajectory up to round t , (p^1, p^2, \dots, p^t) . The sum in (48) can be decomposed, by means of the union bound, as

$$\begin{aligned}
\sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, \mathcal{B}_{k,t}^c \right\} &= \sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, (\|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| > \epsilon/2 \quad \text{or} \quad S_k(t) > t(\sigma_k^2 + \epsilon)) \right\} \\
&= \sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \right\} \\
&\quad + \sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, S_k(t) \geq t(\sigma_k^2 + \epsilon) \right\}, \tag{49}
\end{aligned}$$

where the second sum in the right-hand side satisfies

$$\begin{aligned}
\sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, S_k(t) \geq t(\sigma_k^2 + \epsilon) \right\} &\leq \sum_{t=3}^T \mathbf{P} \left\{ S_k(t) \geq t(\sigma_k^2 + \epsilon) \right\} \\
&\leq \sum_{t=0}^{\infty} e^{-th(\epsilon/\sigma_k^2)} \\
&= \frac{1}{1 - e^{-h(\epsilon/\sigma_k^2)}}. \tag{50}
\end{aligned}$$

On the other hand, because p_k^t is \mathcal{F}_t measurable, meaning that $\sigma(p_k^1, \dots, p_k^t) = \mathcal{P}_{t+1} \subset \mathcal{F}_t$, the first sum in the right-hand side of (49) can be upper bounded as

$$\begin{aligned}
&\sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \right\} \\
&= \mathbb{E} \left\{ \sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \mid \mathcal{P}_{t+1} \right\} \right\} \\
&= \mathbb{E} \left\{ \sum_{i=0}^T \sum_{t=3}^T \mathbf{1} \{z_t \in (i, i+1)\} \mathbf{P} \left\{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \mid \mathcal{P}_{t+1} \right\} \right. \\
&\quad \left. \times \mathbf{P} \left\{ \tilde{k}_t^* = k \mid \mathcal{P}_{t+1}, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \right\} \right\} \\
&= \mathbb{E} \left\{ \sum_{i=0}^T \sum_{t=3}^T \mathbf{1} \{z_t \in (i, i+1)\} e^{-z_t \epsilon^2 / 4\sigma_k^2} \mathbf{P} \left\{ \tilde{k}_t^* = k \mid \mathcal{P}_{t+1}, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \right\} \right\} \\
&\leq \sum_{i=0}^T e^{-i\epsilon^2 / 4\sigma_k^2} \sum_{t=3}^T \mathbb{E} \left\{ \mathbf{1} \{z_t \in (i, i+1)\} \mathbf{P} \left\{ \tilde{k}_t^* = k \mid \mathcal{P}_{t+1}, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \right\} \right\} \\
&= \sum_{i=0}^T e^{-i\epsilon^2 / 4\sigma_k^2} \sum_{t=3}^T \mathbb{E} \left\{ \mathbf{1} \{z_t \in (i, i+1)\} \right. \\
&\quad \left. \times \mathbb{E} \left\{ \mathbf{P} \left\{ \tilde{k}_t^* = k \mid \mathcal{F}_t, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \right\} \mid \mathcal{P}_{t+1} \right\} \right\} \\
&= \sum_{i=0}^T e^{-i\epsilon^2 / 4\sigma_k^2} \mathbb{E} \left\{ \sum_{t=3}^T \mathbf{1} \{z_t \in (i, i+1)\} \mathbb{E} \left\{ p_k^t \mid \mathcal{P}_{t+1}, \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon/2 \right\} \right\} \\
&= \sum_{i=0}^T e^{-i\epsilon^2 / 4\sigma_k^2} \mathbb{E} \left\{ \sum_{t=3}^T \mathbf{1} \{z_t \in (i, i+1)\} p_k^t \right\} \\
&\leq 2 \sum_{i=0}^{\infty} e^{-i\epsilon^2 / 4\sigma_k^2} \\
&= \frac{2}{1 - e^{-\epsilon^2 / 4\sigma_k^2}},
\end{aligned}$$

hence,

$$\begin{aligned}
\sum_{t=3}^T \mathbf{P} \left\{ \tilde{k}_t^* = k, \mathcal{B}_{k,t}^c \right\} &\leq \frac{2}{1 - e^{-\epsilon^2 / 4\sigma_k^2}} + \frac{1}{1 - e^{-h(\epsilon/\sigma_k^2)}} \\
&= \mathcal{O}(\epsilon^{-2}) + \mathcal{O}(\epsilon^{-2}), \tag{51}
\end{aligned}$$

finishing the proof. \blacksquare

Lemma 6. Consider the conditions of Theorem 2. Let $\epsilon > 0$, and define the events

$$\mathcal{A}_t = \{\|\tilde{\boldsymbol{\mu}}^*(t)\| \geq \|\boldsymbol{\mu}_1\| - \epsilon\},$$

for every $t \in \mathbb{N}$. Then,

$$\sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* \neq k^*, \mathcal{A}_t^c \right\} = \mathcal{O}(\epsilon^{-4}), \quad (52)$$

for every arm $k \neq k^*$, where the $\mathcal{O}(\epsilon^{-4})$ term is independent of t . \triangle

Proof. Assume, without of generality, that $k^* = 1$. Let $z_t := \sum_{\ell=1}^t p_1^\ell$ denote the (random) cumulative power on arm 1. Additionally, observe that event $\{\tilde{k}_t^* \neq 1, \mathcal{A}_t^c\}$ is included in the event $\{\|\tilde{\boldsymbol{\mu}}_1(t)\| \leq \|\boldsymbol{\mu}_1\| - \epsilon\}$. Consider the following chain of inequalities:

$$\begin{aligned} \sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* \neq 1, \mathcal{A}_t^c \right\} &= \mathbb{E} \left\{ \sum_{t=3}^T \mathbb{P} \left\{ \tilde{k}_t^* \neq 1, \mathcal{A}_t^c \mid \mathcal{F}_t \right\} \right\} \\ &\leq \mathbb{E} \left\{ \sum_{t=3}^T \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t)\| \leq \|\boldsymbol{\mu}_1\| - \epsilon \mid \mathcal{F}_t \right\} \right\} \\ &\leq \mathbb{E} \left\{ \sum_{t=3}^T \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \right\} \right\} \end{aligned} \quad (53)$$

Now, observe that each of the terms in the above equation satisfy

$$\begin{aligned} &\mathbb{E} \left\{ \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \right\} \right\} \\ &= \mathbb{E} \left\{ \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \right\} \mathbf{1} \left\{ \|\bar{\mathbf{x}}_1(t-1) - \boldsymbol{\mu}_1\| > \epsilon/4 \right\} \right\} \\ &\quad + \mathbb{E} \left\{ \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \right\} \right. \\ &\quad \quad \left. \times \mathbf{1} \left\{ \|\bar{\mathbf{x}}_1(t-1) - \boldsymbol{\mu}_1\| \leq \epsilon/4, S_1(t-1) \geq 2(t-1)\sigma_1^2 \right\} \right\} \\ &\quad + \mathbb{E} \left\{ \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \right\} \right. \\ &\quad \quad \left. \times \mathbf{1} \left\{ \|\bar{\mathbf{x}}_1(t-1) - \boldsymbol{\mu}_1\| \leq \epsilon/4, S_1(t-1) < 2(t-1)\sigma_1^2 \right\} \right\}, \end{aligned} \quad (54)$$

so we now find an upper bound for the above expression by bounding each of its three terms. The first term satisfies

$$\begin{aligned} &\mathbb{E} \left\{ \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \right\} \mathbf{1} \left\{ \|\bar{\mathbf{x}}_1(t-1) - \boldsymbol{\mu}_1\| > \epsilon/4 \right\} \right\} \\ &\leq \mathbb{E} \left\{ \mathbb{E} \left\{ \mathbf{1} \left\{ \|\bar{\mathbf{x}}_1(t-1) - \boldsymbol{\mu}_1\| \geq \epsilon/4 \right\} \mid z_{t-1} \right\} \right\} \\ &\leq \mathbb{E} \left\{ e^{-z_t \epsilon^2 / (16\sigma_1^2)} \right\}, \end{aligned} \quad (55)$$

which follows from applying Lemma 7. Similarly, the second result in Lemma 7 implies that

$$\begin{aligned} &\mathbb{E} \left\{ \mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \right\} \mathbf{1} \left\{ \|\bar{\mathbf{x}}_1(t-1) - \boldsymbol{\mu}_1\| \geq \epsilon/4, S_1(t-1) \geq 2(t-1)\sigma_1^2 \right\} \right\} \\ &\leq \mathbb{P} \left\{ S_1(t-1) \geq 2(t-1)\sigma_1^2 \right\} \\ &\leq e^{-(t-1)h(1)}. \end{aligned} \quad (56)$$

For the third term, we invoke Lemma 2, yielding

$$\begin{aligned}
& \mathbb{P} \{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \} \\
&= \int_{\tilde{\boldsymbol{\mu}}: \|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}_1\| \geq \epsilon} \frac{z_{t-1}(t-3)}{\pi S_1(t-1)} \left(1 + \frac{z_{t-1} \|\tilde{\boldsymbol{\mu}} - \bar{\boldsymbol{x}}_1(t-1)\|^2}{S_1(t-1)} \right)^{-t+2} d\tilde{\boldsymbol{\mu}} \\
&\stackrel{(a)}{\leq} \int_{\tilde{\boldsymbol{\mu}}: \|\tilde{\boldsymbol{\mu}} - \bar{\boldsymbol{x}}_1(t-1)\| \geq \epsilon/2} \frac{z_{t-1}(t-3)}{\pi S_1(t-1)} \left(1 + \frac{z_{t-1} \|\tilde{\boldsymbol{\mu}} - \bar{\boldsymbol{x}}_1(t-1)\|^2}{S_1(t-1)} \right)^{-t+2} d\tilde{\boldsymbol{\mu}} \\
&= \int_0^{2\pi} \int_{\epsilon/2}^{\infty} \frac{z_{t-1}(t-3)}{\pi S_1(t-1)} \left(1 + \frac{z_{t-1} r^2}{S_1(t-1)} \right)^{-t+2} r dr d\theta \\
&= 2 \left(1 + \frac{z_{t-1} \epsilon^2}{4S_1(t-1)} \right)^{-t+3} \\
&\stackrel{(b)}{\leq} 2 \left(1 + \frac{z_{t-1} \epsilon^2}{8(t-1)\sigma_k^2} \right)^{-t+3}, \tag{57}
\end{aligned}$$

where (a) follows from the fact that $\bar{\boldsymbol{x}}_1(t-1)$ is $\epsilon/4$ close to $\boldsymbol{\mu}_1$, meaning that one can always construct a circle of radius $\epsilon/2$ around $\bar{\boldsymbol{x}}_1(t-1)$ (that encircles $\boldsymbol{\mu}_1$) that is completely contained inside $\{\tilde{\boldsymbol{\mu}}: \|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}_1\| \geq \epsilon\}$. The inequality in (b) holds because the sample variance is restricted to $S_1(t-1) \leq 2(t-1)\sigma_k^2$. Putting these three bounds together leads to

$$\begin{aligned}
& \mathbb{E} \{ \mathbb{P} \{ \|\tilde{\boldsymbol{\mu}}_1(t) - \boldsymbol{\mu}_1\| \geq \epsilon \mid \mathcal{F}_t \} \} \\
&\leq \mathbb{E} \left\{ e^{-z_{t-1}\epsilon^2/(16\sigma_1^2)} + e^{-(t-1)h(1)} + 2 \left(1 + \frac{z_{t-1}\epsilon^2}{8(t-1)\sigma_1^2} \right)^{-t+3} \right\}, \tag{58}
\end{aligned}$$

where the expectation is respect to the sequence z_{t-1} . Then, the expected sum in (53) satisfies

$$\begin{aligned}
& \sum_{t=3}^T \mathbb{P} \{ \tilde{k}_t^* \neq k^*, \mathcal{A}_t^c \} \\
&\leq \mathbb{E} \left\{ \sum_{t=3}^{\infty} e^{-z_{t-1}\epsilon^2/(16\sigma_1^2)} + e^{-(t-1)h(1)} + 2 \left(1 + \frac{z_{t-1}\epsilon^2}{8(t-1)\sigma_1^2} \right)^{-t+3} \right\} \\
&= \mathbb{E} \left\{ \sum_{t=2}^{\infty} e^{-z_t\epsilon^2/(16\sigma_1^2)} + e^{-th(1)} + 2 \left(1 + \frac{z_t\epsilon^2}{8t\sigma_1^2} \right)^{-t+2} \right\} \\
&\stackrel{(a)}{\leq} \mathbb{E} \left\{ \sum_{i=0}^{\infty} \sum_{z_t \in (i, i+1)} \left(e^{-z_t\epsilon^2/(16\sigma_1^2)} + 2 \left(1 + \frac{\epsilon^2}{8\sigma_1^2} \right) e^{-z_t\epsilon^2/(\epsilon^2+8\sigma_1^2)} \right) \right\} + \frac{1}{1 - e^{-h(1)}} \\
&\leq \mathbb{E} \left\{ \sum_{i=0}^{\infty} \left(e^{-i\epsilon^2/(16\sigma_1^2)} + 2 \left(1 + \frac{\epsilon^2}{8\sigma_1^2} \right) e^{-i\epsilon^2/(\epsilon^2+8\sigma_1^2)} \right) \sum_{z_t \in (i, i+1)} 1 \right\} + \frac{1}{1 - e^{-h(1)}} \\
&= \sum_{i=0}^{\infty} \left(e^{-i\epsilon^2/(16\sigma_1^2)} + 2 \left(1 + \frac{\epsilon^2}{8\sigma_1^2} \right) e^{-i\epsilon^2/(\epsilon^2+8\sigma_1^2)} \right) \mathbb{E} \left\{ \sum_{z_t \in (i, i+1)} 1 \right\} + \frac{1}{1 - e^{-h(1)}} \\
&\stackrel{(b)}{\leq} \sum_{i=0}^{\infty} \left(e^{-i\epsilon^2/(16\sigma_1^2)} + 2 \left(1 + \frac{\epsilon^2}{8\sigma_1^2} \right) e^{-i\epsilon^2/(\epsilon^2+8\sigma_1^2)} \right) C_{\boldsymbol{\mu}, \boldsymbol{\sigma}} + \frac{1}{1 - e^{-h(1)}} \\
&= C_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \left(\frac{1}{1 - e^{-\epsilon^2/(16\sigma_1^2)}} + 2 \left(1 + \frac{\epsilon^2}{8\sigma_1^2} \right) \frac{1}{1 - e^{-\epsilon^2/(\epsilon^2+8\sigma_1^2)}} \right) + \frac{1}{1 - e^{-h(1)}} \\
&= \mathcal{O}(\epsilon^{-2}) + \mathcal{O}(\epsilon^{-2}) + \mathcal{O}(\epsilon^{-2}) \\
&= \mathcal{O}(\epsilon^{-2}), \tag{59}
\end{aligned}$$

where $h(1) = 1 - \log 2 > 0$. Step (a) is obtained from the inequality $(1+x/y)^{-y} = e^{-y \log(1+x/y)} \leq e^{-y \frac{x/y}{1+x/y}}$, with $x = z_t \epsilon^2 / (8\sigma_1^2)$ and $y = t$, so

$$\exp \left\{ \frac{-t \frac{z_t \epsilon^2}{8\sigma_1^2}}{t + \frac{z_t \epsilon^2}{8\sigma_1^2}} \right\} \leq \exp \left\{ \frac{-t \frac{z_t \epsilon^2}{8\sigma_1^2}}{t \left(1 + \frac{1\epsilon^2}{8\sigma_1^2}\right)} \right\},$$

because $z_t \leq t$. Inequality (b) is a very crucial part of this result that deserves its own statement and proof, which can be found in Lemma 9. The importance of that result is that it allows us to upper bound the expected number of rounds it takes to apply a cumulative power exceeding 1 at arm k^* . Such a bound is problem-dependent constant $C_{\mu, \sigma}$ but it does not depend on t nor ϵ . ■

D.2 Technical Lemmas

In this section we provide three technical lemmas that support the ones stated in Section D.1:

- Lemma 7 provides concentration inequalities (and equalities) for the sufficient statistics $\bar{\mathbf{x}}_k(t)$ and $S_k(t)$, the posterior means $\tilde{\boldsymbol{\mu}}_k(t)$.
- Lemma 8 proves a key property of WTS, in which the cumulative power allocated at every arm diverges to ∞ with probability 1. This is, in fact, a very desirable property since learning the optimal arm involves knowing $\boldsymbol{\mu}$ completely which, by the law of large numbers [32], is only achievable by sampling all arms an infinite number of times or, in our case, with an infinite cumulative power.
- Lemma 9 bounds the minimum number of rounds required to bring the cumulative power at arm 1 from $i \in \mathbb{N}$ to $i+1$ when the power profiles are determined by WTS (see Algorithm 1). Such a bound is provided in the form of a problem-dependent constant, independent of T and ϵ . This lemma is invoked only in the proof of Lemma 6.

Lemma 7. For every $k \in \{1, \dots, K\}$, $t \in \{1, \dots, T\}$, and $\epsilon > 0$, it holds that

$$\mathbb{P} \left\{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| + \epsilon \mid p_k^1, p_k^2, \dots, p_k^t \right\} = e^{-\frac{\sum_{\ell=1}^t p_k^\ell}{\sigma_k^2} \epsilon^2}, \quad (60)$$

and

$$\mathbb{P} \left\{ S_k(t) \geq t(\sigma_k^2 + \epsilon) \right\} \leq e^{-th(\epsilon/\sigma_k^2)}, \quad (61)$$

where $h(x) := x - \log(1+x)$, $\forall x > 0$. Furthermore,

$$\mathbb{P} \left\{ \|\tilde{\boldsymbol{\mu}}_k(t) - \bar{\mathbf{x}}_k(t-1)\| \geq \delta \mid \mathcal{F}_t \right\} = \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell \delta^2}{S_k(t-1)} \right)^{-t+3}. \quad (62)$$

△

Proof. Let \mathcal{P}_{t+1} denote the sigma algebra generated by $p_k^1, p_k^2, \dots, p_k^t$. By Lemma 1, we have that

$$\begin{aligned} \mathbb{P} \left\{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \geq \epsilon \mid \mathcal{P}_{t+1} \right\} &= \int_{\mathbf{x}: \|\mathbf{x} - \boldsymbol{\mu}_k\| \geq \epsilon} \frac{\sum_{\ell=1}^t p_k^\ell}{\pi \sigma_k^2} e^{-\frac{\sum_{\ell=1}^t p_k^\ell}{\sigma_k^2} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2} d\mathbf{x} \\ &= \int_{\mathbf{z}: \|\mathbf{z}\| \geq \epsilon} \frac{\sum_{\ell=1}^t p_k^\ell}{\pi \sigma_k^2} e^{-\frac{\sum_{\ell=1}^t p_k^\ell}{\sigma_k^2} \|\mathbf{z}\|^2} d\mathbf{z} \\ &= \int_0^{2\pi} \int_\epsilon^\infty \frac{\sum_{\ell=1}^t p_k^\ell}{\pi \sigma_k^2} e^{-\frac{\sum_{\ell=1}^t p_k^\ell}{\sigma_k^2} r^2} r dr d\theta \\ &= e^{-\frac{\sum_{\ell=1}^t p_k^\ell}{\sigma_k^2} \epsilon^2}, \end{aligned} \quad (63)$$

which follows from the change of variables $\mathbf{z} := \mathbf{x} - \boldsymbol{\mu}_k$, followed by a polar change of variable.

To derive (61), we relax Chernoff's bound on $S_k(t)/(\sigma_k^2/2) \sim \chi_{2(t-1)}^2$:

$$\begin{aligned}
\mathbb{P}\{S_k(t) \geq t(\sigma_k^2 + \epsilon)\} &\leq e^{\inf_{\lambda < 1/\sigma_k^2} \log \mathbb{E}\{e^{\lambda S_k(t)}\} - \lambda t(\sigma_k^2 + \epsilon)} \\
&= \left(\frac{t}{t-1}\right)^{t-1} \left(1 + \frac{\epsilon}{\sigma_k^2}\right)^{t-1} e^{-(t\epsilon/\sigma_k^2 + 1)} \\
&= \left(\frac{t}{t-1}\right)^{t-1} \left(1 + \frac{\epsilon}{\sigma_k^2}\right)^{-1} e^{t \log(1 + \epsilon^2/\sigma_k^2)} e^{-(t\epsilon/\sigma_k^2 + 1)} \\
&= \left(\frac{t}{t-1}\right)^{t-1} e^{-1} \underbrace{\left(\frac{\sigma_k^2}{\sigma_k^2 + \epsilon^2}\right)}_{\leq 1} e^{-t(\epsilon^2/\sigma_k^2 - \log(1 + \epsilon^2/\sigma_k^2))} \\
&\leq e^{-t h(\epsilon/\sigma_k^2)}, \tag{64}
\end{aligned}$$

since $\left(\frac{t}{t-1}\right)^{t-1} \leq e^1$ for every $t > 1$.

For the last part of the proof, observe that, by Lemma 2, and by changing to polar coordinates, we have that

$$\begin{aligned}
&\mathbb{P}\{\|\tilde{\boldsymbol{\mu}}_k(t) - \bar{\boldsymbol{x}}_k(t-1)\| \geq \delta \mid \mathcal{F}_t\} \\
&= \int_{\tilde{\boldsymbol{\mu}}: \|\tilde{\boldsymbol{\mu}} - \bar{\boldsymbol{x}}_k(t-1)\| > \delta} \frac{\sum_{\ell=1}^{t-1} p_k^\ell(t-3)}{\pi S_k(t-1)} \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell \|\tilde{\boldsymbol{\mu}} - \bar{\boldsymbol{x}}_k(t-1)\|^2}{S_k(t-1)}\right)^{-t+2} d\tilde{\boldsymbol{\mu}} \\
&= \int_{\boldsymbol{z}: \|\boldsymbol{z}\| > \delta} \frac{\sum_{\ell=1}^{t-1} p_k^\ell(t-3)}{\pi S_k(t-1)} \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell \|\boldsymbol{z}\|^2}{S_k(t-1)}\right)^{-t+2} d\boldsymbol{z} \\
&= \int_0^{2\pi} \int_\delta^\infty \frac{\sum_{\ell=1}^{t-1} p_k^\ell(t-3)}{\pi S_k(t-1)} \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell r^2}{S_k(t-1)}\right)^{-t+2} r dr d\theta \\
&= \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell \delta^2}{S_k(t-1)}\right)^{-t+3}. \tag{65}
\end{aligned}$$

finishing the proof. ■

Lemma 7 shows that, as in TS [11], the posterior means are symmetrically distributed around the sample (empirical) mean. It also illustrates that the sample means ($\bar{\boldsymbol{x}}_k(t)$) concentrate around their respective means ($\boldsymbol{\mu}_k$) at different speeds, depending on which algorithm is being used: in TS the sample mean of arm k concentrates at a rate $\propto e^{-N_k(t)\epsilon^2}$, with $N_k(t)$ being the number of times arm k has been sampled up to round t , while in WTS, it concentrates at a rate $\propto e^{-\sum_{\ell=1}^t p_k^\ell \epsilon^2}$ (compare this Lemma 7 to Lemmas 3 and 4 in [11]). On the other hand, the sample variances $S_k(t)$ concentrate exponentially in t around $t\sigma_k^2$, regardless of the algorithm (TS or WTS). The concentration speed of $S_k(t)$ is consequence of two facts: (i) the outcome from arm k provides enough information about σ_k^2 no matter how small p_k^t is, and (ii) albeit small, WTS constantly assigns resources to all arms, as the following lemma shows. It is the imbalance in the concentration rates of the sample mean and variance under WTS what introduces difficulties when trying to prove that WTS is an optimal algorithm.

Lemma 8. *Consider problem (RM) under unknown variance σ . Then, the cumulative power at arm k satisfies*

$$\lim_{T \rightarrow \infty} \sum_{t=1}^T p_k^t = \infty, \quad \text{a.s.} \tag{66}$$

for every $k = 1, \dots, K$. △

Proof. The sequence $(\sum_{\ell=1}^t p_k^\ell)_t$ is non-decreasing, since $p_k^t \geq 0$ for every $t = 1, \dots, T$. Furthermore, p_k^t is never equal to zero since the posterior means have infinite support. The proof is by

contradiction. Fix $k \in \{1, \dots, K\}$ and assume $\lim_{T \rightarrow \infty} \sum_{t=1}^T p_k^t = c < \infty$. Then, by Lemmas 1 and 7, $\bar{\mathbf{x}}_k(t)$ and $S_k(t)/t$ converge to some (random) $\bar{\mathbf{x}}_k := \lim_{t \rightarrow \infty} \bar{\mathbf{x}}_k(t) \sim \mathcal{N}(\boldsymbol{\mu}_k, \sigma_k^2/(2c)\mathbf{I}_2)$ and $\lim_{t \rightarrow \infty} S_k(t)/t = \sigma_k^2$ (almost surely). From Lemma 2, we have that the limit posterior mean distribution is

$$\begin{aligned} \lim_{t \rightarrow \infty} f_{\boldsymbol{\mu}} | \mathcal{F}_t(\tilde{\boldsymbol{\mu}}) &= \lim_{t \rightarrow \infty} \frac{\sum_{\ell=1}^{t-1} (t-3)}{\pi S_k(t-1)} \left(1 + \frac{\sum_{\ell=1}^{t-1} \|\tilde{\boldsymbol{\mu}} - \bar{\mathbf{x}}_k(t-1)\|^2}{S_k(t-1)} \right)^{-t+2} \\ &= \frac{c}{\sigma_k^2 \pi} e^{-c \|\tilde{\boldsymbol{\mu}} - \bar{\mathbf{x}}_k\|^2 / \sigma_k^2}, \quad \text{a.s.} \end{aligned} \quad (67)$$

thus, given that $\bar{\mathbf{x}}_k \neq \tilde{\boldsymbol{\mu}}$ (which has measure zero), we have that $\lim_{t \rightarrow \infty} f_{\boldsymbol{\mu}_k} | \mathcal{F}_t$ defines a non-degenerate distribution of infinite support and, therefore, $\lim_{t \rightarrow \infty} \mathbb{P}\{\cap_{i \neq k} \{\|\tilde{\boldsymbol{\mu}}_i(t)\| \geq \|\tilde{\boldsymbol{\mu}}_i\|\} | \mathcal{F}_t\} = c' \neq 0$, for some $c' = c'(\boldsymbol{\mu}, \boldsymbol{\sigma}) > 0$ (provided such limits exist). We then conclude that $\lim_{t \rightarrow \infty} p_k^t = c' > 0$, implying that the sequence $(p_k^\ell)_\ell$ does not converge, so neither does $\sum_{\ell=1}^t p_k^\ell$, leading to a contradiction. This concludes the proof. \blacksquare

Lemma 9. Consider $t \geq 3$. The expected number of rounds in which the cumulative power at arm $k^* = 1$, $z_t := \sum_{\ell=1}^t p_1^\ell$, belongs to the set $(i, i+1)$ satisfies

$$\mathbb{E} \left\{ \sum_{t: z_t \in (i, i+1)} 1 \right\} \leq C_{\boldsymbol{\mu}, \boldsymbol{\sigma}} := 2 \left(1 - e^{\frac{3}{2K} \left(\frac{\Delta_{\min}}{\sigma_{\max}} \right)^2} \right)^{-2K} \left(1 - e^{-2h(1)} \right)^{-K},$$

where $\Delta_{\min} := \min_k \Delta_k$ and $\sigma_{\max} := \max_k \sigma_k$ that is, $C_{\boldsymbol{\mu}, \boldsymbol{\sigma}}$ is a problem-dependent constant, independent of t . \triangle

Proof. Define $z_k^t := \sum_{\ell=1}^t p_k^\ell$, for every $k = 1, \dots, K$, and observe that $z_t = z_1^t$. Proving this result involves firstly lower bounding the expected value of p_1^{t+1} given z_t . Recall that $p_1^{t+1} = \mathbb{P}\{\tilde{k}_{t+1}^* = 1\}$. Now, observe that the event $\cap_{k=1}^K \{\|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \Delta_{\min}, \|\tilde{\boldsymbol{\mu}}_k(t+1) - \bar{\mathbf{x}}_k(t)\| \leq \Delta_{\min}\}$ implies that $\tilde{k}_{t+1}^* = 1$. In words, the event in which the sample means are close to their respective true means, together with the event in which the posteriors are close to their respective means (the sample means), all of them very likely to happen, imply that $\tilde{k}_{t+1}^* = 1$. Observe that $\tilde{k}_{t+1}^* = 1$ holds regardless of the values of $\{S_k(t)\}_{k,t}$ and, in particular, of the event $\cap_{k=1}^K \{S_k(t) \leq 2t\sigma_k^2\}$ (which is also likely), so it

holds that

$$\begin{aligned}
& \mathbb{E} \{ p_1^{t+1} \mid z_t \} \\
&= \mathbb{E} \left\{ \mathbb{E} \left\{ \mathbb{1} \left\{ \tilde{k}_{t+1}^* = 1 \right\} \mid z_t \right\} \mid \mathcal{F}_t \right\} \\
&= \mathbb{P} \left\{ \tilde{k}_{t+1}^* = 1 \mid z_t \right\} \\
&\geq \mathbb{P} \left\{ \bigcap_{k=1}^K \left\{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \Delta_{\min}, S_k(t) \leq 2t\sigma_k^2, \|\tilde{\boldsymbol{\mu}}_k(t+1) - \bar{\mathbf{x}}_k(t)\| \leq \Delta_{\min} \right\} \mid z_t \right\} \\
&\geq \mathbb{E} \left\{ \mathbb{1} \left\{ \bigcap_{k=1}^K \left\{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \Delta_{\min}, S_k(t) \leq 2t\sigma_k^2, \|\tilde{\boldsymbol{\mu}}_k(t+1) - \bar{\mathbf{x}}_k(t)\| \leq \Delta_{\min} \right\} \right\} \mid z_t \right\} \\
&\stackrel{(a)}{=} \mathbb{E} \left\{ \mathbb{1} \left\{ \bigcap_{k=1}^K \left\{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \Delta_{\min}, S_k(t) \leq 2t\sigma_k^2 \right\} \right\} \prod_{k=1}^K \left(1 - \left(1 + \frac{z_k^t \Delta_{\min}^2}{S_k(t)} \right)^{-t+3} \right) \mid z_t \right\} \\
&\stackrel{(b)}{\geq} \left(1 - e^{-\frac{3}{2K} \left(\frac{\Delta_{\min}}{\sigma_{\max}} \right)^2} \right)^K \mathbb{E} \left\{ \mathbb{1} \left\{ \bigcap_{k=1}^K \left\{ \|\bar{\mathbf{x}}_k(t) - \boldsymbol{\mu}_k\| \leq \Delta_{\min}, S_k(t) \leq 2t\sigma_k^2 \right\} \right\} \mid z_t \right\} \\
&\stackrel{(c)}{=} \left(1 - e^{-\frac{3}{2K} \left(\frac{\Delta_{\min}}{\sigma_{\max}} \right)^2} \right)^K \prod_{k=1}^K \left(1 - e^{-z_t \Delta_{\min}^2 / \sigma_k^2} \right) \left(1 - e^{-th(1)} \right) \\
&\geq \left(1 - e^{-\frac{3}{2K} \left(\frac{\Delta_{\min}}{\sigma_{\max}} \right)^2} \right)^{2K} \left(1 - e^{-2h(1)} \right)^K \\
&=: c_{\boldsymbol{\mu}, \boldsymbol{\sigma}}, \tag{68}
\end{aligned}$$

where the equalities in (a) and (c) follow from the tower property of expectations with $\sigma(z_t) \subset \mathcal{F}_t$ and $\sigma(z_t) \subset \sigma(z_2^t, z_3^t, \dots, z_K^t)$, respectively. The inequalities are a consequence of invoking Lemma 7, and the inequality (b) follows because WTS enforces that the first two power profiles p_1^1, p_2^2 satisfy $p_k^t = K^{-1}$, for every arm $k = 1, \dots, K$ (see Algorithm 1).

Fix $i \in \mathbb{N}$ and let us introduce $\underline{t} := \min\{t \in \mathbb{N} : z_t \geq i\}$ and $\bar{t} := \min\{t > \underline{t} : z_t \geq i+1\}$. Observe that \bar{t} denotes a stopping time in $\sigma(z_1, z_2, \dots, z_t)$. Define also the partial sums $\xi_t := \sum_{n=\underline{t}}^{\bar{t}} p_1^n$. Then,

$$\begin{aligned}
\mathbb{E} \{ \xi_{\bar{t}} \} &= \sum_{t=\underline{t}}^{\infty} \mathbb{E} \{ \xi_t \mathbb{1} \{ \bar{t} = t \} \} \\
&= \sum_{t=\underline{t}}^{\infty} \sum_{n=t}^{\bar{t}} \mathbb{E} \{ p_1^n \mathbb{1} \{ \bar{t} = t \} \} \\
&= \sum_{n=\underline{t}}^{\infty} \mathbb{E} \{ p_1^n \mathbb{1} \{ \bar{t} \geq n \} \} \\
&\stackrel{(a)}{=} \sum_{n=\underline{t}}^{\infty} \mathbb{E} \left\{ \mathbb{P} \left\{ \tilde{k}_n^* = 1 \mid z_{n-1} \right\} \mathbb{1} \{ \bar{t} \geq n \} \right\} \\
&\stackrel{(b)}{\geq} c_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \sum_{n=\underline{t}}^{\infty} \mathbb{P} \{ \bar{t} \geq n \} \\
&= c_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \sum_{n=0}^{\infty} \mathbb{P} \{ \bar{t} - \underline{t} \geq n \} \\
&= c_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \mathbb{E} \{ \bar{t} - \underline{t} \}, \tag{69}
\end{aligned}$$

where in step (a) we used the fact that $\{\bar{t} < n\}^c = \{\bar{t} \geq n\}$ is $\sigma(z_1, z_2, \dots, z_{n-1})$ -measurable, and (b) is a consequence of (68). Now notice that the partial sums satisfy $1 \leq \xi_{\bar{t}} \leq 2$, which allows us

to conclude that $\mathbb{E} \left\{ \sum_{t: z_t \in (i, i+1)} 1 \right\} = \mathbb{E} \{ \bar{t} - t \} \leq c_{\mu, \sigma}^{-1} \mathbb{E} \{ \xi_{\bar{t}} \} \leq 2c_{\mu, \sigma}^{-1} =: C_{\mu, \sigma}$, concluding the proof. \blacksquare

E Approximating ρ^t via Monte Carlo simulations

The goal of this section is to provide with a method that allows us to update ρ^t in Algorithm 1 in a tractable manner. Obtaining $\rho^t = \text{Prob}\{\tilde{k}_t^* = k \mid \mathcal{F}_t\}$ can be carried out indirectly by first finding the distributions $(f_{\mu_k \mid \mathcal{F}_t})_{k=1}^K$, as explained in (6). However, finding a closed form for the mapping from $(f_{\mu_k \mid \mathcal{F}_t})_{k=1}^K$ to ρ^t involves the calculation of K -dimensional very complicated integrals:

$$\rho_k^t = \int_{(\tilde{\mu}_1, \dots, \tilde{\mu}_K) \in \mathbb{R}^{2K}} \mathbb{1} \left\{ \bigcap_{i \neq k} \{ \|\tilde{\mu}_i\| \leq \|\tilde{\mu}_k\| \} \right\} f_{\mu_1 \mid \mathcal{F}_t}(\tilde{\mu}_1) \dots f_{\mu_K \mid \mathcal{F}_t}(\tilde{\mu}_K) d(\tilde{\mu}_1, \dots, \tilde{\mu}_K),$$

for which no closed-form solution exists. We overcome this practical issue by estimating ρ^t through the Monte Carlo method [47] defined in Algorithm 2.

Algorithm 2 Estimation of ρ^t

- 1: Input: $(f_{\mu_1 \mid \mathcal{F}_t}, f_{\mu_2 \mid \mathcal{F}_t}, \dots, f_{\mu_K \mid \mathcal{F}_t})$ (posterior distributions at round t), M (number of samples per arm)
 - 2: Draw M samples $\tilde{\mu}_k^{(1)}, \tilde{\mu}_k^{(2)}, \dots, \tilde{\mu}_k^{(M)} \sim f_{\mu_k \mid \mathcal{F}_t}$, for every $k = 1, \dots, K$
 - 3: **for** $k = 1$ to K **do**
 - 4: Set $\rho_k^t \approx \frac{1}{M} \sum_{m=1}^M \mathbb{1} \left\{ \arg \max_i \{ \|\tilde{\mu}_i^{(m)}\| \} = k \right\}$
 - 5: **end for**
-

Observe that, in Algorithm 2, $\frac{1}{M} \sum_{m=1}^M \mathbb{1} \{ \arg \max_i \{ \|\tilde{\mu}_i^{(m)}\| \} = k \} \xrightarrow{\text{a.s.}} \rho_k^t$ as $M \rightarrow \infty$, so picking M large ensures a good approximation of ρ^t . Observe also that WTS in Algorithm 1 coincides with TS when ρ^t is updated using Algorithm 2 with $M = 1$.

For completeness, we explain here how to sample from $\tilde{\mu}_k(t) \sim f_{\mu_k \mid \mathcal{F}_t}$, for every $k = 1, \dots, K$, and $t = 4, 5, \dots$. Fix k and t , and recall that, from Lemma 2, the posterior means $\tilde{\mu}_k(t)$ satisfy

$$f_{\mu_k \mid \mathcal{F}_t}(\tilde{\mu}) = \frac{\sum_{\ell=1}^{t-1} p_k^\ell (t-3)}{\pi S_k(t-1)} \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell \|\tilde{\mu} - \bar{x}_k(t-1)\|^2}{S_k(t-1)} \right)^{-t+2}.$$

Now, consider the following change of variables into polar coordinates $\tilde{\mu}_k(t) = \bar{x}_k(t-1) + r[\cos \theta \quad \sin \theta]^\top$. Then, the distribution of the new variables (r, θ) is [32]

$$\begin{aligned} f_{r, \theta}(r, \theta) &= r f_{\mu_k \mid \mathcal{F}_t} \left(\bar{x}_k(t-1) + r \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \right) \\ &= \underbrace{\frac{1}{2\pi}}_{=: f_\theta(\theta)} \underbrace{\frac{2r \sum_{\ell=1}^{t-1} p_k^\ell (t-3)}{\pi S_k(t-1)} \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell r^2}{S_k(t-1)} \right)^{-t+2}}_{=: f_r(r)}, \end{aligned}$$

meaning that θ and r are statistically independent random variables that can be sampled separately. Observe that θ is uniformly distributed over $[0, 2\pi]$. On the other hand, sampling r can be achieved by means of the inverse sampling theorem. To this end, let $F(r)$ denote the cumulative density function of r :

$$F(r) = \int_0^r f_r(u) du = 1 - \left(1 + \frac{\sum_{\ell=1}^{t-1} p_k^\ell r^2}{S_k(t-1)} \right)^{-t+3},$$

meaning that the inverse F^{-1} is defined by

$$r = F^{-1}(\xi) = \sqrt{\frac{S_{t-1}}{\sum_{\ell=1}^{t-1} p_k^\ell} \left((1 - \xi)^{\frac{1}{-t+3}} - 1 \right)},$$

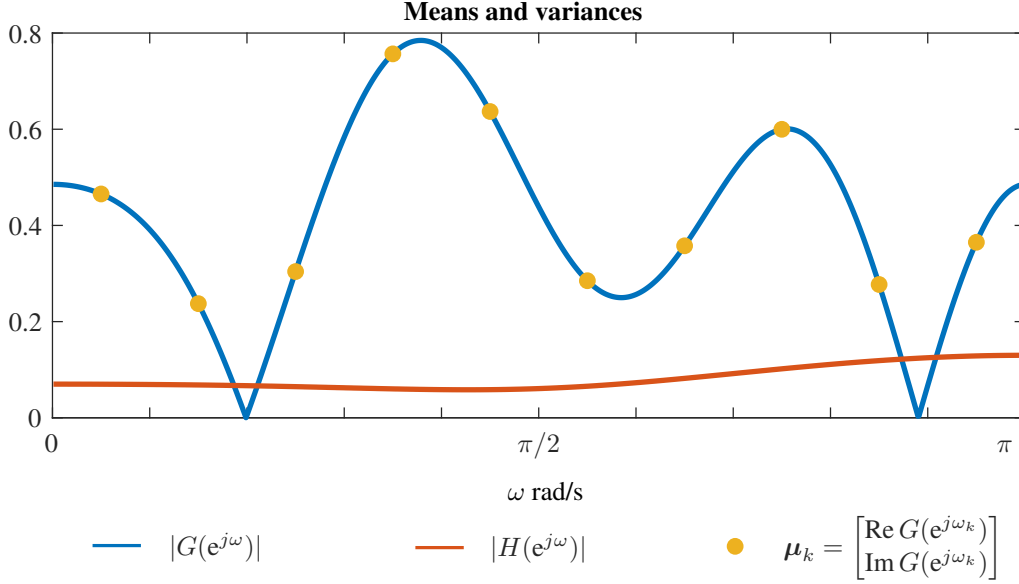


Figure 4: The values of $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ used in the simulation study in Section 4.3 based on the frequency responses of G and H . The two-dimensional means and the variances satisfy $\|\boldsymbol{\mu}_k\| = |G(e^{j\omega_k})|$ while $\sigma_k^2 = |H(e^{j\omega_k})|^2$, respectively.

for every $\xi \in [0, 1]$ because the image of F is $[0, 1]$. Therefore, a sample from $r \sim f_r$ can be obtained by first sampling from $\xi \sim \mathcal{U}([0, 1])$ (where \mathcal{U} stands for uniform distribution) and then obtaining $r = F^{-1}(\xi)$ [32]. In consequence, a sample from $\tilde{\boldsymbol{\mu}}_k(t) \sim f_{\boldsymbol{\mu}_k} | \mathcal{F}_t$ can be obtained as $\tilde{\boldsymbol{\mu}}_k(t) = \boldsymbol{x}_k(t-1) + r[\cos \theta \quad \sin \theta]^\top$.

F Details of simulation study 1 in Section 4.3

In this section we provide more details on the values of $(K, \boldsymbol{\mu}, \boldsymbol{\sigma})$ in the simulation study in Section 4.3. Since the problem formulated in this work is inspired by the gain estimation problem of an LTI system, we use LTI filters G and H to define $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ from their frequency responses in Fig. 4. We choose $K = 10$ in order to have noticeably different lower bounds (see Fig. 1) and where $\boldsymbol{\mu}_k = |G(e^{j\omega_k})|$ for $\omega_k = 2\pi/(2K+1)$, $k = 1, \dots, K$. On the other hand, the variance of the experiments is set to $\sigma_k^2 = |H(e^{j\omega_k})|^2$, for every $k = 1, \dots, K$. The value of ρ^t in WTS (see Algorithm 1) is updated via Algorithm 2 in Appendix E with $M = 500$ Monte Carlo simulations.

Algorithm 3 WTS under known $\boldsymbol{\sigma}$ [10]

- 1: Input: $T, \rho^1 = (1/K, 1/K, \dots, 1/K)$ (prior distribution for \tilde{k}_1^*), λ^2, M
- 2: **for** $t = 1$ to T **do**
- 3: Apply the power profile $p^{t, \text{WTS}} = \rho^1$ and collect the outcome X^t
- 4: Obtain $f_{\boldsymbol{\mu}_k} | \mathcal{F}_{t+1} = \mathcal{N}(\boldsymbol{m}_k^{t+1}, v_k^{t+1} \mathbf{I}_2)$ for every $k = 1, \dots, K$

$$\boldsymbol{m}_k^{t+1} = \frac{\lambda^2 \sum_{\ell=1}^t p_k^\ell X_k^\ell}{\sigma_k^2 + \lambda^2 \sum_{\ell=1}^t p_k^\ell}$$

$$v_k^t = \frac{\sigma_k^2 \lambda^2}{1 + \lambda^2 \sum_{\ell=1}^t p_k^\ell}$$

- 5: Update ρ^{t+1} via Algorithm 2
 - 6: **end for**
-

We implemented four algorithms:

- TS under known σ based on the algorithm introduced in [10],
- WTS under known σ based on the algorithm introduced in [10],
- TS under unknown σ based on the algorithm described in [11], and
- WTS as proposed in Algorithm 1 in this work.

It is important to mention that TS and WTS under known variance are initialized with Gaussian prior distributions for every mean μ_k of the form $f_{\mu_k | \mathcal{F}_1} = \mathcal{N}(\mathbf{0}, \lambda^2 \mathbf{I})$, where $\lambda = 1$ is arbitrarily chosen. Different values of λ may improve (or worsen) the finite-time regrets (or the regret in a particular episode) achieved by TS and WTS under known variance, however, as discussed in [10], the expected value of the regrets asymptotically matches the lower bound predicted by [3] despite the value of $\lambda > 0$. For completeness, we summarize WTS under known σ in Algorithm 3. We only specify how WTS is implemented since TS can be directly recovered by choosing $M = 1$ in Algorithm 3.