

# Forced-exploration free Strategies for Unimodal Bandits

**Hassan Saber**

HASSAN.SABER@INRIA.FR

*SequeL Research Group  
Inria Lille-Nord Europe & CRIStAL  
Villeneuve-d'Ascq, Parc scientifique de la Haute-Borne, France*

**Pierre Ménard**

PIERRE.MENARD@INRIA.FR

*SequeL Research Group  
Inria Lille-Nord Europe & CRIStAL  
Villeneuve-d'Ascq, Parc scientifique de la Haute-Borne, France*

**Odalric-Ambrym Maillard**

ODALRIC.MAILLARD@INRIA.FR

*SequeL Research Group  
Inria Lille-Nord Europe & CRIStAL  
Villeneuve-d'Ascq, Parc scientifique de la Haute-Borne, France*

**Editor:**

## Abstract

We consider a multi-armed bandit problem specified by a set of Gaussian or Bernoulli distributions endowed with a unimodal structure. Although this problem has been addressed in the literature (Combes and Proutiere, 2014), the state-of-the-art algorithms for such structure make appear a forced-exploration mechanism. We introduce IMED-UB, the first forced-exploration free strategy that exploits the unimodal-structure, by adapting to this setting the Indexed Minimum Empirical Divergence (IMED) strategy introduced by Honda and Takemura (2015). This strategy is proven optimal. We then derive KLUCB-UB, a KLUCB version of IMED-UB, which is also proven optimal. Owing to our proof technique, we are further able to provide a concise finite-time analysis of both strategies in an unified way. Numerical experiments show that both IMED-UB and KLUCB-UB perform similarly in practice and outperform the state-of-the-art algorithms.

**Keywords:** Structured Bandits, Indexed Minimum Empirical Divergence, Optimal Strategy

## 1. Introduction

The multi-armed bandit problem is a popular framework to formalize sequential decision making problems. It was first introduced in the context of medical trials (Thompson, 1933, 1935) and later formalized by Robbins (1952): A bandit is specified by a set of unknown probability distributions  $\nu = (\nu_a)_{a \in \mathcal{A}}$  with means  $(\mu_a)_{a \in \mathcal{A}}$ . At each time  $t \in \mathbb{N}$ , the learner chooses an arm  $a_t \in \mathcal{A}$ , based only on the past, the learner then receives and observes a reward  $X_t$ , conditionally independent, sampled according to  $\nu_{a_t}$ . The goal of the learner is to maximize the expected sum of rewards received over time (up to some unknown horizon  $T$ ), or equivalently minimize the *regret* with respect to the strategy constantly receiving the highest mean reward

$$R(\nu, T) = \mathbb{E}_\nu \left[ \sum_{t=1}^T \mu^* - X_t \right] \text{ where } \mu^* = \max_{a \in \mathcal{A}} \mu_a.$$

Both means and distributions are *unknown*, which makes the problem non trivial, and the learner only knows that  $\nu \in \mathcal{D}$  where  $\mathcal{D}$  is a given set of bandit configurations. This problem received increased attention in the middle of the 20<sup>th</sup> century, and the seminal paper [Lai and Robbins \(1985\)](#) established the first lower bound on the cumulative regret, showing that designing a strategy that is optimal uniformly over a given set of configurations  $\mathcal{D}$  comes with a price. The study of the lower performance bounds in multi-armed bandits successfully lead to the development of asymptotically optimal strategies for specific configuration sets, such as the KLUCB strategy ([Lai, 1987](#); [Cappé et al., 2013](#); [Maillard, 2018](#)) for exponential families, or alternatively the DMED and IMED strategies from [Honda and Takemura \(2011, 2015\)](#). The lower bounds from [Lai and Robbins \(1985\)](#), later extended by [Burnetas and Katehakis \(1997\)](#) did not cover all possible configurations, and in particular *structured* configuration sets were not handled until [Agrawal et al. \(1989\)](#) and then [Graves and Lai \(1997\)](#) established generic lower bounds. Here, structure refers to the fact that pulling an arm may reveals information that enables to refine estimation of other arms. Unfortunately, designing numerical efficient strategies that are provably optimal remains a challenge for many structures.

**Structured configurations.** Motivated by the growing popularity of bandits in a number of industrial and societal application domains, the study of *structured configuration sets* has received increasing attention over the last few years: The linear bandit problem is one typical illustration ([Abbasi-Yadkori et al., 2011](#); [Srinivas et al., 2010](#); [Durand et al., 2017](#)), for which the linear structure considerably modifies the achievable lower bound, see [Lattimore and Szepesvari \(2017\)](#). The study of a *unimodal* structure naturally appears in many contexts, e.g. single-peak preference economics, voting theory or wireless communications, and has been first considered in [Yu and Mannor \(2011\)](#) from a bandit perspective, then in [Combes and Proutiere \(2014\)](#) providing an explicit lower bound together with a strategy exploiting this specific structure. Other structures include Lipschitz bandits [Magureanu et al. \(2014\)](#), and we refer to the manuscript [Magureanu \(2018\)](#) for other examples, such as cascading bandits that are useful in the context of recommender systems. In [Combes et al. \(2017\)](#), a generic strategy is introduced called OSB (Optimal Structured Stochastic Bandit), stepping the path towards generic multi-armed bandit strategies that are adaptive to a given structure.

**Unimodal-structure.** In this paper, we provide novel regret minimization results related to the following structure. We assume a *unimodal* structure similar to that considered in [Yu and Mannor \(2011\)](#) and [Combes and Proutiere \(2014\)](#). That is, there exists an undirected graph  $G = (\mathcal{A}, E)$  whose vertices are arms  $\mathcal{A}$ , and whose edges  $E$  characterize a partial order among means  $(\mu_a)_{a \in \mathcal{A}}$ . This partial order is assumed unknown to the learner. We assume that there exists a unique optimal arm  $a^* = \operatorname{argmax}_{a \in \mathcal{A}} \mu_a$  and that for all sub-optimal arm  $a \neq a^*$ , there exists a path  $P_a = (a_1 = a, \dots, a_{\ell_a} = a^*) \in \mathcal{A}^{\ell_a}$  of length  $\ell_a \geq 2$  such that for all  $i \in [1, \ell_a - 1]$ ,  $(a_i, a_{i+1}) \in E$  and  $\mu_{a_i} < \mu_{a_{i+1}}$ . Lastly, we assume that  $\nu \subset \mathcal{P} := \{p(\mu), \mu \in \Theta\}$ , where  $p(\mu)$  is an exponential-family distribution probability with density  $f(\cdot, \mu)$  with respect to some positive measure  $\lambda$  on  $\mathbb{R}$  and mean  $\mu \in \Theta \subset \mathbb{R}$ .  $\mathcal{P}$  is assumed to be known to the learner. Thus, for all  $a \in \mathcal{A}$  we have  $\nu_a = p(\mu_a)$ . We denote by  $\mathcal{D}_{(\mathcal{P}, G)}$  or simply  $\mathcal{D}$  the structured set of such unimodal-bandit distributions characterized by  $(\mathcal{P}, G)$ . In the following, we assume that  $\mathcal{P}$  is either the set of real Gaussian distributions with means in  $\mathbb{R}$  and variance 1 or the set of Bernouilli distributions with means in  $(0, 1)$ .

**Goal.** A key contribution in the study of unimodal bandits is the work [Combes and Proutiere \(2014\)](#), where the authors establish lower confidence bounds on the regret for the unimodal structure, and introduce an asymptotically optimal strategy called OSUB. One may then consider that unimodal

bandits are solved. Unfortunately, a closer look at the proposed approach reveals that the considered strategy forces some arms to be played (this is different than what is called forced exploration in structured bandits; it is rather a forced exploitation scheme). In this paper, our goal is to introduce alternative strategies to OSUB, that do not use any such forcing scheme, but consider variants of the pseudo-index induced by the lower bound analysis. Whether or not forcing mechanisms are desirable features is currently still under debate in the community; by providing the first strategy without any requirement for forcing in a structured bandit setup, we show that such mechanisms are not always required, which we believe opens an interesting avenue of research.

**Contributions.** In this paper, we first revisit the Indexed Minimum Empirical Divergence (IMED) strategy from [Honda and Takemura \(2011\)](#) introduced for unstructured multi-armed bandits, and adapt it to the unimodal-structured setting. We introduce in Section 3 the IMED-UB strategy that is limited to the pulling of the current best arm or their no more than  $d$  nearest arms at each time step, with  $d$  the maximum degree of nodes in  $G$ . Being constructed from IMED, IMED-UB does not require any optimization procedure and does not separate exploration from exploitation rounds. IMED-UB appears to be a *local* strategy. Motivated by practical considerations, under the assumption that  $G$  is a tree, when the number of arms  $|\mathcal{A}|$  becomes large, we further develop d-IMED-UB, an algorithm that behaves like IMED-UB while resorting to a dichotomic second order exploration over all nodes of the graph. This helps quickly identify the best arm  $a^*$  within a large set of arms  $\mathcal{A}$  by empirical considerations. We also introduce for completeness the KLUCB-UB strategy, that is similar to IMED-UB, but inspired from UCB strategies. We prove in Theorem 9 that IMED-UB, d-IMED-UB and KLUCB-UB are asymptotically optimal strategies that do not require forcing scheme. Furthermore, our unified finite time analysis shows that IMED-UB and KLUCB-UB are closely related. Furthermore, these novel strategies significantly outperform OSUB in practice. This is confirmed by numerical illustrations on synthetic data. We believe that the construction of these algorithms together with the proof techniques developed in this paper are of independent interest for the bandit community.

**Notations.** Let  $\nu \in \mathcal{D}$ . Let  $\mu^* = \max_{a \in \mathcal{A}} \mu_a$  be the optimal mean and  $a^* = \operatorname{argmax}_{a \in \mathcal{A}} \mu_a$  be the optimal arm of  $\nu$ . We define for an arm  $a \in \mathcal{A}$  its sub-optimality gap  $\Delta_a = \mu^* - \mu_a$ . Considering an horizon  $T \geq 1$ , thanks to the chain rule we can rewrite the regret as follows:

$$R(\nu, T) = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}_\nu [N_a(T)], \quad (1)$$

where  $N_a(t) = \sum_{s=1}^t \mathbb{1}_{\{a_s=a\}}$  is the number of pulls of arm  $a$  at time  $t$ .

## 2. Regret Lower bound

In this subsection, we recall for completeness the known lower bound on the regret when we assume a unimodal structure. In order to obtain non trivial lower bound we consider strategies that are *consistent* (aka uniformly-good).

**Definition 1 (Consistent strategy)** A strategy is consistent on  $\mathcal{D}$  if for all configuration  $\nu \in \mathcal{D}$ , for all sub-optimal arm  $a$ , for all  $\alpha > 0$ ,

$$\lim_{T \rightarrow \infty} \mathbb{E}_\nu \left[ \frac{N_a(T)}{T^\alpha} \right] = 0.$$

We can derive from the notion of consistency an asymptotic lower bound on the regret, see [Combes and Proutiere \(2014\)](#). To this end, we introduce  $\mathcal{V}_a = \{a' \in \mathcal{A} : (a, a') \in E\}$  to denote the neighbourhood of an arm  $a \in \mathcal{A}$ .

**Proposition 2 (Lower bounds on the regret)** *Let us consider a consistent strategy. Then, for all configuration  $\nu \in \mathcal{D}$ , it must be that*

$$\liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \geq c(\nu) := \sum_{a \in \mathcal{V}_{a^*}} \frac{\Delta_a}{\text{KL}(\mu_a | \mu^*)},$$

where  $\text{KL}(\mu | \mu') = \int_{\mathbb{R}} \log(f(x, \mu) / f(x, \mu')) f(x, \mu) \lambda(dx)$  denotes the Kullback-Leibler divergence between  $\nu = p(\mu)$  and  $\nu' = p(\mu')$ , for  $\mu, \mu' \in \Theta$ .

**Remark 3** *The quantity  $c(\nu)$  is a fully explicit function of  $\nu$  (it does not require solving any optimization problem) for some set of distributions  $\nu$  (see Remark 4). This useful property no longer holds in general for arbitrary structures. Also, it is noticeable that  $c(\nu)$  does not involve all the sub-optimal arms but only the ones in  $\mathcal{V}_{a^*}$ . This indicates that sub-optimal arms outside  $\mathcal{V}_{a^*}$  are sampled  $o(\log(T))$ , which contrasts with the unstructured stochastic multi-armed bandits. See [Combes and Proutiere \(2014\)](#) for further insights.*

**Remark 4** *For Gaussian distributions (variance  $\sigma^2 = 1$ ), we assume  $\lambda$  to be the Lebesgue measure,  $\Theta = \mathbb{R}$ , and for  $\mu \in \mathbb{R}$ ,  $f(\cdot, \mu) =: x \in \mathbb{R} \mapsto (\sqrt{2\pi})^{-1} e^{-(x-\mu)^2/2}$ . Then for all  $\mu, \mu' \in \mathbb{R}$ ,  $\text{KL}(\mu | \mu') = (\mu' - \mu)^2/2$ . For Bernoulli distributions, a possible setting is to assume  $\lambda = \delta_0 + \delta_1$  (with  $\delta_0, \delta_1$  Dirac measures),  $\Theta = (0, 1)$  and for  $\mu \in \Theta$ ,  $f(\cdot, \mu) =: x \in \{0, 1\} \mapsto \mu^x (1 - \mu)^{1-x}$ . Then for all  $\mu, \mu' \in [0, 1]$ ,  $\text{KL}(\mu | \mu') = \text{kl}(\mu | \mu')$ , where*

$$\text{kl}(\mu | \mu') := \begin{cases} 0 & \text{if } \mu = \mu', \\ +\infty & \text{if } \mu < \mu' = 1, \\ \mu \log\left(\frac{\mu}{\mu'}\right) + (1 - \mu) \log\left(\frac{1 - \mu}{1 - \mu'}\right) & \text{otherwise,} \end{cases}$$

with the convention  $0 \times \log(0) = 0$ .

### 3. Forced-exploration free strategies for unimodal-structured bandits

We present in this section three novel strategies that both match the asymptotic lower bound of Proposition 2. Two of these strategies are inspired by the Indexed Minimum Empirical Divergence (IMED) proposed by [Honda and Takemura \(2011\)](#). The other one is based on Kullback–Leibler Upper Confidence Bounds (KLUCB), using insights from IMED. The general idea behind these algorithms is, following the intuition given by the lower bound, to narrow on the current best arm and its neighbourhood for pulling an arm at a given time step.

**Notations.** The empirical mean of the rewards from the arm  $a$  is denoted by  $\hat{\mu}_a(t) = \sum_{s=1}^t \mathbb{I}_{\{a_s=a\}} X_s / N_a(t)$  if  $N_a(t) > 0$ , 0 otherwise. We also denote by  $\hat{\mu}^*(t) = \max_{a \in \mathcal{A}} \hat{\mu}_a(t)$  and  $\hat{\mathcal{A}}^*(t) = \text{argmax}_{a \in \mathcal{A}} \hat{\mu}_a(t)$  respectively the current best mean and the current set of optimal arms.

For convenience, we recall below the 0SUB (Optimal sampling for Unimodal Bandits) strategy from [Combes and Proutiere \(2014\)](#).

---

**Algorithm 1** OSUB

---

Pull an arbitrary arm  $a_1 \in \mathcal{A}$   
**for**  $t = 1 \dots T - 1$  **do**  
  Choose  $\hat{a}_t^* \in \operatorname{argmin}_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} N_{\hat{a}^*}(t)$  (chosen arbitrarily)  
  Pull  $a_{t+1} = \begin{cases} \hat{a}_t^* & \text{if } \frac{L_t(\hat{a}_t^*) - 1}{d+1} \in \mathbb{N} \\ \operatorname{argmax}_{a \in \mathcal{V}_{\hat{a}_t^*}} u_a(t) & \text{else} \end{cases}$   
**end for**

---

In Algorithm 1, for some numerical constant  $c > 0$ , the index computed by OSUB strategy for arm  $a \in \mathcal{A}$  and step  $t \geq 1$  is

$$u_a(t) = \sup \left\{ u \geq \hat{\mu}_a(t) : N_a(t) \operatorname{KL}(\hat{\mu}_a(t) | u) \leq f_c(L_t(\hat{a}_t^*)) \right\},$$

where  $L_t(a) = \sum_{t'=1}^t \mathbb{I}_{\{\hat{a}_{t'}^* = a\}}$  counts how many times arm  $a$  was a leader (best empirical arm),  $d$  is the maximum degree of nodes in  $G$ , and  $f_c(\cdot) = \log(\cdot) + c \log \log(\cdot)$ .

### 3.1 The IMED-UB strategy.

For all arm  $a \in \mathcal{A}$  and time step  $t \geq 1$  we introduce the IMED index

$$I_a(t) = N_a(t) \operatorname{KL}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) + \log(N_a(t)),$$

with the convention  $0 \times \infty = 0$ . This index can be seen as a transportation cost for moving a sub-optimal arm to an optimal one plus an exploration term: the logarithm of the numbers of pulls. When an optimal arm is considered, the transportation cost is null and there is only the exploration part. Note that, as stated in [Honda and Takemura \(2011\)](#),  $I_a(t)$  is an index in the weaker sense since it cannot be determined only by samples from the arm  $a$  but also uses empirical means of current optimal arms. We define IMED-UB (Indexed Minimum Empirical Divergence for Unimodal Bandits), described in Algorithm 2, to be the strategy consisting of pulling an arm  $a_t \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}$  with minimum index at each time step  $t$ , where is  $\hat{a}_t^* \in \operatorname{argmin}_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} N_{\hat{a}^*}(t)$  is a current best arm. This is a natural algorithm since the lower bound on the regret given in Proposition 2 involves only the arms in  $\mathcal{V}_{a^*}$ , the neighbourhood of the arm  $a^*$  of maximal mean.

---

**Algorithm 2** IMED-UB

---

Pull an arbitrary arm  $a_1 \in \mathcal{A}$   
**for**  $t = 1 \dots T - 1$  **do**  
  Choose  $\hat{a}_t^* \in \operatorname{argmin}_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} N_{\hat{a}^*}(t)$  (chosen arbitrarily)  
  Pull  $a_{t+1} \in \operatorname{argmin}_{a \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}} I_a(t)$  (chosen arbitrarily)  
**end for**

---

### 3.2 The KLUCB-UB strategy

For all arm  $a \in \mathcal{A}$  and time step  $t \geq 1$  we introduce the following Upper Confidence Bound

$$U_a(t) = \max \left\{ u \geq \hat{\mu}_a(t) \mid N_a(t) \text{KL}(\hat{\mu}_a(t) | u) + \log(N_a(t)) \leq \log(N_{\hat{a}_t^*}(t)) \right\}$$

with  $\hat{a}_t^* \in \operatorname{argmin}_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} N_{\hat{a}^*}(t)$ .

By convention, we set  $U_a(t) = \hat{\mu}_a(t)$  if for  $a \in \mathcal{A}$ ,  $\log(N_a(t)) > \log(N_{\hat{a}_t^*}(t))$ .

**Remark 5** A classical KLUCB strategy would replace the term  $\log(N_{\hat{a}_t^*}(t)/N_a(t))$  with  $\log(t)$ , and a KLUCB<sup>+</sup> would use  $\log(t/N_a(t))$ . This is a simple yet crucial modification. Indeed, although this makes KLUCB-UB not an index strategy, this enables to get a more intrinsic strategy, to simplify the analysis and get improved numerical results.

As for IMED-UB and IMED,  $U_a(t)$  is an index in a weaker sense since it cannot be determined only by samples from the arm  $a$  but also uses numbers of pulls of current optimal arms. We define KLUCB-UB (Kullback-Leibler Upper Confidence Bounds for Unimodal Bandits) to be the strategy consisting of pulling an arm  $a_t \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}$  with maximum index at each time step  $t$ . This algorithm can be seen as a KLUCB version of the IMED-UB strategy.

---

#### Algorithm 3 KLUCB-UB

---

```

Pull  $a_1 \in \mathcal{A}$  at random.
for  $t = 1 \dots T - 1$  do
    Choose  $\hat{a}_t^* \in \operatorname{argmin}_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} N_{\hat{a}^*}(t)$  (chosen arbitrarily)
    Pull  $a_{t+1} \in \operatorname{argmax}_{a \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}} U_a(t)$  (chosen arbitrarily)
end for

```

---

**Remark 6** IMED-UB does not require solving any optimization problem, unlike OSUB or KLUCB-UB. We believe this feature, inherited from IMED, makes it an especially appealing strategy. KLUCB-UB solves an optimization similar to that of the KLUCB strategy for unstructured bandits, and also related to the optimization used in OSUB from [Combes and Proutiere \(2014\)](#). The difference between KLUCB-UB and OSUB is that it does not use any forced exploitation.

### 3.3 The d-IMED-UB strategy for large set of arms

When the set of arms is large, a bad initialization of IMED-UB (that is, choose arm  $a_1$  far from  $a^*$ ) comes with high initial regret. Indeed, IMED-UB does not allow to explore outside the neighbourhood  $\mathcal{V}_{\hat{a}_t^*}$  of  $\hat{a}_t^*$ . When  $\mathcal{A}$  is large compared to the neighbourhoods, this may generate a large burn-in phase. To overcome this practical limitation, it is natural to explore outside the neighbourhood of the current best arm. However, to be compatible with the lower bound on the regret stated in Proposition 2 such exploration must be asymptotically negligible. We now consider  $\mathcal{G}$  to be a tree, and introduce d-IMED-UB, a strategy that trades-off between these two types of exploration. d-IMED-UB shares with IMED-UB the same exploitation criteria and explores if the index of the current best arm exceeds the indexes of arms in its neighbourhood. However, in exploration phase, d-IMED-UB runs an IMED type

strategy to choose between exploring within *or outside* the neighbourhood of the current best arm. For all time step  $t \geq 1$ , for all arm  $a' \in \mathcal{V}_{\hat{a}_t^*}$ , for all arm  $a \in \hat{G}_{a'}(t)$ , where  $G_{a'}(t)$  denotes the sub-tree containing  $a'$  obtained by cutting edge  $(a', \hat{a}_t^*)$ , we define the second order IMED index relative to  $a'$ , as

$$I_a^{(a')}(t) = N_a(t) \text{KL}^+(\hat{\mu}_a(t) | \hat{\mu}_{a'}(t)) + \log(N_a(t)) ,$$

where  $\text{KL}^+(\mu | \mu') = \text{KL}(\mu | \mu')$  if  $\mu < \mu'$ , 0 otherwise. At each exploration time step, d-IMED-UB pulls an arm in  $\mathcal{S}_t$  with minimal secondary index relative to the arm  $\underline{a}_t$  with current minimal index and belonging to the neighbourhood of the current best arm, where  $\mathcal{S}_t$  is a sub-tree of  $\hat{G}_{\underline{a}_t}(t)$  dichotomously chosen that contains  $\underline{a}_t$ . We illustrate in Appendix E, a way to dynamically choose  $\mathcal{S}_t$ .

**Remark 7** Assuming that  $G$  is a tree ensures that for all  $a' \in \mathcal{V}_{a^*}$ , the nodes of  $G_{a'}$ , the sub-tree containing  $a'$  obtained by cutting edge  $(a', a^*)$ , induce a unimodal bandit configuration with optimal arm  $a'$ . This specific property allows establishing the optimality of d-IMED-UB.

---

**Algorithm 4** d-IMED-UB

---

Pull an arbitrary arm  $a_1 \in \mathcal{A}$   
**for**  $t = 1 \dots T - 1$  **do**  
    Choose  $\hat{a}_t^* \in \underset{\hat{a}^* \in \hat{\mathcal{A}}^*(t)}{\text{argmin}} N_{\hat{a}^*}(t)$  (chosen arbitrarily)  
    Choose  $\underline{a}_t \in \underset{a \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}}{\text{argmin}} I_a(t)$  (chosen arbitrarily)  
    **if**  $\underline{a}_t = \hat{a}_t^*$  **then**  
        Pull  $a_{t+1} = \underline{a}_t$   
    **else**  
        Pull  $a_{t+1} \in \underset{a \in \mathcal{S}_t}{\text{argmin}} I_a^{(\underline{a}_t)}(t)$   
    **end if**  
**end for**

---

### 3.4 Asymptotic optimality of IMED-UB, d-IMED-UB and KLUCB-UB

In this section, we state the main theoretical result of this paper.

**Theorem 8 (Upper bounds)** Let us consider a set of Gaussian or Bernoulli distributions  $\nu \in \mathcal{D}$  and let  $a^*$  its optimal arm. Let  $\mathcal{V}_{a^*}$  be the sub-optimal arms in the neighbourhood of  $a^*$ . Then under IMED-UB and KLUCB-UB strategies for all  $0 < \varepsilon < \varepsilon_\nu$ , for all horizon time  $T \geq 1$ , for all  $a \in \mathcal{V}_{a^*}$ ,

$$\mathbb{E}_\nu[N_a(T)] \leq \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a | \mu_{a^*})} \log(T) + d |\mathcal{A}|^2 C_\varepsilon + 1$$

and, for all  $a \notin \{a^*\} \cup \mathcal{V}_{a^*}$ ,

$$\mathbb{E}_\nu[N_a(T)] \leq d |\mathcal{A}|^2 C_\varepsilon + 1 ,$$

where  $d$  is the maximum degree of nodes in  $G$ ,  $\varepsilon_\nu = \min \{1 - \mu^*, \min_{a \neq a'} |\mu_a - \mu_{a'}|/4\}$ ,  $C_\varepsilon = 34 \log(1/\varepsilon) \varepsilon^{-6}$  and where  $\alpha_\nu(\cdot)$  is a non-negative function depending only on  $\nu$  such that  $\lim_{\varepsilon \rightarrow 0} \alpha_\nu(\varepsilon) = 0$  (see Section 4.1 for more details).

Furthermore, if the considered graph is a tree, then under  $d$ -IMED-UB, for all horizon  $T \geq 1$ , for all  $a \in \mathcal{V}_{a^*}$ ,

$$\mathbb{E}_\nu[N_a(T)] \leq \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a|\mu_{a^*})} \log(T) + d |\mathcal{A}|^2 C_\varepsilon + 1$$

and, for all  $a \notin \{a^*\} \cup \mathcal{V}_{a^*}$ ,

$$\begin{aligned} \mathbb{E}_\nu[N_a(T)] &\leq \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_a|\mu_{\underline{a}})} \log\left(\frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_{\underline{a}}|\mu_{a^*})} \log(T)\right) \\ &\quad + d |\mathcal{A}|^2 C_\varepsilon + 1. \end{aligned}$$

In particular one can note that the arms in the neighbourhood of the optimal one are pulled  $\mathcal{O}(\log(T))$  times while the other sub-optimal arms are pulled a finite number of times under IMED-UB and KLUCB-UB, and  $\mathcal{O}(\log\log(T))$  times under  $d$ -IMED-UB. This is coherent with the lower bound that only involves the neighbourhood of the best arm. More precisely, combining Theorem 8 and the chain rule (1) gives the asymptotic optimality of IMED-UB and KLUCB-UB with respect to the lower bound of Proposition 2.

**Corollary 9 (Asymptotic optimality)** *With the same notations as in Theorem 8, then under IMED-UB and KLUCB-UB strategies*

$$\limsup_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \leq c(\nu) = \sum_{a \in \mathcal{V}_{a^*}} \frac{\Delta_a}{\text{KL}(\mu_a|\mu^*)}.$$

*If the considered graph is a tree, same result holds under  $d$ -IMED-UB strategy.*

See respectively Section 4 and Appendix C for a finite time analysis of IMED-UB,  $d$ -IMED-UB and KLUCB-UB.

#### 4. IMED-UB finite time analysis

At a high level, the key interesting step of the proof is to realize that the considered strategies imply empirical lower and empirical upper bounds on the numbers of pulls (see Lemma 10, Lemma 11 for IMED-UB). Then, based on concentration lemmas (see Section A.1), the strategy-based empirical lower bounds ensure the reliability of the estimators of interest (Lemma 14). This makes use of more classical arguments based on concentration of measure. Then, combining the reliability of these estimators with the obtained strategy-base empirical upper bounds, we obtain upper bounds on the average numbers of pulls (Theorem 8).

In this section, we only detail the finite time analysis of IMED-UB algorithm and defer those of  $d$ -IMED-UB and KLUCB-UB to the appendix, as it follows essentially the same steps. Indeed, we show that KLUCB and  $d$ -IMED-UB strategies imply empirical bounds (Lemmas 19,20, Lemmas 25,26) very similar to IMED-UB strategy. These inequalities are the cornerstone of the analysis. We believe that this general way of proceeding is of independent interest as it simplifies the proof steps.



## 4.1 Notations

Let us consider  $\nu \in \mathcal{D}$  and let us denote by  $a^*$  its best arm. We recall that for all  $a \in \mathcal{A}$ ,  $\mathcal{V}_a = \{a' \in \mathcal{A} : (a, a') \in E\}$  is the neighbourhood of arm  $a$  in graph  $G = (\mathcal{A}, E)$ , and that

$$\widehat{d} = \max_{a \in \mathcal{A}} |\mathcal{V}_a|, \quad \varepsilon_\nu = \min \left\{ 1 - \mu^*, \min_{a \neq a'} \frac{|\mu_a - \mu_{a'}|}{4} \right\}.$$

Then, there exists a function  $\alpha_\nu(\cdot)$  such that for all  $a \neq a'$ , for all  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\frac{\text{kl}(\mu_a | \mu_{a'})}{1 + \alpha_\nu(\varepsilon)} \leq \text{kl}(\mu_a + \varepsilon | \mu_{a'} - \varepsilon) \leq (1 + \alpha_\nu(\varepsilon)) \text{kl}(\mu_a | \mu_{a'})$$

and  $\lim_{\varepsilon \downarrow 0} \alpha_\nu(\varepsilon) = 0$ . For all studied strategy, at each time step  $t \geq 1$ ,  $\widehat{a}_t^*$  is arbitrarily chosen in

$\text{argmin}_{a \in \widehat{\mathcal{A}}^*(t)} N_a(t)$  where  $\widehat{\mathcal{A}}^*(t) = \text{argmax}_{a \in \mathcal{A}} \widehat{\mu}_a(t)$ .

For all arms  $a \in \mathcal{A}$  and  $n \geq 1$ , we introduce the stopping times  $\tau_{a,n} = \inf \{t \geq 1 : N_a(t) = n\}$  and define the empirical means corresponding to local times

$$\widehat{\mu}_a^n = \frac{1}{n} \sum_{m=1}^n X_{\tau_{a,m}}.$$

For a subset of times  $\mathcal{E} \subset \{t \geq 1\}$ , we denote by  $\mathcal{E}^c$  its complementary in  $\{t \geq 1\}$ .

## 4.2 Strategy-based empirical bounds

IMED-UB strategy implies inequalities between the indexes that can be rewritten as inequalities on the numbers of pulls. While lower bounds involving  $\log(t)$  may be expected in view of the asymptotic regret bounds, we show lower bounds on the numbers of pulls involving instead  $\log(N_{a_{t+1}}(t))$ , the logarithm of the number of pulls of the current chosen arm. We also provide upper bounds on  $N_{a_{t+1}}(t)$  involving  $\log(t)$ .

We believe that establishing these empirical lower and upper bounds is a key element of our proof technique, that is of independent interest and not *a priori* restricted to the unimodal structure.

**Lemma 10 (Empirical lower bounds)** *Under IMED-UB, at each step time  $t \geq 1$ , for all  $a \in \mathcal{V}_{\widehat{a}_t^*}$ ,*

$$\log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\widehat{\mu}_a(t) | \widehat{\mu}^*(t)) + \log(N_a(t))$$

and

$$N_{a_{t+1}}(t) \leq N_{\widehat{a}_t^*}(t).$$

**Proof** For  $a \in \mathcal{A}$ , by definition, we have  $I_a(t) = N_a(t) \text{KL}(\widehat{\mu}_a(t) | \widehat{\mu}^*(t)) + \log(N_a(t))$ , hence

$$\log(N_a(t)) \leq I_a(t).$$

This implies, since the arm with minimum index is pulled,  $\log(N_{a_{t+1}}(t)) \leq I_{a_{t+1}}(t) = \min_{a' \in \{\widehat{a}_t^*\} \cup \mathcal{V}_{\widehat{a}_t^*}} I_{a'}(t) \leq$

$I_{\widehat{a}_t^*}(t) = \log(N_{\widehat{a}_t^*}(t))$ . By taking the  $\exp(\cdot)$ , the last inequality allows us to conclude.  $\blacksquare$

**Lemma 11 (Empirical upper bounds)** Under IMED-UB at each step time  $t \geq 1$ ,

$$N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t) | \hat{\mu}^*(t)) \leq \log(t).$$

**Proof** As above, by construction we have

$$I_{a_{t+1}}(t) \leq I_{\hat{a}_t^*}(t).$$

It remains, to conclude, to note that

$$N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t) | \hat{\mu}^*(t)) \leq I_{a_{t+1}}(t),$$

and

$$I_{\hat{a}_t^*}(t) = \log(N_{\hat{a}_t^*}(t)) \leq \log(t).$$

■

### 4.3 Reliable current best arm and means

In this subsection, we consider the subset  $\mathcal{T}_\varepsilon$  of times where everything is well behaved: The current best arm corresponds to the true one and the empirical means of the best arm and the current chosen arm are  $\varepsilon$ -accurate for  $0 < \varepsilon < \varepsilon_\nu$ , that is

$$\mathcal{T}_\varepsilon := \left\{ t \geq 1 : \begin{array}{l} \hat{\mathcal{A}}^*(t) = \{a^*\} \\ \forall a \in \{a^*, a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \end{array} \right\}.$$

We will show that its complementary set is finite on average. In order to prove this we decompose the set  $\mathcal{T}_\varepsilon$  in the following way. Let  $\mathcal{E}_\varepsilon$  be the set of times where the means are well estimated,

$$\mathcal{E}_\varepsilon := \left\{ t \geq 1 : \forall a \in \hat{\mathcal{A}}^*(t) \cup \{a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \right\},$$

and  $\Lambda_\varepsilon$  the set of times where an arm that is not the current optimal neither pulled is underestimated

$$\Lambda_\varepsilon := \left\{ t \geq 1 : \exists a \in \mathcal{V}_{\hat{a}_t^*} \setminus \{a_{t+1}, \hat{a}_t^*\} \text{ s.t. } \hat{\mu}_a(t) < \mu_a - \varepsilon \text{ and } \log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t) | \mu_a - \varepsilon) + \log(N_a(t)) \right\}.$$

Then we prove below the following inclusion.

**Lemma 12 ( Relations between the subsets of times)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathcal{T}_\varepsilon^c \setminus \mathcal{E}_\varepsilon^c \subset \Lambda_\varepsilon. \quad (2)$$

**Proof** Let us consider  $t \in \mathcal{T}_\varepsilon^c \setminus \mathcal{E}_\varepsilon^c$ . Since  $t \in \mathcal{E}_\varepsilon$  and  $\varepsilon < \varepsilon_\nu$  we have

$$\forall a \in \hat{\mathcal{A}}^*(t) \cup \{a_{t+1}\}, \quad |\hat{\mu}_a(t) - \mu_a| < \varepsilon.$$

By triangle inequality this implies, for all  $\hat{a}^* \in \hat{\mathcal{A}}^*(t)$ ,

$$\begin{aligned} |\mu_{\hat{a}_t^*} - \mu_{\hat{a}^*}| - 2\varepsilon &\leq |\mu_{\hat{a}_t^*} - \mu_{\hat{a}^*}| - |\mu_{\hat{a}_t^*} - \hat{\mu}_{\hat{a}_t^*}(t)| - |\hat{\mu}_{\hat{a}_t^*}(t) - \mu_{\hat{a}^*}| \\ &\leq |\hat{\mu}_{\hat{a}_t^*}(t) - \hat{\mu}_{\hat{a}_t^*}(t)| = 0 \end{aligned}$$

and

$$\widehat{\mathcal{A}}^*(t) = \{\widehat{a}_t^*\} .$$

Thus, since  $t \notin \mathcal{T}_\varepsilon$ , we have  $\widehat{a}_t^* \neq a^*$ . In particular, since  $(\mu_a)_{a \in \mathcal{A}}$  is unimodal, there exists  $a \in \mathcal{V}_{\widehat{a}_t^*}$  such that  $\mu_a > \mu_{\widehat{a}_t^*}$ . From Lemma 10 we have the following empirical lower bound

$$\log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\widehat{\mu}_a(t) | \widehat{\mu}^*(t)) + \log(N_a(t)) .$$

Furthermore, since  $t \in \mathcal{E}_\varepsilon$  and  $\varepsilon < \varepsilon_\nu$ , we have

$$\widehat{\mu}_a(t) \leq \widehat{\mu}^*(t) = \widehat{\mu}_{\widehat{a}_t^*}(t) < \mu_{\widehat{a}_t^*} + \varepsilon < \mu_a - \varepsilon .$$

Since  $|\widehat{\mu}_{a_{t+1}}(t) - \mu_{a_{t+1}}| < \varepsilon$ , it indicates in particular that  $a \in \mathcal{V}_{\widehat{a}_t^*} \setminus \{a_{t+1}, \widehat{a}_t^*\}$ . In addition, the monotony of the  $\text{KL}(\cdot | \cdot)$  implies

$$\text{KL}(\widehat{\mu}_a(t) | \widehat{\mu}^*(t)) \leq \text{KL}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) .$$

Therefore for such  $t$  we have  $\widehat{\mu}_a(t) < \mu_a - \varepsilon$  and

$$\log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) + \log(N_a(t)) ,$$

which concludes the proof. ■

We can now resort to classical concentration arguments in order to control the size of these sets, which yields the following upper bounds. We defer the proof to Appendix A.2 as they follow standard arguments.

**Lemma 13 (Bounded subsets of times)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] \leq \frac{10 |\mathcal{A}|^2}{\varepsilon^4} \quad \mathbb{E}_\nu[|\Lambda_\varepsilon|] \leq 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6} ,$$

where  $d$  is the maximum degree of nodes in  $G$ .

Thus combining them with (2) we obtain

$$\begin{aligned} \mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] &\leq \mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] + \mathbb{E}_\nu[|\Lambda_\varepsilon|] \\ &\leq \frac{10 |\mathcal{A}|^2}{\varepsilon^4} + 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6} \\ &\leq 33d |\mathcal{A}|^2 \frac{\log(1/\varepsilon)}{\varepsilon^6} . \end{aligned}$$

Hence, we just proved the following lemma.

**Lemma 14 (Reliable estimators)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] \leq 33d |\mathcal{A}|^2 \frac{\log(1/\varepsilon)}{\varepsilon^6} ,$$

where  $d$  is the maximum degree of nodes in  $G$ .

#### 4.4 Upper bounds on the numbers of pulls of sub-optimal arms

In this section, we now combine the different results of the previous sections to prove Theorem 8.

**Proof** [Proof of Theorem 8.] From Lemma 14, considering the following subset of times

$$\mathcal{T}_\varepsilon := \left\{ t \geq 1 : \begin{array}{l} \widehat{\mathcal{A}}^\star(t) = \{a^\star\} \\ \forall a \in \{a^\star, a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \end{array} \right\}.$$

we have

$$\mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] \leq 33d|\mathcal{A}|^2 \frac{\log(1/\varepsilon)}{\varepsilon^6},$$

$|\mathcal{A}| = 11$   $|\mathcal{A}| = 10^2$   $|\mathcal{A}| = 10^3$   $|\mathcal{A}| = 10^4$  where  $d$  is the maximum degree of nodes in  $G$ . Then, let us consider  $a \neq a^\star$  and a time step  $t \in \mathcal{T}_\varepsilon$  such that  $a_{t+1} = a$ . From Lemma 11 we get

$$N_a(t) \text{KL}(\hat{\mu}_a(t) | \hat{\mu}^\star(t)) \leq \log(t) \leq \log(T).$$

Furthermore, since  $t \in \mathcal{T}_\varepsilon$ , we have

$$\hat{a}_t^\star = a^\star \quad \text{and} \quad |\hat{\mu}_a(t) - \mu_a|, |\hat{\mu}_{a^\star}(t) - \mu_{a^\star}| < \varepsilon.$$

According to the strategy  $a = a_{t+1} \in \mathcal{V}_{a^\star}$  and by construction of  $\alpha_\nu(\cdot)$  (see Section 4.1 Notations)

$$\text{KL}(\hat{\mu}_a(t) | \hat{\mu}^\star(t)) = \text{KL}(\hat{\mu}_a(t) | \hat{\mu}_{a^\star}(t)) \geq \frac{\text{KL}(\mu_a | \mu_{a^\star})}{1 + \alpha_\nu(\varepsilon)}$$

and

$$N_a(t) \leq \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a | \mu_{a^\star})} \log(T).$$

Thus, we have shown that for  $a \neq a^\star$ ,

$$\forall t \in \mathcal{T}_\varepsilon \text{ s.t. } a_{t+1} = a : a \in \mathcal{V}_{a^\star}$$

and

$$N_a(t) \leq \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a | \mu_{a^\star})} \log(T).$$

This implies:

$$N_a(T) \leq \begin{cases} \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a | \mu_{a^\star})} \log(T) + |\mathcal{T}_\varepsilon^c| + 1, & \text{if } a \in \mathcal{V}_{a^\star} \\ |\mathcal{T}_\varepsilon^c| + 1, & \text{otherwise.} \end{cases}$$

Averaging these inequalities allows us to conclude. ■

## 5. Numerical experiments

In this section, we consider Gaussian distributions with variance  $\sigma^2 = 1$  and compare empirically the following strategies introduced beforehand: OSUB described in Algorithm 1, IMED-UB, d-IMED-UB described in Algorithms 2,4, KLUCB-UB described in Algorithm 3 as well as the baseline IMED by Honda and Takemura (2011) that does not exploit the structure and finally the generic OSSB strategy by Combes et al. (2017) that adapts to several structures. We compare these strategies on two setups.

**Fixed configuration** (Figure 1). For the first experiments we consider a small number of arms  $|\mathcal{A}| = 11$  and investigate these strategies over 500 runs on *fixed* Gaussian configuration  $\nu^0 \in \mathcal{D}$  with means  $(\mu_a^0)_{a \in \mathcal{A}} = (0, 0.2, 0.4, 0.6, 0.8, 1, 0.8, 0.6, 0.4, 0.2, 0)$ .

**Random configurations** (Figure 2). In this experiment we consider larger numbers of arms  $|\mathcal{A}| \in \{10^2, 10^3, 10^4\}$  and average regrets over 500 random Gaussian configurations uniformly sampled in  $\{\nu \in \mathcal{D} : (\mu_a)_{a \in \mathcal{A}} \in [0, 1]^{\mathcal{A}}\}$ .

It seems that for a small number of arms IMED-UB and KLUCB-UB perform better than the baseline IMED whereas OSUB performs very poorly for unimodal structure (this may be the price its genericity). Both IMED-UB and KLUCB-UB outperform OSUB significantly. When the set of arms becomes larger, only d-IMED-UB benefits from the unimodal structure and outperforms the baseline IMED.

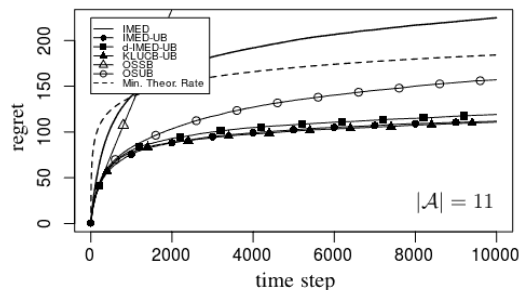


Figure 1: Regret approximated over 500 runs for  $\nu_0 \in \mathcal{D}$ .

**Remark 15** It is generally observed in bandit problems that theoretical asymptotic lower bounds on the regret are larger than the actual regret in finite horizon, as is it in Figure 1.

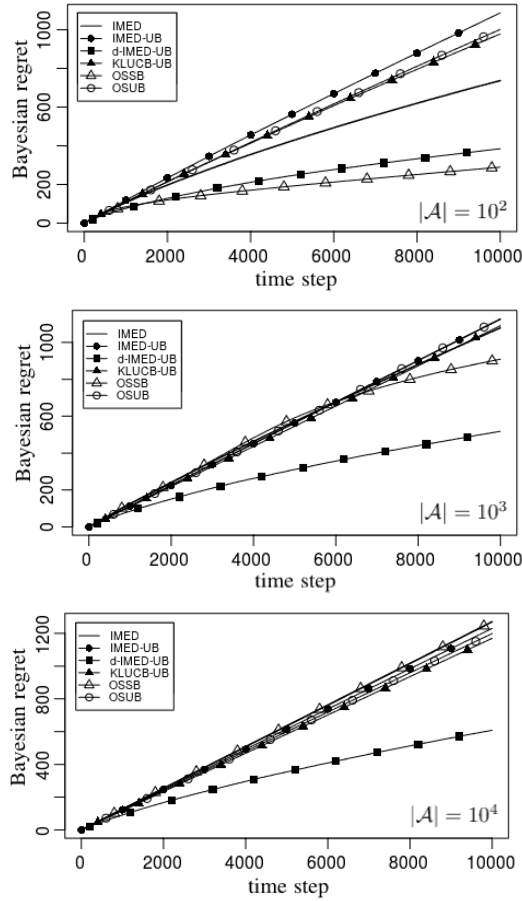


Figure 2: Regret averaged over 500 random configurations in  $\mathcal{D}$ .

## Conclusion

In this paper, we have revisited the setup of unimodal multi-armed bandits: We introduced three novel variants, two based on the IMED strategy and a second one using a KLUCB type index but modified using tools similar to IMED. These strategies do not require forcing to play specific arms (unlike for instance OSUB) on top of the naturally introduced score. Remarkably, the IMED-UB and d-IMED-UB strategies do not require any optimization procedure, which can be interesting for practitioners. We also provided a novel proof strategy (inspired from IMED), in which we make explicit empirical lower and upper bounds, before tackling the handling of bad events by more standard concentration tools. This proof technique greatly simplifies and shortens the analysis of IMED-UB (compared to that of OSUB), and is also employed to analyze KLUCB-UB and d-IMED-UB, in a somewhat unified way. Last, we provided numerical experiments that show the practical advantages of the novel approach over the OSUB strategy.

## Appendix A. IMED-UB finite time analysis

We regroup in this section, for completeness, the proofs of the remaining lemmas used in the analysis of IMED-UB in Section 4.

### A.1 Concentration lemmas

We state two concentration lemmas that do not depend on the followed strategy. Lemma 16 comes from Lemma B.1 in [Combes and Proutiere \(2014\)](#) and Lemma 17 comes from Lemma 14 in [Honda and Takemura \(2015\)](#). Proofs are provided in Appendix B.

**Lemma 16 (Concentration inequalities)** *Independently of the considered strategy, for all set of Gaussian or Bernoulli distributions  $\nu \in \mathcal{D}$ , for all  $0 < \varepsilon \leq 1/2$ , for all  $a, a' \in \mathcal{A}$ , we have*

$$\mathbb{E}_\nu \left[ \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a, N_{a'}(t) \geq N_a(t), |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon\}} \right] \leq \frac{10}{\varepsilon^4}.$$

**Lemma 17 (Large deviation probabilities)** *Let us consider a set of Gaussian or Bernoulli distributions  $\nu \in \mathcal{D}$ . Let  $0 < \varepsilon \leq \min(1 - \mu^*, 1/2)$  and  $a \in \mathcal{A}$ . Let  $\lambda = \mu_a - \varepsilon$ . Then, independently of the considered strategy, we have*

$$\mathbb{E}_\nu \left[ \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda\}} n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda)) \right] \leq \frac{23 \log(1/\varepsilon)}{\varepsilon^6}.$$

### A.2 Proof of Lemma 13 (Bounded subsets of times)

Using Lemma 10 we have

$$\forall t \geq 1, \quad N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t).$$

Since  $\hat{a}_t^* \in \operatorname{argmin}_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} N_{\hat{a}^*}(t)$ , this implies

$$\forall t \geq 1, \forall \hat{a}^* \in \hat{\mathcal{A}}^*(t), \quad N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t) \leq N_{\hat{a}^*}(t).$$

Then, based on the concentration inequalities from Lemma 16, we obtain

$$\begin{aligned} \mathbb{E}_\nu [|\mathcal{E}_\varepsilon^c|] &\leq \sum_{a, a' \in \mathcal{A}} \mathbb{E}_\nu \left[ \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a, N_{a'}(t) \geq N_a(t), |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon\}} \right] \\ &\leq \sum_{a, a' \in \mathcal{A}} \frac{10}{\varepsilon^4} \\ &\leq \frac{10 |\mathcal{A}|^2}{\varepsilon^4}. \end{aligned}$$

Furthermore, for  $t \geq 1$  and  $a \in \mathcal{A}$ , we have

$$\log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t) | \lambda_a) + \log(N_a(t)) \Leftrightarrow N_{a_{t+1}}(t) \leq N_a(t) \exp(N_a(t) \text{KL}(\hat{\mu}_a(t) | \lambda_a)),$$

where  $\lambda_a = \mu_a - \varepsilon$  for all arm  $a \in \mathcal{A}$ .  
Thus, we have

$$\begin{aligned}
|\Lambda_\varepsilon| &\leq \sum_{t \geq 1} \sum_{a \in \mathcal{V}_{\hat{a}_t^*} \setminus \{a_{t+1}, \hat{a}_t^*\}} \mathbb{I}_{\{\hat{\mu}_a(t) < \lambda_a \text{ and } N_{a_{t+1}}(t) \leq N_a(t) \exp(N_a(t) \text{KL}(\hat{\mu}_a(t) | \lambda_a))\}} \\
&= \sum_{t \geq 1} \sum_{\hat{a}^* \in \mathcal{A}} \sum_{a' \in \{\hat{a}^*\} \cup \mathcal{V}_{\hat{a}^*}} \sum_{a \in \mathcal{V}_{\hat{a}^*} \setminus \{a', \hat{a}^*\}} \sum_{n \geq 1} \mathbb{I}_{\{\hat{a}_t^* = \hat{a}^*, a_{t+1} = a', N_a(t) = n\}} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda_a, N_{a'}(t) \leq n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda_a))\}} \\
&\leq \sum_{t \geq 1} \sum_{\hat{a}^* \in \mathcal{A}} \sum_{a' \in \{\hat{a}^*\} \cup \mathcal{V}_{\hat{a}^*}} \sum_{a \in \mathcal{V}_{\hat{a}^*} \setminus \{a', \hat{a}^*\}} \sum_{n \geq 1} \mathbb{I}_{\{a_{t+1} = a'\}} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda_a\}} \mathbb{I}_{\{N_{a'}(t) \leq n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda_a))\}} \\
&= \sum_{\hat{a}^* \in \mathcal{A}} \sum_{a' \in \{\hat{a}^*\} \cup \mathcal{V}_{\hat{a}^*}} \sum_{a \in \mathcal{V}_{\hat{a}^*} \setminus \{a', \hat{a}^*\}} \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda_a\}} \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1} = a' \text{ and } N_{a'}(t) \leq n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda_a))\}} \\
&\leq \sum_{\hat{a}^* \in \mathcal{A}} \sum_{a' \in \{\hat{a}^*\} \cup \mathcal{V}_{\hat{a}^*}} \sum_{a \in \mathcal{V}_{\hat{a}^*} \setminus \{a', \hat{a}^*\}} \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda_a\}} n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda_a)) \\
&\leq \sum_{\hat{a}^* \in \mathcal{A}} \sum_{a' \in \{\hat{a}^*\} \cup \mathcal{V}_{\hat{a}^*}} \sum_{a \in \mathcal{V}_{\hat{a}^*} \setminus \{a', \hat{a}^*\}} \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda_a\}} n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda_a))
\end{aligned}$$

and

$$\mathbb{E}_\nu[|\Lambda_\varepsilon|] \leq \sum_{\hat{a}^* \in \mathcal{A}} \sum_{a' \in \{\hat{a}^*\} \cup \mathcal{V}_{\hat{a}^*}} \sum_{a \in \mathcal{V}_{\hat{a}^*} \setminus \{a', \hat{a}^*\}} \mathbb{E}_\nu \left[ \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda_a\}} n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda_a)) \right].$$

Then, by applying Lemma 17 based on large deviation probabilities, we have

$$\forall a \in \mathcal{A}, \quad \mathbb{E}_\nu \left[ \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda_a\}} n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda_a)) \right] \leq \frac{23 \log(1/\varepsilon)}{\varepsilon^6}.$$

It comes:

$$\mathbb{E}_\nu[|\Lambda_\varepsilon|] \leq \sum_{\hat{a}^* \in \mathcal{A}} \sum_{a' \in \{\hat{a}^*\} \cup \mathcal{V}_{\hat{a}^*}} \sum_{a \in \mathcal{V}_{\hat{a}^*} \setminus \{a', \hat{a}^*\}} \frac{23 \log(1/\varepsilon)}{\varepsilon^6} \leq 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6}.$$

## Appendix B. Concentration lemmas

**Lemma** Independently of the considered strategy, for all set of Gaussian or Bernoulli distributions  $\nu \in \mathcal{D}$ , for all  $0 < \varepsilon \leq 1/2$ , for all  $a, a' \in \mathcal{A}$ , we have

$$\mathbb{E}_\nu \left[ \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1} = a, N_{a'}(t) \geq N_a(t), |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon\}} \right] \leq \frac{10}{\varepsilon^4}.$$

**Proof** Considering the stopping times  $\tau_{a,n} = \inf \{t \geq 1, N_a(t) = n\}$  we will rewrite the sum  $\sum_{t \geq 1} \mathbb{I}_{\{a_{t+1} = a, N_{a'}(t) \geq N_a(t), |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon\}}$  and use an Hoeffding's type argument for distributions with



support included in  $[0, 1]$ .

$$\begin{aligned}
& \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a, N_{a'}(t) \geq N_a(t), |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon\}} \\
&= \sum_{t \geq 1} \sum_{n \geq 2, m \geq 1} \mathbb{I}_{\{\tau_{a,n}=t+1, N_{a'}(t)=m\}} \mathbb{I}_{\{m \geq n-1, |\hat{\mu}_{a'}^m - \mu_{a'}| \geq \varepsilon\}} \\
&= \sum_{m \geq 1} \sum_{n \geq 2} \mathbb{I}_{\{m \geq n-1, |\hat{\mu}_{a'}^m - \mu_{a'}| \geq \varepsilon\}} \sum_{t \geq 1} \mathbb{I}_{\{\tau_{a,n}=t+1, N_{a'}(t)=m\}} \\
&\leq \sum_{m \geq 1} \sum_{n \geq 2} \mathbb{I}_{\{m \geq n-1, |\hat{\mu}_{a'}^m - \mu_{a'}| \geq \varepsilon\}} \sum_{t \geq 1} \mathbb{I}_{\{\tau_{a,n}=t+1\}} \\
&\leq \sum_{m \geq 1} \sum_{n \geq 2} \mathbb{I}_{\{m \geq n-1, |\hat{\mu}_{a'}^m - \mu_{a'}| \geq \varepsilon\}}
\end{aligned}$$

Taking the expectation, it comes

$$\begin{aligned}
& \mathbb{E}_\nu \left[ \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a, N_{a'}(t) \geq N_a(t), |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon\}} \right] \\
&\leq \sum_{m \geq 1} \sum_{n \geq 2} \mathbb{I}_{\{m \geq n-1\}} \Pr_\nu (|\hat{\mu}_{a'}^m - \mu_{a'}| \geq \varepsilon) \\
&\leq \sum_{m \geq 1} \sum_{n \geq 2} \mathbb{I}_{\{m \geq n-1\}} \max \left( 2e^{-2m\varepsilon^2}, 2e^{-m\varepsilon^2/2} \right) \quad (\text{Hoeffding's inequality}) \\
&= \sum_{m \geq 1} \sum_{n \geq 2} \mathbb{I}_{\{m \geq n-1\}} 2e^{-m\varepsilon^2/2} \\
&= 2 \sum_{m \geq 1} m e^{-m\varepsilon^2/2} \\
&= \frac{2e^{-\varepsilon^2/2}}{(1 - e^{-\varepsilon^2/2})^2} = \frac{2e^{\varepsilon^2/2}}{(e^{\varepsilon^2/2} - 1)^2} \leq \frac{8e^{1/8}}{\varepsilon^4} \leq \frac{10}{\varepsilon^4}. \quad (0 < \varepsilon \leq 1/2)
\end{aligned}$$

■

**Lemma** Let us consider a set of Gaussian or Bernoulli distributions  $\nu \in \mathcal{D}$ . Let  $0 < \varepsilon \leq \min(1 - \mu^*, 1/2)$  and  $a \in \mathcal{A}$ . Let  $\lambda = \mu_a - \varepsilon$ . Then, independently of the considered strategy, we have

$$\mathbb{E}_\nu \left[ \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda\}} n \exp(n \text{KL}(\hat{\mu}_a^n | \lambda)) \right] \leq \frac{23 \log(1/\varepsilon)}{\varepsilon^6}.$$

We provide two proofs, one for Gaussian distributions and another for Bernoulli distributions, that can be read separately.

**Proof** [For Gaussian distributions] The proof is based on a Chernoff type inequality and a calculation by measurement change.

Since  $\nu_a \sim \mathcal{N}(\mu_a, 1)$  we have for all  $\varepsilon > 0$ ,

$$\forall n \geq 1, \quad \Pr_\nu(\hat{\mu}_a^n - \mu_a \leq -\varepsilon) \leq e^{-n\varepsilon^2/2}.$$

In addition,  $\forall \mu, \mu' \in \mathbb{R}$ ,  $\text{KL}(\mu|\mu') = \frac{(\mu - \mu')^2}{2}$ . Let  $n \geq 1$ . We have:

$$\begin{aligned}
& \mathbb{E}_\nu \left[ \mathbb{I}_{\{\hat{\mu}_a^n \leq \lambda\}} n e^{n\text{KL}(\hat{\mu}_a^n|\lambda)} \right] \\
&= \int_0^\infty \Pr_\nu \left( \mathbb{I}_{\{\hat{\mu}_a^n \leq \lambda\}} n e^{n\text{KL}(\hat{\mu}_a^n|\lambda)} > x \right) dx \\
&= \int_0^\infty \Pr_\nu \left( n e^{n\text{KL}(\hat{\mu}_a^n|\lambda)} > x, \hat{\mu}_a^n \leq \lambda \right) dx \\
&= \int_{-\infty}^\infty n^2 e^{nu} \Pr_\nu \left( \text{KL}(\hat{\mu}_a^n|\lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du \quad (\text{variable change } x = n e^{nu}, dx = n^2 e^{nu} du) \\
&= \int_{-\infty}^0 n^2 e^{nu} \Pr_\nu \left( \text{KL}(\hat{\mu}_a^n|\lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du + \int_0^\infty n^2 e^{nu} \Pr_\nu \left( \text{KL}(\hat{\mu}_a^n|\lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du \\
&= n^2 \Pr_\nu \left( \hat{\mu}_a^n - \mu_a \leq -\varepsilon \right) \int_{-\infty}^0 e^{nu} du + \int_0^\infty n^2 e^{nu} \Pr_\nu \left( \hat{\mu}_a^n - \mu_a \leq -\varepsilon - \sqrt{2u} \right) du \\
&\leq n^2 e^{-n\varepsilon^2/2} \frac{1}{n} + \int_0^\infty n^2 e^{nu} e^{-n(\varepsilon + \sqrt{2u})^2/2} du \\
&= n e^{-n\varepsilon^2/2} + n^2 e^{-n\varepsilon^2/2} \int_0^\infty e^{-n\varepsilon\sqrt{2u}} du \\
&= n e^{-n\varepsilon^2/2} + n^2 e^{-n\varepsilon^2/2} \int_0^\infty y e^{-n\varepsilon y} dy \quad (\text{variable change } u = \frac{y^2}{2}, du = y dy) \\
&= n e^{-n\varepsilon^2/2} + n^2 e^{-n\varepsilon^2/2} \frac{1}{(n\varepsilon)^2} \\
&= n e^{-n\varepsilon^2/2} + \frac{1}{\varepsilon^2} e^{-n\varepsilon^2/2}
\end{aligned}$$

To ends the proof, we use the following equalities for  $r > 0$

$$\begin{aligned}
\sum_{n \geq 1} e^{-nr} &= \frac{e^{-r}}{1 - e^{-r}} \\
\sum_{n \geq 1} n e^{-nr} &= \frac{e^{-r}}{(1 - e^{-r})^2}
\end{aligned}$$

and obtain

$$\mathbb{E}_\nu \left[ \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda\}} n \exp(n\text{KL}(\hat{\mu}_a^n|\lambda)) \right] \leq \frac{1}{\varepsilon^2} \frac{e^{-\varepsilon^2/2}}{1 - e^{-\varepsilon^2/2}} + \frac{e^{-\varepsilon^2/2}}{(1 - e^{-\varepsilon^2/2})^2} \leq \frac{10}{\varepsilon^4}.$$

■

**Proof** [For Bernoulli distributions] The proof is based on a Chernoff type inequality and a calculation by measurement change.

Since the support of  $\nu_a$  is included in  $[0, 1]$  we have by Chernoff's and Pinsker's inequalities

$$\forall 0 \leq v \leq \mu_a, \forall n \geq 1, \Pr_\nu \left( \hat{\mu}_a^n \leq v \right) \leq e^{-n\text{kl}(v|\mu_a)} \leq e^{-2n(\mu_a - v)^2}.$$

Let  $n \geq 1$ . We have

$$\begin{aligned}
& \mathbb{E}_\nu \left[ \mathbb{I}_{\{\hat{\mu}_a^n \leq \lambda\}} n e^{n \text{kl}(\hat{\mu}_a^n | \lambda)} \right] \\
&= \int_0^\infty \Pr_\nu \left( \mathbb{I}_{\{\hat{\mu}_a^n \leq \lambda\}} n e^{n \text{kl}(\hat{\mu}_a^n | \lambda)} > x \right) dx \\
&= \int_0^\infty \Pr_\nu \left( n e^{n \text{kl}(\hat{\mu}_a^n | \lambda)} > x, \hat{\mu}_a^n \leq \lambda \right) dx \\
&= \int_{-\infty}^\infty n^2 e^{nu} \Pr_\nu \left( \text{kl}(\hat{\mu}_a^n | \lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du \quad (\text{variable change } x = n e^{nu}, dx = n^2 e^{nu} du) \\
&= \int_{-\infty}^0 n^2 e^{nu} \Pr_\nu \left( \text{kl}(\hat{\mu}_a^n | \lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du + \int_0^{\text{kl}(0|\lambda)} n^2 e^{nu} \Pr_\nu \left( \text{kl}(\hat{\mu}_a^n | \lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du.
\end{aligned}$$

On the one hand

$$\begin{aligned}
& \int_{-\infty}^0 n^2 e^{nu} \Pr_\nu \left( \text{kl}(\hat{\mu}_a^n | \lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du \\
&= n^2 \Pr_\nu \left( \hat{\mu}_a^n \leq \lambda \right) \int_{-\infty}^0 e^{nu} du \\
&\leq n^2 e^{-2n(\mu_a - \lambda)^2} \frac{1}{n} \\
&= n e^{-2n\varepsilon^2}.
\end{aligned}$$

On the other hand, using variable change  $u = \text{kl}(v|\lambda)$  and Lemma 18, it comes

$$\begin{aligned}
& \int_0^{\text{kl}(0|\lambda)} n^2 e^{nu} \Pr_\nu \left( \text{kl}(\hat{\mu}_a^n | \lambda) > u, \hat{\mu}_a^n \leq \lambda \right) du \\
&= \int_0^\lambda n^2 e^{n \text{kl}(v|\lambda)} \Pr_\nu \left( \text{kl}(\hat{\mu}_a^n | \lambda) > \text{kl}(v|\lambda), \hat{\mu}_a^n \leq \lambda \right) \left[ -\frac{\partial \text{kl}}{\partial p}(v|\lambda) \right] dv \\
&= \int_0^\lambda n^2 e^{n \text{kl}(v|\lambda)} \Pr_\nu \left( \hat{\mu}_a^n < v \right) \left[ -\frac{\partial \text{kl}}{\partial p}(v|\lambda) \right] dv \\
&\leq \int_0^\lambda n^2 e^{n \text{kl}(v|\lambda)} e^{-n \text{kl}(v|\mu_a)} \left[ -\frac{\partial \text{kl}}{\partial p}(v|\lambda) \right] dv \\
&= \int_0^\lambda n^2 e^{-n(\text{kl}(v|\mu_a) - \text{kl}(v|\lambda))} \left[ -\frac{\partial \text{kl}}{\partial p}(v|\lambda) \right] dv \\
&\leq \int_0^\lambda n^2 e^{-n \frac{(\mu_a - \lambda)^2}{2}} \left[ -\frac{\partial \text{kl}}{\partial p}(v|\lambda) \right] dv \quad (\text{Lemma 18}) \\
&= \text{kl}(0|\lambda) n^2 e^{-n\varepsilon^2/2} \\
&\leq (-\log(1 - \mu^*)) n^2 e^{-n\varepsilon^2/2} \\
&\leq \log(1/\varepsilon) n^2 e^{-n\varepsilon^2/2},
\end{aligned}$$

where  $\frac{\partial \text{kl}}{\partial p}$  corresponds to the derivative of the  $\text{kl}(\cdot|\cdot)$  according to the first variable. Thus we have

$$\mathbb{E}_\nu \left[ \mathbb{I}_{\{\hat{\mu}_a^n \leq \lambda\}} n e^{n \text{kl}(\hat{\mu}_a^n | \lambda)} \right] \leq n e^{-2n\varepsilon^2} + \log(1/\varepsilon) n^2 e^{-n\varepsilon^2/2}.$$

To ends the proof, we use the following equalities for  $r > 0$

$$\begin{aligned}\sum_{n \geq 1} e^{-nr} &= \frac{e^{-r}}{1 - e^{-r}} = \frac{1}{e^r - 1} \\ \sum_{n \geq 1} n e^{-nr} &= \frac{e^r}{(e^r - 1)^2} = \frac{1}{e^r - 1} + \frac{1}{(e^r - 1)^2} \\ \sum_{n \geq 1} n^2 e^{-nr} &= \frac{e^r}{(e^r - 1)^2} + \frac{2e^r}{(e^r - 1)^3} = \frac{e^{2r} + e^r}{(e^r - 1)^3}\end{aligned}$$

and obtain

$$\mathbb{E}_\nu \left[ \sum_{n \geq 1} \mathbb{I}_{\{\hat{\mu}_a^n < \lambda\}} n \exp(n \text{kl}(\hat{\mu}_a^n | \lambda)) \right] \leq \frac{e^{2\varepsilon^2}}{(e^{2\varepsilon^2} - 1)^2} + \frac{\log(1/\varepsilon)e^{\varepsilon^2} + e^{\varepsilon^2/2}}{(e^{\varepsilon^2/2} - 1)^3} \leq \frac{23 \log(1/\varepsilon)}{\varepsilon^6}.$$

■

**Lemma 18** For all  $0 \leq v \leq \lambda < \mu < 1$  we have

$$\text{kl}(v|\mu) - \text{kl}(v|\lambda) \geq \frac{(\mu - \lambda)^2}{2}.$$

**Proof** Using monotony of the  $\text{kl}(\cdot|\cdot)$  we get

$$\text{kl}(v|\mu) - \text{kl}(v|\lambda) \geq \text{kl}(v|\mu) - \text{kl}\left(v \left| \frac{\mu + \lambda}{2} \right.\right).$$

Using convexity of the  $\text{kl}(\cdot|\cdot)$  we get

$$\frac{\text{kl}(v|\mu) - \text{kl}\left(v \left| \frac{\mu + \lambda}{2} \right.\right)}{\mu - \frac{\mu + \lambda}{2}} \geq \frac{\partial \text{kl}}{\partial q}\left(v \left| \frac{\mu + \lambda}{2} \right.\right) \geq \frac{\partial \text{kl}}{\partial q}\left(\lambda \left| \frac{\mu + \lambda}{2} \right.\right),$$

where  $\frac{\partial \text{kl}}{\partial q}$  corresponds to the derivative of the  $\text{kl}(\cdot|\cdot)$  according to the second variable. From Lemma B.4 in [Combes and Proutiere \(2014\)](#) we have

$$\left(\frac{\mu + \lambda}{2} - \lambda\right) \frac{\partial \text{kl}}{\partial q}\left(\lambda \left| \frac{\mu + \lambda}{2} \right.\right) \geq \text{kl}\left(\lambda \left| \frac{\mu + \lambda}{2} \right.\right).$$

Then Pinsker's inequality implies

$$\left(\frac{\mu + \lambda}{2} - \lambda\right) \frac{\partial \text{kl}}{\partial q}\left(\lambda \left| \frac{\mu + \lambda}{2} \right.\right) \geq 2 \left(\frac{\mu + \lambda}{2} - \lambda\right)^2 = \frac{(\mu - \lambda)^2}{2},$$

which ends the proof. ■

## Appendix C. KLUCB-UB finite time analysis

KLUCB-UB strategy implies similar lower bounds and empirical upper bounds on the numbers of pulls as IMED-UB strategy. An additional random process  $(\gamma_t)_{t \geq 1} \in \{0, 1\}$  appears in the empirical lower bounds induced by KLUCB-UB strategy (Lemma 19). When  $\gamma_t = 1$ , the empirical bounds are the same of the ones induced by IMED-UB strategy. And we show that the process  $(\gamma_t)_{t \geq 1}$  reaches zero only a finite number of times for which a use of the empirical bounds is needed. Then, similar reasoning as the one developed in Section 4 can be re-used and gives similar finite time analysis.

### C.1 Notations

Please, refer to Section 4.1.

### C.2 Strategy-based empirical bounds

In this subsection, we provide empirical bounds very similar to the ones induced by IMED-UB strategy. We first establish preliminary results on the indexes.

It is noticeable that for all time step  $t \geq 1$ ,

$$\forall \hat{a}^* \in \hat{\mathcal{A}}^*(t), \quad U_{\hat{a}^*}(t) = \hat{\mu}_{\hat{a}^*}(t) = \hat{\mu}^*(t). \quad (3)$$

In addition for  $a \notin \hat{\mathcal{A}}^*(t)$ ,

$$\text{if } N_a(t) \geq N_{\hat{a}_t^*}(t), \quad U_a(t) = \hat{\mu}_a(t) \text{ and } U_a(t) < U_{\hat{a}_t^*}(t), \quad (4)$$

$$\text{if } N_a(t) < N_{\hat{a}_t^*}(t), \quad N_a(t) \text{KL}(\hat{\mu}_a(t)|U_a(t)) + \log(N_a(t)) = \log(N_{\hat{a}_t^*}(t)). \quad (5)$$

In particular we have

$$N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t) \quad (6)$$

$$\text{and } N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) + \log(N_{a_{t+1}}(t)) = \log(N_{\hat{a}_t^*}(t)). \quad (7)$$

**Lemma 19 (Empirical lower bounds)** *Under KLUCB-UB, at each step time  $t \geq 1$ ,*

1.  $\forall a \in \mathcal{V}_{\hat{a}_t^*}, \quad \gamma_t \log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t)) + \log(N_a(t))$
2.  $N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t),$

where  $\gamma_t = \mathbb{I}_{\{a_{t+1} \in \hat{\mathcal{A}}^*(t)\}} + \mathbb{I}_{\{a_{t+1} \notin \hat{\mathcal{A}}^*(t) \text{ and } \log(N_{a_{t+1}}(t)) \leq N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t))\}} \in \{0, 1\}$ .

**Proof** Let  $t \geq 1$ . We have already seen that  $N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t)$  in (6). This corresponds to point 2. In the following, we prove point 1.

For  $a \in \mathcal{V}_{\hat{a}_t^*}$  such that  $a = a_{t+1}$  or  $N_a(t) \geq N_{a_{t+1}}(t)$ , point 1. is naturally satisfied.

Let  $a \in \mathcal{V}_{\hat{a}_t^*}$  such that  $a \neq a_{t+1}$  and  $N_a(t) < N_{a_{t+1}}(t)$ .

**Case 1** :  $a_{t+1} = \hat{a}_t^*$

Then  $N_a(t) < N_{\hat{a}_t^*}(t) \leq N_{a_{t+1}}(t)$  and from equations (5) and (6) it comes:

$$\log(N_{a_{t+1}}(t)) = \log(N_{\hat{a}_t^*}(t)) \quad N_a(t) \text{KL}(\hat{\mu}_a(t)|U_a(t)) + \log(N_a(t)) = \log(N_{a_{t+1}}(t)).$$

According to the followed strategy and equation 3

$$\hat{\mu}_a(t) \leq U_a(t) \leq U_{a_{t+1}}(t) \text{ and } U_{a_{t+1}}(t) = \hat{\mu}^*(t).$$

Since  $a_{t+1} = \hat{a}_t^*$ , this implies

$$\hat{\mu}_a(t) \leq U_a(t) \leq \hat{\mu}^*(t).$$

Then the monotony of the KL implies

$$\text{KL}(\hat{\mu}_a(t)|U_a(t)) \leq \text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t))$$

and

$$\begin{aligned} \log(N_{a_{t+1}}(t)) &= N_a(t)\text{KL}(\hat{\mu}_a(t)|U_a(t)) + \log(N_a(t)) \\ &\leq N_a(t)\text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t)) + \log(N_a(t)). \end{aligned}$$

**Case 2** :  $a_{t+1} \neq \hat{a}^*(t)$

From equations (5) and (7), we get

$$N_{a_{t+1}}(t)\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) + \log(N_{a_{t+1}}(t)) = N_a(t)\text{KL}(\hat{\mu}_a(t)|U_a(t)) + \log(N_a(t)). \quad (8)$$

Since  $N_a(t) < N_{a_{t+1}}(t)$ , this implies

$$\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) < \text{KL}(\hat{\mu}_a(t)|U_a(t)).$$

According to the followed strategy we have  $\hat{\mu}_a(t) \leq U_a(t) \leq U_{a_{t+1}}(t)$ . Then, the monotony of the KL implies

$$\text{KL}(\hat{\mu}_a(t)|U_a(t)) \leq \text{KL}(\hat{\mu}_a(t)|U_{a_{t+1}}(t)).$$

Thus, we have:

$$\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) < \text{KL}(\hat{\mu}_a(t)|U_{a_{t+1}}(t)) \text{ and } \hat{\mu}_a(t) \leq U_{a_{t+1}}(t).$$

Since  $\hat{\mu}_{a_{t+1}}(t) \leq U_{a_{t+1}}(t)$ , the monotony of the KL implies

$$\hat{\mu}_a(t) < \hat{\mu}_{a_{t+1}}(t).$$

Then from equation (8) we deduce

$$N_{a_{t+1}}(t)\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) \leq N_a(t)\text{KL}(\hat{\mu}_a(t)|U_{a_{t+1}}(t)) + \log(N_a(t))$$

and

$$N_{a_{t+1}}(t) \leq \frac{\text{KL}(\hat{\mu}_a(t)|U_{a_{t+1}}(t))}{\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t))} N_a(t) + \frac{\log(N_a(t))}{\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t))}.$$

Similarly, since  $\hat{\mu}_{a_{t+1}}(t) \leq \hat{\mu}^*(t) = U_{\hat{a}_t^*}(t) \leq U_{a_{t+1}}(t)$ , the monotony of the KL implies

$$\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) \geq \text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t)).$$

This implies

$$N_{a_{t+1}}(t) \leq \frac{\text{KL}(\hat{\mu}_a(t)|U_{a_{t+1}}(t))}{\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t))} N_a(t) + \frac{\log(N_a(t))}{\text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t))}.$$

Since  $\hat{\mu}_a(t) \leq \hat{\mu}_{a_{t+1}}(t) \leq \hat{\mu}^*(t) \leq U_{a_{t+1}}(t)$ , we have from Lemma 21:

$$\frac{\text{KL}(\hat{\mu}_a(t)|U_{a_{t+1}}(t))}{\text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t))} \leq \frac{\text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t))}{\text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t))}.$$

This implies

$$N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t)) + \log(N_a(t)).$$

■

**Lemma 20 (Empirical upper bounds)** *Under KLUCB-UB at each step time  $t \geq 1$ ,*

$$N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t)) \leq \log(t).$$

**Proof** From equation (7) we deduce

$$\begin{aligned} N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) &\leq N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) + \log(N_{a_{t+1}}(t)) \\ &= \log(N_{\hat{a}_t^*}(t)) \\ &\leq \log(t). \end{aligned}$$

Furthermore, according to the followed strategy, we have

$$\hat{\mu}_{a_{t+1}}(t) \leq \hat{\mu}^*(t) = U_{\hat{a}_t^*}(t) \leq U_{a_{t+1}}(t).$$

Then the monotony of the KL implies

$$\text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t)) \leq \text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t))$$

and

$$N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|\hat{\mu}^*(t)) \leq N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t)|U_{a_{t+1}}(t)) \leq \log(N_{\hat{a}_t^*}(t)) \leq \log(t).$$

■

**Lemma 21** *Let  $0 \leq \mu \leq \mu' \leq \mu'' \leq 1$ . We have:*

$$\forall u \in [\mu'', 1], \quad \frac{\text{KL}(\mu|u)}{\text{KL}(\mu'|u)} \leq \frac{\text{KL}(\mu|\mu'')}{\text{KL}(\mu'|\mu'')}.$$

We prove this lemma only for Bernoulli distributions when  $\text{KL}(\cdot|\cdot) = \text{kl}(\cdot|\cdot)$ . The proof is simpler for Gaussian distributions.

**Proof** [For Bernoulli distributions] We denote by  $\frac{\partial \text{kl}}{\partial p}(\cdot|\cdot)$  and  $\frac{\partial \text{kl}}{\partial q}(\cdot|\cdot)$  the derivatives of  $\text{kl}(\cdot|\cdot)$  respectively according to the first and second variables. Let us consider  $0 \leq \mu \leq \mu' \leq \mu'' < 1$  and  $f: u \in (\mu'', 1) \mapsto \text{kl}(\mu'|\mu'') \text{kl}(\mu|u) - \text{kl}(\mu|\mu'') \text{kl}(\mu'|u)$ .  $f$  is a C-1 function and for  $u \in (\mu'', 1)$ ,

$$f'(u) = \text{kl}(\mu'|\mu'') \frac{\partial \text{kl}}{\partial q}(\mu|u) - \text{kl}(\mu|\mu'') \frac{\partial \text{kl}}{\partial q}(\mu'|u) = \frac{\text{kl}(\mu'|\mu'') (u - \mu) - \text{kl}(\mu|\mu'') (u - \mu')}{u(1 - u)}.$$

Let us introduce  $g: u \in (\mu'', 1) \mapsto \text{kl}(\mu'|\mu'') (u - \mu) - \text{kl}(\mu|\mu'') (u - \mu')$ .  $g$  is a C-1 function and for  $u \in (\mu'', 1)$ ,

$$g'(u) = \text{kl}(\mu'|\mu'') - \text{kl}(\mu|\mu'').$$

Since  $\mu \leq \mu' \leq \mu''$ , the monotony of the kl implies  $g'(u) \leq 0$ . Then  $g$  is a non-increasing function. In addition

$$g(\mu'') = \text{kl}(\mu'|\mu'') (\mu'' - \mu) - \text{kl}(\mu|\mu'') (\mu'' - \mu') = (\mu'' - \mu) \times \left( \text{kl}(\mu'|\mu'') - \frac{\text{kl}(\mu|\mu'')}{(\mu'' - \mu)} (\mu'' - \mu') \right).$$

Lastly, let us consider  $h: p \in [0, \mu''] \mapsto \frac{\text{kl}(p|\mu'')}{\mu'' - p}$ .  $h$  is a C-1 function or  $p \in [0, \mu'']$ ,

$$h'(p) = \frac{(\mu'' - p) \frac{\partial \text{kl}}{\partial p}(p|\mu'') + \text{kl}(p|\mu'')}{(\mu'' - p)^2} = \frac{-\text{kl}(\mu''|p)}{(\mu'' - p)^2} \leq 0.$$

Then  $h$  is a non-increasing function. In particular, since  $\mu \leq \mu'$ ,

$$\frac{\text{kl}(\mu|\mu'')}{(\mu'' - \mu)} \geq \frac{\text{kl}(\mu'|\mu'')}{(\mu'' - \mu')}.$$

This implies  $g(\mu'') \leq 0$  and  $g \leq 0$ , since  $g$  is a non-increasing function. Then  $f' \leq 0$  and  $f$  is a non-increasing function. Since  $f(\mu'') = 0$ , this implies  $f \leq 0$ , which ends the proof.  $\blacksquare$

### C.3 Reliable current best arm and means

As in IMED-UB analysis, we consider the subset  $\mathcal{T}_\varepsilon$  of times where everything is well behaved, that is: the current best arm corresponds to the true one and the empirical means of the best arm and the current chosen arm are  $\varepsilon$ -accurate for  $0 < \varepsilon < \varepsilon_\nu$ , i.e.

$$\mathcal{T}_\varepsilon := \left\{ t \geq 1 : \widehat{\mathcal{A}}^*(t) = \{a^*\} \text{ and } \forall a \in \{a^*, a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \right\}.$$

We will show that its complementary set is finite on average. In order to prove this we decompose the set  $\mathcal{T}_\varepsilon$  in the following way. Let  $\mathcal{E}_\varepsilon$  be the set of times where the means are well estimated,

$$\mathcal{E}_\varepsilon := \left\{ t \geq 1 : \forall a \in \widehat{\mathcal{A}}^*(t) \cup \{a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \right\},$$

and  $\Lambda_\varepsilon$  the set of times where an arm that is not the current optimal neither pulled is underestimated

$$\Lambda_\varepsilon := \left\{ t \geq 1 : \exists a \in \mathcal{V}_{\widehat{a}_t^*} \setminus \{a_{t+1}, \widehat{a}_t^*\} \text{ s.t. } \left\{ \begin{array}{l} \hat{\mu}_a(t) < \mu_a - \varepsilon \\ \log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon) + \log((N_a(t))) \end{array} \right\} \right\}.$$

Then, in the same way as for Lemma 12, we can prove the following inclusion.



**Lemma 22 (Relations between the subsets of times)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathcal{T}_\varepsilon^c \cap \{t \geq 1 : \gamma_t = 1\} \cap \mathcal{E}_\varepsilon \subset \Lambda_\varepsilon. \quad (9)$$

We can now resort to classical concentration arguments in order to control the size of these sets, which yields the following upper bounds.

**Lemma 23 (Bounded subsets of times)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] \leq \frac{10(d+1)|\mathcal{A}|}{\varepsilon^4} \quad \mathbb{E}_\nu[|\Lambda_\varepsilon|] \leq 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6} \quad \mathbb{E}_\nu[|\{t \geq 1 : \gamma_t = 0\} \cap \mathcal{E}_\varepsilon|] \leq \frac{|\mathcal{A}|}{\varepsilon^4},$$

where  $d$  is the maximum degree of nodes in  $G$ .

**Proof** Refer to Lemma 13 to prove  $\mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] \leq \frac{10(d+1)|\mathcal{A}|}{\varepsilon^4}$ ,  $\mathbb{E}_\nu[|\Lambda_\varepsilon|] \leq 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6}$ . It is exactly the same proof.

Let  $t \in \{t \geq 1 : \gamma_t = 0\} \cap \mathcal{E}_\varepsilon$ . Then  $\gamma_t = 0$ . This implies

$$a_{t+1} \neq \hat{a}_t^* \text{ and } \log(N_{a_{t+1}}(t)) > N_{a_{t+1}}(t) \text{KL}(\hat{\mu}_{a_{t+1}}(t) | \hat{\mu}^*(t)).$$

Since  $t \in \mathcal{E}_\varepsilon$  and  $\varepsilon < \varepsilon_\nu$ , we have

$$|\hat{\mu}_{\hat{a}_t^*}(t) - \mu_{\hat{a}_t^*}| < \varepsilon \text{ and } \hat{\mu}_a(t) \leq \hat{\mu}^*(t) = \hat{\mu}_{\hat{a}_t^*}(t) < \mu_{\hat{a}_t^*} + \varepsilon < \mu_a - \varepsilon.$$

This implies by Pinsker's inequality

$$\text{KL}(\hat{\mu}_{a_{t+1}}(t) | \hat{\mu}^*(t)) \geq \min\left(2(\hat{\mu}^*(t) - \hat{\mu}_{a_{t+1}}(t))^2, \frac{(\hat{\mu}^*(t) - \hat{\mu}_{a_{t+1}}(t))^2}{2}\right) > (2\varepsilon)^2/2 = 4\varepsilon^2.$$

In addition, for all  $N \geq 1$ ,  $\log(N) \leq 2\sqrt{N}$ . Thus we have

$$2\sqrt{N_{a_{t+1}}(t)} > 4N_{a_{t+1}}(t)\varepsilon^2 \quad \text{i.e.} \quad N_{a_{t+1}}(t) < \frac{1}{4\varepsilon^4} < \frac{1}{\varepsilon^4}.$$

This implies

$$\begin{aligned} |\{t \geq 1 : \gamma_t = 0\} \cap \mathcal{E}_\varepsilon| &\leq \sum_{t \geq 1} \mathbb{I}_{\{N_{a_{t+1}}(t) < 1/\varepsilon^4\}} \\ &= \sum_{a \in \mathcal{A}} \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1} = a \text{ and } N_a(t) < 1/\varepsilon^4\}} \\ &\leq \sum_{a \in \mathcal{A}} \frac{1}{\varepsilon^4} = \frac{|\mathcal{A}|}{\varepsilon^4}. \end{aligned}$$

■

Thus combining them with (9) we obtain

$$\begin{aligned} \mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] &\leq \mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] + \mathbb{E}_\nu[|\Lambda_\varepsilon|] + \mathbb{E}_\nu[|\{t \geq 1 : \gamma_t = 0\} \cap \mathcal{E}_\varepsilon|] \\ &\leq \frac{10(d+1)|\mathcal{A}|}{\varepsilon^4} + 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6} + \frac{|\mathcal{A}|}{\varepsilon^4} \\ &\leq 34d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6}. \end{aligned}$$

Indeed, we have

$$\mathcal{T}_\varepsilon^c \subset (\mathcal{T}_\varepsilon^c \cap \{t \geq 1 : \gamma_t = 1\} \cap \mathcal{E}_\varepsilon) \cup (\{t \geq 1 : \gamma_t = 0\} \cap \mathcal{E}_\varepsilon) \cup \mathcal{E}_\varepsilon^c.$$

Hence, we just proved the following lemma.

**Lemma 24 (Reliable estimators)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] \leq 34d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6},$$

where  $d$  is the maximum degree of nodes in  $G$ .

#### C.4 Upper bounds on the numbers of pulls of sub-optimal arms

In this section, we now combine the different results of the previous sub-sections to prove Theorem 8.

**Proof** [Proof of Theorem 8.] Please refer to Section 4.4. It is exactly the same proof.  $\blacksquare$

### Appendix D. IMED-UB finite time analysis

In this section we assume that  $G$  is a tree.  $d$ -IMED-UB behaves as IMED-UB except during second order exploration phases. Thus,  $d$ -IMED-UB strategy implies the same lower bounds and empirical upper bounds on the numbers of pulls as IMED-UB strategy most of times. Then similar guaranties as those obtained under IMED-UB can be established based on the same reasoning for  $d$ -IMED-UB. These guaranties involve the numbers of pulls of arms in  $\mathcal{V}_{a^*}$  which are shown to be of order  $\mathcal{O}(\log(T))$ , and the assumption that  $G$  is a tree ensures the best arms  $(\underline{a}_t)_{t \geq 1}$  of the sub-trees  $(\hat{G}_{\underline{a}_t}(t))_{t \geq 1}$  belong to  $\mathcal{V}_{a^*}$  most of times. Then, since  $\mathcal{S}_t$  is built as a sub-tree of  $\hat{G}_{\underline{a}_t}(t)$  that contains  $\underline{a}_t$  for all time step  $t \geq 1$ , the IMED type strategy followed during the second order exploration phases implies that exploration outside  $\mathcal{V}_{a^*}$  is of order  $\mathcal{O}(\log(\mathcal{O}(\log(T)))) = \mathcal{O}(\log \log(T))$ .

#### D.1 Notations

Please, refer to Section 4.1.

#### D.2 Strategy-based empirical bounds

In this subsection, we provide empirical bounds very similar to the ones induced by IMED-UB strategy.

**Lemma 25 (Empirical lower bounds)** Under  $d$ -IMED-UB, at each step time  $t \geq 1$ ,

$$1. \quad \forall a \in \mathcal{V}_{\hat{a}_t^*}, \log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) + \log(N_a(t)) \quad \text{and} \quad N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t).$$

Furthermore, if  $a_{t+1} \notin \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}$ , we have

$$2. \quad \forall a \in \mathcal{S}_t, \log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}^+(\hat{\mu}_a(t) | \hat{\mu}_{\underline{a}_t}(t)) + \log(N_a(t)) \quad \text{and} \quad N_{a_{t+1}}(t) \leq N_{\underline{a}_t}(t) \leq N_{\hat{a}_t^*}(t).$$

**Proof**

**Case 1** :  $a_{t+1} \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}$ .

This means there is no second order exploration at time  $t$  and d-IMED-UB behaves as IMED-UB. Then point 1. is satisfied according to Lemma 10.

**Case 2** :  $a_{t+1} \notin \{\widehat{a}_t^*\} \cup \mathcal{V}_{\widehat{a}_t^*}$ .

This means  $\underline{a}_t \neq \widehat{a}_t^*$  and according to d-IMED-UB strategy

$$\forall a \in \mathcal{S}_t, \quad \log(N_{a_{t+1}}(t)) \leq I_{a_{t+1}}^{(\underline{a}_t)}(t) \leq I_a^{(\underline{a}_t)}(t) = N_a(t) \text{KL}^+(\hat{\mu}_a(t) | \hat{\mu}_{\underline{a}_t}(t)) + \log(N_a(t)) .$$

Since  $\underline{a}_t \in \mathcal{S}_t$  and  $I_{\underline{a}_t}^{(\underline{a}_t)}(t) = \log(N_{\underline{a}_t}(t))$ , by taking the  $\exp(\cdot)$  we get  $N_{a_{t+1}}(t) \leq N_{\underline{a}_t}(t)$  and prove point 2. . Furthermore, still according to d-IMED-UB strategy, we have

$$\forall a \in \{\widehat{a}_t^*\} \cup \mathcal{V}_{\widehat{a}_t^*}, \quad \log(N_{a_{t+1}}(t)) \leq \log(N_{\underline{a}_t}(t)) \leq I_{\underline{a}_t}(t) \leq I_a(t) = N_a(t) \text{KL}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) + \log(N_a(t)) .$$

Since  $I_{\widehat{a}_t^*}(t) = \log(N_{\widehat{a}_t^*}(t))$ , by taking the  $\exp(\cdot)$  we get in particular  $N_{a_{t+1}}(t) \leq N_{\widehat{a}_t^*}(t)$  and prove point 1. . ■

**Lemma 26 (Empirical upper bounds)** Under d-IMED-UB at each step time  $t \geq 1$ ,

$$1. \quad N_{\underline{a}_t}(t) \text{KL}(\hat{\mu}_{\underline{a}_t}(t) | \hat{\mu}^*(t)) \leq \log(t) .$$

Furthermore, if  $a_{t+1} \notin \{\widehat{a}_t^*\} \cup \mathcal{V}_{\widehat{a}_t^*}$ , we have

$$2. \quad N_{a_{t+1}}(t) \text{KL}^+(\hat{\mu}_{a_{t+1}}(t) | \hat{\mu}_{\underline{a}_t}(t)) \leq \log\left(\frac{\log(t)}{\text{KL}(\hat{\mu}_{\underline{a}_t}(t) | \hat{\mu}^*(t))}\right) .$$

**Proof** 1. According to the followed strategy, we have

$$I_{\underline{a}_t}(t) \leq I_{\widehat{a}_t^*}(t) .$$

It remains, to conclude, to note that

$$N_{\underline{a}_t}(t) \text{KL}(\hat{\mu}_{\underline{a}_t}(t) | \hat{\mu}^*(t)) \leq I_{\underline{a}_t}(t) ,$$

and

$$I_{\widehat{a}_t^*}(t) = \log(N_{\widehat{a}_t^*}(t)) \leq \log(t) .$$

2. We assume that  $a_{t+1} \notin \{\widehat{a}_t^*\} \cup \mathcal{V}_{\widehat{a}_t^*}$ . According to the followed strategy, we have

$$I_{a_{t+1}}^{(\underline{a}_t)}(t) \leq I_{\underline{a}_t}^{(\underline{a}_t)}(t) .$$

Furthermore, by definition of the second order IMED indexes we have

$$N_{a_{t+1}}(t) \text{KL}^+(\hat{\mu}_{a_{t+1}}(t) | \hat{\mu}_{\underline{a}_t}(t)) \leq I_{a_{t+1}}^{(\underline{a}_t)}(t) ,$$

and

$$I_{\underline{a}_t}^{(\underline{a}_t)}(t) = \log(N_{\underline{a}_t}(t)) .$$

We conclude the proof using point 1. we just proved. ■

### D.3 Reliable current best arm and means

As in IMED-UB analysis, we consider the subset  $\mathcal{T}_\varepsilon$  of times where everything is well behaved, that is: the current best arm corresponds to the true one and the empirical means of the best arm, the arm with minimal current index and the current chosen arm are  $\varepsilon$ -accurate for  $0 < \varepsilon < \varepsilon_\nu$ , i.e.

$$\mathcal{T}_\varepsilon := \left\{ t \geq 1 : \hat{\mathcal{A}}^\star(t) = \{a^\star\} \text{ and } \forall a \in \{a^\star, \underline{a}_t, a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \right\}.$$

We will show that its complementary set is finite on average. In order to prove this we decompose the set  $\mathcal{T}_\varepsilon$  in the following way. Let  $\mathcal{E}_\varepsilon$  be the set of times where the means are well estimated,

$$\mathcal{E}_\varepsilon := \left\{ t \geq 1 : \forall a \in \hat{\mathcal{A}}^\star(t) \cup \{\underline{a}_t, a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \right\},$$

and  $\Lambda_\varepsilon$  the set of times where an arm that is not the current optimal neither pulled is underestimated

$$\Lambda_\varepsilon := \left\{ t \geq 1 : \exists a \in \mathcal{V}_{\hat{a}_t^\star} \setminus \{a_{t+1}, \hat{a}_t^\star\} \text{ s.t. } \begin{cases} \hat{\mu}_a(t) < \mu_a - \varepsilon \\ \log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t) | \mu_a - \varepsilon) + \log(N_a(t)) \end{cases} \right\}.$$

Then we get the same relation between these sets as for IMED-UB strategy.

**Lemma 27 (Relations between the subsets of times)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathcal{T}_\varepsilon^c \setminus \mathcal{E}_\varepsilon^c \subset \Lambda_\varepsilon. \quad (10)$$

**Proof** The proof is exactly the same as for Lemma 12. ■

We can now resort to classical concentration arguments in order to control the size of these sets, which yields the following upper bounds.

**Lemma 28 (Bounded subsets of times)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] \leq \frac{10 |\mathcal{A}|^2}{\varepsilon^4} \quad \mathbb{E}_\nu[|\Lambda_\varepsilon|] \leq 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6},$$

where  $d$  is the maximum degree of nodes in  $G$ .

**Proof** Refer to Lemma 13 to prove  $\mathbb{E}_\nu[|\Lambda_\varepsilon|] \leq 23d^2 |\mathcal{A}| \frac{\log(1/\varepsilon)}{\varepsilon^6}$ . It is exactly the same proof.

Using Lemma 25 we have

$$\forall t \geq 1, \quad N_{a_{t+1}}(t) \leq N_{\underline{a}_t} \leq N_{\hat{a}_t^\star}(t).$$

Since  $\hat{a}_t^\star \in \operatorname{argmin}_{\hat{a}^\star \in \hat{\mathcal{A}}^\star(t)} N_{\hat{a}^\star}(t)$ , this implies

$$\hat{a}^\star \in \hat{\mathcal{A}}^\star(t)$$

$$\forall t \geq 1, \forall \hat{a}^\star \in \hat{\mathcal{A}}^\star(t), \quad N_{a_{t+1}}(t) \leq N_{\hat{a}_t^\star}(t) \leq N_{\hat{a}^\star}(t).$$

Then, based on the concentration inequalities from Lemma 16, we obtain

$$\begin{aligned} \mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] &\leq \sum_{a, a' \in \mathcal{A}} \mathbb{E}_\nu \left[ \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a, N_{a'}(t) \geq N_a(t), |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon\}} \right] \\ &\leq \sum_{a, a' \in \mathcal{A}} \frac{10}{\varepsilon^4} \\ &\leq \frac{10 |\mathcal{A}|^2}{\varepsilon^4}. \end{aligned}$$

■

Thus combining them with (10) we obtain

$$\begin{aligned} \mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] &\leq \mathbb{E}_\nu[|\mathcal{E}_\varepsilon^c|] + \mathbb{E}_\nu[|\Lambda_\varepsilon|] \\ &\leq \frac{10|\mathcal{A}|^2}{\varepsilon^4} + 23d^2|\mathcal{A}|\frac{\log(1/\varepsilon)}{\varepsilon^6} \\ &\leq 33d|\mathcal{A}|^2\frac{\log(1/\varepsilon)}{\varepsilon^6}. \end{aligned}$$

Indeed, we have

$$\mathcal{T}_\varepsilon^c \subset (\mathcal{T}_\varepsilon^c \setminus \mathcal{E}_\varepsilon^c) \cup \mathcal{E}_\varepsilon^c.$$

Hence, we just proved the following lemma.

**Lemma 29 (Reliable estimators)** For  $0 < \varepsilon < \varepsilon_\nu$ ,

$$\mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] \leq 33d|\mathcal{A}|^2\frac{\log(1/\varepsilon)}{\varepsilon^6},$$

where  $d$  is the maximum degree of nodes in  $G$ .

#### D.4 Upper bounds on the numbers of pulls of sub-optimal arms

In this section, we now combine the different results of the previous sub-sections to prove Theorem 8.

**Proof** [Proof of Theorem 8.] From Lemma 29, considering the following subset of times

$$\mathcal{T}_\varepsilon := \left\{ t \geq 1 : \begin{array}{l} \widehat{\mathcal{A}}^*(t) = \{a^*\} \\ \forall a \in \{a^*, \underline{a}_t, a_{t+1}\}, |\hat{\mu}_a(t) - \mu_a| < \varepsilon \end{array} \right\}.$$

we have

$$\mathbb{E}_\nu[|\mathcal{T}_\varepsilon^c|] \leq 33d|\mathcal{A}|^2\frac{\log(1/\varepsilon)}{\varepsilon^6},$$

where  $d$  is the maximum degree of nodes in  $G$ . Then, let us consider  $a \neq a^*$  and a time step  $t \in \mathcal{T}_\varepsilon$  such that  $a_{t+1} = a$ . Since  $t \in \mathcal{T}_\varepsilon$ , we have

$$\widehat{a}_t^* = a^* \text{ and } |\hat{\mu}_a(t) - \mu_a|, |\hat{\mu}_{\underline{a}_t}(t) - \mu_{\underline{a}_t}|, |\hat{\mu}_{a^*}(t) - \mu_{a^*}| < \varepsilon.$$

Then  $\underline{a}_t \neq a^*$  and, by construction of  $\alpha_\nu(\cdot)$  (see Section 4.1 Notations),

$$\text{KL}(\hat{\mu}_{\underline{a}_t}(t)|\hat{\mu}^*(t)) = \text{KL}(\hat{\mu}_{\underline{a}_t}(t)|\hat{\mu}_{a^*}(t)) \geq \frac{\text{KL}(\mu_{\underline{a}_t}|\mu_{a^*})}{1 + \alpha_\nu(\varepsilon)}.$$

**Case 1** :  $a_{t+1} \in \{\widehat{a}_t^*\} \cup \mathcal{V}_{\widehat{a}_t^*}$ , that is  $a = \underline{a}_t \in \mathcal{V}_{a^*}$

Then from Lemma 26 we get

$$N_a(t)\text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t)) \leq \log(t) \leq \log(T),$$

This implies

$$N_a(t) \leq \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a|\mu_{a^*})} \log(T).$$

**Case 2** :  $a_{t+1} \notin \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}$ , that is  $a \in G_{\underline{a}_t} \setminus \{a_t\}$  and  $\underline{a}_t \in \mathcal{V}_{a^*}$

Then from Lemma 26 we get

$$N_a(t) \text{KL}^+(\hat{\mu}_a(t) | \hat{\mu}_{\underline{a}_t}(t)) \leq \log \left( \frac{\log(t)}{\text{KL}(\hat{\mu}_{\underline{a}_t}(t) | \hat{\mu}^*(t))} \right) \leq \log \left( \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_{\underline{a}_t} | \mu_{a^*})} \log(T) \right) \leq \log \left( \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_{\underline{a}} | \mu_{a^*})} \log(T) \right).$$

Since  $G$  is a tree and  $a \in G_{\underline{a}_t} \setminus \{a_t\}$ , we have  $\mu_a < \mu_{\underline{a}_t}$ . Since  $\varepsilon < \varepsilon_\nu$ , we have  $\hat{\mu}_a(t) < \mu_a + \varepsilon < \mu_{\underline{a}_t} - \varepsilon < \hat{\mu}_{\underline{a}_t}(t)$  and  $\text{KL}^+(\hat{\mu}_a(t) | \hat{\mu}_{\underline{a}_t}(t)) = \text{KL}(\hat{\mu}_a(t) | \hat{\mu}_{\underline{a}_t}(t))$ . By construction of  $\alpha_\nu(\cdot)$ , it comes

$$\text{KL}(\hat{\mu}_a(t) | \hat{\mu}_{\underline{a}_t}(t)) = \text{KL}(\hat{\mu}_a(t) | \hat{\mu}_{\underline{a}_t}(t)) \geq \frac{\text{KL}(\mu_a | \mu_{\underline{a}_t})}{1 + \alpha_\nu(\varepsilon)} \geq \frac{1}{1 + \alpha_\nu(\varepsilon)} \min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_a | \mu_{\underline{a}}).$$

Then we have

$$N_a(t) \leq \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_a | \mu_{\underline{a}})} \log \left( \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_{\underline{a}} | \mu_{a^*})} \log(T) \right).$$

Thus, we have shown that for  $a \neq a^*$ , for all  $t \in \mathcal{T}_\varepsilon$  such that  $a_{t+1} = a$ ,

$$N_a(T) \leq \begin{cases} \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a | \mu_{a^*})} \log(T) & , \text{ if } a \in \mathcal{V}_{a^*} \\ \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_a | \mu_{\underline{a}})} \log \left( \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_{\underline{a}} | \mu_{a^*})} \log(T) \right) & , \text{ otherwise.} \end{cases}$$

This implies:

$$N_a(T) \leq \begin{cases} \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a | \mu_{a^*})} \log(T) + |\mathcal{T}_\varepsilon^c| + 1 & , \text{ if } a \in \mathcal{V}_{a^*} \\ \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_a | \mu_{\underline{a}})} \log \left( \frac{1 + \alpha_\nu(\varepsilon)}{\min_{\underline{a} \in \mathcal{V}_{a^*}} \text{KL}(\mu_{\underline{a}} | \mu_{a^*})} \log(T) \right) + |\mathcal{T}_\varepsilon^c| + 1 & , \text{ otherwise.} \end{cases}$$

Averaging these inequalities allows us to conclude. ■

## Appendix E. Details on numerical experiments

In this section we briefly describe how the subsets  $(\mathcal{S}_t)_{t \geq 1}$  used in d-IMED-UB are dynamically chosen for the experiments. We assume in this section that  $\mathcal{A} = \llbracket 1, A \rrbracket$  with  $A \geq 2$ .

Let us introduce the function  $d(\cdot)$  that extracts dichotomously from an interval  $\llbracket a, a' \rrbracket$  a subset of arms from their extreme values to its median.

---

**Algorithm 5** Dichotomous function  $d(\cdot)$ 

---

**input**  $\llbracket a, a' \rrbracket \subset \mathcal{A}$ , where  $a < a'$   
**if**  $a' - a < 4$  **then**  
    **return**  $\llbracket a, a' \rrbracket$   
**else**  
    **return**  $\{a, a'\} \cup d(\llbracket a + \lfloor (a' - a)/4 \rfloor, a' - \lfloor (a' - a)/4 \rfloor)$   
**end if**

---

Using function  $d(\cdot)$  we dynamically build a sequence of subsets  $(\tilde{\mathcal{S}}_t)_{t \geq 1}$  as described in Algorithm 6, where:

- $median(\cdot)$  returns the median of input subset,
  - $list(\cdot)$  creates the list (indexed from 1) of input elements,
  - $index(e, L)$  returns the index of element  $e$  in list  $L$ ,
  - $element(I, L)$  returns the elements of list  $L$  with indexes in  $I$ ,
  - $distance(a, \mathcal{S}) = \min_{a' \in \mathcal{S}} |a' - a|$ , for all arm  $a \in \mathcal{A}$  and all subset of arms  $\mathcal{S} \subset \mathcal{A}$ ,
  - $append(e, L)$  returns list  $L$  to which is added element  $e$ .
- 

**Algorithm 6** Dynamic sequence of subsets  $(\tilde{\mathcal{S}}_t)_{t \geq 1}$ 

---

$\tilde{\mathcal{S}}_1 \leftarrow d(\mathcal{A})$   
 $\tilde{a}_1^* \leftarrow median(\mathcal{S}_1)$   
 $List_{\tilde{\mathcal{S}}} \leftarrow list(\tilde{\mathcal{S}}_1)$   
 $List_{\tilde{a}^*} \leftarrow list(\tilde{a}_1^*)$   
**for**  $t = 2 \dots T$  **do**  
    **if**  $\hat{a}_t^* \in \tilde{\mathcal{S}}_{t-1}$  **then**  
         $\tilde{a}_t^* \leftarrow \hat{a}_t^*$   
        **if**  $\tilde{a}_t^* \in List_{\tilde{a}^*}$  **then**  
             $i \leftarrow index(\hat{a}_t^*, List_{\tilde{a}^*})$   
             $\tilde{\mathcal{S}}_t \leftarrow element(i, List_{\tilde{\mathcal{S}}})$   
  
             $List_{\tilde{\mathcal{S}}} \leftarrow element(\llbracket 1, i \rrbracket, List_{\tilde{\mathcal{S}}})$   
             $List_{\tilde{a}^*} \leftarrow element(\llbracket 1, i \rrbracket, List_{\tilde{a}^*})$   
        **else**  
             $\Delta \leftarrow distance(\tilde{a}_t^*, \tilde{\mathcal{S}}_{t-1} \setminus \{\tilde{a}_t^*\})$   
             $\tilde{\mathcal{S}}_t \leftarrow \tilde{\mathcal{S}}_{t-1} \cup d(\llbracket \tilde{a}_t^* - \Delta, \tilde{a}_t^* + \Delta \rrbracket)$   
  
             $List_{\tilde{\mathcal{S}}} \leftarrow append(\tilde{\mathcal{S}}_t, List_{\tilde{\mathcal{S}}})$   
             $List_{\tilde{a}^*} \leftarrow append(\tilde{a}_t^*, List_{\tilde{a}^*})$   
        **end if**  
    **end if**  
**end for**

---

Then we build the sequence of subsets  $(\mathcal{S}_t)_{t \geq 1}$  as follows:

$$\forall t \geq 1, \quad \mathcal{S}_t = \begin{cases} \{\underline{a}_t\} \cup \{a \in \tilde{\mathcal{S}}_t : a < \underline{a}_t\} & \text{if } \underline{a}_t < \hat{a}_t^*, \\ \{\underline{a}_t\} \cup \{a \in \tilde{\mathcal{S}}_t : a > \underline{a}_t\} & \text{if } \underline{a}_t > \hat{a}_t^*. \end{cases}$$

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- Rajeev Agrawal, Demosthenis Teneketzis, and Venkatachalam Anantharam. Asymptotically efficient adaptive allocation schemes for controlled iid processes: Finite parameter space. *IEEE Transactions on Automatic Control*, 34(3), 1989.
- Apostolos N. Burnetas and Michael N. Katehakis. Optimal adaptive policies for Markov decision processes. *Mathematics of Operations Research*, 22(1):222–255, 1997.
- Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, and Gilles Stoltz. Kullback–Leibler upper confidence bounds for optimal sequential allocation. *Annals of Statistics*, 41(3):1516–1541, 2013.
- Richard Combes and Alexandre Proutiere. Unimodal bandits: Regret lower bounds and optimal algorithms. In *International Conference on Machine Learning*, 2014.
- Richard Combes, Stefan Magureanu, and Alexandre Proutiere. Minimal exploration in structured stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 1763–1771, 2017.
- Audrey Durand, Odalric-Ambrym Maillard, and Joelle Pineau. Streaming kernel regression with provably adaptive mean, variance, and regularization. *arXiv preprint arXiv:1708.00768*, 2017.
- Todd L Graves and Tze Leung Lai. Asymptotically efficient adaptive choice of control laws in controlled markov chains. *SIAM journal on control and optimization*, 35(3):715–743, 1997.
- Junya Honda and Akimichi Takemura. An asymptotically optimal policy for finite support models in the multiarmed bandit problem. *Machine Learning*, 85(3):361–391, 2011.
- Junya Honda and Akimichi Takemura. Non-asymptotic analysis of a new bandit algorithm for semi-bounded rewards. *Machine Learning*, 16:3721–3756, 2015.
- Tze Leung Lai. Adaptive treatment allocation and the multi-armed bandit problem. *The Annals of Statistics*, pages 1091–1114, 1987.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- Tor Lattimore and Csaba Szepesvari. The end of optimism? an asymptotic analysis of finite-armed linear bandits. In *Artificial Intelligence and Statistics*, pages 728–737, 2017.
- Stefan Magureanu. *Efficient Online Learning under Bandit Feedback*. PhD thesis, KTH Royal Institute of Technology, 2018.
- Stefan Magureanu, Richard Combes, and Alexandre Proutiere. Lipschitz bandits: Regret lower bounds and optimal algorithms. *Machine Learning*, 35:1–25, 2014.



- O-A Maillard. Boundary crossing probabilities for general exponential families. *Mathematical Methods of Statistics*, 27(1):1–31, 2018.
- H. Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematics Society*, 58:527–535, 1952.
- Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: no regret and experimental design. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 1015–1022. Omnipress, 2010.
- William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- William R Thompson. On a criterion for the rejection of observations and the distribution of the ratio of deviation to sample standard deviation. *The Annals of Mathematical Statistics*, 6(4):214–219, 1935.
- Jia Yuan Yu and Shie Mannor. Unimodal bandits. In *ICML*, pages 41–48. Citeseer, 2011.