



**HAL**  
open science

## Structure adaptation in bandit theory

Hassan Saber

► **To cite this version:**

Hassan Saber. Structure adaptation in bandit theory. Artificial Intelligence [cs.AI]. Université de Lille, 2022. English. NNT : 2022ULILB049 . tel-04143097v2

**HAL Id: tel-04143097**

**<https://theses.hal.science/tel-04143097v2>**

Submitted on 27 Jun 2023

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

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

Université de Lille  
École Doctorale MADIS

# Thèse de Doctorat

Spécialité **Informatique**

présentée par

**Hassan SABER**

---

## Structure Adaptation in Bandit Theory

---

Adaptation à la Structure en Théorie des Bandits

sous la direction de **Odalric-Ambrym Maillard**

---

Soutenue publiquement le **19 décembre 2022** à **Villeneuve-d'Ascq**, devant le jury composé de

Gilles <b>Stoltz</b>	Directeur de recherche, CNRS	Rapporteur
Junya <b>Honda</b>	Professeur associé, Université de Kyoto	Rapporteur
Arnak <b>Dalalyan</b>	Professeur affilié, ENSAE/CREST	Examineur & Président
Alexandre <b>Proutière</b>	Professeur, KTH Royal Institute of Technology	Examineur
Claire <b>Vernade</b>	Chargée de recherche, DeepMind London	Examinatrice
Odalric-Ambrym <b>Maillard</b>	Chargé de recherche (HDR), Inria	Directeur de thèse
Richard <b>Combes</b>	Maître de conférence, CentraleSupélec	Invité

Centre de Recherche en Informatique, Signal et Automatique de Lille (CRIStAL),  
UMR 9189 Équipe Scool, 59650, Villeneuve d'Ascq, France





## Acknowledgements

I thank my advisor Odalric-Ambrym Maillard, my colleagues and managers at Inria for allowing and offering me the best conditions to pursue my PhD. I thank all those who encouraged me, helped me, motivated me, loved me. Those who knew how to correct me when it was necessary, those who knew how to be patient when it was necessary. I thank the staff of the university of Lille. I thank the organizers of conferences, workshops and summer schools in which I was allowed to participate and share my work. Finally, I thank the members of the jury for offering me the opportunity to defend my thesis and compete for the title of doctor.

## Remerciements

Je réserve cette section à mon encadrant et directeur de thèse pour des remerciements que je souhaite en français. Odalric est un homme avec beaucoup d'humanité et dont je pense avoir beaucoup appris. Bien que nos relations soient essentiellement académiques et professionnelles, il me laisse voir en sa personnalité une personne passionnée et talentueuse, un personnage robuste et raisonné. C'est pour moi un bel exemple d'accomplissement professionnel. Un accomplissement qui donne sens au métier de chercheur, un accomplissement en devenir. Pour tout cela, Odalric, je te remercie et te souhaite le meilleur des jours à venir.

# Abstract

Understanding the dynamics of complex systems, and how to optimally act in them impacts all aspects of human societies where a careful management of natural, energetic, human and computational resources is required. To overcome the limitations of human capabilities to process large amounts of data, researchers from the field of machine learning and mathematical statistics for sequential decision making pursue the long-term goal of developing an optimal and automatic method that can, from partial observations and sequential interactions with a complex system, learn an optimal behavior. While optimal control considers the dynamics of the system is assumed to be known, reinforcement learning is interested in the case when the dynamics is unknown and must be learned from observations only. A key difficulty to design a solution for these problems is that typically, when a decision is made, one only gets to see a noisy effect of that decision, and little about the effect of other alternatives. This gives rise to the study of the fundamental exploration- exploitation trade-off: Shall we follow a algorithm that has been used a lot in the past and has empirically proven good until now (exploitation), or shall we explore a less known but potentially promising algorithm (exploration)? Addressing this trade-off does not yield the same approach depending on the underlying structure of the dynamical system, where structure can be understood in several ways. Since not taking advantage of this (possibly hidden) structure is prone to obtaining loose algorithms, it is crucial to investigate how to build learning algorithms that can be adaptive to it. In this thesis, we want to better understand how the notion of structure modifies the learning guarantees and suggests novel improved algorithms in the context of bandits. We give special attention to the cases of unimodal, multimodal and graph structure. We introduce an algorithm for each of these structures, respectively  $\text{IMED-UB}$ ,  $\text{IMED-MB}$  and  $\text{IMED-GS}$ . These algorithms are extensions of the popular Indexed Minimum Empirical Divergence ( $\text{IMED}$ ) algorithm from [Honda and Takemura \(2015\)](#) to the considered structures. We provide a finite time analysis for each of them and prove their asymptotic optimality. In particular, we considered new exploring mechanisms (second order exploration allowed by  $\text{IMED}$  approach) and developed new tools (concentrations inequalities) we think that are of independent interest for the bandit community. Furthermore, these novel algorithms perform well in practice. This is confirmed by numerical illustrations on synthetic data.

# Résumé

Dans nos sociétés modernes, une gestion minutieuse des ressources naturelles, énergétiques, humaines et informatiques est nécessaire. Comprendre les dynamiques de systèmes complexes et gérer les interactions de manière optimale constituent un enjeu majeur. Pour surmonter la limite des capacités humaines à traiter de grandes quantités de données, les chercheurs en apprentissage automatique et statistiques mathématiques pour la prise de décision séquentielle ont pour objectif de développer une méthode automatique et optimale pouvant, à partir d'observations partielles et d'interactions séquentielles avec un système complexe, apprendre un comportement optimal. Alors que le contrôle optimal considère que la dynamique du système est supposée connue, l'apprentissage par renforcement s'intéresse au cas où la dynamique est inconnue et doit être apprise uniquement à partir d'observations. Une difficulté clé pour concevoir une solution à ces problèmes est que, lorsqu'une décision est prise, seulement un effet bruité de cette décision est observée. Et cela renseigne peu sur les décisions alternatives. Cela donne lieu à l'étude du compromis fondamental entre exploration et exploitation: doit-on suivre un algorithme déjà beaucoup utilisé par le passé et empiriquement viable (exploitation) ou doit-on explorer un algorithme moins connue mais potentiellement meilleur (exploration)? Aborder ce compromis donne lieu à diverses approches en fonction de la structure sous-jacente au système dynamique, la notion de structure pouvant être comprise de différentes façons. Puisque ne pas tirer avantage de la structure (éventuellement cachée) est susceptible d'occasionner des algorithmes largement sous-optimaux, il est crucial d'examiner la manière de concevoir des algorithmes d'apprentissage qui peuvent s'adapter à celle-ci. Dans cette thèse, nous souhaitons mieux comprendre comment la notion de structure modifie les garanties d'apprentissage et suggère de nouveaux algorithmes plus performants dans le contexte des bandits. Notre attention sera particulièrement portée sur les structures dites unimodal, multimodal et de graphe. Pour chacune de ces structures nous introduisons un algorithme, respectivement  $IMED-UB$ ,  $IMED-MB$  et  $IMED-GS$ , dont nous prouvons l'optimalité asymptotique et dont nous étudions les performances à horizon fini. Ces algorithmes vise à adapter l'algorithme  $IMED$  introduit par [Honda and Takemura \(2015\)](#) pour le cas des bandits non-structurés au cas des bandits unimodaux, multimodaux et structure de graphe. En outre, l'étude de ces algorithmes, nous a amené à développer des outils nouveaux (inégalités de concentration) et considérer de nouveaux mécanismes d'exploration (exploration de second ordre permise par l'approche  $IMED$ ) dont nous pensons qu'ils bénéficieront à la communauté. Enfin, nous illustrons par des simulations l'efficacité pratique de ces nouveaux algorithmes.

# Contents

<b>1</b>	<b>Introduction</b>	<b>10</b>
1.1	The stochastic multi-armed bandit problem	10
1.2	Optimal algorithms for non-structured bandit problems	12
1.3	Structured configurations	13
1.4	Publications	16
1.5	Contributions	17
<b>2</b>	<b>Unimodal Bandits</b>	<b>20</b>
2.1	Introduction	20
2.2	Regret lower bound	21
2.3	Optimal algorithm for unimodal bandits	21
2.3.1	OSUB algorithm	21
2.3.2	IMED-UB algorithm	22
2.3.3	Asymptotic optimality of IMED-UB	22
2.3.4	Asymptotic optimality of IMED	23
2.4	IMED-UB finite time analysis	23
2.4.1	Notations	24
2.4.2	Algorithm-based empirical bounds	24
2.4.3	Non-reliable current best arm	25
2.4.4	Reliable current means and current best arm	26
2.4.5	Upper bounds on the numbers of pulls of sub-optimal arms	27
2.5	Numerical experiments	29
<b>3</b>	<b>Multimodal Bandits</b>	<b>30</b>
3.1	Introduction	30
3.2	Regret lower bound	31
3.3	Numerically efficient algorithm for multimodal bandits	33
3.3.1	Notations	33
3.3.2	IMED-MB algorithm.	34
3.4	Numerical experiments	36
<b>4</b>	<b>Graph-Structured Bandits</b>	<b>40</b>
4.1	Introduction	40
4.2	Regret lower bound	41
4.3	Optimal algorithm for graph-structured bandits	42
4.3.1	IMED-GS algorithm	43
4.3.2	Asymptotic optimality of IMED-GS	46
4.4	Numerical experiments	47

4.5	Finite time properties of IMED-GS algorithm . . . . .	49
4.5.1	Additional assumptions . . . . .	49
4.5.2	Algorithm-based empirical bounds . . . . .	50
4.5.3	Non-reliable current means . . . . .	53
4.5.4	Non-reliable current best arm and current informative sets of arms . . . . .	57
4.5.5	Reliable current means of current informative sets of Arms . . . . .	60
4.6	Upper bounds under IMED-GS algorithm . . . . .	62
4.6.1	Upper bounds on the optimal numbers of pulls . . . . .	62
4.6.2	Upper bounds on the numbers of pulls . . . . .	65
4.6.3	Almost-sure upper bound on the cumulative regret . . . . .	67
4.6.4	Upper bound on the regret . . . . .	70
<b>5</b>	<b>Concentration Inequalities for Structured Bandits</b>	<b>72</b>
<b>6</b>	<b>Routine Bandits</b>	<b>78</b>
6.1	The Routine Bandit Setting . . . . .	78
6.2	The KLUCB-RB Strategy . . . . .	79
6.3	Sketch of Proof . . . . .	82
6.4	Numerical Experiments . . . . .	84
6.4.1	More Arms than Bandits: A Beneficial Case . . . . .	84
6.4.2	Increasing the Number of Bandit Instances . . . . .	85
6.4.3	Critical Settings . . . . .	86
6.5	Conclusion . . . . .	87
<b>7</b>	<b>Conclusion</b>	<b>89</b>
7.1	Summary of the presented contributions . . . . .	89
7.2	Some personal satisfactions . . . . .	89
7.3	Future work . . . . .	90
7.3.1	Algorithms with computational efficiency . . . . .	90
7.3.2	IMED for Markov Decision Processes with known transition probabilities . . . . .	90
7.3.3	From routine bandit to non-stationary bandit problem . . . . .	90
<b>A</b>	<b>Graph-Structured Bandits: Complements</b>	<b>96</b>
A.1	Structures Unimodal, Lipschitz and Aggregate of Bandits . . . . .	96
A.2	Proof related to the regret lower bound . . . . .	97
A.2.1	Proof of Proposition 7 . . . . .	98
A.2.2	Proof of Proposition 3 . . . . .	102
A.3	Technical results . . . . .	103
A.4	Additional experiments . . . . .	105
A.4.1	PO-IMED-GS algorithm . . . . .	105
A.4.2	Regrets Averaged on Random Structured Configurations . . . . .	106
<b>B</b>	<b>Generic Tools</b>	<b>108</b>
B.1	Non-reliable current means . . . . .	108
B.2	Concentration of measure . . . . .	111



<b>C</b>	<b>Routine Bandits : Proof and Additional Experiments</b>	<b>112</b>
C.1	Proof of Proposition 5	112
C.2	Proof of Theorem 5	113
C.2.1	Proof of Lemma 18	114
C.2.2	Proof of Lemma 19	116
C.2.3	Proof of Lemma 20	119
C.2.4	Proof of Lemma 21	121
C.2.5	Proof of Proposition 6	122
C.2.6	Tools from Concentration of Measure	124
C.3	Additional Experiments: Ideal Cases for which Bandits are Close Enough on the Subset of Optimal Arms	124
C.3.1	A Single Instance	125
C.3.2	Similarity of Different Instances on the Optimal Subset $\mathcal{A}^*$	125
C.3.3	Complement of Sections 6.4.2 and 6.4.3	126
<b>D</b>	<b>Bandits with Groups of Similar Arms</b>	<b>128</b>
D.1	Introduction	128
D.2	A regret lower bound with combinatorial and non-combinatorial parts	129
D.3	Information Minimization for bandits with equivalence class	131
D.4	Regret analysis	132
D.5	Experiments	133
D.6	Conclusion	135

# List of Symbols

$\mathcal{A}$	Set of arms
$\mathcal{D}$	Set of one-dimensional exponential family distributions
$\nu$	Distribution in $\mathcal{D}$
$I$	Interval in $\mathbb{R}$
$\mu$	Means of distribution $\nu$ with values in $\mathbb{I}^A$
$\mu^*$	Maximum of means $\mu$
$\Delta_a$	Gap between maximal mean $\mu^*$ and mean $\mu_a$ of arm $a \in \mathcal{A}$
$\mathbb{P}_\nu(\cdot)$	Probability measure under distribution $\nu$
$\mathbb{E}_\nu[\cdot]$	Expectation under distribution $\nu$
$T$	Horizon
$N_a(t)$	Number of pulls of arm $a \in \mathcal{A}$ at time step $t \geq 1$
$R(\nu, T)$	Averaged regret
$\mathfrak{C}_{\mathcal{D}}(\mu)$	Constant solution of the constrained linear-optimization problem from the regret lower bound
$\text{KL}(\cdot \cdot)$	Kullback-Leibler divergence between one-dimensional exponential family distributions
$\text{kl}(\cdot \cdot)$	Kullback-Leibler divergence between Bernoulli distributions
$\hat{\mu}_a(t)$	Empirical mean of arm $a \in \mathcal{A}$ at time step $t \geq 1$
$\hat{\mu}^*(t)$	Maximum of empirical means at time step $t \geq 1$
$I_a(t)$	IMED index of arm $a \in \mathcal{A}$ at time step $t \geq 1$
$\mathcal{E}_{a,a'}^-(\varepsilon)$	Set of time steps where the current mean of arm $a$ $\varepsilon$ -deviates from below while $N_a(t) \geq N_{a'}(t)$
$\mathcal{K}_{a,a'}^-(\varepsilon)$	Set of time steps where couple of arms $(a, a')$ shows $\varepsilon^-$ -KL-log deviation
$\mathcal{D}_{\text{uni}}$	Set of distributions for the unimodal structure
$\mathcal{D}_{\text{M-modal}}$	Set of distributions for the multimodal structure
$\Theta$	Set of means in $(0, 1)^A$
$\mathcal{D}_\Theta$	Set of Bernoulli distributions for graph structure

# Chapter 1

## Introduction

### 1.1 The stochastic multi-armed bandit problem

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 instance is specified by a configuration, that is a set of unknown probability distributions,  $\nu = (\nu_a)_{a \in \mathcal{A}}$  with means  $(\mu_a(\nu))_{a \in \mathcal{A}}$ . When there is no possible confusion, the means are simply denoted  $(\mu_a)_{a \in \mathcal{A}}$ . 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 algorithm 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 an algorithm 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 algorithms for specific configuration sets, such as KLUCB algorithm (Lai, 1987; Cappé et al., 2013; Maillard, 2018) for exponential families, or alternatively DMED and IMED algorithms from Honda and Takemura (2011, 2015). We refer to Lattimore and Szepesvári (2020) for a recent survey. In this regard, it should be highlighted two another main approaches to solve optimally the stochastic bandit problem: Bayesian algorithm (Thompson, 1933) and algorithms based on re-sampling methods, such as RB-SDA introduced in a recent work by Baudry et al. (2020).

Both to propose an effective presentation of some relevant algorithms and start introducing the framework within which our work is included, we assume one-dimensional exponential family distributions and then explain the known lower bounds on the regret from Lai and Robbins (1985) under this assumption.

**Assumption 1** (One-dimensional exponential family distributions). *For all  $\nu \in \mathcal{D}$ ,  $\nu \subset \mathcal{P}_1 := \{p(\mu), \mu \in \mathbf{I}\}$ , where  $p(\mu)$  is a canonical exponential-family distribution probability with parameter  $\eta(\mu)$  and density  $f(\cdot, \mu)$  with respect to some positive measure  $\lambda$  on  $\mathbb{R}$  and mean  $\mu \in \mathbf{I} \subset \mathbb{R}$ .  $f(\cdot, \mu)$  has the following shape:*

$$f(\cdot, \mu) : \quad x \mapsto h(x) \exp(\eta(\mu) T(x) - A(\mu)),$$

where  $h \in \mathbb{R}_+^{\mathbb{R}}$ ,  $T \in \mathbb{R}^{\mathbb{R}}$  and  $A(\mu) = \log \int h(x) \exp(\eta(\mu) T(x)) \lambda(dx)$  are such that  $|A(\mu)| < \infty$ .

We consider algorithms that are *consistent* in order to obtain non trivial lower bound on the regret.

**Definition 1** (Consistent algorithm). *An algorithm is consistent on  $\mathcal{D}$  if for all configuration  $\nu \in \mathcal{D}$ , for all sub-optimal arm  $a \notin \mathcal{A}^* := \operatorname{argmax}_{a \in \mathcal{A}} \mu_a$ , for all  $\alpha > 0$ ,*

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

where  $N_a(t) = \sum_{s=1}^t \mathbb{I}_{\{a_s=a\}}$  is the number of pulls of arm  $a$  at time  $t \geq 0$ .

In particular, for  $\alpha = 1$ , the number of pulls of a sub-optimal arm is sub-linear under a consistent algorithm.

Thanks to the *tower rule*, we have

$$R(\nu, T) = T\mu^* - \mathbb{E}_\nu \left[ \sum_{t=1}^T X_t \right] = T\mu^* - \mathbb{E}_\nu \left[ \sum_{t=1}^T \mu_{a_t} \right].$$

We can then rewrite the regret as follows.

$$R(\nu, T) = \mathbb{E}_\nu \left[ \sum_{t=1}^T \mu^* - \mu_{a_t} \right] = \mathbb{E}_\nu \left[ \sum_{t=1}^T \sum_{a \in \mathcal{A}} \mathbb{I}_{\{a_t=a\}} (\mu^* - \mu_a) \right] = \sum_{a \in \mathcal{A}} \mathbb{E}_\nu [N_a(T)] \Delta_a. \quad (1.1)$$

When  $\mathcal{D} = \mathcal{P}_1^{\mathcal{A}}$  and under a consistent algorithm, from [Lai and Robbins \(1985\)](#) we have the following lower bounds on the numbers of pulls:

$$\forall a \notin \mathcal{A}^*, \quad \liminf_{T \rightarrow \infty} \frac{\mathbb{E}_\nu [N_a(T)] \operatorname{KL}(\mu_a | \mu^*)}{\log(T)} \geq 1. \quad (1.2)$$

This implies

$$\liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \geq \mathfrak{C}_1(\mu) := \sum_{a \in \mathcal{A}} \frac{\Delta_a}{\operatorname{KL}(\mu_a | \mu^*)}, \quad (1.3)$$

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

We refer to the proof of Theorem 1 from [Garivier et al. \(2016\)](#) for a proof of Equation (1.2) that provides lower bounds on the numbers of pulls when assuming non-structured bandits and consistent algorithms. Note that we adapt this same proof in Proposition 2 to obtain a lower bound on the regret when assuming multimodal structure (see Chapter 3) and consistent algorithms.

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

**Remark 2.** *We allow and reserve the notation  $\operatorname{kl}(\cdot | \cdot)$  for the Kullback-Leibler divergence between Bernoulli distributions.*

## 1.2 Optimal algorithms for non-structured bandit problems

We briefly present in this section some of the main (*asymptotically*) optimal algorithms (under Assumption 1) for the classical multi-armed bandit problem. Here, optimal means that  $\liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \leq \mathfrak{C}_1(\mu)$  under these algorithms.

**Kullback-Leibler Upper Confidence Bounds (KLUCB).** For an arm  $a \in \mathcal{A}$  and a time step  $t \geq 1$ , the upper confidence bound is defined as follows:

$$U_a(t) = \sup \left\{ \lambda \geq \hat{\mu}_a(t) : \text{KL}(\hat{\mu}_a(t) | \lambda) \leq \frac{f(t)}{N_a(t)} \right\}$$

where  $f(t) = \log(t) + 3 \log \log(t)$  if  $t \geq 3$ , 0 otherwise, with the convention  $0/0 = 0$ . This upper bound is motivated by the following concentration inequality:

$$\mathbb{P}_\nu(U_a(t) < \mu_a) \leq e \lceil f(t) \log(t) \rceil \exp(-f(t)).$$

We refer to Theorem 10 of [Garivier and Cappé \(2011\)](#) for a more general formulation and a proof of this concentration inequality. KLUCB algorithm then consists in pulling an arm with maximal upper bound at each time step,  $a_{t+1} = \operatorname{argmax}_{a \in \mathcal{A}} U_a(t)$ . KLUCB algorithm is summarized in Algorithm 1.

---

### Algorithm 1 KLUCB

---

**for**  $t = 0 \dots T - 1$  **do**

Pull  $a_{t+1} \in \operatorname{argmax}_{a \in \mathcal{A}} U_a(t)$  (chosen arbitrarily)

**end for**

---

**Indexed Minimum Empirical Divergence (IMED).** For an arm  $a \in \mathcal{A}$  and a time step  $t \geq 1$ , the IMED index is defined as follows:

$$I_a(t) = N_a(t) \text{KL}(\hat{\mu}_a(t), \hat{\mu}_*(t)) + \log(N_a(t)).$$

This quantity can be seen as a transportation cost for “moving” a sub-optimal arm to an optimal one, plus exploration terms (the logarithms of the numbers of pulls). When an optimal arm is considered, the transportation cost is null and it remains 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 pair  $a$  but also uses empirical means of current optimal arms. IMED is the algorithm that consists in pulling an arm with minimal index at each time step,  $a_{t+1} = \operatorname{argmin}_{a \in \mathcal{A}} I_a(t)$ . Finite time guarantees for IMED algorithm and its optimality are established in Corollary 2 from next Chapter 2. IMED algorithm is summarized in Algorithm 2.

---

### Algorithm 2 IMED

---

**for**  $t = 0 \dots T - 1$  **do**

Pull  $a_{t+1} \in \operatorname{argmin}_{a \in \mathcal{A}} I_a(t)$  (chosen arbitrarily)

**end for**

---

**Thomson Sampling (TS).** After pulling each arm once, we denote by  $\pi_a(t)$  the posterior distribution of arm  $a \in \mathcal{A}$  at time  $t \geq |\mathcal{A}|$  and sample  $\tilde{\mu}_a(t)$  from  $\pi_a(t)$ , what we note  $\tilde{\mu}_a(t) \sim \pi_a(t)$ . TS algorithm consists then in pulling the arm with the maximal sample,  $a_{t+1} = \operatorname{argmax}_{a \in \mathcal{A}} \tilde{\mu}_a(t)$ . Note the choice of the posteriors is crucial.

TS algorithm is proven to be asymptotically optimal for one-dimensional exponential family distributions in [Korda et al. \(2013\)](#) when the posteriors are obtained from Jeffreys priors. TS algorithm is summarized in [Algorithm 3](#).

---

**Algorithm 3** TS

---

Pull each arm once  
**for**  $t = |\mathcal{A}| \dots T - 1$  **do**  
  **for**  $a \in \mathcal{A}$  **do**  
     $\tilde{\mu}_a(t) \sim \pi_a(t)$   
  **end for**  
  Pull  $a_{t+1} \in \operatorname{argmax}_{a \in \mathcal{A}} \tilde{\mu}_a(t)$  (chosen arbitrarily)  
**end for**

---

**Best Empirical Sampled Average (BESA).** We restrain ourselves to the case of two arms when  $\mathcal{A} = \{1, 2\}$ . We refer to [Baransi et al. \(2014\)](#) and [Baudry et al. \(2020\)](#) for an optimal generalization beyond the case of two arms of the following approach. After pulling each arm once, we consider at each time  $t \geq 2$ , the sample size  $n(t) = \min \{N_1(t), N_2(t)\}$  and choose arbitrarily two samples  $\mathcal{S}_1(t)$  and  $\mathcal{S}_2(t)$  of rewards of same size  $n(t)$ , respectively from arm 1 and arm 2:  $\mathcal{S}_a(t) \subset \{X_s : a_s = a, 1 \leq s \leq t\}$ ,  $|\mathcal{S}_a(t)| = n(t)$ ,  $a \in \{1, 2\}$ . The empirical sub-sampled averages can be then computed as follows:

$$\hat{m}_a(t) = \sum_{X \in \mathcal{S}_a(t)} X/n(t).$$

BESA algorithm consists in pulling an arm with maximal empirical sample average,  $a_{t+1} = \operatorname{argmax}_{i \in \{1,2\}} \hat{m}_i(t)$ .

BESA algorithm is summarized in [Algorithm 4](#).

---

**Algorithm 4** BESA

---

Pull each arm once  
**for**  $a \in \{1, 2\}$  **do**  
  Choose an arbitrarily  $n(t)$ -sized sample  $\mathcal{S}_a(t)$  of rewards from arm  $a$   
  Compute the empirical sampled average  $\hat{m}_a(t)$   
**end for**  
  Pull  $a_{t+1} \in \operatorname{argmax}_{a \in \{1,2\}} \hat{m}_a(t)$  (chosen arbitrarily)

---

### 1.3 Structured configurations

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 algorithms that are provably optimal remains a challenge for many structures.

The study of specific *structured configuration sets*  $\mathcal{D}$  has received increasing attention over the last few years, motivated by the growing popularity of bandits in a number of industrial and societal application domains: The study of 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 \(2014a\)](#); [Trinh et al. \(2020\)](#); [Saber et al. \(2021a\)](#) providing an explicit lower bound and corresponding algorithms. Note that in [Saber et al. \(2021a\)](#), we adapt IMED algorithm to the Unimodal structure simply by narrowing on the current best arm and its neighbourhood for pulling an arm at a given time step. See also [Kunne et al. \(2020\)](#); [Gao et al. \(2019\)](#). Lipschitz bandits were studied in [Magureanu et al. \(2014\)](#); [Wang et al. \(2020\)](#); [Lu et al. \(2019\)](#) while combinatorial structures have been studied e.g. in [Kveton et al. \(2015\)](#); [Magureanu \(2018\)](#), and more recently [Cuvelier et al. \(2021b\)](#). The Linear bandit problem is also one typical illustration ([Abbasi-Yadkori et al. \(2011\)](#); [Srinivas et al. \(2010\)](#); [Durand et al. \(2017\)](#); [Kveton et al. \(2020\)](#)), see [Lattimore and Szepesvari \(2017\)](#) for a study of the lower bound (and [Degegne et al. \(2020a\)](#) for the related pure-exploration setup). Another body of work focuses on proving asymptotic minimax optimality in the worst-case setting rather than instance-dependent performance bounds, also targeting order optimal rather than exact optimal regret bounds. This is the case for example in [Kleinberg et al. \(2008\)](#) and [Bubeck et al. \(2008\)](#), respectively introducing ZOOMING and HOO algorithms. In particular the provided bounds on the regret are not instance-dependent and instance-dependent optimality is not established for these algorithms. Such a worst-case setting is out of the scope of this thesis.

For structured bandit problems and one-dimensional exponential family distributions (Assumption 1), the lower bound on the regret takes the generic form

$$\liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \geq \mathfrak{C}_{\mathcal{D}}(\mu), \quad (1.4)$$

where  $\mathfrak{C}_{\mathcal{D}}(\mu)$  is a constant, solution of a constrained linear-optimization problem. When there is no structure,  $\mathcal{D} = \mathcal{P}_1^A$  and then  $\mathfrak{C}_{\mathcal{P}_1^A}(\mu) = \mathfrak{C}_1(\mu)$ . We note that

$$\begin{aligned} \mathfrak{C}_1(\mu) &:= \min_{n \in \mathbb{R}_+^A} \sum_{a \in A} n_a (\max \mu - \mu_a) \\ &\text{s.t. } \forall a \notin \text{argmax } \mu, \quad n_a \text{KL}(\mu_a | \max \mu) \geq 1. \end{aligned} \quad (1.5)$$

The  $(n_a)_{a \in A}$  from (1.5) then appear as the normalized numbers of pulls  $(\mathbb{E}_{\nu}[N_a(T)] / \log(T))_{a \in A}$ .

In [Graves and Lai \(1997\)](#) a generic algorithm was proposed to solve any structured bandit problems, with however *prohibitive* computational complexity, including the requirement to compute a version of  $\mathfrak{C}_{\mathcal{D}}(\mu)$  at each time step. In [Combes et al. \(2017\)](#), another generic algorithm called OSSB (Optimal Structured Stochastic Bandit) is introduced, stepping the path towards generic structure-adaptive bandit algorithms. Although asymptotically optimal, the algorithm often suffers from poor finite-time numerical performances (as the asymptotic regime kicks-in possibly late), and still high computational cost. Inspired by combinatorial structures for which computing  $\mathfrak{C}_{\mathcal{D}}(\mu)$  is simply not feasible, a relaxation of the generic constrained optimization problem was recently proposed in [Cuvelier et al. \(2021a\)](#), however at the price of trading-off regret optimality for computational efficiency. In [Degegne et al. \(2020b\)](#), the authors explore an adaptation of KLUCB algorithm to structured  $\mathcal{D}$ . In [Van Parys and Golrezaei and \(2020\)](#), the authors propose an approach base on convex duality. Motivated by these issues, we make specific efforts to build algorithms that are regret efficient in practice.

We will address several structures such as Unimodal and Lipschitz structures which makes relevant to focus on the following generic algorithm: Optimal Sampling for Structured Bandits (OSSB) from [Combes et al. \(2017\)](#). We will naturally compare the algorithms we introduce to this known algorithm.

**Optimal Sampling for Structured Bandits (OSSB).** When assuming there is a unique optimal arm  $a^*$  (which is a very common assumption), that is  $\mathcal{A}^* = \{a^*\}$ , and under the assumptions of Theorem 1 from Combes et al. (2017), it is shown that  $\mathfrak{C}_{\mathcal{D}}(\mu)$  can be chosen so that

$$\begin{aligned} \mathfrak{C}_{\mathcal{D}}(\mu) &:= \min_{n \in \mathbb{R}_+^{\mathcal{A}}} \sum_{a \in \mathcal{A}} n_a (\max \mu - \mu_a) \\ \text{s.t. } \forall \lambda \in \Lambda(\mu), & \sum_{a \in \mathcal{A}} n_a \text{KL}(\mu_a | \lambda_a) \geq 1, \end{aligned} \quad (1.6)$$

where

$$\Lambda(\mu) = \left\{ \begin{array}{l} (1) \ \nu' \in \mathcal{D} \\ \mu(\nu') \in \mathbb{I}^{\mathcal{A}} : \begin{array}{l} (2) \ \mu_{a^*}(\nu') = \mu^* \\ (3) \ \max \mu(\nu') > \mu^* \end{array} \end{array} \right\} \quad (1.7)$$

is the set of ‘‘confusing’’ means for means  $\mu$  of configuration  $\nu \in \mathcal{D}$ .

Let  $(n_a(\mu))_{a \in \mathcal{A}}$  be the solution of the previous optimization problem. Then  $n_a(\mu) \log(T)$  indicates the asymptotic number of times sub-optimal arm  $a \neq a^*$  should be played under efficient and consistent algorithms.

Note that for all arm  $a \in \text{argmax} \mu$ , we have  $n_a(\mu) = 0$ .

The key idea of OSSB algorithm is then to compute *at each time step*  $(n_a(\hat{\mu}(t)))_{a \in \mathcal{A}}$ , the current empirical solution of  $\mathfrak{C}_{\mathcal{D}}(\hat{\mu}(t))$ , for  $t \geq 1$ , and explore so that the number of pull of sub-optimal arm  $a \neq a^*$ ,  $N_a(t)$  grows at speed  $n_a(\hat{\mu}(t)) \log(T)$ , that is so that  $N_a(t) \approx n_a(\hat{\mu}(t)) \log(T)$ .

In detail, OSSB algorithm takes theoretically two parameters  $\varepsilon, \gamma > 0$  in input. Then, at each time step  $t \geq 1$  the algorithm alternate between three phases: exploitation, estimation and exploration.

If for all arm  $a \in \mathcal{A}$ , we have  $N_a(t) \geq n_a(\hat{\mu}(t))(1 + \gamma) \log(t)$ , then we consider there is enough information to consider the structure is well estimated and we enter an exploitation phase: we pull arbitrarily an arm with maximal current mean, that is  $a_{t+1} \in \text{argmax}_{a \in \mathcal{A}} \hat{\mu}_a(t)$ .

Otherwise, we consider more information is needed to identify the optimal arm. Thus, either we enter an estimation phase and pull arbitrarily an arm with a minimal current number of pulls, that is  $a_{t+1} \in \text{argmin}_{a \in \mathcal{A}} N_a(t)$ ,

or we enter an exploration phase and pull arbitrarily an arm with a minimal current ratio, where current ratio of arm  $a \in \mathcal{A}$  is  $N_a(t)/n_a(\hat{\mu}(t))$  if  $n_a(\hat{\mu}(t)) > 0$ ,  $+\infty$  otherwise. Note that the current ratios measure how much the desired constraints  $N_a(t) \gtrsim n_a(\hat{\mu}(t)) \log(t)$ , for  $a \in \mathcal{A}$ , are satisfied.

More precisely, we consider a counter  $s(t)$  of the number of times we have not entered an exploitation phase. If  $\min_{a \in \mathcal{A}} N_a(t) \leq \varepsilon s(t)$ , we enter an estimation phase, otherwise we enter an exploration phase.

Finally, we note that parameters  $\varepsilon$  and  $\gamma$  are set equal to 0 in practice for the numerical experiments.

OSSB algorithm is summarized in Algorithm 5.



---

**Algorithm 5** OSSB

---

**Input:**  $\varepsilon, \gamma > 0$ Pull  $a_1 \in \mathcal{A}$  (arbitrarily chosen) $s(1) \leftarrow 0$ **for**  $t = 1 \dots T - 1$  **do**    Compute  $(n_a(\widehat{\mu}(t)))_{a \in \mathcal{A}}$ , the solution of current optimization problem  $\mathfrak{C}_{\mathcal{D}}(\widehat{\mu}(t))$     **if**  $N_a(t) \geq n_a(\widehat{\mu}(t))(1 + \gamma) \log(t)$ , for all  $a \in \mathcal{A}$  **then**         $s(t + 1) \leftarrow s(t)$         Pull  $a_{t+1} \in \operatorname{argmax}_{a \in \mathcal{A}} \widehat{\mu}_a(t)$  (arbitrarily chosen) ▷ Exploitation    **else**         $s(t + 1) \leftarrow s(t) + 1$         **if**  $\min_{a \in \mathcal{A}} N_a(t) \leq \varepsilon s(t)$  **then**            Pull  $a_{t+1} \in \operatorname{argmin}_{a \in \mathcal{A}} N_a(t)$  (arbitrarily chosen) ▷ Estimation        **else**            Pull  $a_{t+1} \in \operatorname{argmin}_{a \in \mathcal{A}} N_a(t) / n_a(\widehat{\mu}(t))$  (arbitrarily chosen) ▷ Exploration        **end if**    **end if****end for**

---

**Comment.** OSSB algorithm from [Combes et al. \(2017\)](#) is a quite simple, natural and generic algorithm for structured bandit problems, which makes it a reference in the field. However, OSSB presents this major constraint of solving current optimization problem  $\mathfrak{C}_{\mathcal{D}}(\widehat{\mu}(t))$  at each time step  $t \geq 1$  that limits its practical applications. Furthermore, parameters  $\varepsilon$  and  $\gamma$  are set equal to 0 in practice without analysis of the algorithm for this choice of setting. That being said, we will retain the main result from [Combes et al. \(2017\)](#) about the asymptotic optimality of OSSB algorithm that is reproduced in the following theorem. Note that OSSB algorithm is the first algorithm that optimally solve genetic structured bandit problems.

**Theorem 1** (OSSB optimality). *Let us consider a configuration  $\nu \in \mathcal{D}$ . If the assumptions of Theorem 2 from [Combes et al. \(2017\)](#) hold, then under OSSB algorithm, we have for all  $0 < \varepsilon < |\mathcal{A}|^{-1}$ ,*

$$\limsup_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \leq \mathfrak{C}_{\mathcal{D}}(\mu) F(\mu, \varepsilon, \gamma),$$

with  $F$  a function such that  $F(\mu, \varepsilon, \gamma) \rightarrow 1$  as  $\varepsilon \rightarrow 0$  and  $\gamma \rightarrow 0$ .

## 1.4 Publications

In this section, we detail the papers published in journals and conferences.

Saber, H., Ménard, P., and Maillard, O.-A. (2021a). Indexed minimum empirical divergence for unimodal bandits. *International Conference on Neural Information Processing Systems (NeurIPS)*

Pesquerel, F., Saber, H., and Maillard, O.-A. (2021). Stochastic bandits with groups of similar arms. *International Conference on Neural Information Processing Systems (NeurIPS)*

Saber, H., Saci, L., Maillard, O.-A., and Durand, A. (2021b). Routine bandits: Minimizing regret on recurring problems. *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PPKD)*

## 1.5 Contributions

We summarize our contributions for structured bandit problems.

Unimodal, multimodal and graph structures are respectively the subject of Chapter 2, 3 and 4. Our work on routine bandits is the object of Chapter 6. These chapters have been written to be largely independent and thus facilitate the reading of the manuscript. Our work on bandits with groups of similar arms is included in Appendix D. The complements from [Saber et al. \(2020\)](#) relating to unimodal bandits are not reproduced in the manuscript.

**Unimodal bandits.** In the unimodal bandit problem, a graph  $G$  supports the structure and allows a notion of neighbourhood between the arms. An arm  $a \in \mathcal{A}$  is then considered as a local maximum if its mean  $\mu_a$  is greater than the means of the arms in its neighbourhood  $\mathcal{V}_a \subset \mathcal{A} \setminus \{a\}$ , that is  $\mu_a > \mu_{a'}$ , for all  $a' \in \mathcal{V}_a$ . In the unimodal bandit problem, we only consider distributions  $\nu$  with a unique local maximum. This is formalized in Chapter 2, where we provide novel regret minimization results related to the unimodal structure. We first revisit Indexed Minimum Empirical Divergence (IMED) algorithm from [Honda and Takemura \(2015\)](#) introduced for unstructured multi-armed bandits, and adapt it to the unimodal setting. We introduce in Section 2.3 IMED-UB algorithm 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* algorithm. We prove in Theorem 1 that IMED-UB is asymptotically optimal. Furthermore, this novel algorithm competes with the state-of-the-art algorithms in practice. This is confirmed by numerical illustrations on synthetic data. This is the subject of paper [Saber et al. \(2021a\)](#). In [Saber et al. \(2020\)](#), we introduced and studied KLUCB-UB algorithm, a KLUCB version of IMED-UB. Our KLUCB-UB finite time analysis highlighted strong links between KLUCB and IMED approaches. We further developed  $d$ -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 within a large set of arms  $\mathcal{A}$  by empirical considerations. This second order exploration appears as a new mechanism of exploring sub-optimal arms (and the underlying structure) and seems to be specific to IMED approach. This mechanism operates in IMED-MB algorithm introduced in Chapter 3, where we define the multimodal structure that generalises the unimodal bandit problem. Note that our work from [Saber et al. \(2020\)](#) on KLUCB-UB and  $d$ -IMED-UB algorithms is not reproduced in the manuscript.

**Multimodal bandits.** In the multimodal bandit problem, a graph  $G$  supports the structure and allows a notion of neighbourhood between the arms. An arm  $a \in \mathcal{A}$  is then considered as a local maximum if its mean  $\mu_a$  is greater than the means of the arms in its neighbourhood  $\mathcal{V}_a \subset \mathcal{A} \setminus \{a\}$ , that is  $\mu_a > \mu_{a'}$ , for all  $a' \in \mathcal{V}_a$ . In the multimodal bandit problem, we only consider distributions  $\nu$  with a fixed number of local maximums that is known in advance to the learner. This is formalized in Chapter 3, where we generalize the unimodal structure and precisely define the multimodal bandit problem. We provide regret minimization results related to the multimodal structure. We again revisit the Indexed Minimum Empirical Divergence (IMED) algorithm from [Honda and Takemura \(2015\)](#) introduced for unstructured multi-armed bandits, and adapt it to the multimodal-structured setting. We introduce in Section 3.3 IMED-MB algorithm that is limited to the pulling of the arms with locally maximal empirical means or their no more than  $d$  nearest arms at each time step, with  $d$  the maximum degree of nodes in  $G$ . We do not provide finite time analysis of IMED-MB algorithm, but it could be the focus of future work. This novel algorithm performs well in practice and competes with the state-of-the-art algorithms when assuming unimodal structure. This is confirmed by numerical illustrations on synthetic data. We believe that the construction of this algorithm and the use of IMED type indexes to test the local

maximality of arms with locally maximal empirical means (Equation 3.17) are of independent interest for the bandit community.

**Bandits with groups of similar arms.** We consider a variant of the stochastic multi-armed bandit problem where arms are known to be organized into different groups having the same mean. The groups are unknown but a lower bound  $q$  on their size is known. This situation typically appears when each arm can be described with a list of categorical attributes, and the (unknown) mean reward function only depends on a subset of them, the others being redundant. In this case,  $q$  is linked naturally to the number of attributes considered redundant, and the number of categories of each attribute. For this structured problem of practical relevance, we first derive the asymptotic regret lower bound and corresponding constrained optimization problem. They reveal the achievable regret can be substantially reduced when compared to the unstructured setup, possibly by a factor  $q$ . However, solving exactly the exact constrained optimization problem involves a combinatorial problem. Owing to this key insight, we introduce  $\text{IMED-EC}$ , an adaptation of  $\text{IMED}$  algorithm from [Honda and Takemura \(2015\)](#) to the considered structured set of bandits. One advantage of  $\text{IMED}$  over a  $\text{KLUCB}$  alternative is its reduced complexity, which translates to the equivalence class setup. At each time step, the complexity of computing the next arm to be pulled by  $\text{IMED-EC}$  is no more than the one of sorting a list of  $|\mathcal{A}|$  elements once the  $\text{IMED}$  indexes have been computed, which is only  $\log(|\mathcal{A}|)$  times larger than looking for the minimal  $\text{IMED}$  index. We prove that  $\text{IMED-EC}$  achieves a controlled asymptotic regret that matches the non-combinatorial part of the lower bound and is at most (less than) a factor of 2 times the optimal regret bound. Last, we illustrate the benefit of the  $\text{IMED-EC}$  over its unstructured version, where it shows a substantial improvement. This work is the subject of [Pesquerel et al. \(2021\)](#) and has been initiated by Fabien Pesquerel.

**Graph-structured bandits.** In Chapter 4, we consider a structured variant of the multi-armed bandit problem when the difference of means between any pair of arms is constrained not to exceed some value. This graph structure is introduced to encompass some classical structures such as Unimodal and Lipschitz. We assume Bernoulli distributions although most of this work can be generalised to one-dimensional exponential family distributions. In Section 4.2, we first provide an instance-dependent lower bound on the regret for generic graph-structured bandits, see Proposition 3. The proof follows standard change-of-measure techniques and reveals interesting sets of “confusing” bandit instances. In Section 4.3, building on  $\text{IMED}$  algorithm originally introduced for unstructured bandits by [Honda and Takemura \(2015\)](#), we introduce  $\text{IMED-GS}$  algorithm, to which we incorporate some key modifications in order to adapt it to a graph-structure. This  $\text{IMED}$  type approach directly builds on the lower bounds. We detail  $\text{IMED-GS}$  in Section 4.3.1 and Algorithm 9. The main theoretical result is Theorem 3 that provides a *finite-time* upper bound on the regret under  $\text{IMED-GS}$ , which then implies asymptotic optimality as shown in Corollary 3. The proof of this result, detailed in Section 4.5, involves a technique exploiting *empirical lower bounds*. It also involves a novel concentration inequality of independent interest that combines stochastic orderings with approximations. In Section 4.4, we illustrate and discuss the benefit of this approach over the state-of-the-art regarding instance-dependent optimality on some specific structures, including Unimodal and Lipschitz bandits.

**Concentration inequalities.** In Chapter 5, we highlight new concentration inequalities introduced to provide a finite time analysis of  $\text{IMED-GS}$  algorithm and prove its asymptotic optimality for graph-structured bandits (Chapter 4). These concentration inequalities enable to obtain more precise control of the concentration terms compared to alternative tools present in the literature like in Theorem 2 from [Magureanu et al. \(2014\)](#).

**Routine bandits.** In the online recommender system problem items are recommended to users. Items can be seen as arms and users as bandit instances. When a recommender system is deployed on multiple users, one does not typically assume that the best recommendation is the same for all users. The naive algorithm in this

situation is to consider each user as being a different bandit instance and learning from scratch for each user. When users can be recognized (e.g., characterized by features), this information can be leveraged to speed up the learning process by sharing observations across users. The resulting setting is known as *contextual bandit* (Langford and Zhang, 2007; Lu et al., 2010). We tackle the case where users cannot be or do not want to be identified (e.g., for privacy reasons), but where we assume that there exists a (unknown) finite set of possible user profiles (bandit instances), such that information may be shared between the current user and some previously encountered users. We call the resulting setting as *routine bandit* (Saber et al., 2021b). We establish lower bounds on the achievable cumulative regret that adapt the bound from Lai and Robbins (1985) to the routine setting. We then extend KLUCB algorithm, known to be optimal under the classical stochastic bandit setting, into a new algorithm called KLUCB-RB that leverages the information obtained on previously encountered bandits. We provide a theoretical analysis of KLUCB-RB and investigate the performance of the algorithm using extensive numerical experiments. These results highlight the empirical conditions required so that past information can be efficiently leveraged to speed up the learning process. The main contributions of this work are 1) the newly proposed routine bandit setting, 2) KLUCB-RB algorithm that solves this problem with asymptotically optimal regret minimization guarantees, and 3) an empirical illustration of the conditions for past information to be beneficial to the learning agent. This work is the subject of Saber et al. (2021b) and has been in collaboration with Léo Saci and Audrey Durand. It is the subject of Chapter 6.

# Chapter 2

## Unimodal Bandits

We consider a multi-armed bandit problem specified by a set of one-dimensional family exponential distributions endowed with a unimodal structure. We introduce `IMED-UB`, an algorithm that optimally exploits the unimodal-structure, by adapting to this setting Indexed Minimum Empirical Divergence (`IMED`) algorithm introduced by [Honda and Takemura \(2015\)](#). Owing to our proof technique, we are able to provide a concise finite-time analysis of `IMED-UB` algorithm. Numerical experiments show that `IMED-UB` competes with the state-of-the-art algorithms. This chapter reproduces the contents from [Saber et al. \(2021a\)](#) published at NeurIPS conference.

### 2.1 Introduction

We assume a *unimodal* structure similar to that considered in [Yu and Mannor \(2011\)](#) and [Combes and Proutiere \(2014a\)](#). 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  $\mu_{a_i} < \mu_{a_{i+1}}$ . We denote by  $\mathcal{V}_a := \{a' \in \mathcal{A} \setminus \{a\} : (a, a') \in E\}$  the neighbourhood of arm  $a \in \mathcal{A}$  in graph  $G$ . Lastly, we assume that  $\nu \subset \mathcal{P} := \{p(\mu), \mu \in \mathbb{I}\}$ , 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  $\mu \in \mathbb{I} \subset \mathbb{R}$ .  $\mathcal{P}$  is assumed to be known to the learner. Thus, for all  $a \in \mathcal{A}$  the distribution  $\nu_a$  is fully specified by its mean  $\mu_a$ , and we have  $\nu_a = p(\mu_a)$ . We denote by  $\mathcal{D}_{(\mathcal{P}, G)}$  or  $\mathcal{D}_{\text{uni}}$  (or simply  $\mathcal{D}$  when there is no confusion) the structured set of such unimodal-bandit distributions characterized by  $(\mathcal{P}, G)$ . In the following, we assume that  $\mathcal{P}$  is a set of one-dimensional exponential family distributions.

**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 tower rule we can rewrite the regret as follows:

$$R(\nu, T) = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}_\nu [N_a(T)], \quad (2.1)$$

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

## 2.2 Regret lower bound

In this section, 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 algorithms that are *consistent* (Definition 1). We can derive from the notion of consistency an asymptotic lower bound on the regret, see [Combes and Proutiere \(2014a\)](#).

**Proposition 1** (Lower bounds on the regret). *Let us consider a consistent algorithm. Then, for all configuration  $\nu \in \mathcal{D}_{uni}$  with means  $\mu \in \mathbf{I}^A$ , it must be that*

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

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  $\mu, \mu' \in \mathbf{I}$ , and where  $\mathcal{V}_{a^*} \subset \mathcal{A} \setminus \{a^*\}$  is the neighbourhood of optimal arm  $a^*$ .

**Remark 3.**  $\mathfrak{C}_{\mathcal{D}_{uni}}(\mu)$  is simply denoted  $\mathfrak{C}_{uni}(\mu)$ .

**Remark 4.** The quantity  $\mathfrak{C}_{uni}(\mu)$  is a fully explicit function of  $\mu$  (it does not require solving any optimization problem) for some set of distributions  $\nu$  (see Remark 1). This useful property no longer holds in general for arbitrary structures. Also, it is noticeable that  $\mathfrak{C}_{uni}(\mu)$  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 \(2014a\)](#) for further insights.

## 2.3 Optimal algorithm for unimodal bandits

We present in this section a novel algorithm that matches the asymptotic lower bound of Proposition 1. This algorithm is inspired by the Indexed Minimum Empirical Divergence (IMED) proposed by [Honda and Take-mura \(2011\)](#). The general idea behind this algorithm 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  $a$  at time  $t$  is denoted by  $\hat{\mu}_a(t) = \frac{\sum_{s=1}^t \mathbb{I}_{\{a_s=a\}} X_s}{N_a(t)}$  if  $N_a(t) > 0$ , 0 otherwise. We also denote by  $\hat{\mu}^*(t) = \max_{a \in \mathcal{A}} \hat{\mu}_a(t)$  and  $\hat{\mathcal{A}}^*(t) = \operatorname{argmax}_{a \in \mathcal{A}} \hat{\mu}_a(t)$  respectively the current best mean and the current set of optimal arms.

### 2.3.1 OSUB algorithm

For completeness, we recall in this subsection OSUB (Optimal Sampling for Unimodal Bandits) algorithm from [Combes and Proutiere \(2014a\)](#).

---

#### Algorithm 6 OSUB

---

Pull each arm once

**for**  $t = |\mathcal{A}| \dots T - 1$  **do**

Choose  $\hat{a}_t^* \in \operatorname{argmin}_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} N_{\hat{a}^*}(t)$  (chosen arbitrarily) so that  $\hat{a}_t^* \in \hat{\mathcal{A}}^*(t)$  and  $N_{\hat{a}_t^*}(t) \leq N_{\hat{a}^*}(t)$ ,  $\forall \hat{a}^* \in \hat{\mathcal{A}}^*(t)$

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 6, for some numerical constant  $c > 0$ , the index computed by OSUB algorithm for arm  $a \in \mathcal{A}$  and step  $t \geq |\mathcal{A}|$  is

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

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),  $d = \max_{a \in \mathcal{A}} |\mathcal{V}_a|$  is the maximum degree of nodes in  $G$ , and  $f_c(\cdot) = \log(\cdot) + c \log \log(\cdot)$ . We set  $c = 1$  for the numerical experiments.

### 2.3.2 IMED-UB algorithm

We first pull each arm once. For all arm  $a \in \mathcal{A}$  and time step  $t \geq 1$  we introduce the IMED index

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

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 number 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 the empirical mean of the current optimal arm. We define IMED-UB (Indexed Minimum Empirical Divergence for Unimodal Bandits), described in Algorithm 7, to be the algorithm 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  $\hat{a}_t^* \in \underset{\hat{a}^* \in \hat{\mathcal{A}}^*(t)}{\text{argmin}} N_{\hat{a}^*}(t)$  is a current best arm in  $\hat{\mathcal{A}}^*(t)$  such that for all current best arm  $\hat{a}^* \in \hat{\mathcal{A}}^*(t)$ ,  $N_{\hat{a}^*}(t) \geq N_{\hat{a}_t^*}(t)$ . This is a natural algorithm since the lower bound on the regret given in Proposition 1 involves only the arms in  $\mathcal{V}_{a^*}$ , the neighbourhood of the arm  $a^*$  of maximal mean.

---

#### Algorithm 7 IMED-UB

---

Pull each arm once

**for**  $t = |\mathcal{A}| \dots T - 1$  **do**

    Choose  $\hat{a}_t^* \in \underset{\hat{a}^* \in \hat{\mathcal{A}}^*(t)}{\text{argmin}} N_{\hat{a}^*}(t)$  (chosen arbitrarily) so that  $\hat{a}_t^* \in \hat{\mathcal{A}}^*(t)$  and  $N_{\hat{a}_t^*}(t) \leq N_{\hat{a}^*}(t)$ ,  $\forall \hat{a}^* \in \hat{\mathcal{A}}^*(t)$

    Pull  $a_{t+1} \in \underset{a \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}}{\text{argmin}} I_a(t)$  (chosen arbitrarily)

**end for**

---

### 2.3.3 Asymptotic optimality of IMED-UB

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

**Theorem 2** (Upper bounds). *Let us consider a set of distributions  $\nu \in \mathcal{D}_{\text{uni}}$  with means  $\mu \in \mathcal{I}^A$  and let  $a^*$  its optimal arm. Let  $\mathcal{V}_{a^*}$  be the sub-optimal arms in the neighbourhood of  $a^*$ . Then under IMED-UB algorithm for all  $0 < \varepsilon < \varepsilon_\nu$ , for all horizon time  $T \geq 1$ , for all  $a \in \mathcal{V}_{a^*}$ ,*

$$\begin{aligned} \mathbb{E}_\nu[N_a(T)] &\leq \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a|\mu_{a^*})} \log(T) + 2 \min\{|\mathcal{A}|, d(d+1)\} C_\varepsilon \sqrt{\log(c_\varepsilon T)} \\ &\quad + \min\{|\mathcal{A}|, d(d+1)\} (1 + c_{\varepsilon_\nu}^{-1}) + \min\{3|\mathcal{A}|, (d+1)(d+3)\} \frac{2\sigma_{\varepsilon_\nu}^2 e^{\varepsilon_\nu^2/2\sigma_\varepsilon^2}}{\varepsilon^2} + 1 \end{aligned}$$

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

$$\begin{aligned} \mathbb{E}_\nu[N_a(T)] &\leq 2 \min\{|\mathcal{A}|, d(d+1)\} C_\varepsilon \sqrt{\log(c_\varepsilon T)} \\ &\quad + \min\{|\mathcal{A}|, d(d+1)\} (1 + c_{\varepsilon_\nu}^{-1}) + \min\{3|\mathcal{A}|, (d+1)(d+3)\} \frac{2\sigma_{\varepsilon_\nu}^2 e^{\varepsilon_\nu^2/2\sigma_\varepsilon^2}}{\varepsilon^2} + 1, \end{aligned}$$

where  $d$  is the maximum degree of nodes in  $G$ ,  $\varepsilon_\nu = \frac{1}{2} \min_{a \neq a^*} \max_{a' \in \mathcal{V}_a} \mu_{a'} - \mu_a$ ,

$\sigma_\varepsilon^2 = \max_{a \in \mathcal{A}} \left\{ \mathbb{V}_{X \sim p(\mu')} (X) : \mu' \in [\mu_a - \varepsilon, \mu_a] \right\}$  and  $c_\varepsilon, C_\varepsilon > 0$  are the constants involved in Theorem 6.  $\alpha_\nu(\cdot)$  is a non-negative function depending only on  $\nu$  such that  $\lim_{\varepsilon \rightarrow 0} \alpha_\nu(\varepsilon) = 0$  (see Section 2.4.1 for more details).

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  $\mathcal{O}\left(\sqrt{\log(T)}\right)$  of times under IMED-UB. This is coherent with the lower bound that only involves the neighbourhood of the best arm. More precisely, combining Theorem 2 and the tower rule (2.1) gives the asymptotic optimality of IMED-UB with respect to the lower bound of Proposition 1.

**Corollary 1** (Asymptotic optimality). *With the same notations as in Theorem 2, then under IMED-UB algorithm*

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

A finite time analysis of IMED-UB is provided in following Section 2.4.

### 2.3.4 Asymptotic optimality of IMED

The non-structured multi-armed bandit problem can be seen as a particular case of unimodal bandits. Indeed, if we consider  $E = \{(a, a') : a \neq a'\}$  then  $G$  is the fully connected graph and  $\mathcal{V}_a = \mathcal{A} \setminus \{a\}$  for all arm  $a \in \mathcal{A}$ . In particular, the path from any sub-optimal arm  $a$  to  $a^*$ , the optimal one, is simply  $(a, a^*) \in E$ . Furthermore, in such a case, IMED and IMED-UB algorithms, respectively summarized in Algorithms 2 and 7, coincide perfectly. Thus, a straightforward application of Theorem 2 gives us optimal finite time guarantees for IMED algorithm. Note that this natural shift from the structured unimodal case to the unstructured case from both an algorithmic and analytical point of view makes the IMED approach an innovative approach to deal with structured bandit problems.

**Corollary 2** (Upper bounds under IMED algorithm). *Let us consider a set of distributions  $\nu \in \mathcal{D}$  with means  $\mu \in \mathbb{I}^A$  and let  $a^*$  its optimal arm. Then under IMED algorithm for all  $0 < \varepsilon < \varepsilon_\nu$ , for all horizon time  $T \geq 1$ , for all sub-optimal arm  $a \neq a^*$ ,*

$$\mathbb{E}_\nu[N_a(T)] \leq \frac{1 + \alpha_\nu(\varepsilon)}{\text{KL}(\mu_a | \mu_{a^*})} \log(T) + 2 |\mathcal{A}| C_\varepsilon \sqrt{\log(c_\varepsilon T)} + |\mathcal{A}| (1 + c_\varepsilon^{-1}) + 3 |\mathcal{A}| \frac{2\sigma_{\varepsilon_\nu}^2 e^{\varepsilon_\nu^2/2\sigma_\varepsilon^2}}{\varepsilon^2} + 1$$

where  $\varepsilon_\nu = \frac{1}{2} \min_{a \neq a^*} \mu_{a^*} - \mu_a$ ,  $\sigma_\varepsilon^2 = \max_{a \in \mathcal{A}} \left\{ \mathbb{V}_{X \sim p(\mu')} (X) : \mu' \in [\mu_a - \varepsilon, \mu_a] \right\}$  and  $c_\varepsilon, C_\varepsilon > 0$  are the constants involved in Theorem 6.  $\alpha_\nu(\cdot)$  is a non-negative function depending only on  $\nu$  such that  $\lim_{\varepsilon \rightarrow 0} \alpha_\nu(\varepsilon) = 0$  (see Section 2.4.1 for more details). Thus under IMED algorithm, we have

$$\limsup_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \leq \mathfrak{C}_1(\mu) = \sum_{a \neq a^*} \frac{\Delta_a}{\text{KL}(\mu_a | \mu_{a^*})}.$$

## 2.4 IMED-UB finite time analysis

At a high level, the key interesting step of the proof is to realize that the considered algorithm implies empirical lower and empirical upper bounds on the numbers of pulls (see Lemma 1, Lemma 2). Then, based on



concentration lemmas (see Section B.2), the algorithm-based empirical lower bounds ensure the reliability of the estimators of interest (Lemma 4). Interestingly, this makes use of arguments based on recent concentration of measure that enable to control the concentration without adding some log log bonus (such a bonus was required for example in the initial analysis of KLUCB algorithm from Cappé et al. (2013)). Then, combining the reliability of these estimators with the obtained algorithm-based empirical upper bounds, we obtain upper bounds on the average numbers of pulls (Theorem 2). The proof is concise to fit mostly in the next few pages.

### 2.4.1 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

$$d = \max_{a \in \mathcal{A}} |\mathcal{V}_a|, \quad \varepsilon_\nu = \frac{1}{2} \min_{a \neq a^*} \max_{a' \in \mathcal{V}_a} \mu_{a'} - \mu_a. \quad (2.2)$$

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

$$\text{KL}(\mu_a + \varepsilon | \mu^* - \varepsilon) \geq (1 + \alpha_\nu(\varepsilon))^{-1} \text{KL}(\mu_a | \mu^*) \quad (2.3)$$

and  $\lim_{\varepsilon \downarrow 0} \alpha_\nu(\varepsilon) = 0$ . Indeed, since  $\text{KL}(\cdot | \cdot)$  is a convex function, it is a continuous function within the interior of its domain of definition. At each time step  $t \geq 1$ ,  $\hat{a}_t^*$  is arbitrarily chosen in  $\underset{a \in \hat{\mathcal{A}}^*(t)}{\text{argmin}} N_a(t)$  where

$\hat{\mathcal{A}}^*(t) = \underset{a \in \mathcal{A}}{\text{argmax}} \hat{\mu}_a(t)$ , that is

$$\begin{cases} \hat{a}_t^* \in \hat{\mathcal{A}}^*(t) \\ \forall \hat{a}^* \in \hat{\mathcal{A}}^*(t), \quad N_{\hat{a}^*}(t) \geq N_{\hat{a}_t^*}(t). \end{cases}$$

### 2.4.2 Algorithm-based empirical bounds

IMED-UB algorithm 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 1** (Empirical lower bounds). *Under IMED-UB, at each step time  $t \geq |\mathcal{A}|$ , for all  $a \in \mathcal{V}_{\hat{a}_t^*}$ ,*

$$\log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) + \log(N_a(t)) \quad (2.4)$$

and

$$N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t). \quad (2.5)$$

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

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

This implies, since the arm with minimum index is pulled,  $\log(N_{a_{t+1}}(t)) \leq I_{a_{t+1}}(t) = \min_{a' \in \{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}} I_{a'}(t) \leq I_{\hat{a}_t^*}(t) = \log(N_{\hat{a}_t^*}(t))$ . By taking the  $\log^{-1}(\cdot)$ , the last inequality allows us to conclude.  $\square$

**Lemma 2** (Empirical upper bounds). *Under IMED-UB at each step time  $t \geq |\mathcal{A}|$ ,*

$$N_{a_{t+1}}(t) \text{KL}(\widehat{\mu}_{a_{t+1}}(t) | \widehat{\mu}^*(t)) \leq \log(t). \quad (2.6)$$

*Proof.* As above, by construction we have

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

It remains, to conclude, to note that

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

and

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

□

### 2.4.3 Non-reliable current best arm

Before going further in the analysis, we inform the reader that sets  $\mathcal{E}_{a,a'}^+(\varepsilon)$ ,  $\mathcal{E}_{a,a'}^-(\varepsilon)$ ,  $\mathcal{K}_{a,a'}^-(\varepsilon)$  for  $a, a' \in \mathcal{A}$ ,  $\varepsilon > 0$ , used in this section are introduced and study in Section B.1.

For accuracy  $\varepsilon > 0$ , let  $\mathcal{M}^*(\varepsilon)$  be the set of times  $t \geq 1$  that do not belong to  $\mathcal{E}_{\widehat{a}_t^*, a_{t+1}}^+(\varepsilon)$  and where some of the current best arm  $\widehat{a}_t^*$  differs from  $a^*$ ,

$$\mathcal{M}^*(\varepsilon) := \left\{ t \geq |\mathcal{A}| : \begin{array}{l} (1) \ t \notin \mathcal{E}_{\widehat{a}_t^*, a_{t+1}}^+(\varepsilon) \\ (2) \ \widehat{a}_t^* \neq a^* \end{array} \right\}. \quad (2.7)$$

**Lemma 3** (Relation between subsets of times). *Under IMED-UB, for all accuracy  $0 < \varepsilon < \varepsilon_\nu$ ,*

$$\mathcal{M}^*(\varepsilon) \subset \bigcup_{\substack{t \geq 1 \\ a \in \mathcal{V}_{\widehat{a}_t^*}}} \mathcal{K}_{a, a_{t+1}}^-(\varepsilon_\nu), \quad (2.8)$$

where  $\varepsilon_\nu = \frac{1}{2} \min_{a \neq a^*} \max_{a' \in \mathcal{V}_a} \mu_{a'} - \mu_a$ .

*Proof.* Let us consider  $t \in \mathcal{M}^*(\varepsilon)$ . Since  $\widehat{a}_t^* \neq a^*$ , there exists  $a \in \underset{a' \in \mathcal{V}_{\widehat{a}_t^*}}{\text{argmax}} \mu_{a'}$  such that

$$\mu_a > \mu_{\widehat{a}_t^*}. \quad (2.9)$$

Then, since  $\widehat{a}_t^* \in \underset{a \in \mathcal{A}}{\text{argmax}} \widehat{\mu}_a(t)$ , we have

$$\widehat{\mu}_{\widehat{a}_t^*}(t) = \widehat{\mu}^*(t) \geq \widehat{\mu}_a(t). \quad (2.10)$$

Since  $t \in \mathcal{M}^*(\varepsilon)$ ,  $t \notin \mathcal{E}_{\widehat{a}_t^*, a_{t+1}}^+(\varepsilon)$ , where  $\mathcal{E}_{\widehat{a}_t^*, a_{t+1}}^+(\varepsilon) = \{t \in [1, T-1] : N_{a_{t+1}}(t) \leq N_{\widehat{a}_t^*}(t), \widehat{\mu}_{\widehat{a}_t^*}(t) \geq \mu_{\widehat{a}_t^*} + \varepsilon\}$  according to Equation (B.1). From empirical lower bounds (2.5), we have  $N_{a_{t+1}}(t) \leq N_{\widehat{a}_t^*}(t)$ . This implies

$$\mu_{\widehat{a}_t^*} + \varepsilon \geq \widehat{\mu}_{\widehat{a}_t^*}(t). \quad (2.11)$$

By combining Equations (2.10) and (2.11), it comes

$$\mu_{\widehat{a}_t^*} + \varepsilon \geq \widehat{\mu}^*(t) \geq \widehat{\mu}_a(t). \quad (2.12)$$

Since  $\varepsilon < \varepsilon_\nu \leq |\mu_a - \mu_{\hat{a}_t^*}|/2$ , Equation (2.9) and previous Equation (2.12) imply

$$\mu_a - \varepsilon_\nu > \hat{\mu}_{\hat{a}_t^*}(t) \geq \hat{\mu}_a(t). \quad (2.13)$$

Since  $a \in \mathcal{V}_{\hat{a}_t^*}$ , empirical lower bounds (2.4) imply

$$\log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t)) + \log(N_a(t)). \quad (2.14)$$

The classical monotonic properties of  $\text{KL}(\cdot|\cdot)$  and Equation (2.13) imply

$$\begin{cases} \hat{\mu}_a(t) < \mu_a - \varepsilon_\nu \\ \text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t)) \leq \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon_\nu). \end{cases} \quad (2.15)$$

Combining Equations (2.14) and (2.15), we get

$$\begin{cases} \hat{\mu}_a(t) < \mu_a - \varepsilon_\nu \\ \log(N_{a_{t+1}}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon_\nu) + \log(N_a(t)), \end{cases} \quad (2.16)$$

which means  $t \in \mathcal{K}_{a, a_{t+1}}^-(\varepsilon_\nu)$ .  $\square$

#### 2.4.4 Reliable current means and current best arm

In this subsection, we characterize subsets of times where both the mean of current pulled arm and the optimal mean are well-estimated.

Let us consider for  $0 < \varepsilon < \varepsilon_\nu$ , for  $a \neq a^*$ ,

$$\mathcal{U}_a(\varepsilon) = \{t \geq |\mathcal{A}| : a_{t+1} = a\} \cap \left( \bigcup_{t \geq 1} \mathcal{E}_{a_{t+1}, a_{t+1}}^+(\varepsilon) \cup \mathcal{E}_{\hat{a}_t^*, a_{t+1}}^-(\varepsilon) \cup \mathcal{E}_{\hat{a}_t^*, a_{t+1}}^+(\varepsilon) \cup \mathcal{M}^*(\varepsilon) \right). \quad (2.17)$$

According to IMED-UB algorithm (summarized in Algorithm 7), for all  $t \geq |\mathcal{A}|$ , current pulled arm  $a_{t+1}$  is chosen in  $\{\hat{a}_t^*\} \cup \mathcal{V}_{\hat{a}_t^*}$  so that  $\hat{a}_t^* \in \{a_{t+1}\} \cup \mathcal{V}_{a_{t+1}}$ . Furthermore, from Lemma 3 we have  $\mathcal{M}^*(\varepsilon) \subset \bigcup_{\substack{t \geq 1 \\ a \in \mathcal{V}_{\hat{a}_t^*}}} \mathcal{K}_{a, a_{t+1}}^-(\varepsilon_\nu)$ .

This implies

$$\mathcal{U}_a(\varepsilon) \subset \bigcup_{\substack{a' \in \{a\} \cup \mathcal{V}_a \\ a'' \in \mathcal{V}_{a'}}} \mathcal{E}_{a', a}^-(\varepsilon) \cup \mathcal{E}_{a', a}^+(\varepsilon) \cup \mathcal{K}_{a'', a}^-(\varepsilon_\nu) \subset \bigcup_{a''' \in \mathcal{A}} \mathcal{E}_{a''', a}^-(\varepsilon) \cup \mathcal{E}_{a''', a}^+(\varepsilon) \cup \mathcal{K}_{a''', a}^-(\varepsilon_\nu), \quad (2.18)$$

where dummy variables  $a_{t+1}$  and  $\hat{a}_t^*$  have been respectively replaced by  $a$  and  $a'$ .

In particular, from Lemma 30 and previous Equation (2.18) we have

$$\begin{aligned} \mathbb{E}_\nu[\mathcal{U}_a(\varepsilon)] &\leq 2 \min\{|\mathcal{A}|, (d+1)\} \frac{2\sigma_\varepsilon^2 e^{\varepsilon^2/2\sigma_\varepsilon^2}}{\varepsilon^2} \\ &\quad + \min\{|\mathcal{A}|, d(d+1)\} \left( \frac{2\sigma_{\varepsilon_\nu}^2 e^{\varepsilon_\nu^2/2\sigma_{\varepsilon_\nu}^2}}{\varepsilon_\nu^2} + 1 + c_{\varepsilon_\nu}^{-1} + 2C_\varepsilon \sqrt{\log(c_\varepsilon T)} \right) \\ &\leq \min\{3|\mathcal{A}|, (d+1)(d+3)\} \frac{2\sigma_{\varepsilon_\nu}^2 e^{\varepsilon_\nu^2/2\sigma_{\varepsilon_\nu}^2}}{\varepsilon^2} \\ &\quad + \min\{|\mathcal{A}|, d(d+1)\} \left( 1 + c_{\varepsilon_\nu}^{-1} + 2C_\varepsilon \sqrt{\log(c_\varepsilon T)} \right), \end{aligned} \quad (2.19)$$

where  $d = \max_{a \in \mathcal{A}} |\mathcal{V}_a|$  is the maximum degree of nodes in graph  $\mathcal{G}$ .

**Lemma 4 (Reliable current means).** *Under IMED-UB, for all accuracy  $0 < \varepsilon < \varepsilon_\nu$ , for all sub-optimal arm  $a \neq a^*$ , for all time step  $t \notin \mathcal{U}_a(\varepsilon)$ ,  $t \geq |\mathcal{A}|$ , such that  $a_{t+1} = a$ ,*

$$\begin{cases} \hat{a}_t^* = a^* \\ \hat{\mu}^*(t) \geq \mu^* - \varepsilon \\ \hat{\mu}_a(t) \leq \mu_a + \varepsilon. \end{cases}$$

*Proof.* For  $0 < \varepsilon < \varepsilon_\nu = \min_{a \neq a'} |\mu_a - \mu_{a'}|/2$ , for  $a \neq a^*$ , let us consider a time step  $t \notin \mathcal{U}_a(\varepsilon)$ ,  $t \geq |\mathcal{A}|$  such that  $a_{t+1} = a$ .

Since  $a_{t+1} = a$  and  $t \notin \mathcal{U}_{a_{t+1}}(\varepsilon)$  then  $t \notin \mathcal{E}_{a_{t+1}, a_{t+1}}^+(\varepsilon)$ , that is  $\hat{\mu}_{a_{t+1}}(t) < \mu_{a_{t+1}} + \varepsilon$  or  $\hat{\mu}_a(t) < \mu_a + \varepsilon$  (since  $a_{t+1} = a$ ).

Since  $a_{t+1} = a$  and  $t \notin \mathcal{U}_{a_{t+1}}(\varepsilon)$  then  $t \notin \mathcal{E}_{\hat{a}_t^*, a_{t+1}}^-(\varepsilon)$ , where

$$\mathcal{E}_{\hat{a}_t^*, a_{t+1}}^-(\varepsilon) = \{t \in \llbracket 1, T-1 \rrbracket : N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t), \hat{\mu}_{\hat{a}_t^*}(t) \leq \mu_{\hat{a}_t^*} - \varepsilon\}$$

according to Equation (B.2). From empirical lower bounds (2.5), we have  $N_{a_{t+1}}(t) \leq N_{\hat{a}_t^*}(t)$ . This implies

$$\hat{\mu}^*(t) = \hat{\mu}_{\hat{a}_t^*}(t) > \mu_{\hat{a}_t^*} - \varepsilon. \quad (2.20)$$

Since  $a_{t+1} = a$  and  $t \notin \mathcal{U}_{a_{t+1}}(\varepsilon)$  then  $t \notin \mathcal{E}_{\hat{a}_t^*, a_{t+1}}^+(\varepsilon) \cup \mathcal{M}^*(\varepsilon)$ . From Equation (2.7), this implies

$$\hat{a}_t^* = a^*. \quad (2.21)$$

By combining Equations (2.20) and (2.21), we get

$$\hat{\mu}^*(t) > \mu_{a^*} - \varepsilon = \mu^* - \varepsilon. \quad (2.22)$$

□

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

In this subsection, we now combine the different results of the previous subsections to prove Theorem 2.

*Proof of Theorem 2.* For  $0 < \varepsilon < \varepsilon_\nu$ , for  $a \neq a^*$ , let us consider  $t \notin \mathcal{U}_a(\varepsilon)$ ,  $t \geq |\mathcal{A}|$ , such that  $a_{t+1} = a$ . From empirical upper bounds (2.6), we have

$$N_a(t) \text{KL}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) \leq \log(t). \quad (2.23)$$

From Lemma 4 and Algorithm 7, we have  $a \in \mathcal{V}_{a^*}$  and  $\hat{\mu}_a(t) \leq \mu_a + \varepsilon < \mu^* - \varepsilon \leq \hat{\mu}^*(t)$ . From classical monotonic properties of  $\text{KL}(\cdot | \cdot)$  and Equation (2.3), we have  $\text{KL}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) \geq \text{KL}(\mu_a + \varepsilon | \mu^* - \varepsilon) \geq (1 + \alpha_\nu(\varepsilon))^{-1} \text{KL}(\mu_a | \mu^*)$ . In view of Equation (2.23), this implies

$$\forall t \notin \mathcal{U}_a(\varepsilon), t \geq |\mathcal{A}|, \text{ such that } a_{t+1} = a, \quad \begin{cases} a \in \mathcal{V}_{a^*} \\ N_a(t) \leq \frac{(1 + \alpha_\nu(\varepsilon)) \log(t)}{\text{KL}(\mu_a | \mu^*)}. \end{cases} \quad (2.24)$$

For all arm  $a \in \mathcal{A}$ , for all time step  $t \geq |\mathcal{A}|$ , we denote by

$$\tau_a(t) = \max \{t' \in \llbracket |\mathcal{A}|, t \rrbracket : a_{t'+1} = a \text{ and } t' \notin \mathcal{U}_a(\varepsilon)\} \quad (2.25)$$

the last time step before time step  $t$  that does not belong to  $\mathcal{U}_a(\varepsilon)$  such that we pull arm  $a$ .

Then, from Equations (2.24) and (2.25) we have

$$\begin{aligned}
\forall a \neq a^*, \forall t \geq 1, \quad N_a(t) &= N_a(|\mathcal{A}|) + \sum_{t' \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{t'+1}=a\}} \\
&\leq 1 + \sum_{t' \geq 1}^{t-1} \mathbb{I}_{\{a_{t'+1}=a, t' \in \mathcal{U}_a(\varepsilon)\}} + \sum_{t' \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{t'+1}=a, t' \notin \mathcal{U}_a(\varepsilon)\}} \\
&\leq 1 + |\mathcal{U}_a(\varepsilon)| + \sum_{t' \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{t'+1}=a, t' \notin \mathcal{U}_a(\varepsilon)\}} \\
&\leq 1 + |\mathcal{U}_a(\varepsilon)| + \mathbb{I}_{\{a \notin \mathcal{V}_{a^*}\}} \times 0 + \mathbb{I}_{\{a \in \mathcal{V}_{a^*}\}} \times N_a(\tau_a(t)) \\
&\leq 1 + |\mathcal{U}_a(\varepsilon)| + \mathbb{I}_{\{a \in \mathcal{V}_{a^*}\}} \frac{(1 + \alpha_\nu(\varepsilon)) \log(\tau_a(t))}{\text{KL}(\mu_a | \mu^*)} \\
&\leq 1 + |\mathcal{U}_a(\varepsilon)| + \mathbb{I}_{\{a \in \mathcal{V}_{a^*}\}} \frac{(1 + \alpha_\nu(\varepsilon)) \log(t)}{\text{KL}(\mu_a | \mu^*)}.
\end{aligned}$$

This implies

$$\forall a \neq a^*, \forall t \geq 1, \quad N_a(t) \leq \begin{cases} \frac{(1 + \alpha_\nu(\varepsilon)) \log(t)}{\text{KL}(\mu_a | \mu^*)} + |\mathcal{U}_a(\varepsilon)| + 1 & \text{if } a \in \mathcal{V}_{a^*} \\ |\mathcal{U}_a(\varepsilon)| + 1 & \text{if } a \notin \mathcal{V}_{a^*}. \end{cases} \quad (2.26)$$

From Equation (2.19), averaging these inequalities allows us to conclude.  $\square$

## 2.5 Numerical experiments

In this section, we compare empirically the following algorithms : OSUB, UTS (Combes and Proutiere, 2014a; Trinh et al., 2020) and IMED-UB described in Algorithm 7. We illustrate how performs IMED-UB algorithm under Bernoulli, Gaussian (variance  $\sigma^2 = 0.25$ ) or Exponential distribution assumption. For the experiments we consider a graph  $\mathcal{G}$  with maximal degree  $d = 2$  and the unimodal unimodal vectors of means  $\mu = (0.05, 0.10, 0.15, 0.20, 0.25, 0.20, 0.15, 0.10, 0.05)$ , and average regrets over 500 runs for each distribution family. Based on these experiments (Figure 2.1), it seems that IMED-UB competes with OSUB and UTS.

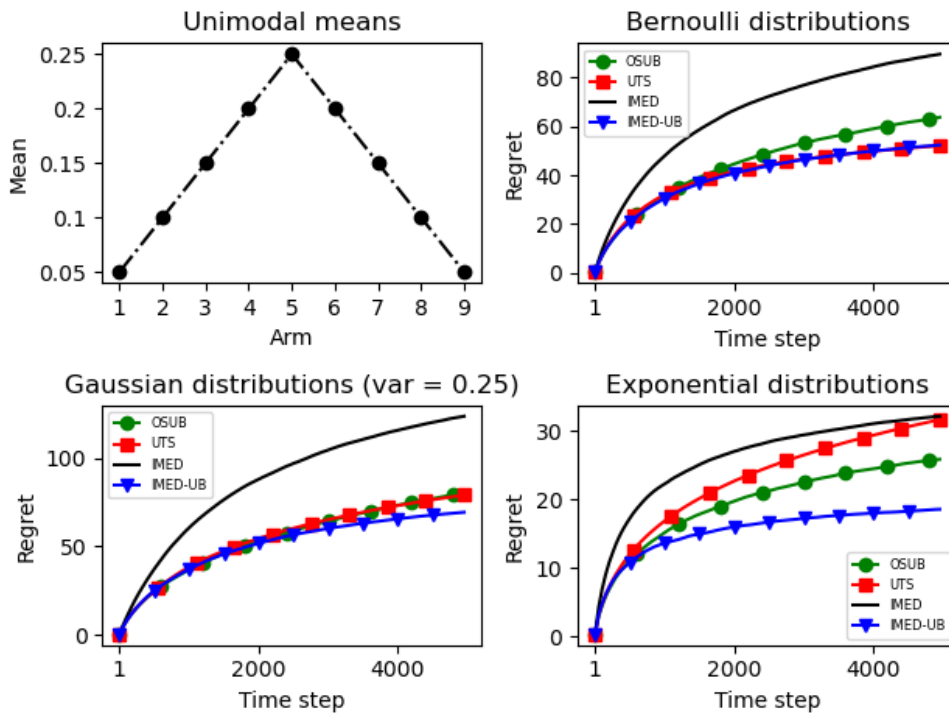


Figure 2.1: Cumulative regrets averaged over 500 runs.

# Chapter 3

## Multimodal Bandits

We consider a multi-armed bandit problem specified by a set of one-dimensional family exponential distributions endowed with a multimodal structure. The multimodal structure naturally extends the unimodal structure studied in Chapter 2. We introduce  $\text{IMED-MB}$ , an algorithm that exploits the multimodal structure, by adapting to this setting the popular Indexed Minimum Empirical Divergence ( $\text{IMED}$ ) algorithm. Numerical experiments show that  $\text{IMED-MB}$  performs well in practice and competes with the state-of-the-art algorithms when assuming unimodal structure.

### 3.1 Introduction

We assume 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 denote by  $\mathcal{V}_a = \{a' \neq a : (a, a') \in E\}$  the neighbours of arm  $a \in \mathcal{A}$  in graph  $G = (\mathcal{A}, E)$  and by  $\mathcal{A}^+(\nu) = \{a \in \mathcal{A} : \forall a' \in \mathcal{V}_a, \mu'_a < \mu_a\}$  the set of arms with locally maximal means. When there is no possible confusion  $\mathcal{A}^+(\nu)$  is simply denoted  $\mathcal{A}^+$ . We assume that  $|\mathcal{A}^+(\nu)|$ , the size of subset  $\mathcal{A}^+(\nu)$ , is equal to  $M$ , the number of local maximums. We assume that  $M$  is known to the learner (Assumption 2). Lastly, we assume that  $\nu \subset \mathcal{P} := \{p(\mu), \mu \in \mathbb{I}\}$ , 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  $\mu \in \mathbb{I} \subset \mathbb{R}$ .  $\mathcal{P}$  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  $\mathcal{D}_{\text{M-modal}}$  (or simply  $\mathcal{D}$  when there is no confusion) the structured set of such unimodal-bandit distributions characterized by  $(\mathcal{P}, G)$ . In the following, we assume that  $\mathcal{P}$  is a set of one-dimensional exponential family distributions. For  $\nu \subset \mathcal{P}$ , we denote by  $\mathcal{A}^*(\nu) = \operatorname{argmax}_{a \in \mathcal{A}} \mu_a$  the set of optimal arms of  $\nu$ . When there is no possible confusion  $\mathcal{A}^*(\nu)$  is simply denoted  $\mathcal{A}^*$ . In particular, we have

$$\mathcal{A}^* \subset \mathcal{A}^+. \quad (3.1)$$

**Assumption 2** (Local maximums). *The number  $M = |\mathcal{A}^+|$  of arms with locally maximal means is known to the learner.*

**Assumption 3** (Unique optimal arm). *We assume there exists  $a^* \in \mathcal{A}$  such that  $\mathcal{A}^* = \{a^*\}$ .*

**Unimodal Structure.** When  $\mathcal{A}^+ = \mathcal{A}^* = \{a^*\}$ , the introduced *multimodal* coincides with the *unimodal* structure (Chapter 2) that has been first considered in [Yu and Mannor \(2011\)](#) from a bandit perspective. The study of unimodal structure naturally appears in many contexts, e.g. single-peak preference economics, voting theory or wireless communications, and [Combes and Proutiere \(2014a\)](#); [Trinh et al. \(2020\)](#) provide an explicit lower bound and optimal corresponding algorithms, respectively  $\text{OSUB}$ ,  $\text{UTS}$ .

**Notations.** Let  $\nu \in \mathcal{D}$ . Let  $\mu^* = \max_{a \in \mathcal{A}} \mu_a$  be the optimal mean. 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 tower rule we can rewrite the regret as follows:

$$R(\nu, T) = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}_\nu [N_a(T)], \quad (3.2)$$

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

## 3.2 Regret lower bound

In order to obtain non trivial lower bound on the regret we consider algorithms that are *consistent* (Definition 1). We can derive from the notion of consistency an asymptotic lower bound on the regret.

**Proposition 2** (Lower bounds on the regret). *Let us consider a consistent algorithm. Let us consider a configuration  $\nu \in \mathcal{D}_{M\text{-modal}}$  with means  $\mu \in \mathbf{I}^A$  such that for all  $a_1^+, a_2^+ \in \mathcal{A}^+(\nu)$ ,*

$$a_1^+ \neq a_2^+ \quad \Rightarrow \quad \mathcal{V}_{a_1^+} \cap \mathcal{V}_{a_2^+} = \emptyset. \quad (3.3)$$

Then it must be that

$$\liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \geq \mathfrak{C}_{\text{multi}}(\mu) := \sum_{a^+ \in \mathcal{A}^+} \sum_{\substack{a \in \{a^+\} \cup \mathcal{V}_{a^+} \\ \mu_a \neq \mu^*}} \frac{\Delta_a}{\text{KL}(\mu_a | \mu^*)},$$

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  $\mu, \mu' \in \mathbf{I}$ , and where  $\mathcal{V}_a$  is the neighbourhood of arm  $a \in \mathcal{A}$ .

A proof of Proposition 2 is provided after the discussion below.

**Remark 5.**  $\mathfrak{C}_{\mathcal{D}_{M\text{-modal}}}(\mu)$  is simply denoted  $\mathfrak{C}_{\text{multi}}(\mu)$ .

**Remark 6.** The quantity  $\mathfrak{C}_{\text{multi}}(\mu)$  is a fully explicit function of  $\mu$  (it does not require solving any optimization problem) for some set of distributions  $\nu$  (see Remark 1). This useful property no longer holds in general for arbitrary structures. Also, it is noticeable that  $\mathfrak{C}_{\text{multi}}(\mu)$  does not involve all the sub-optimal arms but only the ones in  $\cup_{a^+ \in \mathcal{A}^+} \{a^+\} \cup \mathcal{V}_{a^+}$ . This indicates that sub-optimal arms outside  $\cup_{a^+ \in \mathcal{A}^+} \{a^+\} \cup \mathcal{V}_{a^+}$  are sampled  $o(\log(T))$ , which contrasts with the unstructured stochastic multi-armed bandits.

**Discussion.** In Proposition 2, the conditions on  $\mathcal{A}^+(\nu)$  detailed in Equation 3.3 ensure that for all local maximum  $a^+ \in \mathcal{A}^+$ , for all sub-optimal  $a \in \mathcal{V}_{a^+}$  in its neighbourhood, all “most confusing”<sup>1</sup> configuration  $\nu^{(a)}$  (defined in Equation (3.4)) still remains in  $\mathcal{D}_{M\text{-modal}}$ . This allows us to use the consistency on  $\mathcal{D}_{M\text{-modal}}$  of the considered algorithm. Without the conditions detailed in Equation (3.3), a most confusing configuration  $\nu^{(a)}$  could have only  $|\mathcal{A}^+(\nu)| - 1 = M - 1$  local maximums. However, these conditions are not very restrictive and only imply that two local maximums cannot follow each other in graph  $G$ . For clarity, let us illustrate in Figure 3.1 the situation in dimension one when the initial configuration has 2-modal means but may have a most confusing configuration with unimodal means.

<sup>1</sup>This notion of “most confusing” refers to the generic proof technique used to derive regret lower bounds. It involves a change-of-measure argument, from the initial configuration in which the arm is sub-optimal to another one chosen to make it optimal.



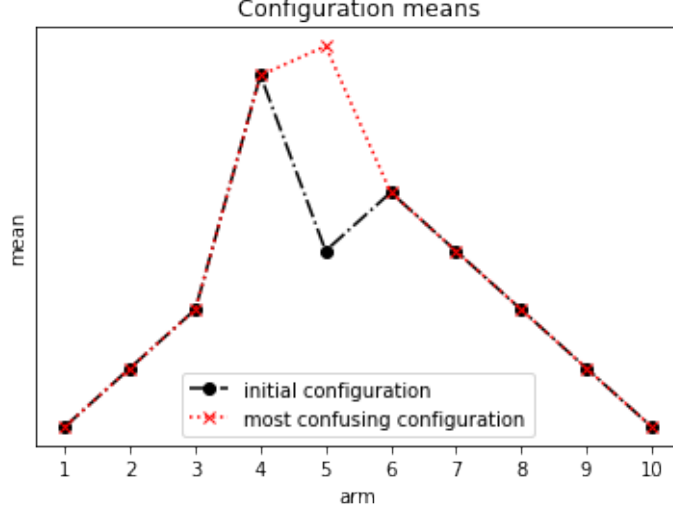


Figure 3.1: Configuration with 1-dimensional 2-modal means admitting unimodal most confusing configuration. Here,  $\mathcal{A} = \llbracket 1, 10 \rrbracket$  and  $\mathcal{V}_a = \{a-1, a+1\} \cap \mathcal{A}$  for all  $a \in \mathcal{A}$  so we have  $\mathcal{A}^+ = \{4, 6\}$  and  $\mathcal{V}_4 \cap \mathcal{V}_6 = \{5\} \neq \emptyset$ .

*Proof.* We consider a configuration  $\nu \in \mathcal{D}$  satisfying Equation (3.3), that is such that for all  $a_1^+, a_2^+ \in \mathcal{A}^+(\nu)$ ,

$$a_1^+ \neq a_2^+ \quad \Rightarrow \quad \mathcal{V}_{a_1^+} \cap \mathcal{V}_{a_2^+} = \emptyset,$$

where  $\mathcal{V}_a$  is the neighbourhood of arm  $a \in \mathcal{A}$ .

Let us consider a sub-optimal  $a \in \cup_{a^+ \in \mathcal{A}^+} \{a^+\} \cup \mathcal{V}_{a^+} \setminus \mathcal{A}^+(\nu)$ . The proof consists in using Lemma 5 below from Garivier et al. (2016) with configuration  $\nu$  and the most confusing configuration  $\nu^{(a)}(\varepsilon)$  for  $\varepsilon > 0$ , with means  $\mu^{(a)}(\varepsilon)$ , where

$$\forall a' \in \mathcal{A}, \quad \mu_{a'}^{(a)}(\varepsilon) = \begin{cases} \mu_{a'} & \text{if } a' \neq a \\ \mu^* + \varepsilon & \text{if } a' = a. \end{cases} \quad (3.4)$$

Note that the set of optimal arms for the most confusing configuration  $\nu^{(a)}$  reduces to the singleton  $\mathcal{A}^*(\nu^{(a)}) = \{a\}$  and that, due to Equation (3.3), the most confusing configuration  $\nu^{(a)}(\varepsilon)$  still belongs to  $\mathcal{D}_{\text{M-modal}}$ , that is  $\mu^{(a)}(\varepsilon)$  has exactly  $|\mathcal{A}^+(\nu)| = M$  local maximums.

Let us consider the random variable  $Z_T = N_a(T)/T \in [0, 1]$ . Then previous Lemma 5 implies

$$\sum_{a' \in \mathcal{A}} \mathbb{E}_\nu[N_{a'}(T)] \text{KL}(\mu_{a'} | \mu_{a'}^{(a)}(\varepsilon)) \geq \text{kl}(\mathbb{E}_\nu[Z_T] | \mathbb{E}_{\nu^{(a)}(\varepsilon)}[Z_T]). \quad (3.5)$$

Since for all  $a' \neq a$  we have the equality of means  $\mu_{a'} = \mu_{a'}^{(a)}(\varepsilon)$  and since  $\mu_a^{(a)}(\varepsilon) = \mu^* + \varepsilon$ , previous Equation (3.5) rewrites

$$\mathbb{E}_\nu[N_a(T)] \text{KL}(\mu_a | \mu^* + \varepsilon) \geq \text{kl}(\mathbb{E}_\nu[Z_T] | \mathbb{E}_{\nu^{(a)}(\varepsilon)}[Z_T]). \quad (3.6)$$

From there, what remains of the proof is classic. For instance, the reader can refer to the proof of Theorem 1 in Garivier et al. (2016).

Since we consider a consistent algorithm on  $\mathcal{D}_{\text{M-modal}}$  and  $\begin{cases} \nu \in \mathcal{D}_{\text{M-modal}} \\ a \notin \mathcal{A}^*(\nu) \end{cases}$ , the averaged number of pulls of arm  $a$  for configuration  $\nu$  is sub-linear and

$$\lim_{T \rightarrow \infty} \mathbb{E}_\nu[Z_T] = \lim_{T \rightarrow 0} \mathbb{E}_\nu[N_a(T)]/T = 0. \quad (3.7)$$

Since we consider a consistent algorithm on  $\mathcal{D}_{\text{M-modal}}$  and  $\begin{cases} \nu^{(a)} \in \mathcal{D}_{\text{M-modal}} \\ \{a\} = \mathcal{A}^*(\nu^{(a)}(\varepsilon)) \end{cases}$ , the averaged number of pulls of arm  $a$  for configuration  $\nu^{(a)}$  is linear and

$$\lim_{T \rightarrow \infty} \mathbb{E}_{\nu^{(a)}(\varepsilon)}[Z_T] = \lim_{T \rightarrow 0} \mathbb{E}_{\nu^{(a)}(\varepsilon)}[N_a(T)]/T = 1. \quad (3.8)$$

By combining Equation (3.7) and (3.8), we have in particular when  $T$  tends to  $\infty$  that

$$\text{kl}(\mathbb{E}_\nu[Z_T] | \mathbb{E}_{\nu^{(a)}(\varepsilon)}[Z_T]) \underset{T \rightarrow \infty}{\sim} \log \left( \frac{1}{1 - \mathbb{E}_{\nu^{(a)}(\varepsilon)}[Z_T]} \right). \quad (3.9)$$

Note that the right term of the last equation can be rewritten as follows,

$$\log \left( \frac{1}{1 - \mathbb{E}_{\nu^{(a)}(\varepsilon)}[Z_T]} \right) = \log \left( \frac{T}{\sum_{a' \notin \mathcal{A}^*(\nu^{(a)}(\varepsilon))} \mathbb{E}_{\nu^{(a)}(\varepsilon)}[N_{a'}(T)]} \right) = \log \left( \frac{T}{O(T^\alpha)} \right), \quad \forall \alpha > 0. \quad (3.10)$$

In particular, by combining previous Equation (3.10) and Equation (3.9) we get the following asymptotic result,

$$\lim_{T \rightarrow \infty} \frac{\text{kl}(\mathbb{E}_\nu[Z_T] | \mathbb{E}_{\nu^{(a)}(\varepsilon)}[Z_T])}{\log(T)} = 1. \quad (3.11)$$

We prove Proposition 2 by combining this last Equation (3.11) with Equation (3.6).  $\square$

**Lemma 5** (Fundamental inequality). *Let us consider a consistent algorithm on  $\mathcal{D}$ . Then for all configurations  $\nu, \nu' \in \mathcal{D}$  with means  $\mu, \mu' \in \mathbb{I}^A$ , for all horizon  $T \geq 1$ , for random variable  $Z_T$  with values in  $[0, 1]$ ,*

$$\sum_{a \in \mathcal{A}} \mathbb{E}_\nu[N_a(T)] \text{KL}(\mu_a | \mu'_a) \geq \text{kl}(\mathbb{E}_\nu[Z_T] | \mathbb{E}_{\nu'}[Z_T]),$$

where  $\text{kl}(p|q) = p \log(\frac{p}{q}) + (1-p) \log(\frac{1-p}{1-q})$  for  $p, q \in [0, 1]$ .

### 3.3 Numerically efficient algorithm for multimodal bandits

We start this section by considering some useful notations before introducing and defining IMED-MB algorithm.

#### 3.3.1 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. We denote by  $\hat{a}_t^*$  an arm arbitrarily chosen in  $\hat{\mathcal{A}}^*(t)$ . We denote by  $\bar{\mathcal{A}}(t) := \{a \in \mathcal{A} : \forall a' \in \mathcal{V}_a, \hat{\mu}_{a'}(t) \leq \hat{\mu}_a(t)\}$  the set of arms with locally maximal empirical means and by

$\widehat{\mathcal{A}}^+(t) \subset \overline{\mathcal{A}}(t)$  the set of no more than  $|\mathcal{A}^+|$  arms with locally maximal empirical means (ties are broken arbitrarily),

$$\begin{cases} \text{if } |\overline{\mathcal{A}}(t)| < |\mathcal{A}^+| : & \widehat{\mathcal{A}}^+(t) = \overline{\mathcal{A}}(t) \\ \text{if } |\overline{\mathcal{A}}(t)| \geq |\mathcal{A}^+| : & \widehat{\mathcal{A}}^+(t) \in \underset{a_1, \dots, a_{|\mathcal{A}^+|} \in \overline{\mathcal{A}}(t)}{\operatorname{argmax}} \sum_{a \in \{a_1, \dots, a_{|\mathcal{A}^+|}\}} \widehat{\mu}_a(t). \end{cases}$$

Note that this choice of  $\widehat{\mathcal{A}}^+(t)$  gives priority to local maximums with greater current empirical means.

### 3.3.2 IMED-MB algorithm.

For all arm  $a \in \mathcal{A}$  and time step  $t \geq 1$  we introduce the IMED index from [Honda and Takemura \(2015\)](#),

$$I_a(t) = N_a(t) \operatorname{KL}(\widehat{\mu}_a(t) | \widehat{\mu}^*(t)) + \log(N_a(t)), \quad (3.12)$$

with the convention  $0 \times \infty = 0$  and  $\log(0) = -\infty$ , and denote by  $\bar{a}_t$  the arm with minimal IMED index

$$\bar{a}_t \in \underset{a \in \mathcal{A}}{\operatorname{argmin}} I_a(t) \quad (\text{arbitrarily chosen}). \quad (3.13)$$

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 \(2015\)](#),  $I_a(t)$  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. If multimodal structure is not considered, arm  $\bar{a}_t$  may be seen as the current most informative arm. However, regarding the lower bound for multimodal structure ([Proposition 2](#)), the current most informative arm may be

$$\bar{\bar{a}}_t \in \underset{\substack{a \in \widehat{\mathcal{A}}^+(t) \\ \bar{\bar{a}}_t \in \widehat{\mathcal{A}}^+(t)}}{\operatorname{argmin}} I_a(t) \quad (\text{arbitrarily chosen}). \quad (3.14)$$

When  $|\widehat{\mathcal{A}}^+(t)| < |\mathcal{A}^+|$ , obviously the underlying multimodal structure is poorly estimated and arm  $\bar{a}_t$  is used to deal with the trade-off exploitation versus exploration: if  $\bar{a}_t$  coincides with the current best arm, that is  $\bar{a}_t = \widehat{a}_t^*$ , we exploit and then pull best arm  $\widehat{a}_t^*$ , otherwise we explore in order to better estimate the underlying multimodal structure and then pull the arm with the minimal number of pull in exploration phases

$$\underline{\bar{a}}_t \in \underset{a \in \mathcal{A}}{\operatorname{argmin}} N_a(t) \quad (\text{arbitrarily chosen}). \quad (3.15)$$

When  $|\widehat{\mathcal{A}}^+(t)| = |\mathcal{A}^+|$ , arm  $\bar{\bar{a}}_t$  is now used to deal with the trade-off exploitation versus exploration. We do not directly pull  $\bar{\bar{a}}_t$  during exploration phases but consider second order IMED type indexes ([Eq. 3.17](#)) to ensure better estimations of the arms with locally maximal empirical means, that is the arms from set  $\widehat{\mathcal{A}}^+(t)$ . We explore in order to improve confidence in arms with locally maximal empirical means and then pull an arm with minimal second order index ([Eq. 3.18](#)). Let us consider

$$\widehat{\bar{a}}_t^+ \in \underset{a \in \{\bar{\bar{a}}_t\} \cup \mathcal{V}_{\bar{\bar{a}}_t}}{\operatorname{argmax}} \widehat{\mu}_a(t) \quad (\text{arbitrarily chosen}). \quad (3.16)$$

In particular, we have  $\widehat{a}_t^+ \in \widehat{\mathcal{A}}^+(t)$ . We then consider the second order index for arm  $a \in \mathcal{A}$ ,

$$\underline{I}_a(t) = \begin{cases} \text{if } \widehat{\mu}_a(t) < \widehat{\mu}_{\widehat{a}_t^+}(t) : \\ \quad N_a(t) \text{KL}\left(\widehat{\mu}_a(t) \middle| \widehat{\mu}_{\widehat{a}_t^+}(t)\right) + \log(N_a(t)) \\ \text{if } \widehat{\mu}_a(t) \geq \widehat{\mu}_{\widehat{a}_t^+}(t) : \\ \quad \log(N_a(t)) \end{cases} \quad (3.17)$$

and

$$\underline{a}_t \in \underset{a \in \{\widehat{a}_t^+\} \cup \mathcal{V}_{\widehat{a}_t^+}}{\text{argmin}} \underline{I}_a(t) \quad (\text{arbitrarily chosen}). \quad (3.18)$$

We speak of second order exploration because of the following inequalities on the indexes. At each time step  $t \geq 1$ , we have

$$\underline{I}_{\underline{a}_t}(t) \leq \log\left(N_{\widehat{a}_t^+}(t)\right) \leq I_{\widehat{a}_t^+}(t). \quad (3.19)$$

Note that when  $\widehat{a}_t^+ = \widehat{a}_t^*$ , then both IMED type indexes coincide, that is  $\underline{I}_{\underline{a}_t}(t) = I_{\underline{a}_t}(t)$ . However, when  $\widehat{a}_t^+ \neq \widehat{a}_t^*$ , we expect arm  $\widehat{a}_t^+$  to be a sub-optimal one and so to be pulled a logarithmic number of times. In that particular case and according to previous Equation (3.19), we then expect  $\underline{I}_{\underline{a}_t}(t)$  to be a  $O(\log \log t)$ .

IMED-MB algorithm is summarized in Algorithm 8.

---

**Algorithm 8** IMED-MB

---

```

for  $t = 1 \dots T - 1$  do
  if  $|\widehat{\mathcal{A}}^+(t)| < |\mathcal{A}^+|$  then                                 $\triangleright \triangleright \triangleright$  MULTIMODAL STRUCTURE BADLY ESTIMATED
    if  $\bar{a}_t = \widehat{a}_t^*$  then                                        $\triangleright$  Exploitation
      Pull arm  $a_{t+1} = \bar{a}_t$  (Eq. 3.13)
    else                                                          $\triangleright$  Exploration
      Pull arm  $a_{t+1} = \underline{a}_t$  (Eq. 3.15)
    end if
  else                                                          $\triangleright \triangleright \triangleright$  MULTIMODAL STRUCTURE WELL ESTIMATED
    if  $\underline{a}_t = \widehat{a}_t^+$  then                                        $\triangleright$  Reliable local maximum  $\widehat{a}_t^+$ : exploration-exploitation
      Pull arm  $a_{t+1} = \bar{a}_t$ 
    else                                                          $\triangleright \triangleright$  Non-reliable local maximum  $\widehat{a}_t^+$ : second order exploration
      Pull arm  $a_{t+1} = \underline{a}_t$  (Eq. 3.18)
    end if
  end if
end for

```

---

### 3.4 Numerical experiments

In this section, we compare empirically the following algorithms : OSUB, UTS (Combes and Proutiere, 2014; Trinh et al., 2020) and IMED-MB described in Algorithm 8. We illustrate how perform IMED-MB algorithm under Bernoulli, Gaussian (variance  $\sigma^2 = 0.25$ ) or Exponential distribution assumption. For the experiments we consider a graph  $\mathcal{G}$  with maximal degree  $d = 8$  and the multimodal vectors of means is of dimension 2, and average regrets over 100 runs for each distribution family. Additional details are provided in the paragraph below. Based on these experiments (Figures 3.2, 3.3), it seems that IMED-MB competes with OSUB and UTS for the unimodal structure.

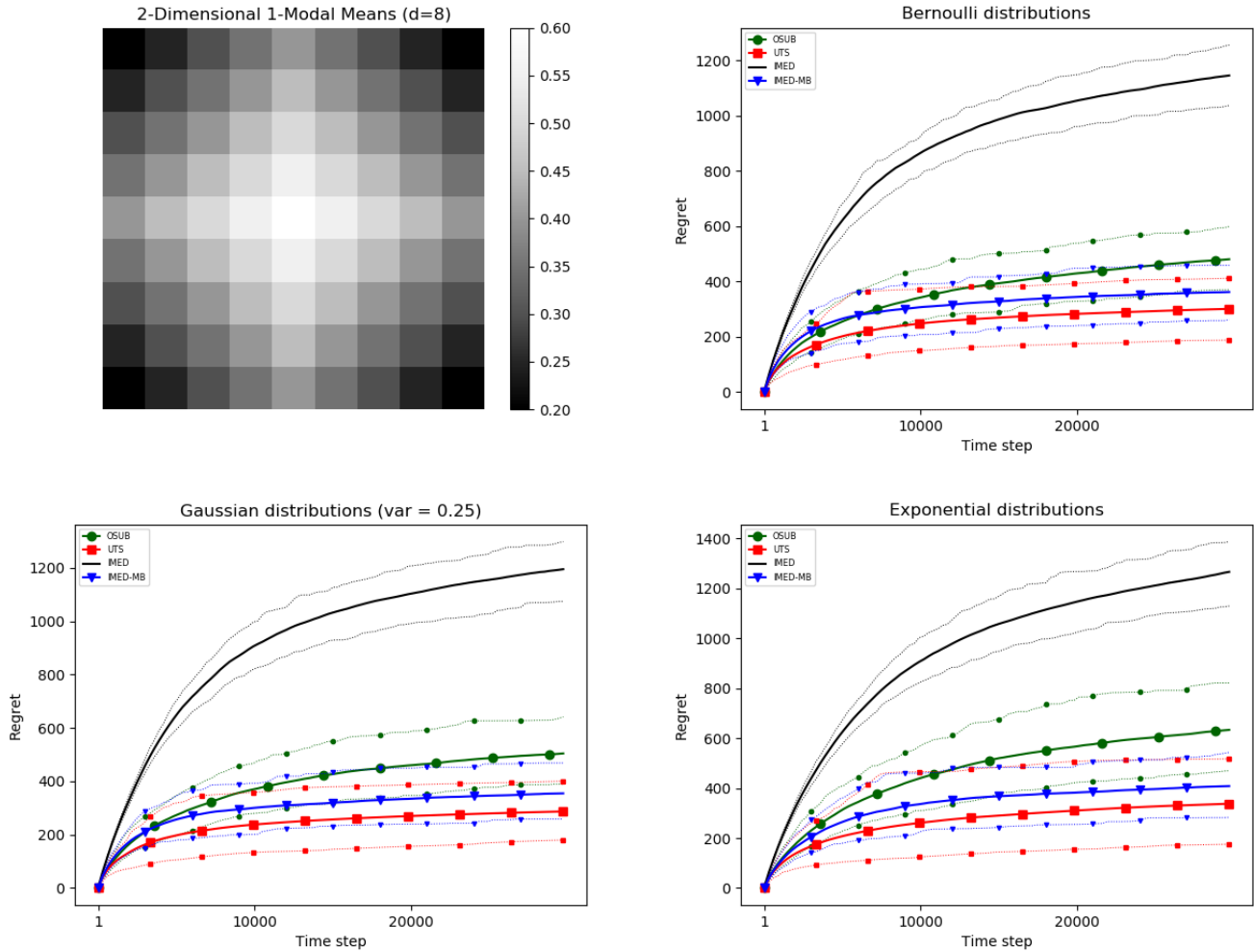


Figure 3.2: Cumulative regrets averaged over 100 runs for  $\nu$  with 2-dimensional 1-modal means.

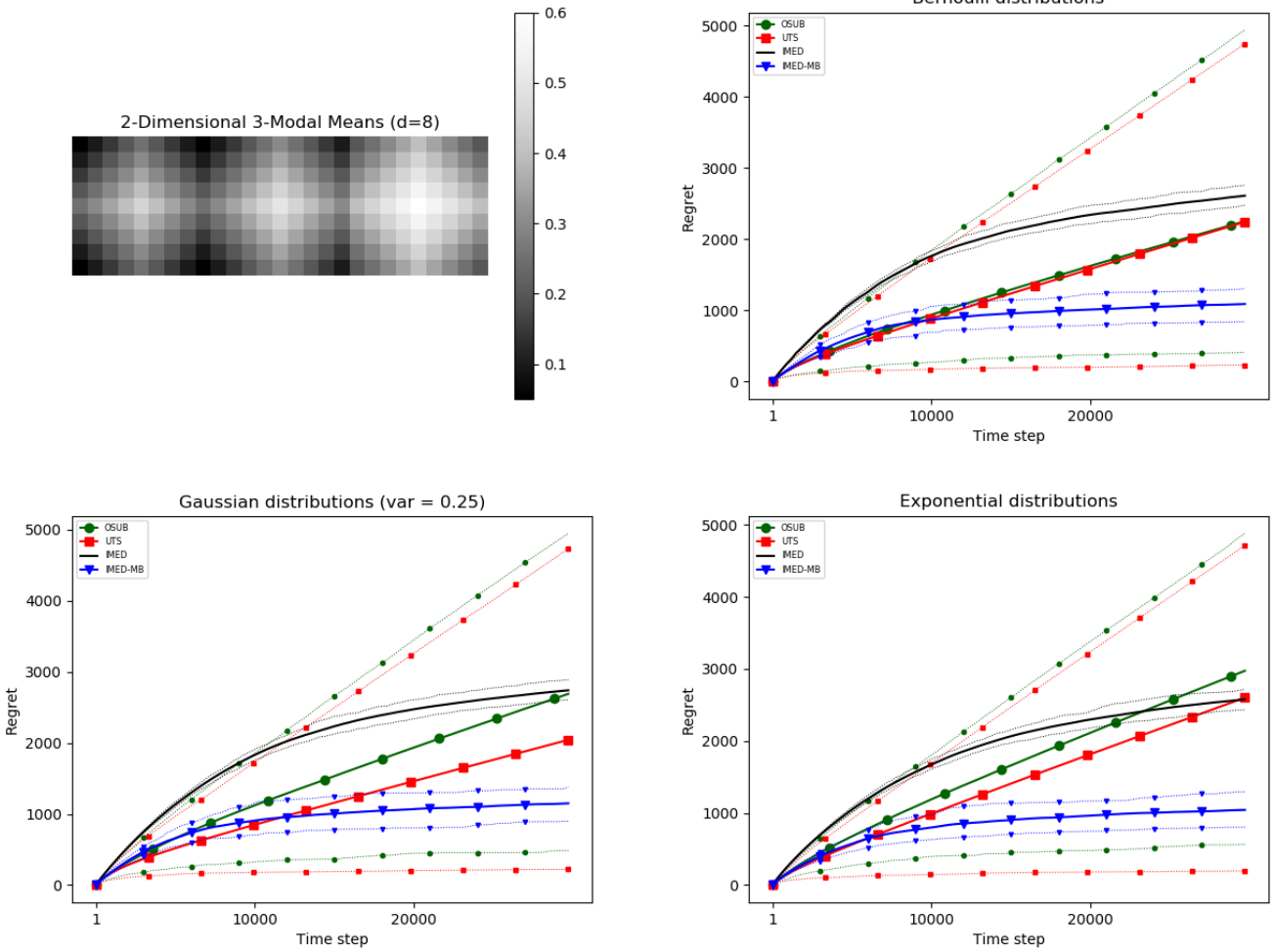


Figure 3.3: Cumulative regrets averaged over 100 runs for  $\nu$  with 2-dimensional 3-modal means.

**Additional details.** In the following we provide additional details about experiments summarized in Figures 3.2 and 3.3.

The parameter of OSUB algorithm is set equal to  $\gamma_{\text{OSUB}} = d$ , where  $d = \max_{a \in \mathcal{A}} |\mathcal{V}_a|$  is the maximal degree of nodes. The parameter of UTS algorithm is set equal to  $\gamma_{\text{UTS}} = 2$ . These parameters are those recommended in Combes and Proutiere (2014a) and Trinh et al. (2020).

For all the algorithms, we do not consider initial phase consisting in pulling each arm once. Indeed, we are interested in algorithm that do not necessarily explore all the arms. This is motivated by practical considerations: since the set of arms  $\mathcal{A}$  can be large, we focus on algorithms that a priori do not explore all the arms. We show in Figures 3.4 and 3.5 that IMED-MB indeed concentrate the pulls around the arms with local maximal means. In both figures we represent the averaged numbers of pulls of sub-optimal arms under IMED-MB algorithm used to plot the averaged regrets in Figures 3.2 and Figures 3.3. Note that for readability purpose, we removed the pulls of optimal arms in both Figures.

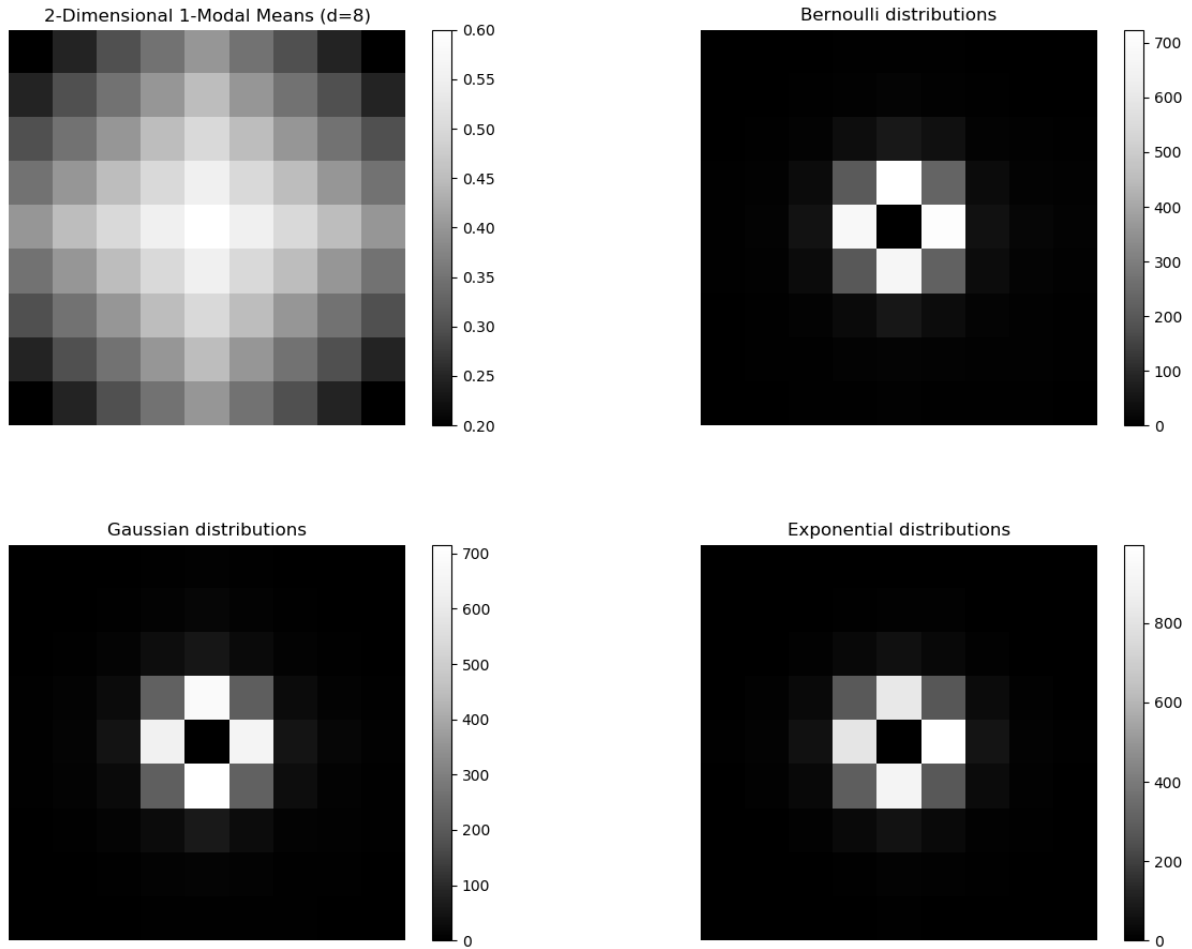


Figure 3.4: Number of pulls of sub-optimal arms under IMED-MB used in the unimodal case of Figure 3.2.

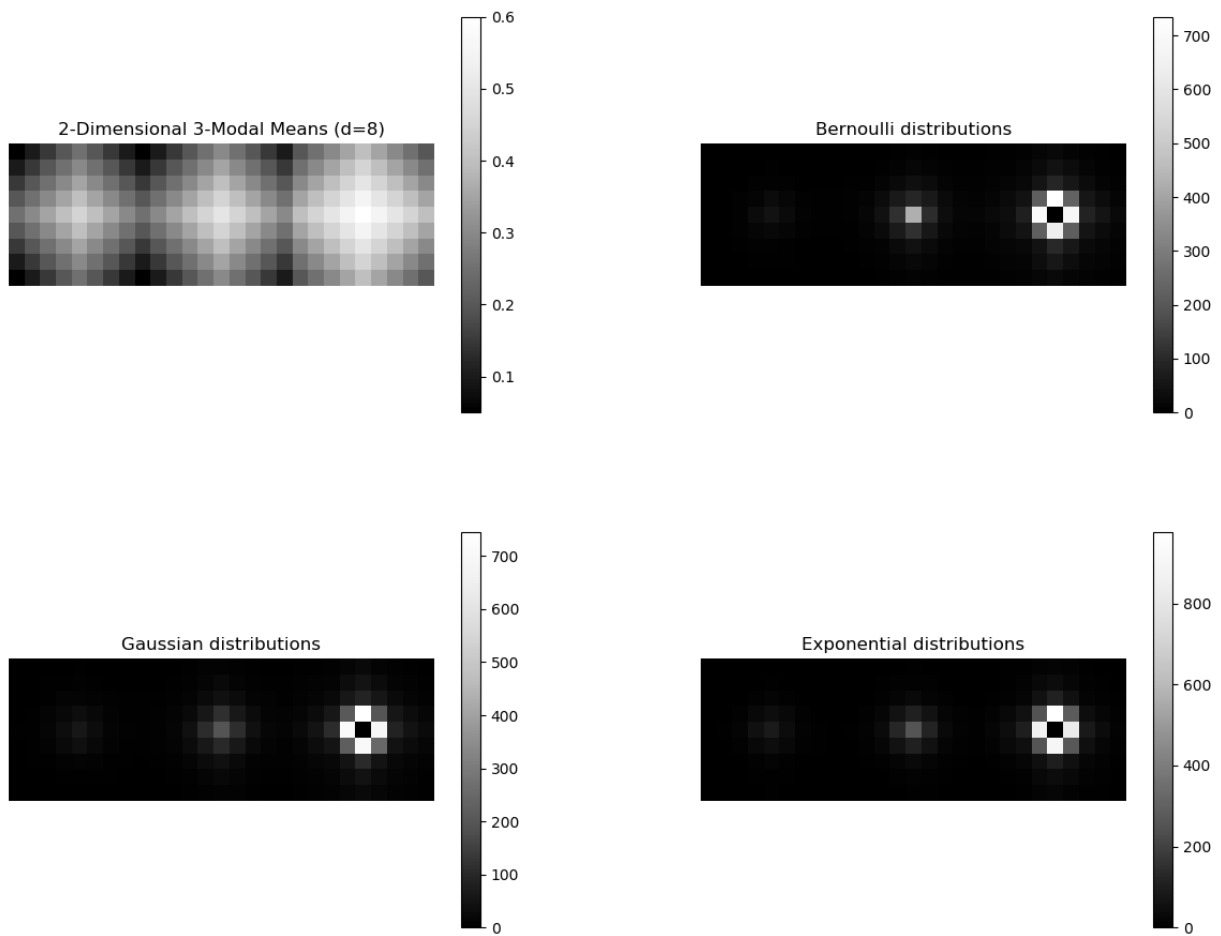


Figure 3.5: Numbers of pulls of sub-optimal arms under IMED-MB used in the 3-modal case of Figure 3.3.



# Chapter 4

## Graph-Structured Bandits

We consider a structured variant of the multi-armed bandit problem when the difference of means between any pair of arms is constrained not to exceed some value. This graph structure is introduced to encompass as special cases the classical structures Unimodal and Lipschitz. We derive the asymptotic lower bound on the cumulative regret for this structure, and introduce `IMED-GS` an extension of the popular Indexed Minimum Empirical Divergence (`IMED`) algorithm to such structured configurations. In order to show asymptotic optimality in a graph-structured scenario, we further add a tracking step to the `IMED` approach whose aim is to ensure the sub-optimal arms are played with correct asymptotic frequencies. Interestingly, this tracking step that requires solving an optimization problem is only triggered when the `IMED` index suggests exploration, which provably happens no more than  $O(\log(T))$  times within  $T$  rounds. We further carefully handle the rounds when the structure cannot be exploited due to large uncertainty (e.g. initial rounds), by combining structured and unstructured versions of `IMED`. Our analysis enables an explicit finite-time regret bound emphasizing the role of the second-order terms. Last, we illustrate the benefit of `IMED-GS` over alternative structured bandit algorithms on numerical experiments.

### 4.1 Introduction

For a given *closed set of relationship matrices*  $\Theta \subset [-1, 1]^{\mathcal{A}^2}$  we introduce

$$\mathcal{D}_\Theta := \left\{ \nu \in \mathcal{B} : \exists \theta \in \Theta, \forall a, a' \in \mathcal{A}, \mu_a - \mu_{a'} \leq \theta_{a,a'} \right\}, \quad (4.1)$$

where  $\mathcal{B}$  is the set of Bernoulli distributions with means in  $(0, 1)$ . We call this a *graph structure*, as it constraints pairs of means. We further impose a pseudo-metric property on each matrix  $\theta \in \Theta$ .

**Remark 7.**  $\mathcal{D}_\Theta$  may be denoted  $\mathcal{D}$  when there is no possible confusion.

**Assumption 4 (Pseudo-metric).** For all  $\theta \in \Theta$ , for all arms  $a, a', a'' \in \mathcal{A}$ , one has  $\theta_{a,a} = 0$  (definiteness) and  $\theta_{a,a''} \leq \theta_{a,a'} + \theta_{a',a''}$  (triangular inequality).

This formulation is flexible enough to capture as special cases popular structures such as Unimodal, Lipschitz or Aggregate of bandits (see details in Section A.1), which makes it appealing to study. Note that it can also naturally interpolates between a fully structured and fully unstructured case.

**Remark 8.** Assumption 4 does not require  $\theta_{a,a'}$  to be non-negative for  $a, a' \in \mathcal{A}$ .

**Remark 9.** The notion of Lipschitz bandit we refer to is the one used in Magureanu et al. (2014). When  $\Theta = \{\theta\}$  is a singleton and  $\theta$  is a symmetric relationship matrix, that is  $\theta_{a,a'} = \theta_{a',a}$  for  $a, a' \in \mathcal{A}$ , we then recover a generic notion of Lipschitz structure as a particular case of graph structure. However, the Unimodal structure, for example, cannot be described with Lipschitz properties but still remains a particular case of graph structure (see details in Section A.1).

**Goal.** First, we target the design of a single algorithm able to achieve, for any given graph structure, instance-dependent asymptotic optimality against any instance. We also want this algorithm to be computationally efficient, in the sense it avoids computing (a version of)  $\mathfrak{C}_{\mathcal{D}}(\mu)$  when estimations are "too noisy", and only solves this costly problem provably scarcely. Second, we want to obtain finite-time analysis with explicit terms using the same generic proof technique for each structure. Last, although our algorithm applies to other structures, we want it to be competitive against both structure-dependent (like e.g. OSUB, CKL-UCB respectively for Unimodal and Lipschitz) and generic (like OSSB) state-of-the-art algorithms in practice on classical structures. We now address these challenges.

## 4.2 Regret lower bound

In this section, we specify the regret lower bound when assuming a graph structure. To this end, we follow the classical approach from [Lai and Robbins \(1985\)](#); [Graves and Lai \(1997\)](#). In order to obtain non trivial lower bound we consider algorithms that are consistent on  $\mathcal{D}_{\Theta}$  (Definition 1). Using this notion, asymptotic regret lower bounds are proved using change-of-measure argument by considering *most confusing* bandit configurations. We identify them first introducing, for each sub-optimal arm  $a \notin \mathcal{A}^*(\nu) = \operatorname{argmax}_{a \in \mathcal{A}} \mu_a$ , the following distribution-dependent subset of relationship matrices

$$\Theta_a(\nu) := \left\{ \theta \in \Theta^*(\nu) : \begin{array}{l} (1) \forall a^* \in \mathcal{A}^*(\nu), \theta_{a,a^*} > 0 \\ (2) \forall a' \notin \mathcal{A}^*(\nu), \theta_{a,a'} \geq 0 \end{array} \right\}, \quad (4.2)$$

where  $\Theta^*(\nu) = \{ \theta \in \Theta : \nu \in \mathcal{D}_{\{\theta\}} \}$ , and then the corresponding set of *informative* sub-optimal arms

$$\mathcal{A}_a(\nu, \theta) := \{ a' \in \mathcal{A} : \mu_{a'} \leq \mu^* - \theta_{a,a'} \}, \quad \theta \in \Theta_a(\nu). \quad (4.3)$$

When "moving" sub-optimal arm  $a$  to make it optimal in a most confusing configuration,  $\mathcal{A}_a(\nu, \theta)$  represents the set of sub-optimal arms which must also be "moved" to ensure the "most confusing"<sup>1</sup> bandit for sub-optimal  $a$  belongs to the structure  $\mathcal{D}_{\theta}$ . Establishing lower bounds is also a key component towards building an efficient algorithm. To avoid technical issues, we make the mild assumption that there exists most confusing configurations in  $\mathcal{D}_{\Theta}$  for each sub-optimal arm.

**Assumption 5.** Each sub-optimal arm  $a \notin \mathcal{A}^*(\nu)$  admits a most confusing instance, that is  $\Theta_a(\nu) \neq \emptyset$ .

**Proposition 3** (Lower bound on the regret). *Let us consider a consistent algorithm. Then, for all configuration  $\nu \in \mathcal{D}_{\Theta}$  with means  $\mu = (\mu_a)_{a \in \mathcal{A}}$ , under Assumptions 4 and 5 it must be that*

$$\begin{aligned} \liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} &\geq \mathfrak{C}_{\Theta}(\mu) := \min_{\eta \in \mathbb{R}_+^{\mathcal{A}}} \sum_{a \in \mathcal{A}} \eta_a (\max(\mu) - \mu_a) \\ \text{s.t. } &\forall a \notin \operatorname{argmax}(\mu), \\ &\min_{\theta \in \bar{\Theta}_a(\nu)} \sum_{\substack{a' \in \mathcal{A} \\ \mu_{a'} \leq \max(\mu) - \theta_{a,a'}}} \operatorname{kl}(\mu_{a'} | \max(\mu) - \theta_{a,a'}) \eta_{a'} \geq 1, \end{aligned} \quad (4.4)$$

where  $\bar{\Theta}_a(\nu) = \{ \theta \in \Theta^*(\nu) : \forall a' \in \mathcal{A}, \theta_{a,a'} \geq 0 \}$  denotes the closure of  $\Theta_a(\nu)$  in  $\Theta$ , for all  $a \notin \mathcal{A}^*(\nu)$ .

<sup>1</sup>These notions of "moving" and "most confusing" refer to the generic proof technique used to derive regret lower bounds. It involves a change-of-measure argument, from the initial configuration in which the arm is sub-optimal to another one chosen to make it optimal.

The proof of this result uses a change-of-measure argument and follows classical proof techniques from the literature, see [Lai and Robbins \(1985\)](#); [Agrawal et al. \(1989\)](#); [Graves and Lai \(1997\)](#); [Cappé et al. \(2013\)](#). We detail it in Section [A.2](#) for completeness. Intuitively, the vector  $\eta$  represents an asymptotic number of pulls of arms rescaled by  $\log(T)$ .

**Remark 10.**  $\mathfrak{C}_{\mathcal{D}_\Theta}(\mu)$  is simply denoted  $\mathfrak{C}_\Theta(\mu)$ .

**Remark 11.** *Proposition 3 is a specification to the graph structure of the known lower bounds on the regret for generic structures. Assumption 5 is satisfied by both Unimodal and Lipschitz structures (Section [A.1](#)). By straightforward calculations, we recover for these structures the regret lower bounds respectively established in [Combes and Proutiere \(2014a\)](#) and [Magureanu et al. \(2014\)](#). The regret lower bound formulation from Proposition 3 motivates the shape of the indexes for IMED-GS algorithm (Section [4.3.1](#)).*

### 4.3 Optimal algorithm for graph-structured bandits

In this section, we introduce our main algorithm IMED-GS to handle graph-structured bandits. It is primarily based on an IMED-type approach (see [Honda and Takemura \(2015\)](#)). IMED has been introduced for the case of *unstructured* bandit configurations but is interesting for two reasons. First, it has been shown to be asymptotically optimal, in the sense of matching the asymptotic (unstructured) instance-dependent lower bounds like KLUCB or Thompson-sampling algorithms. Then, its index is directly derived from the analysis of the lower bound, that does not requires an optimization procedure. For this reason, it constitutes an interesting basis in order to build regret efficient algorithms for structured configurations. When extended to the structured case, the IMED approach naturally targets satisfying the constraints of the optimization problem (which is sometimes called Pareto-optimality), but not solving the optimization problem, which offers an interesting advantage for the practitioner. In order to go beyond Pareto-optimality, IMED-GS employs a weak form of (data dependent) forcing mechanisms called tracking, in which a solution to the constrained optimization problem is computed only when Pareto-optimality is already ensured. Also, in initial rounds when the structure is "hidden" by the large uncertainty about the distributions, IMED-GS effectively reduces to the unstructured IMED. These innovations yield provable optimality for any graph-structure at the price of a slight but controlled increase in exploration.

One idea behind IMED-GS algorithm is, following the intuition given by the lower bound, to define structured IMED-type indexes, which requires being able to estimate  $\bar{\Theta}_a(\nu)$  for sub-optimal arms  $a \notin \mathcal{A}^*(\nu)$  and  $\Theta^*(\nu)$ , where  $\nu$  is unknown. The following assumptions provide a convenient setting to do so and capture Unimodal and Lipschitz structures. Assumption 6 implies in particular that either  $\Theta^*(\nu) = \Theta$ , or  $\Theta^*(\nu)$  can be well-estimated when the best arm is well-identified.

**Assumption 6** (Arm-supported structure).  $\Theta$  is either a singleton or the structure is supported by the best arm, that is  $\Theta = \cup_{a \in \mathcal{A}} \{\theta^{(a)}\}$  and for  $a \in \mathcal{A}$ ,  $\nu \in \mathcal{D}_{\{\theta^{(a)}\}} \Leftrightarrow a = a^*$ .

In the following Section [4.3.1](#), we introduce and detail the IMED-GS algorithm. It is summarized in Algorithm [9](#) (all ties are broken arbitrarily). We discuss its optimality properties in Section [4.3.2](#).

---

**Algorithm 9** IMED-GS
 

---

**Input:** Structure  $\Theta$ , sequences  $(\gamma_t)_{t \geq 1}$ , positive constant  $\Gamma$ ,  $\xi$ , integer  $d$ .

Pull each arm once

**for**  $t = |\mathcal{A}| \dots T - 1$  **do**

  Compute structured IMED choice  $\bar{a}_t \in \operatorname{argmin}_{a \in \mathcal{A}} \bar{I}_a^{(d)}(t)$  (Eq. 4.9)

**if**  $\bar{a}_t \in \hat{\mathcal{A}}^*(t)$  **then** ▷ Exploitation

    Pull  $a_{t+1} = \bar{a}_t$

**else** ▷ Exploration

    Compute  $n^{\text{opt}}(t)$  to form  $N^{\text{opt}}(t)$ , and compute  $\bar{\mathcal{A}}_{\bar{a}_t}(t)$ .

    Compute  $a_t^{\text{opt}} \in \operatorname{argmax}_{a \in \bar{\mathcal{A}}_{\bar{a}_t}(t)} N_a^{\text{opt}}(t) - N_a(t)$  and  $\underline{a}_t \in \operatorname{argmin}_{a \in \mathcal{A}} \underline{I}_a(t)$  (Eq. 4.13)

**if**  $\underline{a}_t == a_t^{\text{opt}}$  **then** ▷ Reliable opt.

      Pull  $a_{t+1} = a_t^{\text{opt}}$  (Eq. 4.13)

**else** ▷ Unreliable opt.

**if**  $\bar{I}_{\bar{a}_t}^{(d)}(t) \leq \Gamma \cdot I_{\dot{a}_t}(t)$  **then**

        Pull  $a_{t+1} = \bar{a}_t$  (Eq. 4.9)

**else**

        Pull  $a_{t+1} = \dot{a}_t$  (Eq. 4.6)

**end for**

---

**Unstructured IMED.** For convenience, we introduce for  $\xi \geq 0$ , the function

$$f_\xi : x \geq 1 \mapsto \log(x) + \xi \log(1 \vee \log(x)) \geq 0. \quad (4.5)$$

We note that  $f_\xi$  is an increasing continuous function on  $[1, +\infty[$ . Using this notation, we redefine in this chapter the IMED indexes for unstructured bandits: for each  $a \in \mathcal{A}$ , for  $t \geq 1$ .

(Unstructured IMED)

$$I_a(t) = N_a(t) \operatorname{kl}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) + f_\xi(N_a(t)). \quad (4.6)$$

**Remark 12.** The classical IMED indexes from [Honda and Takemura \(2015\)](#) are defined using  $\xi = 0$ . We introduce this minor extension to simplify the analysis of IMED-GS detailed in Section 4.5. In particular, later defining the structured index  $\bar{I}_a(t)$  in (4.9), the following straightforward property occurs:

$$\forall t \geq |\mathcal{A}|, \forall a \in \mathcal{A}, \quad f_\xi(N_a(t)) \leq I_a(t) \leq \bar{I}_a(t). \quad (4.7)$$

### 4.3.1 IMED-GS algorithm

At time  $t$ , we denote  $\hat{\mu}^*(t) = \max_{a \in \mathcal{A}} \hat{\mu}_a(t)$  the current (empirical) best mean,  $\hat{a}_t^* \in \hat{\mathcal{A}}^*(t) = \operatorname{argmax}_{a \in \mathcal{A}} \hat{\mu}_a(t)$  the current set of optimal arms and finally  $\hat{\theta}^*(t) = \theta$  if  $\Theta = \{\theta\}$ ,  $\theta^{(\hat{a}_t^*)}$  otherwise.

**Remark 13.** The definition of  $\hat{\theta}^*(t)$ , for  $t \geq 1$ , is inspired from Assumption 6, which is satisfied for the Unimodal, Lipschitz structures and the Aggregate of bandits (defined in Section A.1).

We then define the current set of arms informative about arm  $a \notin \hat{\mathcal{A}}^*(t)$  as

$$\hat{\mathcal{A}}_a(t) := \left\{ a' \notin \hat{\mathcal{A}}^*(t) : \hat{\mu}_{a'}(t) \leq \hat{\mu}^*(t) - \hat{\theta}_{a,a'}^*(t) \right\}. \quad (4.8)$$

In the initial rounds, not all informative arms are considered to compute the index. More precisely, (only) when it holds that  $|\widehat{\mathcal{A}}_a(t)| > d+1$  and  $\sum_{a' \in \widehat{\mathcal{A}}_a(t)} 2(\widehat{\mu}^*(t) - \widehat{\mu}_{a'}(t))^2 N_{a'}(t) < \Phi(2|\widehat{\mathcal{A}}_a(t)|+1)$ , for  $a \notin \widehat{\mathcal{A}}^*(t)$ , where

$d \in \{1, \dots, |\mathcal{A}|\}$  is a parameter and  $\Phi(x) = x \log(x)$ , we replace  $\widehat{\mathcal{A}}_a(t)$  with the set  $\widehat{\mathcal{A}}_a^{(d)}(t) \subset \widehat{\mathcal{A}}_a(t)$  consisting of  $\{a\}$  plus the  $d$ -th most pulled arms from  $\widehat{\mathcal{A}}_a(t) \setminus \{a\}$ . It is justified by the fact that in the beginning, no structure can reasonably be exploited due to the poor estimates. We note that  $a \in \widehat{\mathcal{A}}_a(t)$  since  $\widehat{\theta}_{a,a}^*(t) = 0$ , for all current sub-optimal arm  $a \notin \widehat{\mathcal{A}}^*(t)$ .

**Graph-structured index.** We are now ready to detail IMED-GS algorithm. Guided by the lower bound established in Proposition 3 we generalize in (4.9) the IMED index from Honda and Takemura (2011) to take into account the graph structure as follows. For convenience, we consider  $\widehat{\mathcal{A}}_{\widehat{a}^*}(t) = \widehat{\mathcal{A}}_{\widehat{a}^*}^{(d)}(t) = \{\widehat{a}^*\}$  for all current optimal arm  $\widehat{a}^* \in \widehat{\mathcal{A}}^*(t)$ . Then, for all arm  $a \in \mathcal{A}$  and for all time step  $t \geq |\mathcal{A}|$ , we define

$$\begin{aligned} \bar{I}_a(t) := & \sum_{a' \in \widehat{\mathcal{A}}_a(t)} N_{a'}(t) \text{kl}(\widehat{\mu}_{a'}(t) \parallel \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)) \\ & + \log \left( \sum_{a' \in \widehat{\mathcal{A}}_a(t)} N_{a'}(t) \right), \end{aligned} \quad (4.9)$$

and its reduced version  $\bar{I}_a^{(d)}(t)$  defined using  $\widehat{\mathcal{A}}_a^{(d)}(t)$  in lieu of  $\widehat{\mathcal{A}}_a(t)$  (this coincides with  $\bar{I}_a(t)$  when the condition for using  $\widehat{\mathcal{A}}_a^{(d)}(t)$  does not hold). Note that

$$\forall \widehat{a}^* \in \widehat{\mathcal{A}}^*(t), \quad \bar{I}_{\widehat{a}^*}(t) = \bar{I}_{\widehat{a}^*}^{(d)}(t) = f_\xi(N_{\widehat{a}^*}(t)).$$

This generalized index can be seen as a transportation cost for “moving” a sub-optimal arm to an optimal one, plus an exploration term (the term  $f_\xi(\cdot)$ ). When a current optimal arm is considered, the transportation cost is null and only the exploration part remains. Then, the most informative arm at time step  $t \geq |\mathcal{A}|$  is defined as

$$\bar{a}_t \in \operatorname{argmin}_{a \in \mathcal{A}} \bar{I}_a^{(d)}(t). \quad (4.10)$$

Pulling this arm intuitively ensures the constraints in the optimization problem of the regret lower bound are asymptotically satisfied (Pareto optimality). In order to reach optimality, we need to ensure we pull arms in a way that asymptotically matches  $\mathfrak{C}_\Theta(\mu)$ . In order to avoid repeatedly solving such an optimization problem too often, we now detail an innovative approach that may look rather intricate at first sight, but is carefully designed to avoid unnecessary computations. It is based on an interplay between exploration/exploitation phases and reliability tests that we explain precisely below.

**Exploitation.** In case the most informative arm is currently optimal, that is  $\bar{a}_t \in \widehat{\mathcal{A}}^*(t)$ , we exploit, that is the algorithm simply pulls this arm:  $a_{t+1} = \bar{a}_t$ . We show later that this happens asymptotically often, as the other case detailed below is only triggered about  $O(\log(T))$  times out of  $T$  steps.

**Exploration.** In the other cases, we explore by trading off low regret and information gathering. To this end, we introduce the pseudo-counts  $\left( N_a^{\text{opt}}(t) \right)_{a \in \mathcal{A}} = \left( n_a^{\text{opt}}(t) \bar{I}_{\bar{a}_t}(t) \right)_{a \in \mathcal{A}}$ , where  $n^{\text{opt}}(t)$  is a solution of  $\mathfrak{C}_\Theta(\widehat{\mu}(t))$ ,

the empirical version of (4.4), that is

$$\begin{aligned}
n^{\text{opt}}(t) &\in \operatorname{argmin}_{n \in \mathbb{R}_+^{\mathcal{A}}} \sum_{a \notin \widehat{\mathcal{A}}^*(t)} \widehat{\Delta}_a(t) n_a \\
\text{s.t. } &\forall a \notin \widehat{\mathcal{A}}^*(t), \\
&\sum_{a' \in \widehat{\mathcal{A}}_a(t)} \operatorname{kl}\left(\widehat{\mu}_{a'}(t) \middle| \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)\right) n_{a'} \geq 1.
\end{aligned} \tag{4.11}$$

Note that the pseudo-counts rescale these proportions using  $\bar{I}_{\bar{a}_t}(t)$  and not  $\log(t)$ .  $\bar{I}_{\bar{a}_t}(t)$  behaves asymptotically as  $\log(t)$  but is less conservative and more natural in finite-time. We then *track* the current optimal numbers of pulls by computing

$$a_t^{\text{opt}} \in \operatorname{argmax}_{a' \in \widehat{\mathcal{A}}_{\bar{a}_t}(t)} N_{a'}^{\text{opt}}(t) - N_{a'}(t). \tag{4.12}$$

Solving (4.11) is intuitively only useful asymptotically, provided that the set of informative arms are well-estimated. In order to have meaningful current subsets of informative arms when computing current optimization problem (4.11), for all time step  $t \geq |\mathcal{A}|$ , for all current sub-optimal arm  $a' \notin \widehat{\mathcal{A}}^*(t)$  we introduce the index

$$I_{a'}(t) := \begin{cases} f_\xi(N_{a_t^{\text{opt}}}(t)) & \text{if } a' = a_t^{\text{opt}} \\ N_{a'}(t) \cdot 2\gamma_t^2(\widehat{\Delta}_{a'}(t))^2 + f_\xi(N_{a'}(t)) & \text{if } a' \neq a_t^{\text{opt}}, \end{cases} \tag{4.13}$$

where  $\widehat{\Delta}_{a'}(t) = \widehat{\mu}^*(t) - \widehat{\mu}_{a'}(t)$  and  $(\gamma_t)_{t \geq 1}$  is a decreasing sequence such that  $0 < \gamma_t < \gamma_1$  for all  $t \geq 1$ , with  $\gamma_1 := 1/(7 \vee \sqrt{\Phi(2|\mathcal{A}| + 1)})$ , where  $\Phi(x) = x \log(x)$ . We further compute

$$\underline{a}_t \in \operatorname{argmin}_{a' \notin \widehat{\mathcal{A}}^*(t)} I_{a'}(t). \tag{4.14}$$

This trades-off between estimating the graph structure and achieving low regret. However, rather than directly pulling this arm, we use it as a test. Indeed, pulling this arm only makes sense when the structure is already well estimated.

**Reliability test of current optimization problem.** If  $\underline{a}_t = a_t^{\text{opt}}$ , we consider the set of informative arms is well estimated hence we pull the current *tracked* arm  $a_{t+1} = a_t^{\text{opt}}$ . Otherwise, the current optimization problem is considered to be unreliable and thus either we explore the current most informative arm  $\bar{a}_t$  or we explore disregarding the considered structure. To trade-off between these two options, we compare  $\bar{I}_{\bar{a}_t}^{(d)}(t)$  to the minimum index of an unstructured bandit. Let  $\Gamma > 1$ . If  $\bar{I}_{\bar{a}_t}^{(d)}(t) \leq \Gamma \cdot I_{\underline{a}_t}(t)$ , we explore the current most informative arm, that is we choose  $a_{t+1} = \bar{a}_t$ . Otherwise, we explore according to IMED indexes for unstructured bandit, that is we choose

$$a_{t+1} = \hat{a}_t \in \operatorname{argmin}_{a \in \mathcal{A}} I_a(t). \tag{4.15}$$

**Comment.** When  $\gamma_t$  is set equal to zero, pulling arm  $a_t^{\text{opt}}$  then corresponds to pulling the least pulled current sub-optimal arm in exploration phases, which would a priori lead to a non-optimal asymptotic behavior. Now when  $\gamma_t$  is set close to 1, pulling arm  $a_t^{\text{opt}}$  would asymptotically yields a numbers of pulls of current sub-optimal arms larger than  $O(f_\xi(N_{a_t^{\text{opt}}}(t))) = O(f_\xi(f_\xi(t)))$ , which seems insufficient to properly estimate the current optimization problem. The sequence  $(\gamma_t)_{t \geq 1}$  is introduced in order to manage all the situations between these two borderline cases. Now  $\Gamma$  trades-off between gathering information about the structure and

reliability perspectives. When the structure is not informative, that is when  $\bar{I}_a(t) = I_a(t)$  for all  $a \in \mathcal{A}$ , the ratio  $\bar{I}_{\bar{a}_t}(t)/I_{\bar{a}_t}(t)$  is all equal to 1. A ratio greater than 1 for some current sub-optimal arm means that we gather information about its interactions with other current sub-optimal arms (from a structure perspective) at the price of a poorer estimation of its means. Reliable current means ensures reliable current optimization problems. Lastly, the re-scaling of  $n^{\text{opt}}(t)$ , solution of  $\mathfrak{C}_\Theta(\hat{\mu}(t))$ , by the index  $\bar{I}_{\bar{a}_t}(t)$  and not  $\log(t)$  results in better performance in practice since  $\bar{I}_{\bar{a}_t}(t) \leq f_\xi \left( \max_{a \in \mathcal{A}} N_a(t) \right)$  is upper bounded by the counts  $(N_a(t))_{a \in \mathcal{A}}$  rather than current time step  $t \geq 1$ .

### 4.3.2 Asymptotic optimality of IMED-GS

In this section, we state the main theoretical result of this chapter about the instance-dependent regret bound of our algorithm for graph-structured bandits. To state this result, we first introduce a few technical notations and consider some technical assumptions detailed in Section 4.5.1. We define for the bandit configuration  $\nu$ , its minimal optimality gap as  $\Delta_{\min} = \min_{a \notin \mathcal{A}^*(\nu)} \Delta_a$ . Let us also denote  $\Sigma_{\Delta^{-1}} = \sum_{a \notin \mathcal{A}^*(\nu)} \Delta_a^{-1}$ . With regard to

Assumption 6 and 8, we introduce for convenience the quantity  $\varepsilon_\nu = \frac{\delta_{\min}}{4} \wedge (1 - \mu^*)$ , where  $\delta_{\min} = \Delta_{\min}$  if  $\Theta$  is a singleton, and  $\delta_{\min} = \min_{a \neq a'} \{ \Delta_{\min}, \theta_{a,a'}^* - (\mu_a - \mu_{a'}) \}$  otherwise. In particular, these problem-dependent quantities

are away from 0. We further introduce  $\forall a \in \mathcal{A}, \forall t \geq 1$ , the quantity  $\varepsilon_a(t) = \frac{3}{2} \cdot \frac{\gamma_t}{1 - \gamma_t} \cdot \Delta_a$ , and note that  $\varepsilon_a(t) < \Delta_a/4$  provided that  $\gamma_t < 1/7$ . Last, we recall that  $\gamma_1 = 1/(7 \vee \sqrt{\Phi(2|\mathcal{A}|+1)})$ , where  $\Phi(x) = x \log(x)$ . Our main result is a precise finite-time bound on the regret of IMED-GS, in which we meticulously fill out the details, highlighting the first-order, second-order and constant terms.

**Theorem 3 (Regret upper bound).** *Let us consider a configuration  $\nu \in \mathcal{D}_\Theta$  with means  $\mu = (\mu_a)_{a \in \mathcal{A}} \in (0, 1)^\mathcal{A}$ . Under Assumptions 4-5-7-6-8-9, for all positive decreasing sequence  $(\gamma_t)_{t \geq 1} < \gamma_1$ , for all  $\Gamma \geq 1$ , for all  $\xi > 1$ , for all  $d \geq 0$ , for all accuracy  $0 < \varepsilon < \varepsilon_\nu$ , for all time horizon  $T \geq |\mathcal{A}|$ ,*

$$R(\nu, T) \leq (A) + (B) + C_{\xi, d, \varepsilon}, \text{ with}$$

$$(A) = \inf_{\tau \in [1, T]} \left( \lambda_\tau(\mu, \varepsilon) \cdot \mathfrak{C}_\Theta(\mu(\varepsilon, \tau)) \cdot f_\xi(T) + \sum_{a \in \mathcal{A}} \Delta_a \cdot \tau \right)$$

$$(B) = \frac{\sum_{a \notin \mathcal{A}^*} \Delta_a^{-1} \Gamma \gamma_T^{-2}}{2(\mu^* - \varepsilon)(1 - \mu^* - \varepsilon)} \left[ f_\xi \left( \frac{98 \Sigma_{\Delta^{-1}}}{\delta_{\min}} f_\xi(T) + 1 \right) + f_\xi(C_{\xi, d, \varepsilon}) + 1 \right],$$

where for the first term (A), we have introduced  $\forall t \geq 1, \mu(\varepsilon, t) = (\mu_a(\varepsilon, t))_{a \in \mathcal{A}}$  with

$$\mu_a(\varepsilon, t) = \begin{cases} \mu_a + \varepsilon + \varepsilon_a(t) & , \text{ if } a \notin \mathcal{A}^* \\ \mu_a - \varepsilon & , \text{ if } a \in \mathcal{A}^* \end{cases}$$

$$\text{and } \lambda_t(\mu, \varepsilon) = \max_{a \notin \mathcal{A}^*} \left( \frac{\Delta_a + 2\varepsilon + \varepsilon_a(t)}{\Delta_{a'} - 2\varepsilon - \varepsilon_a(t)} \right)^2.$$

Furthermore,  $(\lambda_t(\mu, \varepsilon))_{t \geq 1}$  and  $(\mathfrak{C}_\Theta(\mu(\varepsilon, t)))_{t \geq 1}$  are non-increasing sequences such that for  $t \geq 1$ ,

$$1 \leq \lambda_t(\mu, \varepsilon) \leq 49 \quad \mathfrak{C}_\Theta(\mu) \leq \mathfrak{C}_\Theta(\mu(\varepsilon, t)) \leq 2\Sigma_{\Delta^{-1}}$$

and, provided that  $\lim_{t \rightarrow \infty} \gamma_t = 0$ ,

$$\lim_{\varepsilon \rightarrow 0} \lim_{t \rightarrow \infty} \lambda_t(\mu, \varepsilon) = 1 \quad \lim_{\varepsilon \rightarrow 0} \lim_{t \rightarrow \infty} \mathfrak{C}_\Theta(\mu(\varepsilon, t)) = \mathfrak{C}_\Theta(\mu).$$

Finally, the term  $C_{\xi, d, \varepsilon} = O(\varepsilon^{-2})$  does not depend on  $T$  and is made explicit in Section 4.6.4.

Asymptotically when  $T \rightarrow \infty$ , provided that  $\lim_{t \rightarrow \infty} \gamma_t = 0$ , the first term is the leading term, and scales with  $\mathfrak{C}_\Theta(\mu) \log(T)$ . Indeed, the infimum involved in the regret bound is bounded above by  $\mathfrak{C}_\Theta(\mu) \log(T) + o(\log(T))$  when  $\tau = \sqrt{\log(T)}$  for instance. Note that besides the fact the terms  $\lambda_\tau(\mu, \varepsilon)$  and  $\mathfrak{C}_\Theta(\mu(\varepsilon, \tau))$  before  $f_\xi(T)$  are asymptotically optimal, these terms are decreasing and provably upper-bounded. The second term scales with  $O(\Gamma \gamma_T^{-2} \log \log(T))$ , and the third one is a constant. Further,  $\varepsilon$  can be chosen such that  $\varepsilon = o(\log(T)^{-1/2})$  in order to obtain an asymptotically optimal bound on the regret that only depends on horizon  $T$ , parameters  $\xi, d, \Gamma$  and  $(\gamma_t)_{t \geq 1}$  of IMED-GS and the intrinsic characteristics of the considered structured bandit. In particular, we deduce the asymptotic optimality of IMED-GS, provided that  $\lim_{t \rightarrow \infty} \gamma_t \sqrt{\frac{\log(t)}{\log \log(t)}} = \infty$ .

**Corollary 3** (Asymptotic optimality). *Let us consider a set of Bernoulli distributions  $\nu \in \mathcal{D}_\Theta$ . Then under IMED-GS algorithm, for  $(\gamma_t)_{t \geq 1}$  such that  $\lim_{t \rightarrow \infty} \gamma_t = 0$  and  $\lim_{t \rightarrow \infty} \gamma_t \sqrt{\frac{\log(t)}{\log \log(t)}} = \infty$ , it holds*

$$\limsup_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} \leq \mathfrak{C}_\Theta(\mu).$$

**Discussion.** Let us remark that Theorem 3 is a non-asymptotic result, were all terms are fully explicit (see Section 4.6 for the precise value of the remaining constant term). This is not very common in structured bandits and contrasts with regret bounds previously obtained for alternative algorithms in the literature. Further, we highlight that we only require the parameter  $\xi$  to exceed 1 thanks to a refinement of concentration inequalities that is of independent interest. We detail this result in Theorem 4 (Chapter 5) as we believe it can benefit other regret analysis. The main innovation is to take into account the fact that for small values of the number of observations, the regret can be handled with other tools than concentration. Hence, concentration only needs to be handled after some burn-in phase and not for all time steps. Now, the use of the reduced indexes with parameter  $d$  is interesting as it enables to decrease to value of the constant  $C_{\xi, d, \varepsilon}$ . Indeed, this enables to make appear an exponential dependency in the parameter  $d$  instead of an exponential dependency of the number of arms  $|\mathcal{A}|$ . In contrast, in Magureanu et al. (2014) the authors required  $\xi > 3|\mathcal{A}|$  (and so in Degenne et al. (2020b), as they employ same concentration results): there exists  $t_0 \geq 1$  such that for all  $t \geq t_0$ ,

$$\mathbb{P}_\nu \left( \sum_{a \in \mathcal{A}} N_a(t) \text{kl}(\hat{\mu}_a(t) | \mu_a) \geq f_{3|\mathcal{A}|+1}(t) \right) \leq \frac{1}{t \log(t)}.$$

A second point is that our concentration result (Theorem 4) does not directly come from stochastic orderings alone due to the presence of the random term  $N_{\mathcal{A}'}(t) = \sum_{a \in \mathcal{A}'} N_a(t)$  for  $\mathcal{A}' \subset \mathcal{A}$ . This random term causes the analysis in Magureanu et al. (2014) to break, hence, we needed to derive a novel analysis to handle this difficulty. Lastly, some terms, e.g. that  $C_{\xi, d, \varepsilon}$ , are reminiscent of other analysis, such that appearing in the optimality bounds for KLUCB or TS (Cappé et al. (2013); Kaufmann et al. (2012)). Indeed these terms are due to concentration inequalities, hence it is expected that they appear here as well. We tighten their control with respect to the previous work. Indeed it is sometimes argued that the non-first order terms are large, making the bound unpractical even for large  $T$ . While this phenomenon cannot be completely avoided, we made specific efforts to control the terms tightly, which is of independent interest.

## 4.4 Numerical experiments

In this section, we illustrate the performance of IMED-GS algorithm when specialized to Unimodal, Lipschitz structures and Aggregates of bandits. These examples enable to compare the regret performance of this algorithm to existing state-of-the-art. We detail these structures in Section A.1.



In all considered bandit configuration, we naturally compare each time IMED-GS to IMED algorithm for unstructured bandits. We also compare IMED-GS to OSSB algorithm for generic structured bandits. For Unimodal structure we add specific comparison with OSUB, UTS and IMED-UB that are specialized to this structure, and for Lipschitz structure we add numerical comparison with CKL-UCB. We further report the IMED-GS run with setting  $d = |\mathcal{A}| - 1$  (hence without the downsizing in the burn-in phase), in order to show that this parameter can be chosen large in practice without hindering much performance (the parameter  $d = 3$  is suggested by theory). Here the parameter  $\Gamma$  is set to  $|\mathcal{A}|^{1.5}$  and  $\xi = 1$ . In each experiment, we considered a time horizon of  $T = 3000$ , and results averaged over 300 independent experiments, reporting 10% and 90% quantiles of the regret on top of its average. The results are reported in Figure 4.1, and show the potential benefit of the algorithm. The distributions and the structures used for Figure 4.1 are provided below. Additional details and complementary experiments are provided in Section A.4.

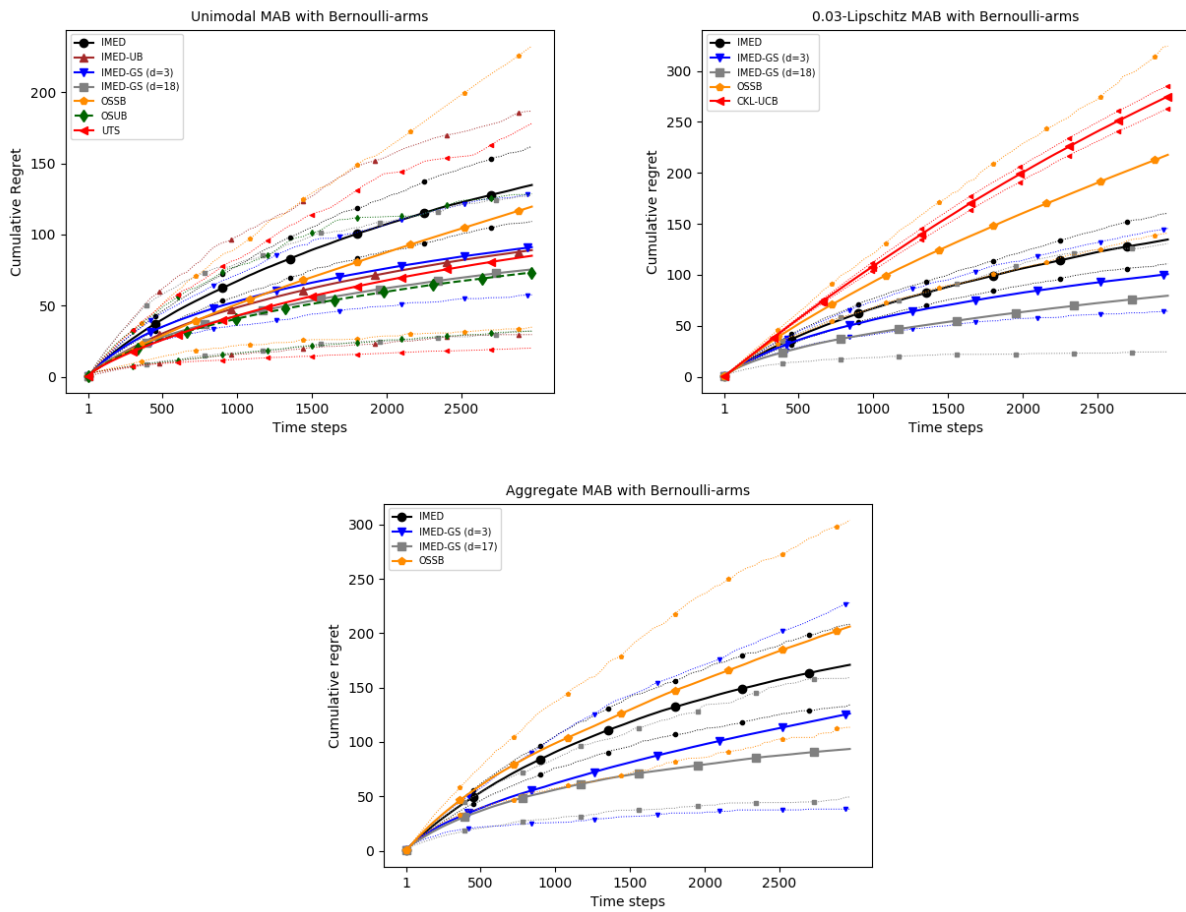


Figure 4.1: Comparison of IMED-GS to other algorithms on several structured bandit instances.

In Figure 4.2 below, we successively represent the vectors of means used for the experiments of Figure 4.1 to respectively illustrate the Unimodal and Lipschitz structures and the Aggregate of bandits. Note that for the Lipschitz structure, the vector of means is 0.03-Lipschitz as assumed in the corresponding experiment. For the Aggregate of bandits, the set of arms is decomposed as  $\mathcal{A} = \mathcal{X} \times \mathcal{K}$  and the means  $\mu = (\mu_{x,k})_{x \in \mathcal{X}, k \in \mathcal{K}}$  are represented in lexicographical order, with  $\mathcal{X} = \llbracket 1, 6 \rrbracket$  and  $\mathcal{K} = \llbracket 1, 3 \rrbracket$ . This means we assume an aggregate of 6 bandits with 3 arms each. We set the relationship matrix equal to  $\omega_{x,x'} = 0.07 |x - x'|$  for  $x, x' \in \mathcal{X}$  and each

arm is represented with a specific marker in Figure 4.2. For the Aggregate of Bandits, the first three arms are the arms of the first bandit, the arms  $\{4, 5, 6\}$  are the arms of the second bandit, and so on.

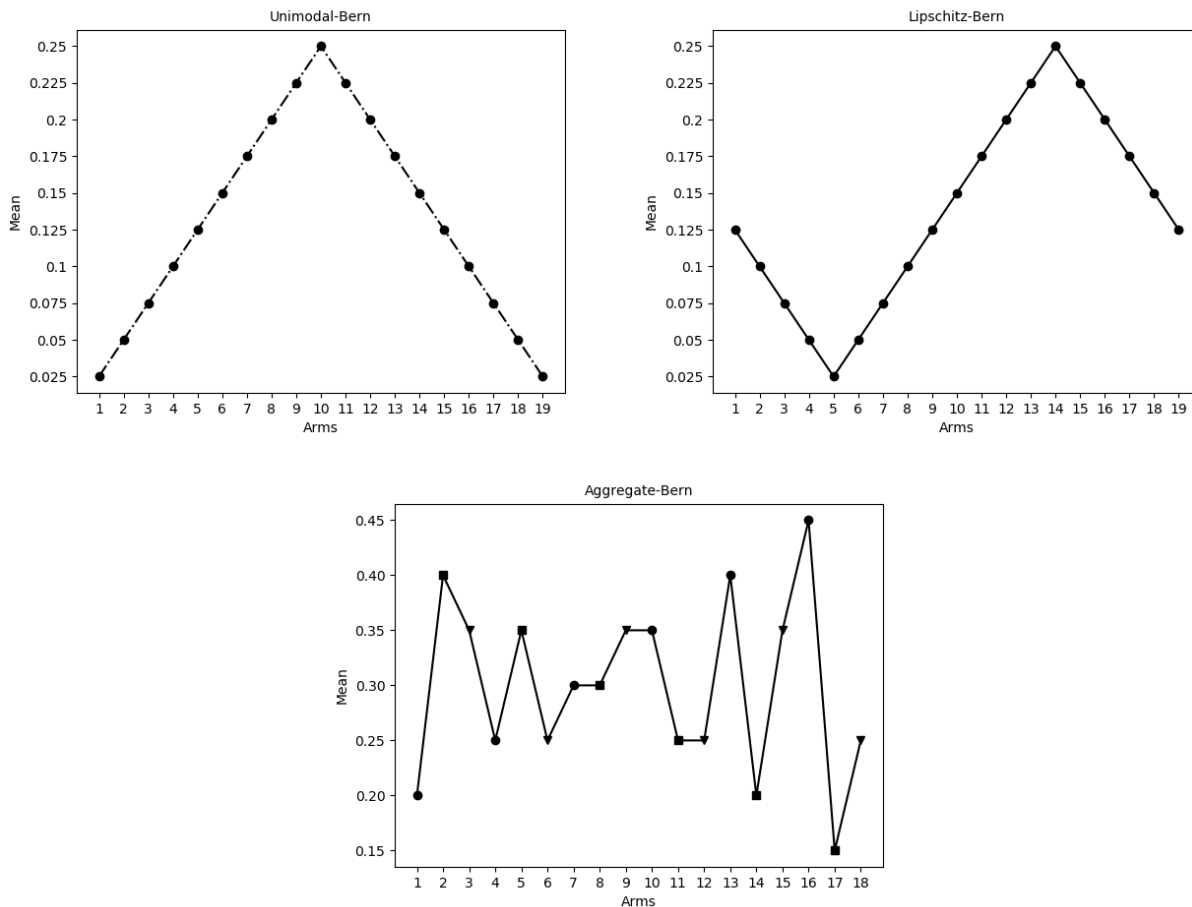


Figure 4.2: Means used for the experiments of Figure 4.1 to illustrate the Unimodal and Lipschitz structures and the Aggregate of bandits

## 4.5 Finite time properties of IMED-GS algorithm

In Section 4.5.1, we state the additional assumptions used in Theorem 3. We prove in Section 4.5.2 algorithm-based empirical bounds on the numbers of pulls. Then we introduce in Section 4.5.3 several notions of deviation of the empirical means and define the subsets of times where these deviations occur before providing upper bounds on the size of these subsets. Then in Section 4.5.4, we consider subset of times where undesirable events occur and establish, using empirical bounds from Section 4.5.2, relations between these subsets and the subsets of times where the empirical means deviate (subsets previously introduced in Section 4.5.3). Lastly, in Section 4.5.5 we show that we have nice reliability properties during most of time steps.

### 4.5.1 Additional assumptions

In this subsection, we state the technical assumptions allowing a simultaneous finite time analysis of the considered structures under the formalism of the graph structure. All these assumptions are satisfied the structure

First, we assume that the instance  $\nu$  has a unique optimal arm, which is a folklore assumption.

**Assumption 7** (Unique optimal arm). *The instance  $\nu$  has a unique optimal arm  $a^* \in \mathcal{A}$ , that is  $\mathcal{A}^* = \{a^*\}$ .*

The next assumption is motivated by technical concerns and is only used when  $\Theta$  is not a singleton. For Unimodal structure, Assumption 8 is trivially satisfied as it reads ( $a'' \in \llbracket a', a \rrbracket$  or  $a'' \in \llbracket a', a \rrbracket$ ) is equivalent to ( $a'' \in \llbracket a, a' \rrbracket$  or  $a'' \in \llbracket a, a' \rrbracket$ ) and that for all  $a < a' \leq a^*$  or  $a^* \leq a' < a$ ,  $\mu_a \neq \mu_{a'}$ . Please refer to Section A.1 for more details.

**Assumption 8** (Symmetry, No adherence). *If  $\Theta$  is not a singleton, then for all  $a, a', a'' \in \mathcal{A}$ ,  $\theta_{a', a''}^{(a)} = \theta_{a, a''}^{(a')}$ .*

**Assumption 9** (No adherence). *For all arms  $a \neq a'$ ,  $\mu_a - \mu_{a'} \neq \theta_{a, a'}^*$ .*

## 4.5.2 Algorithm-based empirical bounds

**Lemma 6** (Empirical lower bounds). *Under IMED-GS, at each step time  $t \geq |\mathcal{A}|$ , for all current sub-optimal arm  $a \notin \hat{\mathcal{A}}^*(t)$ ,*

$$f_\xi(N_{a_{t+1}}(t)) \leq \sum_{a' \in \hat{\mathcal{A}}_a^{(d)}(t)} N_{a'}(t) \text{kl}(\hat{\mu}_{a'}(t) \mid \hat{\mu}^*(t) - \hat{\theta}_{a, a'}^*(t)) + f_\xi(N_{\hat{\mathcal{A}}_a^{(d)}(t)}(t)), \quad (4.16)$$

and, for all current optimal arm  $\hat{a}^* \in \hat{\mathcal{A}}^*(t)$ ,

$$N_{a_{t+1}}(t) \leq N_{\hat{a}^*}(t). \quad (4.17)$$

Furthermore, if time step  $t \geq |\mathcal{A}|$  corresponds to exploration, that is  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$ , for all current sub-optimal arm  $a' \in \mathcal{A}$ ,

$$f_\xi(N_{\bar{a}_t}(t)) \leq N_{a'}(t) \hat{\mathbf{k}}_{a'}(t) + f_\xi(N_{a'}(t)). \quad (4.18)$$

*Proof.* First we note that

$$I_a(t) \leq \bar{I}_a^{(d)}(t), \quad \forall a \in \mathcal{A}, \quad (4.19)$$

implies

$$\min_{a \in \mathcal{A}} I_a(t) \leq \min_{a \in \mathcal{A}} \bar{I}_a^{(d)}(t), \quad (4.20)$$

that is,

$$I_{\hat{a}_t}(t) \leq \bar{I}_{\hat{a}_t}^{(d)}(t) \quad (4.21)$$

From Equations (4.9), (4.6), (4.10) and (4.15) defining indexes  $\left(\bar{I}_a^{(d)}(t)\right)_{a \in \mathcal{A}}$ ,  $(I_a(t))_{a \in \mathcal{A}}$ ,  $\bar{a}_t$  and  $\hat{a}_t$ , we have

$$\begin{aligned} \forall a \notin \hat{\mathcal{A}}^*(t), \forall a' \in \bar{\mathcal{A}}_a(t), \quad f_\xi(N_a(t)) \leq I_a(t) \leq \bar{I}_a^{(d)}(t) \\ f_\xi(N_{a'}(t)) \leq \bar{I}_a^{(d)}(t) \end{aligned} \quad (4.22)$$

$$\forall \hat{a}^* \in \hat{\mathcal{A}}^*(t), \quad f_\xi(N_{\hat{a}^*}(t)) = I_{\hat{a}^*}(t) = \bar{I}_{\hat{a}^*}^{(d)}(t),$$

and

$$\forall a \notin \hat{\mathcal{A}}^*(t), \quad \bar{I}_{\bar{a}_t}^{(d)}(t) \leq \bar{I}_a^{(d)}(t) \leq \sum_{a' \in \hat{\mathcal{A}}_a^{(d)}(t)} N_{a'}(t) \text{kl}(\hat{\mu}_{a'}(t) \mid \hat{\mu}^*(t) - \hat{\theta}_{a, a'}^*(t)) + f_\xi(N_{\hat{\mathcal{A}}_a^{(d)}(t)}(t)), \quad (4.23)$$

$$\forall \hat{a}^* \in \hat{\mathcal{A}}^*(t), \quad \bar{I}_{\hat{a}_t}^{(d)}(t) \leq \bar{I}_{\hat{a}^*}^{(d)}(t) = f_\xi(N_{\hat{a}^*}(t)).$$

Combining Equations (4.21), (4.22) and (4.23) yields

$$\forall a'' \in \overline{\mathcal{A}}_{\bar{a}_t}(t) \cup \{\bar{a}_t\}, \forall a \notin \widehat{\mathcal{A}}^*(t),$$

$$\mathbf{f}_\xi(N_{\hat{a}_t}(t)), \mathbf{f}_\xi(N_{\bar{a}_t}(t)), \mathbf{f}_\xi(N_{a''}(t)) \leq \sum_{a' \in \widehat{\mathcal{A}}_a^{(d)}(t)} N_{a'}(t) \text{kl}(\widehat{\mu}_{a'}(t) | \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)) + \mathbf{f}_\xi\left(N_{\widehat{\mathcal{A}}_a^{(d)}(t)}(t)\right), \quad (4.24)$$

$$\forall a'' \in \overline{\mathcal{A}}_{\bar{a}_t}(t) \cup \{\bar{a}_t\}, \forall \widehat{a}^* \in \widehat{\mathcal{A}}^*(t),$$

$$N_{\hat{a}_t}(t), N_{\bar{a}_t}(t), N_{a''}(t) \leq N_{\widehat{a}^*}(t),$$

where, by convention,  $\overline{\mathcal{A}}_{\widehat{a}^*}(t) = \emptyset$  for all  $\widehat{a}^* \in \widehat{\mathcal{A}}^*(t)$ . We note that  $\{\bar{a}_t\} \subset \widehat{\mathcal{A}}_{\bar{a}_t}(t)$  when  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$ . In particular, by considering Equation (4.12) that defines  $a_t^{\text{opt}}$ , when  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$  we have

$$a_t^{\text{opt}} \in \overline{\mathcal{A}}_{\bar{a}_t}(t)$$

and Equation (4.24) then implies

$$\forall a \notin \widehat{\mathcal{A}}^*(t),$$

$$\begin{aligned} & \mathbf{f}_\xi(N_{\hat{a}_t}(t)), \mathbf{f}_\xi(N_{\bar{a}_t}(t)), \mathbb{I}_{\{\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)\}} \mathbf{f}_\xi\left(N_{a_t^{\text{opt}}}(t)\right) \\ & \leq \sum_{a' \in \widehat{\mathcal{A}}_a^{(d)}(t)} N_{a'}(t) \text{kl}(\widehat{\mu}_{a'}(t) | \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)) + \mathbf{f}_\xi\left(N_{\widehat{\mathcal{A}}_a^{(d)}(t)}(t)\right), \end{aligned} \quad (4.25)$$

$$\forall \widehat{a}^* \in \widehat{\mathcal{A}}^*(t),$$

$$N_{\hat{a}_t}(t), N_{\bar{a}_t}(t), \mathbb{I}_{\{\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)\}} N_{a_t^{\text{opt}}}(t) \leq N_{\widehat{a}^*}(t).$$

If  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$ , from Equations (4.13) and (4.14) defining indexes  $(\underline{I}_{a'}(t))_{a' \in \mathcal{A}}$  and  $\underline{a}_t$ , we have

$$\forall a' \notin \widehat{\mathcal{A}}^*(t), \quad \mathbf{f}_\xi(N_{a'}(t)) \leq \underline{I}_{a'}(t), \quad (4.26)$$

and

$$\forall a' \notin \widehat{\mathcal{A}}^*(t), \quad \underline{I}_{\underline{a}_t}(t) \leq \underline{I}_{a'}(t) \leq N_{a'}(t) \widehat{\mathbf{k}}_{a'}^*(t) + \mathbf{f}_\xi(N_{a'}(t)). \quad (4.27)$$

□

**Lemma 7** (Empirical upper bounds). *Under IMED-GS, at each step time  $t \geq |\mathcal{A}|$  such that  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$ ,*

$$N_{\hat{a}_t}(t) \leq \frac{\mathbf{f}_\xi(t)}{\text{kl}(\widehat{\mu}_{\hat{a}_t}(t) | \widehat{\mu}^*(t))}, \quad (4.28)$$

$$N_{\bar{a}_t}(t) \leq \frac{\mathbf{f}_\xi(t)}{\text{kl}(\widehat{\mu}_{\bar{a}_t}(t) | \widehat{\mu}^*(t))}, \quad (4.29)$$

$$N_{a_t^{\text{opt}}}(t) \leq N_{a_t^{\text{opt}}}^{\text{opt}}(t), \quad (4.30)$$

$$\mathbb{I}_{\{\bar{a}_t \neq a_t^{\text{opt}}\}} N_{\hat{a}_t}(t) \leq \frac{\mathbf{f}_\xi\left(N_{a_t^{\text{opt}}}(t)\right)}{(\gamma_t)^2 \widehat{\mu}^*(t) (1 - \widehat{\mu}^*(t)) \text{kl}(\widehat{\mu}_{\hat{a}_t}(t) | \widehat{\mu}^*(t))}, \quad (4.31)$$

$$\mathbb{I}_{\{\underline{a}_t \neq \underline{a}_t^{\text{opt}}\}} N_{\bar{a}_t}(t) \leq \frac{\bar{I}_{\bar{a}_t}^{(d)}(t)}{I_{\hat{a}_t}(t)} \frac{f_\xi(N_{\underline{a}_t^{\text{opt}}}(t)) + 1}{(\gamma_t)^2 \hat{\mu}^*(t) (1 - \hat{\mu}^*(t)) \text{kl}(\hat{\mu}_{\bar{a}_t}(t) | \hat{\mu}^*(t))}, \quad (4.32)$$

$$N_{\underline{a}_t}(t) \leq \max \left\{ N_{\underline{a}_t}^{\text{opt}}(t), \frac{f_\xi(N_{\underline{a}_t^{\text{opt}}}(t))}{\hat{\mathbf{k}}_{\underline{a}_t}^*(t)} \right\} \leq N_{\underline{a}_t}^{\text{opt}}(t) + \max_{a \notin \hat{\mathcal{A}}^*(t)} \frac{f_\xi(N_a(t))}{\hat{\mathbf{k}}_{a_t}^*(t)}, \quad (4.33)$$

with the convention  $0/0=0$  and  $I/0=\infty$  for  $I>0$ .

*Proof.* From Equations (4.9), (4.10) and Equations (4.6), (4.15) defining  $(\bar{I}_a(t))_{a \in \mathcal{A}}$ ,  $\bar{a}_t$  and  $(I_a(t))_{a \in \mathcal{A}}$ ,  $\hat{a}_t$ , we have for  $\hat{a}^* \in \hat{\mathcal{A}}^*(t)$ ,

$$N_{\hat{a}_t}(t) \text{kl}(\hat{\mu}_{\hat{a}_t}(t) | \hat{\mu}^*(t)) \leq I_{\hat{a}_t}(t) \leq I_{\hat{a}^*}(t) \leq f_\xi(t) \quad (4.34)$$

and

$$N_{\bar{a}_t}(t) \text{kl}(\hat{\mu}_{\bar{a}_t}(t) | \hat{\mu}^*(t)) \leq \bar{I}_{\bar{a}_t}^{(d)}(t) \leq \bar{I}_{\hat{a}^*}(t) \leq f_\xi(t). \quad (4.35)$$

Thus, we deduce Equations (4.28) and (4.29) from Equations (4.34) and (4.35).

From Equation (4.12) that defines  $a_t^{\text{opt}}$ , proving Equation (4.30) from Lemma 7 amounts to prove

$$\max_{a' \in \bar{\mathcal{A}}_{\bar{a}_t}(t)} N_{a'}^{\text{opt}}(t) - N_{a'}(t) \geq 0.$$

Since  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$ , from Equation (4.11) that defines  $(n_a^{\text{opt}}(t))_{a \in \mathcal{A}}$  as solution of current minimization problem, we must have

$$\sum_{a' \in \bar{\mathcal{A}}_{\bar{a}_t}(t)} \text{kl}(\hat{\mu}_{a'}(t) | \hat{\mu}^*(t) - \hat{\theta}_{\bar{a}_t, a'}^*(t)) n_{a'}^{\text{opt}}(t) \geq 1. \quad (4.36)$$

Since  $(N_a^{\text{opt}}(t))_{a \in \mathcal{A}} = (n_a^{\text{opt}}(t) \bar{I}_{\bar{a}_t}(t))_{a \in \mathcal{A}}$ , from Equation (4.36) we have

$$\sum_{a' \in \bar{\mathcal{A}}_{\bar{a}_t}(t)} \text{kl}(\hat{\mu}_{a'}(t) | \hat{\mu}^*(t) - \hat{\theta}_{\bar{a}_t, a'}^*(t)) N_{a'}^{\text{opt}}(t) \geq \bar{I}_{\bar{a}_t}(t). \quad (4.37)$$

Since  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$ , from Equation (4.9) that explains  $\bar{I}_{\bar{a}_t}(t)$  we have

$$\bar{I}_{\bar{a}_t}(t) \geq \sum_{a' \in \bar{\mathcal{A}}_{\bar{a}_t}(t)} \text{kl}(\hat{\mu}_{a'}(t) | \hat{\mu}^*(t) - \hat{\theta}_{\bar{a}_t, a'}^*(t)) N_{a'}(t). \quad (4.38)$$

By combining Equations (4.37) and (4.38), we get

$$\sum_{a' \in \bar{\mathcal{A}}_{\bar{a}_t}(t)} \text{kl}(\hat{\mu}_{a'}(t) | \hat{\mu}^*(t) - \hat{\theta}_{\bar{a}_t, a'}^*(t)) (N_{a'}^{\text{opt}}(t) - N_{a'}(t)) \geq 0. \quad (4.39)$$

Since  $\text{kl}(\hat{\mu}_{a'}(t) | \hat{\mu}^*(t) - \hat{\theta}_{\bar{a}_t, a'}^*(t)) \geq 0$  for all  $a' \in \bar{\mathcal{A}}_{\bar{a}_t}(t)$ , Equation (4.39) implies

$$\max_{a' \in \bar{\mathcal{A}}_{\bar{a}_t}(t)} N_{a'}^{\text{opt}}(t) - N_{a'}(t) \geq 0$$

which ends the proof of Equation (4.30).

From Lemma 23, we have

$$(\gamma_t)^2 \widehat{\mu}^*(t)(1 - \widehat{\mu}^*(t)) \text{kl}(\widehat{\mu}_a(t) | \widehat{\mu}^*(t)) \leq \widehat{k}_a(t), \quad \forall a \in \mathcal{A}. \quad (4.40)$$

From previous Equation (4.40) and Equations (4.13) and (4.6) defining  $(\underline{I}_a(t))_{a \in \mathcal{A}}$  and  $(I_a(t))_{a \in \mathcal{A}}$  we have

$$(\gamma_t)^2 \widehat{\mu}^*(t)(1 - \widehat{\mu}^*(t)) I_{a'}(t) \leq \underline{I}_{a'}(t), \quad \forall a' \neq a_t^{\text{opt}}. \quad (4.41)$$

Since  $I_{\hat{a}_t}(t) = \min_{a \in \mathcal{A}} I_a(t)$ , previous Equation (4.41) implies

$$(\gamma_t)^2 \widehat{\mu}^*(t)(1 - \widehat{\mu}^*(t)) I_{\hat{a}_t}(t) \leq \min_{a' \neq a_t^{\text{opt}}} \underline{I}_{a'}(t). \quad (4.42)$$

Since  $\underline{I}_{a_t}(t) = \min_{a' \in \mathcal{A}} \underline{I}_{a'}(t)$ , previous Equation (4.42) implies

$$\mathbb{I}_{\{a_t \neq a_t^{\text{opt}}\}} (\gamma_t)^2 \widehat{\mu}^*(t)(1 - \widehat{\mu}^*(t)) I_{\hat{a}_t}(t) \leq \mathbb{I}_{\{a_t \neq a_t^{\text{opt}}\}} \underline{I}_{a_t}(t) \leq \underline{I}_{a_t^{\text{opt}}}(t) = \mathbf{f}_\xi(N_{a_t^{\text{opt}}}(t)). \quad (4.43)$$

Combining Equation (4.34) and (4.43), we get

$$\mathbb{I}_{\{a_t \neq a_t^{\text{opt}}\}} (\gamma_t)^2 \widehat{\mu}^*(t)(1 - \widehat{\mu}^*(t)) N_{\hat{a}_t}(t) \text{kl}(\widehat{\mu}_{\hat{a}_t}(t) | \widehat{\mu}^*(t)) \leq \mathbf{f}_\xi(N_{a_t^{\text{opt}}}(t)). \quad (4.44)$$

We deduce Equation (4.31) from previous Equation (4.44).

We note that

$$\bar{I}_{\bar{a}_t}^{(d)}(t) \leq \mathbb{I}_{\{I_{\hat{a}_t}(t) \neq 0\}} \frac{\bar{I}_{\bar{a}_t}^{(d)}(t)}{I_{\hat{a}_t}(t)} I_{\hat{a}_t}(t) + \mathbb{I}_{\{I_{\hat{a}_t}(t) = 0\}} \bar{I}_{\bar{a}_t}^{(d)}(t). \quad (4.45)$$

Combining previous Equation (4.45) and Equations (4.35), (4.43), we get

$$\mathbb{I}_{\{a_t \neq a_t^{\text{opt}}\}} (\gamma_t)^2 \widehat{\mu}^*(t)(1 - \widehat{\mu}^*(t)) N_{\bar{a}_t}(t) \text{kl}(\widehat{\mu}_{\bar{a}_t}(t) | \widehat{\mu}^*(t)) \leq \frac{\bar{I}_{\bar{a}_t}^{(d)}(t)}{I_{\hat{a}_t}(t)} \left( \mathbf{f}_\xi(N_{a_t^{\text{opt}}}(t)) + 1 \right). \quad (4.46)$$

We deduce Equation (4.32) from previous Equation (4.46).  $\square$

### 4.5.3 Non-reliable current means

For all arms  $a, a' \in \mathcal{A}$  and for all accuracy  $\varepsilon > 0$ , let  $\mathcal{E}_{a,a'}^+(\varepsilon)$  be the set of times where the current mean of arm  $a$   $\varepsilon$ -deviates from above while arm  $a$  has more pulls than the current pulled arm  $a'$ ,

$$\mathcal{E}_{a,a'}^+(\varepsilon) := \{t \geq 1 : a_{t+1} = a', N_{a'}(t) \leq N_a(t), \widehat{\mu}_a(t) \geq \mu_a + \varepsilon\}. \quad (4.47)$$

We similarly define

$$\mathcal{E}_{a,a'}^-(\varepsilon) := \{t \geq 1 : a_{t+1} = a', N_{a'}(t) \leq N_a(t), \widehat{\mu}_a(t) \leq \mu_a - \varepsilon\}. \quad (4.48)$$

We also define

$$\mathcal{E}_{a,a'}(\varepsilon) = \mathcal{E}_{a,a'}^+(\varepsilon) \cup \mathcal{E}_{a,a'}^-(\varepsilon). \quad (4.49)$$

**Definition 2** (kl- $f_\xi$  deviation). For  $\varepsilon > 0$ , the couple of arms  $(a, a') \in \mathcal{A}^2$  shows  $\varepsilon^+$ -kl- $f_\xi$  deviation at time step  $t \geq 1$  if the following conditions are satisfied,

- (1)  $a_{t+1} = a'$
- (2)  $\widehat{\mu}_a(t) \geq \mu_a + \varepsilon$
- (3)  $f_\xi(N_{a'}(t)) \leq N_a(t) \text{kl}(1 - \widehat{\mu}_a(t) | 1 - \mu_a - \varepsilon) + f_\xi(N_a(t))$ .

For  $\varepsilon > 0$ , the couple of arms  $(a, a') \in \mathcal{A}^2$  shows  $\varepsilon^-$ -kl- $f_\xi$  deviation at time step  $t \geq 1$  if the following conditions are satisfied,

- (1)  $a_{t+1} = a'$
- (2)  $\widehat{\mu}_a(t) \leq \mu_a - \varepsilon$
- (3)  $f_\xi(N_{a'}(t)) \leq N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) + f_\xi(N_a(t))$ .

For all couple of arms  $(a, a') \in \mathcal{A}^2$  and for all accuracy  $\varepsilon > 0$ , let  $\mathcal{K}_{a,a'}^+(\varepsilon)$  be the set of times where couple of arms  $(a, a')$  shows  $\varepsilon^+$ -kl- $f_\xi$  deviation, that is

$$\mathcal{K}_{a,a'}^+(\varepsilon) := \left\{ t \geq 1 : \begin{array}{l} (1) \ a_{t+1} = a' \\ (2) \ \widehat{\mu}_a(t) \geq \mu_a + \varepsilon \\ (3) \ f_\xi(N_{a'}(t)) \leq N_a(t) \text{kl}(1 - \widehat{\mu}_a(t) | 1 - \mu_a - \varepsilon) + f_\xi(N_a(t)) \end{array} \right\}. \quad (4.50)$$

For all couple of arms  $(a, a') \in \mathcal{A}^2$  and for all accuracy  $\varepsilon > 0$ , let  $\mathcal{K}_{a,a'}^-(\varepsilon)$  be the set of times where couple of arms  $(a, a')$  shows  $\varepsilon^-$ -kl- $f_\xi$  deviation, that is

$$\mathcal{K}_{a,a'}^-(\varepsilon) := \left\{ t \geq 1 : \begin{array}{l} (1) \ a_{t+1} = a' \\ (2) \ \widehat{\mu}_a(t) \leq \mu_a - \varepsilon \\ (3) \ f_\xi(N_{a'}(t)) \leq N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) + f_\xi(N_a(t)) \end{array} \right\}. \quad (4.51)$$

We note that,

$$\mathcal{E}_{a,a'}^+(\varepsilon) \subset \mathcal{K}_{a,a'}^+(\varepsilon) \quad \mathcal{E}_{a,a'}^-(\varepsilon) \subset \mathcal{K}_{a,a'}^-(\varepsilon).$$

We also define

$$\mathcal{K}_{a,a'}(\varepsilon) = \mathcal{K}_{a,a'}^+(\varepsilon) \cup \mathcal{K}_{a,a'}^-(\varepsilon). \quad (4.52)$$

Finally, for accuracy  $\varepsilon > 0$ , we denote by

$$\mathcal{K}^*(\varepsilon) := \left\{ t \geq 1 : \begin{array}{l} (1) \ \widehat{\mu}_a(t) \leq \mu_a - \varepsilon, \quad \forall a \in \widehat{\mathcal{A}}_{a^*}(t) \\ (2) \ 1 \leq N_{\widehat{\mathcal{A}}_{a^*}(t)}(t) \leq |\widehat{\mathcal{A}}_{a^*}(t)| f_\xi(N_{a_{t+1}}(t)) / 2\varepsilon^2 \\ (3) \ \sum_{a \in \widehat{\mathcal{A}}_{a^*}(t)} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) \\ \geq \mathbb{I}_{\{|\widehat{\mathcal{A}}_{a^*}(t)| > d+1\}} \Phi(2|\widehat{\mathcal{A}}_{a^*}(t)| + 1) \vee (f_\xi(N_{a_{t+1}}(t)) - f_\xi(N_{\widehat{\mathcal{A}}_{a^*}(t)}(t))) \\ + \mathbb{I}_{\{|\widehat{\mathcal{A}}_{a^*}(t)| \leq d+1\}} (d+1) \vee (f_\xi(N_{a_{t+1}}(t)) - f_\xi(N_{\widehat{\mathcal{A}}_{a^*}(t)}(t))) \end{array} \right\} \quad (4.53)$$

$$\mathcal{K}_{a',\mathcal{A}'}^*(\varepsilon) = \left\{ t \in \mathcal{K}^*(\varepsilon) : a_{t+1} = a', \widehat{\mathcal{A}}_{a^*}(t) = \mathcal{A}' \right\}, \quad a' \in \mathcal{A}, \mathcal{A}' \subset \mathcal{A} \quad (4.54)$$

and

$$\mathcal{L}^{(d)}(\varepsilon) = \bigcup_{a \in \mathcal{A}} \mathcal{L}_a^{(d)}(\varepsilon) \quad \mathcal{L}_a^{(d)}(\varepsilon) = \left\{ t \geq 1 : \begin{array}{l} (1) \ \widehat{\mu}_a(t) \leq \mu_a - \varepsilon \\ (2) \ N_a(t) \geq \frac{f_\xi(N_{a_{t+1}}(t))}{2\varepsilon^2} \wedge \frac{e^{-d-1}N_{a_{t+1}}(t)}{d+1} \end{array} \right\} \quad (4.55)$$

**Lemma 8** (Bounded subsets of times). For  $\varepsilon > 0$ , for  $(a, a') \in \mathcal{A}^2$ ,

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^+(\varepsilon)|], \mathbb{E}_\nu [|\mathcal{E}_{a,a'}^-(\varepsilon)|] \leq \frac{e^{2\varepsilon^2}}{2\varepsilon^2} \quad (4.56)$$

$$\mathbb{E}_\nu [|\mathcal{L}^{(d)}(\varepsilon)|] \leq |\mathcal{A}| \frac{(d+1)e^{d+1}}{2\varepsilon^2} \left( f_\xi \left( \frac{(d+1)e^{d+1}}{2\varepsilon^2} \right) \right)^2 + \sum_{n \geq 3} \frac{|\mathcal{A}|^2}{n (\log(n))^\xi} \quad (4.57)$$

$$\mathbb{E}_\nu [|\mathcal{K}_{a,a'}^+(\varepsilon)|], \mathbb{E}_\nu [|\mathcal{K}_{a,a'}^-(\varepsilon)|] \leq K_{\xi,1} + \sum_{n \geq 3} \frac{1}{n (\log(n))^\xi} \quad (4.58)$$

$$\mathbb{E}_\nu [|\mathcal{K}^*(\varepsilon)|] \leq K_\xi^*, \quad (4.59)$$

where for  $\xi > 1$ ,

$$\begin{aligned} K_\xi^* &= |\mathcal{A}| \left( e^2 \sqrt{\frac{|\mathcal{A}|}{2\varepsilon^2}} \right) \left( f_\xi \left( e^2 \sqrt{\frac{|\mathcal{A}|}{2\varepsilon^2}} \right) \right)^2 \\ &+ \sum_{n \geq 3} \sum_{\mathcal{A}' \subset \mathcal{A}} \frac{e^{|\mathcal{A}'|+2}}{|\mathcal{A}'|^{|\mathcal{A}'|}} \log(n)^{|\mathcal{A}'|+1} [\Phi(2|\mathcal{A}'|+1) \vee f_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|+1} e^{-[\Phi(2|\mathcal{A}'|+1) \vee f_\xi(n) + 2\varepsilon^2 \underline{N}]} \\ &+ \sum_{n \geq 3} \sum_{\substack{\mathcal{A}' \subset \mathcal{A} \\ |\mathcal{A}'| \leq d+1}} \frac{e^{|\mathcal{A}'|+2}}{|\mathcal{A}'|^{|\mathcal{A}'|}} \log(n)^{|\mathcal{A}'|+1} [(d+1) \vee f_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|+1} e^{-[(d+1) \vee f_\xi(n) + 2\varepsilon^2 \underline{N}]}, \\ K_{\xi,1} &= \left( e^2 \sqrt{\frac{|\mathcal{A}|}{2\varepsilon^2}} \right) \left( f_\xi \left( e^2 \sqrt{\frac{|\mathcal{A}|}{2\varepsilon^2}} \right) \right)^2 + \sum_{n \geq 3} e^3 \log(n)^3 [f_\xi(n) + 2\varepsilon^2 \underline{N}]^3 e^{-[f_\xi(n) + 2\varepsilon^2 \underline{N}]}, \end{aligned}$$

in which we introduce  $\underline{N} = 1 \vee \frac{f_\xi(n)}{1 - \max_{a \in \mathcal{A}'} \mu_a}$ .

*Proof.* We start by proving  $\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^-(\varepsilon)|] \leq e^{2\varepsilon^2}/2\varepsilon^2$ . The proof that  $\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^+(\varepsilon)|] \leq e^{2\varepsilon^2}/2\varepsilon^2$  is similar.

We write

$$|\mathcal{E}_{a,a'}^-(\varepsilon)| = \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a', N_{a'}(t) \leq N_a(t), \mu_a - \hat{\mu}_a(t) \geq \varepsilon\}}. \quad (4.60)$$

Considering the stopped stopping times  $\tau_n = \inf \{t \geq 1, N_{a'}(t) = n\}$  we will rewrite the previous sum of indicators and use Lemma 32.

$$\begin{aligned} |\mathcal{E}_{a,a'}^+(\varepsilon)| &\leq \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a', N_{a'}(t) \leq N_a(t), \mu_a - \hat{\mu}_a(t) \geq \varepsilon\}} \\ &\leq \sum_{n \geq 1} \mathbb{I}_{\{n-1 \leq N_a(\tau_n-1), \mu_a - \hat{\mu}_a(\tau_n-1) \geq \varepsilon\}} \\ &\leq 1 + \sum_{n \geq 2} \mathbb{I}_{\{n-1 \leq N_a(\tau_n-1), \mu_a - \hat{\mu}_a(\tau_n-1) \geq \varepsilon\}}. \end{aligned} \quad (4.61)$$

Taking the expectation of Equation (4.61), it comes

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^+(\varepsilon)|] \leq 1 + \sum_{n \geq 1} \mathbb{P}_\nu \left( \bigcup_{\substack{t \geq |\mathcal{A}| \\ \hat{\mu}_a(t) < \mu_a \\ N_a(t) \geq n}} \mu_a - \hat{\mu}_a(t) \geq \varepsilon \right). \quad (4.62)$$



By Pinsker's inequality, previous Equation (4.62) implies

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^+(\varepsilon)|] \leq 1 + \sum_{n \geq 1} \mathbb{P}_\nu \left( \bigcup_{\substack{t \geq |\mathcal{A}| \\ \hat{\mu}_a(t) < \mu_a \\ N_a(t) \geq n}} \text{kl}(\hat{\mu}_a(t) | \mu_a) \geq 2\varepsilon^2 \right). \quad (4.63)$$

From Lemma 32, previous Equation (4.63) implies

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^+(\varepsilon)|] \leq \sum_{n \geq 0} \exp(-2n\varepsilon^2) = \frac{1}{1 - e^{-2\varepsilon^2}}. \quad (4.64)$$

Finally we note that

$$\frac{1}{1 - e^{-2\varepsilon^2}} = \frac{e^{2\varepsilon^2}}{e^{2\varepsilon^2} - 1} \leq \frac{e^{2\varepsilon^2}}{2\varepsilon^2},$$

which ends the proof.

We then prove Equation (4.57). By using in particular Pinsker's inequality, we have

$$\begin{aligned} |\mathcal{L}^{(d)}(\varepsilon)| &\leq \left| \left\{ t \geq 1 : N_{a_{t+1}}(t) < \frac{(d+1)e^{d+1}}{2\varepsilon^2} \left( f_\xi \left( \frac{(d+1)e^{d+1}}{2\varepsilon^2} \right) \right)^2 \right\} \right| \\ &\quad + \sum_{a \in \mathcal{A}} |\{t \geq 1 : N_a(t) \geq f_\xi(N_{a_{t+1}}(t))/2\varepsilon^2, \hat{\mu}_a(t) < \mu_a, \text{kl}(\hat{\mu}_a(t) | \mu_a) \geq 2\varepsilon^2\}| \\ &\leq |\mathcal{A}| \frac{(d+1)e^{d+1}}{2\varepsilon^2} \left( f_\xi \left( \frac{(d+1)e^{d+1}}{2\varepsilon^2} \right) \right)^2 \\ &\quad + \sum_{a \in \mathcal{A}} |\mathcal{A}| \sum_{n \geq 3} |\{t \geq 1 : N_a(t) \geq f_\xi(n)/2\varepsilon^2, \hat{\mu}_a(t) < \mu_a, \text{kl}(\hat{\mu}_a(t) | \mu_a) \geq 2\varepsilon^2\}|. \end{aligned}$$

By taking the expectation on both sides of previous inequality, it comes

$$\begin{aligned} \mathbb{E}_\nu [|\mathcal{L}^{(d)}(\varepsilon)|] &\leq |\mathcal{A}| \frac{(d+1)e^{d+1}}{2\varepsilon^2} \left( f_\xi \left( \frac{(d+1)e^{d+1}}{2\varepsilon^2} \right) \right)^2 \\ &\quad + \sum_{a \in \mathcal{A}} |\mathcal{A}| \sum_{n \geq 3} \mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ \hat{\mu}_a(t) < \mu_a \\ N_a(t) \geq f_\xi(n)/2\varepsilon^2}} \text{kl}(\hat{\mu}_a(t) | \mu_a) \geq 2\varepsilon^2 \right). \end{aligned} \quad (4.65)$$

From Lemma 32, previous Equation (4.65) implies

$$\begin{aligned} \mathbb{E}_\nu [|\mathcal{L}^{(d)}(\varepsilon)|] &\leq |\mathcal{A}| \frac{(d+1)e^{d+1}}{2\varepsilon^2} \left( f_\xi \left( \frac{(d+1)e^{d+1}}{2\varepsilon^2} \right) \right)^2 \\ &\quad + \sum_{a \in \mathcal{A}} |\mathcal{A}| \sum_{n \geq 3} e^{-f_\xi(n)}, \end{aligned} \quad (4.66)$$

which ends the proof.

We now prove Equation (4.59). The remaining inequalities are proven similarly. We simply note that

$$\begin{aligned}
& |\mathcal{K}^*(\varepsilon)| \\
& \leq \sum_{\substack{a' \in \mathcal{A} \\ \mathcal{A}' \subset \mathcal{A}}} |\mathcal{K}_{a', \mathcal{A}'}^*(\varepsilon)| \\
& \leq \sum_{a' \in \mathcal{A}} n_{\xi, \mathcal{A}, \varepsilon} \\
& \quad + \sum_{\substack{a' \in \mathcal{A} \\ \mathcal{A}' \subset \mathcal{A} \\ n \geq n_{\xi, \mathcal{A}, \varepsilon}}} \left| \left\{ t \geq 1 : \begin{aligned} & \hat{\mu}_a(t) < \mu_a - \varepsilon, \quad \forall a \in \mathcal{A}' \\ & 1 \leq N_{\mathcal{A}'}(t) \leq |\mathcal{A}'| f_\xi(n) / 2\varepsilon^2 \\ & \sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\hat{\mu}_a(t) | \mu_a - \varepsilon) \geq \Phi(2|\mathcal{A}'| + 1) \vee (f_\xi(n) - f_\xi(N_{\mathcal{A}'}(t))) \end{aligned} \right\} \right| \\
& \quad + \sum_{\substack{a' \in \mathcal{A} \\ \mathcal{A}' \subset \mathcal{A}, |\mathcal{A}'| \leq d+1 \\ n \geq n_{\xi, \mathcal{A}, \varepsilon}}} \left| \left\{ t \geq 1 : \begin{aligned} & \hat{\mu}_a(t) < \mu_a - \varepsilon, \quad \forall a \in \mathcal{A}' \\ & 1 \leq N_{\mathcal{A}'}(t) \leq |\mathcal{A}'| f_\xi(n) / 2\varepsilon^2 \\ & \sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\hat{\mu}_a(t) | \mu_a - \varepsilon) \geq (d+1) \vee (f_\xi(n) - f_\xi(N_{\mathcal{A}'}(t))) \end{aligned} \right\} \right|,
\end{aligned}$$

where  $n_{\xi, \mathcal{A}, \varepsilon} = \left( e^2 \vee \frac{|\mathcal{A}|}{2\varepsilon^2} \right) \left( f_\xi \left( e^2 \vee \frac{|\mathcal{A}|}{2\varepsilon^2} \right) \right)^2$ . The proof ends by taking the expectation on both sides of previous inequality and by applying Theorem 4.  $\square$

#### 4.5.4 Non-reliable current best arm and current informative sets of arms

For accuracy  $\varepsilon > 0$ , let  $\mathcal{M}^*(\varepsilon)$  be the set of times  $t \geq 1$  that do not belong to  $\cup_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} \mathcal{E}_{\hat{a}^*, a_{t+1}}^+(\varepsilon)$  and where some of the current best arms do not belong to  $\mathcal{A}^*$ ,

$$\mathcal{M}^*(\varepsilon) := \left\{ t \geq 1 : \begin{aligned} & (1) \quad t \notin \bigcup_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} \mathcal{E}_{\hat{a}^*, a_{t+1}}^+(\varepsilon) \\ & (2) \quad \hat{\mathcal{A}}^*(t) \neq \{a^*\} \end{aligned} \right\}. \quad (4.67)$$

For accuracy  $\varepsilon > 0$ , let  $\mathcal{M}(\varepsilon)$  be the set of times  $t \geq 1$  during exploration phases where the current means of current informative arms are not well  $\varepsilon$ -estimated,

$$\mathcal{M}(\varepsilon) := \left\{ t \geq 1 : \begin{aligned} & (1) \quad \bar{a}_t \notin \hat{\mathcal{A}}^*(t), \quad a_{t+1} = \underline{a}_t = a_t^{\text{opt}} \\ & (2) \quad \exists a' \notin \hat{\mathcal{A}}^*(t), \quad |\hat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon + \hat{\varepsilon}_{a'}^*(t) \end{aligned} \right\}. \quad (4.68)$$

**Lemma 9** (Relation between subsets of times). *For accuracy  $0 < \varepsilon < \varepsilon_0 := \delta_{\min}/3$ ,*

$$\mathcal{M}^*(\varepsilon) \subset \mathcal{K}^*(\varepsilon_0) \cup \mathcal{L}^{(d)}(\varepsilon_0), \quad (4.69)$$

$$\mathcal{M}(\varepsilon) \subset \bigcup_{a \in \mathcal{A}, t \geq 1} \mathcal{K}_{a, a_{t+1}}(\varepsilon). \quad (4.70)$$

*Proof.* We start by proving Equation (4.69).

Let us consider  $t \in \mathcal{M}^*(\varepsilon)$ . Then there exists  $\hat{a}^* \in \hat{\mathcal{A}}^*(t) \setminus \{a^*\}$  and

$$\hat{\mu}_{\hat{a}^*}(t) = \hat{\mu}^*(t) \geq \hat{\mu}_{a^*}(t). \quad (4.71)$$

Since  $t \in \mathcal{M}^*(\varepsilon)$ ,  $t \notin \mathcal{E}_{\hat{a}^*, a_{t+1}}^+(\varepsilon)$ . By considering empirical lower bounds (4.17) and Equation (4.47), we have

$$\mu_{\hat{a}^*} + \varepsilon \geq \hat{\mu}_{\hat{a}^*}(t). \quad (4.72)$$

By combining Equations (4.71) and (4.72), it comes

$$\mu_{\hat{a}^*} + \varepsilon \geq \hat{\mu}_{a^*}(t). \quad (4.73)$$

Since  $\hat{a}^* \notin \mathcal{A}^*$ ,  $\varepsilon < \varepsilon_0 \leq \Delta_{\hat{a}^*}/2$  and

$$\mu^* - \varepsilon_0 > \mu_{\hat{a}^*} + \varepsilon \geq \hat{\mu}_{\hat{a}^*}(t). \quad (4.74)$$

Then Equation (4.73) implies

$$\mu_{a^*} - \varepsilon_0 > \hat{\mu}_{a^*}(t). \quad (4.75)$$

Since  $t \in \mathcal{M}^*(\varepsilon)$ ,  $t \notin \cup_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} \mathcal{E}_{\hat{a}^*, a_{t+1}}^+(\varepsilon)$ . By considering empirical lower bounds (4.17) and Equation (4.47), Equation (4.75) implies

$$a^* \notin \hat{\mathcal{A}}^*(t). \quad (4.76)$$

Let us then consider  $a \in \hat{\mathcal{A}}_{a^*}(t)$ . From Equation (4.8), this means

$$\begin{aligned} \hat{\mu}_a(t) &\leq \hat{\mu}^*(t) - \hat{\theta}_{a^*, a}^* \\ &= \hat{\mu}_{\hat{a}^*}(t) - \theta_{a^*, a}^{(\hat{a}^*)} \end{aligned} \quad (4.77)$$

From Assumption 8, previous Equation (4.77) implies

$$\hat{\mu}_a(t) \leq \hat{\mu}_{\hat{a}^*}(t) - \theta_{a^*, a}^{(a^*)} = \hat{\mu}_{\hat{a}^*}(t) - \theta_{a^*, a}^* \quad (4.78)$$

By combining Equations (4.73) and (4.78), it comes

$$\hat{\mu}_a(t) \leq \hat{\mu}^*(t) - \hat{\theta}_{a^*, a}^* < \mu_{\hat{a}^*} - \theta_{a^*, a}^* + \varepsilon. \quad (4.79)$$

From Equation (4.1) and Assumption 8, we have

$$\mu_{\hat{a}^*} - \theta_{a^*, a}^* \leq \mu_a - 2\varepsilon_0. \quad (4.80)$$

By combining Equations (4.79) and (4.80), we show that

$$\hat{\mu}_a(t) \leq \hat{\mu}^*(t) - \hat{\theta}_{a^*, a}^* \leq \mu_a - \varepsilon_0, \quad \forall a \in \hat{\mathcal{A}}_{a^*}(t) \supset \hat{\mathcal{A}}_{a^*}^{(d)}(t). \quad (4.81)$$

Since  $a^* \notin \hat{\mathcal{A}}^*(t)$ , from empirical lower bounds (4.16) we have

$$\mathbf{f}_\xi(N_{a_{t+1}}(t)) \leq \sum_{a \in \hat{\mathcal{A}}_{a^*}^{(d)}(t)} N_a(t) \text{kl}(\hat{\mu}_a(t) | \hat{\mu}^*(t) - \theta_{a^*, a}^*) + \mathbf{f}_\xi\left(N_{\hat{\mathcal{A}}_{a^*}^{(d)}(t)}(t)\right). \quad (4.82)$$

The monotony of  $\text{kl}(\hat{\mu}_a(t) | \cdot)$ , for  $a \in \hat{\mathcal{A}}_{a^*}^{(d)}(t)$  and Equation (4.81) imply

$$\text{kl}(\hat{\mu}_a(t) | \hat{\mu}^*(t) - \theta_{a^*, a}^*) \leq \text{kl}(\hat{\mu}_a(t) | \mu_a - \varepsilon_0), \quad \forall a \in \hat{\mathcal{A}}_{a^*}^{(d)}(t). \quad (4.83)$$

By combining Equations (4.82) and (4.83), we get

$$\mathbf{f}_\xi(N_{a_{t+1}}(t)) \leq \sum_{a \in \hat{\mathcal{A}}_{a^*}^{(d)}(t)} N_a(t) \text{kl}(\hat{\mu}_a(t) | \mu_a - \varepsilon_0) + \mathbf{f}_\xi\left(N_{\hat{\mathcal{A}}_{a^*}^{(d)}(t)}(t)\right) \quad (4.84)$$

or equivalently

$$\sum_{a \in \widehat{\mathcal{A}}_{a^*}^{(d)}(t)} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon_0) \geq f_\xi(N_{a_{t+1}}(t)) - f_\xi(N_{\widehat{\mathcal{A}}_{a^*}^{(d)}(t)}(t)). \quad (4.85)$$

According to the definition of  $\widehat{\mathcal{A}}_{a^*}^{(d)}(t)$ , if  $|\widehat{\mathcal{A}}_{a^*}^{(d)}(t)| > d+1$  we have  $\sum_{a \in \widehat{\mathcal{A}}_{a^*}^{(d)}(t)} 2(\widehat{\mu}^*(t) - \widehat{\theta}_{a^*,a}^* - \widehat{\mu}_a(t))^2 N_a(t) \geq \Phi(2|\widehat{\mathcal{A}}_{a^*}^{(d)}(t)| + 1)$ . Then, from Equations (4.81), (4.85) and Pinsker's inequality (Lemma 24), this implies  $\sum_{a \in \widehat{\mathcal{A}}_{a^*}^{(d)}(t)} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon_0) \geq \Phi(2|\widehat{\mathcal{A}}_{a^*}^{(d)}(t)| + 1) \sqrt{f_\xi(N_{a_{t+1}}(t)) - f_\xi(N_{\widehat{\mathcal{A}}_{a^*}^{(d)}(t)}(t))}$ . Furthermore, if  $|\widehat{\mathcal{A}}_{a^*}^{(d)}(t)| \leq d+1$  and  $N_{\widehat{\mathcal{A}}_{a^*}^{(d)}(t)}(t) \leq e^{-d-1} N_{a_{t+1}}(t)$ , we have  $\sum_{a \in \widehat{\mathcal{A}}_{a^*}^{(d)}(t)} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon_0) \geq f_\xi(N_{a_{t+1}}(t)) - f_\xi(N_{\widehat{\mathcal{A}}_{a^*}^{(d)}(t)}(t)) \geq \log(N_{a_{t+1}}(t)) - \log(N_{\widehat{\mathcal{A}}_{a^*}^{(d)}(t)}(t)) \geq d+1$ .

Now we prove Equation (4.70).

Since  $t \in \mathcal{M}(\varepsilon)$ , there exists  $a' \notin \widehat{\mathcal{A}}^*(t)$ , such that

$$|\widehat{\mu}_{a'}(t) - \mu_{a'}| \geq \varepsilon + \widehat{\varepsilon}_{a'}^*(t). \quad (4.86)$$

From Equation (4.86), by triangle inequality we have

$$|\widehat{\mu}_{a'}(t) - \mu_{a'} - \varepsilon| \geq \widehat{\varepsilon}_{a'}^*(t). \quad (4.87)$$

Let us introduce

$$\widehat{\mathbf{k}}_{a'}(\varepsilon, t) := \mathbb{I}_{\{\widehat{\mu}_a(t) < \mu_{a'} - \varepsilon\}} \text{kl}(\widehat{\mu}_{a'}(t) | \mu_{a'} - \varepsilon) + \mathbb{I}_{\{1 - \widehat{\mu}_{a'}(t) < 1 - \mu_{a'} - \varepsilon\}} \text{kl}(1 - \widehat{\mu}_{a'}(t) | 1 - \mu_{a'} - \varepsilon) \quad (4.88)$$

$$\widehat{\mathbf{k}}_{a'}(t) := 2(\widehat{\varepsilon}_{a'}^*(t))^2. \quad (4.89)$$

From Equation (4.88), by Pinsker's inequality we have

$$\widehat{\mathbf{k}}_{a'}(\varepsilon, t) \geq 2(\widehat{\mu}_{a'}(t) - \mu_{a'} - \varepsilon)^2. \quad (4.90)$$

By combining Equations (4.87) and (4.90) and Equation (4.89) that defines  $\widehat{\mathbf{k}}_{a'}(t)$ , we obtain

$$\widehat{\mathbf{k}}_{a'}(\varepsilon, t) \geq 2(\widehat{\varepsilon}_{a'}^*(t))^2 = \widehat{\mathbf{k}}_{a'}^*(t). \quad (4.91)$$

Since  $a' \in \widehat{\mathcal{A}}^*(t)$ , by combining empirical lower bounds (4.18) and Equation (4.91), we have

$$f_\xi(N_{a_t}(t)) \leq N_{a'}(t) \widehat{\mathbf{k}}_{a'}(\varepsilon, t) + f_\xi(N_{a'}(t)). \quad (4.92)$$

Since  $t \in \mathcal{M}(\varepsilon)$ ,  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$  and

$$a_{t+1} = \bar{a}_t. \quad (4.93)$$

By combining Equations (4.88), (4.92) and (4.93), we have

$$t \in \mathcal{K}_{a', a_{t+1}}^+(\varepsilon) \cup \mathcal{K}_{a', a_{t+1}}^-(\varepsilon) = \mathcal{K}_{a', a_{t+1}}(\varepsilon).$$

□

### 4.5.5 Reliable current means of current informative sets of Arms

For accuracy  $0 < \varepsilon < \varepsilon_0 = \delta_{\min}/3$ , we define the union of the previous considered subsets of times

$$\mathcal{U}(\varepsilon) := \mathcal{M}(\varepsilon) \cup \mathcal{M}^*(\varepsilon) \cup \mathcal{K}^*(\varepsilon_0) \cup \mathcal{L}^{(d)}(\varepsilon_0) \cup \bigcup_{a \in \mathcal{A}, t \geq 1} \mathcal{K}_{a, a_{t+1}}(\varepsilon) \cup \bigcup_{a \in \mathcal{A}, t \geq 1} \mathcal{E}_{a, a_{t+1}}(\varepsilon). \quad (4.94)$$

From Lemma 9, we have

$$\mathcal{U}(\varepsilon) = \mathcal{K}^*(\varepsilon_0) \cup \mathcal{L}^{(d)}(\varepsilon_0) \cup \bigcup_{a \in \mathcal{A}, t \geq 1} \mathcal{K}_{a, a_{t+1}}(\varepsilon) \cup \bigcup_{a \in \mathcal{A}, t \geq 1} \mathcal{E}_{a, a_{t+1}}(\varepsilon). \quad (4.95)$$

**Lemma 10** (Reliable current mean of current pulled arm). *For accuracy  $0 < \varepsilon < \varepsilon_0 = \delta_{\min}/3$ , for all  $t \notin \mathcal{U}(\varepsilon)$ ,*

$$|\hat{\mu}_{a_{t+1}}(t) - \mu_{a_{t+1}}| < \varepsilon.$$

*Proof.* Let us consider  $t \notin \mathcal{U}(\varepsilon)$ . Then  $t \notin \mathcal{E}_{a_{t+1}, a_{t+1}}(\varepsilon)$ . □

**Lemma 11** (Reliable current best mean). *For accuracy  $0 < \varepsilon < \varepsilon_0 = \delta_{\min}/3$ , for all  $t \notin \mathcal{U}(\varepsilon)$ ,*

$$\hat{\mathcal{A}}^*(t) = \{a^*\},$$

$$|\hat{\mu}^*(t) - \mu^*| < \varepsilon.$$

*Proof.* Let us consider  $t \notin \mathcal{U}(\varepsilon)$ . Then  $t \notin \mathcal{M}^*(\varepsilon) \cup \bigcup_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} \mathcal{E}_{\hat{a}^*, a_{t+1}}^+(\varepsilon)$  and Equation (4.67) implies

$$\hat{\mathcal{A}}^*(t) = \{a^*\}. \quad (4.96)$$

Since  $t \notin \mathcal{U}(\varepsilon)$ ,  $t \notin \bigcup_{\hat{a}^* \in \hat{\mathcal{A}}^*(t)} \mathcal{E}_{\hat{a}^*, a_{t+1}}(\varepsilon)$ . By considering empirical lower bounds (4.17) and Equation (4.48), we have

$$|\hat{\mu}_{\hat{a}^*}(t) - \mu_{\hat{a}^*}| < \varepsilon, \quad \forall \hat{a}^* \in \hat{\mathcal{A}}^*(t). \quad (4.97)$$

Since  $\hat{\mathcal{A}}^*(t) = \operatorname{argmax}_{a \in \mathcal{A}} \hat{\mu}_a(t)$ , by combining Equations (4.96) and (4.97), we have

$$|\hat{\mu}^*(t) - \mu^*| < \varepsilon. \quad (4.98)$$

□

**Lemma 12** (Reliable current means). *For accuracy  $0 < \varepsilon < \varepsilon_0 = \delta_{\min}/3$ , for all time step  $t \notin \mathcal{U}(\varepsilon)$  that corresponds to an exploration phase (that is  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$ ) such that  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ , for all current sub-optimal arm  $a' \notin \hat{\mathcal{A}}^*(t)$ ,*

$$|\hat{\mu}_{a'}(t) - \mu_{a'}| \leq \varepsilon + \hat{\varepsilon}_{a'}^*(t)$$

and

$$\hat{\varepsilon}_{a'}^*(t) \leq \varepsilon_{a'}(t),$$

where

$$\hat{\varepsilon}_{a'}^*(t) \leq \frac{5}{3} \cdot \frac{\gamma_t}{1 - \gamma_t} \cdot \Delta_{a'}$$

$$\varepsilon_{a'}(t) = \frac{5}{3} \cdot \frac{\gamma_t}{1 - \gamma_t} \cdot \Delta_{a'} \leq \frac{\Delta_{a'}}{3}.$$

*Proof.* Let us consider  $t \notin \mathcal{U}(\varepsilon)$  such that  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$  and  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ . Let us consider a current sub-optimal arm  $a' \notin \hat{\mathcal{A}}^*(t)$ . Then  $t \notin \mathcal{M}(\varepsilon)$ ,  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$  and Equation (4.68) that defines  $\mathcal{M}(\varepsilon)$  implies

$$|\hat{\mu}_{a'}(t) - \mu_{a'}| < \varepsilon + \hat{\varepsilon}_{a'}^*(t). \quad (4.99)$$

Since  $t \notin \mathcal{U}(\varepsilon)$ , from Lemma 11 we have

$$|\hat{\mu}^*(t) - \mu^*| < \varepsilon. \quad (4.100)$$

We have by definition

$$\hat{\varepsilon}_{a'}^*(t) = \gamma_t \cdot (\hat{\mu}^*(t) - \hat{\mu}_{a'}(t)). \quad (4.101)$$

By combining Equations (4.99), (4.100) and (4.101), we have

$$\hat{\varepsilon}_{a'}^*(t) \leq \gamma_t \cdot (\mu^* - \mu_{a'} + 2\varepsilon + \hat{\varepsilon}_{a'}^*(t)) = \gamma_t \cdot (\Delta_{a'} + 2\varepsilon + \hat{\varepsilon}_{a'}^*(t)). \quad (4.102)$$

Equation (4.102) implies

$$\hat{\varepsilon}_{a'}^*(t) \leq \frac{\gamma_t}{1 - \gamma_t} \cdot (\Delta_{a'} + 2\varepsilon). \quad (4.103)$$

Since  $0 < \varepsilon < \varepsilon_0 = \delta_{\min}/3$ , Equation (4.103) implies

$$\hat{\varepsilon}_{a'}^*(t) \leq \frac{5}{3} \cdot \frac{\gamma_t}{1 - \gamma_t} \cdot \Delta_{a'}. \quad (4.104)$$

Finally we note that  $\gamma_t$  is set so that

$$\frac{\gamma_t}{1 - \gamma_t} < 1/5. \quad (4.105)$$

□

**Lemma 13** (Reliable current informative sets of arms). *For accuracy  $0 < \varepsilon < \varepsilon_0 = \delta_{\min}/3$ , for all time step  $t \notin \mathcal{U}(\varepsilon)$  that corresponds to an exploration phase, that is  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$ , such that  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ , for all current sub-optimal arm  $a \notin \hat{\mathcal{A}}^*(t)$ ,*

$$\{a\} \subset \mathcal{A}_a(\varepsilon, t) \subset \hat{\mathcal{A}}_a(t),$$

where

$$\mathcal{A}_a(\varepsilon, t) := \left\{ a' \notin \hat{\mathcal{A}}^*(t) : 2\varepsilon + \varepsilon_{a'}(t) \leq \delta_{a,a'}(\theta^*) \right\}.$$

*Proof.* Since  $t \notin \mathcal{U}(\varepsilon)$ ,  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$  and  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ , from Lemma 12 we have

$$\varepsilon_a(t) \leq \frac{\Delta_a}{3}. \quad (4.106)$$

Since  $\varepsilon < \varepsilon_0 \leq \Delta_a/3$ , from Equation (4.106) we have

$$2\varepsilon + \varepsilon_a(t) < \Delta_a. \quad (4.107)$$

Since  $\delta_{a,a}(\theta^*) = \Delta_a$ , Equation (4.107) implies  $a \in \mathcal{A}_a(\varepsilon, t) \neq \emptyset$ .

Let us consider  $a' \in \mathcal{A}_a(\varepsilon, t)$ . Then we have

$$2\varepsilon + \varepsilon_{a'}(t) < \delta_{a,a'}(\theta) = \mu^* - \mu_{a'} - \theta_{a,a'}^*(t), \quad (4.108)$$

that is

$$\mu_{a'} + \varepsilon + \varepsilon_{a'}(t) < \mu^* - \varepsilon - \theta_{a,a'}^*(t). \quad (4.109)$$

Since  $t \notin \mathcal{U}(\varepsilon)$ , from Lemma 11 we have

$$\begin{aligned}\theta^* &= \theta^{(a^*)} = \widehat{\theta}^*(t) \\ \mu^* - \varepsilon &< \widehat{\mu}^*(t).\end{aligned}\tag{4.110}$$

Since  $t \notin \mathcal{U}(\varepsilon)$ ,  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$  and  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ . From Lemma 12 then we have

$$\mu_{a'} + \varepsilon + \varepsilon_{a'}(t) > \widehat{\mu}_{a'}(t).\tag{4.111}$$

By combining Equations (4.109), (4.110) and (4.111), we have

$$\widehat{\mu}_{a'}(t) \leq \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t),\tag{4.112}$$

that is  $a' \in \widehat{\mathcal{A}}_a(t)$ . □

## 4.6 Upper bounds under IMED-GS algorithm

In this section, we now provide the final results related to the control of the numbers of pulls of sub-optimal arms. We start by providing in Section 4.6.1 an upper-bound on the number of pulls  $N_a^{\text{opt}}(t)$  solution to the empirical optimization problem. Then, in Section 4.6.2, we provide an almost sure upper bound on the number of pulls of sub-optimal arms, making appear the size of the random set of times  $\mathcal{U}(\varepsilon)$  introduced in the previous section. This is further used in Section 4.6.3 to provide a fully explicit, finite-time upper bound on the regret valid almost surely. Last, in Section 4.6.4, we handle the size of the random set thanks to the results of Section 4.5 (in particular, Lemma 8), and derive the main regret bound.

### 4.6.1 Upper bounds on the optimal numbers of pulls

**Lemma 14** (Upper bounds on optimal numbers of pulls). *For accuracy  $0 < \varepsilon < \varepsilon_0 = \delta_{\min}/3$ , for all time step  $t \notin \mathcal{U}(\varepsilon)$  that corresponds to an exploration phase, that is  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$ , such that  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ ,*

$$\begin{aligned}\mathfrak{C}_{\{\widehat{\theta}^*(t)\}}(\widehat{\mu}(t)) &\leq \max_{a \notin \mathcal{A}^*} \frac{\Delta_a + 2\varepsilon + \varepsilon_a(t)}{\Delta_a - 2\varepsilon - \varepsilon_a(t)} \cdot \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t)), \\ \sum_{a \notin \mathcal{A}^*} N_a^{\text{opt}}(t) \Delta_a &\leq \max_{a \notin \mathcal{A}^*} \left( \frac{\Delta_a + 2\varepsilon + \varepsilon_a(t)}{\Delta_a - 2\varepsilon - \varepsilon_a(t)} \right)^2 \cdot \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t)) \cdot f_\xi(t),\end{aligned}$$

where for all  $a \in \mathcal{A}$ ,

$$\mu_a(\varepsilon, t) = \begin{cases} \mu_a + \varepsilon + \varepsilon_a(t) & , \text{ if } a \notin \mathcal{A}^* \\ \mu_a - \varepsilon & , \text{ if } a \in \mathcal{A}^*. \end{cases}$$

*Proof.* Let us consider  $t \notin \mathcal{U}(\varepsilon)$  such that  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$ ,  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$  and  $n \in \mathbb{R}_+^{\mathcal{A}}$  such that

$$\forall a \notin \text{argmax}(\mu(\varepsilon, t)), \quad \sum_{\substack{a' \in \mathcal{A} \\ \mu_{a'}(\varepsilon, t) \leq \max((\mu(\varepsilon, t)) - \theta_{a,a'}^*)}} \text{kl}(\mu_{a'}(\varepsilon, t) | \max(\mu(\varepsilon, t)) - \theta_{a,a'}^*) n_{a'} \geq 1.\tag{4.113}$$

Since  $t \notin \mathcal{U}(\varepsilon)$ , from Lemma 11 and Assumptions 7-6 we have

$$\begin{aligned}\widehat{\mathcal{A}}^*(t) &= \mathcal{A}^* = \{a^*\} \\ \theta^* &= \theta^{(a^*)} = \widehat{\theta}^*(t)\end{aligned}\tag{4.114}$$

$$|\widehat{\mu}^*(t) - \mu^*| < \varepsilon. \quad (4.115)$$

Since  $t \notin \mathcal{U}(\varepsilon)$  is such that  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$  and  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ , from Lemma 12 and Equation (4.114) we have

$$|\widehat{\mu}_a(t) - \mu_a| \leq \varepsilon + \varepsilon_a(t), \quad \forall a \notin \mathcal{A}^*, \quad (4.116)$$

$$\varepsilon_a(t) \leq \frac{\Delta_a}{3}, \quad \forall a \notin \mathcal{A}^*. \quad (4.117)$$

From Equations (4.114), (4.115) and (4.117) we have

$$\operatorname{argmax}(\mu(\varepsilon, t)) = \mathcal{A}^* = \{a^*\}. \quad (4.118)$$

Since  $t \notin \mathcal{U}(\varepsilon)$  is such that  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$  and  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ , from Lemma 13 and Equation (4.114) we have for all sub-optimal arm  $a \notin \widehat{\mathcal{A}}^*(t)$ ,

$$\{a\} \subset \mathcal{A}_a(\varepsilon, t) \subset \widehat{\mathcal{A}}_a(t), \quad (4.119)$$

where

$$\begin{aligned} \mathcal{A}_a(\varepsilon, t) &= \{a' \notin \mathcal{A}^* : 2\varepsilon + \varepsilon_{a'}(t) \leq \delta_{a,a'}(\theta^*)\} \\ &= \left\{ a' \notin \mathcal{A}^* : \mu_{a'}(\varepsilon, t) \leq \max((\mu(\varepsilon, t)) - \widehat{\theta}_{a,a'}^*(t)) \right\}. \end{aligned} \quad (4.120)$$

Then, Equation (4.113) rewrites

$$\forall a \notin \mathcal{A}^*, \quad \sum_{a' \in \mathcal{A}_a(\varepsilon, t)} \operatorname{kl}(\mu_{a'} + \varepsilon + \varepsilon_{a'}(t) | \mu^* - \varepsilon - \widehat{\theta}_{a,a'}^*(t)) n_{a'} \geq 1. \quad (4.121)$$

From Equations (4.114), (4.116) and standard monotonic properties of  $\operatorname{kl}(\cdot | \cdot)$ , we have

$$\begin{aligned} \forall a \notin \mathcal{A}^*, \forall a' \in \mathcal{A}_a(\varepsilon, t), \\ \operatorname{kl}(\mu_{a'} + \varepsilon + \varepsilon_{a'}(t) | \mu^* - \varepsilon - \widehat{\theta}_{a,a'}^*(t)) \leq \operatorname{kl}(\widehat{\mu}_{a'}(t) | \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)). \end{aligned} \quad (4.122)$$

By combining Equations (4.121) and (4.122) we have

$$\forall a \notin \mathcal{A}^*, \quad \sum_{a' \in \mathcal{A}_a(\varepsilon, t)} \operatorname{kl}(\widehat{\mu}_{a'}(t) | \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)) n_{a'} \geq 1. \quad (4.123)$$

Then, by combining Equations (4.119) and (4.123) we have

$$\forall a \notin \mathcal{A}^*, \quad \sum_{a' \in \widehat{\mathcal{A}}_a(t)} \operatorname{kl}(\widehat{\mu}_{a'}(t) | \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)) n_{a'} \geq 1. \quad (4.124)$$

Since

$$\begin{aligned} \mathfrak{C}_{\{\widehat{\theta}^*(t)\}}(\widehat{\mu}(t)) &= \min_{n \in \mathbb{R}_+^{\mathcal{A}}} \sum_{a \in \mathcal{A}} n_a (\widehat{\mu}^*(t) - \widehat{\mu}_a(t)) \\ \text{s.t.} \quad \forall a \notin \widehat{\mathcal{A}}^*(t) = \mathcal{A}^*, \quad &\sum_{a' \in \widehat{\mathcal{A}}_a(t)} \operatorname{kl}(\widehat{\mu}_{a'}(t) | \widehat{\mu}^*(t) - \widehat{\theta}_{a,a'}^*(t)) n_{a'} \geq 1, \end{aligned} \quad (4.125)$$



Equation (4.124) implies

$$\mathfrak{C}_{\{\hat{\theta}^*(t)\}}(\hat{\mu}(t)) \leq \sum_{a \in \mathcal{A}} n_a (\hat{\mu}^*(t) - \hat{\mu}_a(t)). \quad (4.126)$$

By combining Equations (4.114), (4.115), (4.116), (4.117) and (4.126) we have

$$\mathfrak{C}_{\{\hat{\theta}^*(t)\}}(\hat{\mu}(t)) \leq \sum_{a \notin \mathcal{A}^*} n_a (\mu^* - \mu_a + 2\varepsilon + \varepsilon_a(t)) = \sum_{a \notin \mathcal{A}^*} n_a (\Delta_a + 2\varepsilon + \varepsilon_a(t)). \quad (4.127)$$

Then Equations (4.118) and (4.127) imply

$$\mathfrak{C}_{\{\hat{\theta}^*(t)\}}(\hat{\mu}(t)) \leq \max_{a \notin \mathcal{A}^*} \frac{\Delta_a + 2\varepsilon + \varepsilon_a(t)}{\Delta_a - 2\varepsilon - \varepsilon_a(t)} \cdot \sum_{a \in \mathcal{A}} n_a (\max(\mu(\varepsilon, t)) - \mu_a(\varepsilon, t)). \quad (4.128)$$

Since previous Equation (4.128) is satisfied for all  $n \in \mathbb{R}_+^{\mathcal{A}}$  satisfying Equation (4.113), we have

$$\mathfrak{C}_{\{\hat{\theta}^*(t)\}}(\hat{\mu}(t)) \leq \max_{a \notin \mathcal{A}^*} \frac{\Delta_a + 2\varepsilon + \varepsilon_a(t)}{\Delta_a - 2\varepsilon - \varepsilon_a(t)} \cdot \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t)). \quad (4.129)$$

Furthermore, by combining Equations (4.114), (4.115), (4.116) and (4.117) we have

$$\begin{aligned} \sum_{a \notin \mathcal{A}^*} N_a^{\text{opt}}(t) \Delta_a &\leq \max_{a \notin \mathcal{A}^*} \left( \frac{\Delta_a}{\Delta_a - 2\varepsilon - \varepsilon_a(t)} \right) \sum_{a \notin \hat{\mathcal{A}}^*(t)} N_a^{\text{opt}}(t) \hat{\Delta}_a(t) \\ &\leq \max_{a \notin \mathcal{A}^*} \left( \frac{\Delta_a + 2\varepsilon + \varepsilon_a(t)}{\Delta_a - 2\varepsilon - \varepsilon_a(t)} \right) \sum_{a \notin \hat{\mathcal{A}}^*(t)} N_a^{\text{opt}}(t) \hat{\Delta}_a(t). \end{aligned} \quad (4.130)$$

Since  $(N_a^{\text{opt}}(t))_{a \in \mathcal{A}} = (\bar{I}_{a_i}(t) n_a^{\text{opt}}(t))_{a \in \mathcal{A}}$ , we have

$$\sum_{a \notin \hat{\mathcal{A}}^*(t)} N_a^{\text{opt}}(t) \hat{\Delta}_a(t) = \min_{a \in \mathcal{A}} \bar{I}_a(t) \sum_{a \notin \hat{\mathcal{A}}^*(t)} n_a^{\text{opt}}(t) \hat{\Delta}_a(t). \quad (4.131)$$

From Equation (4.9) that defines indexes  $(\bar{I}_a(t))_{a \in \mathcal{A}}$  and previous Equation (4.131) we have

$$\sum_{a \notin \hat{\mathcal{A}}^*(t)} N_a^{\text{opt}}(t) \hat{\Delta}_a(t) \leq f_\xi(t) \sum_{a \notin \hat{\mathcal{A}}^*(t)} n_a^{\text{opt}}(t) \hat{\Delta}_a(t). \quad (4.132)$$

From Equation (4.11) that defines  $(n_a^{\text{opt}}(t))_{a \in \mathcal{A}}$  we have

$$\sum_{a \notin \hat{\mathcal{A}}^*(t)} n_a^{\text{opt}}(t) \hat{\Delta}_a(t) = \mathfrak{C}_{\{\hat{\theta}^*(t)\}}(\hat{\mu}(t)). \quad (4.133)$$

By combining Equations (4.132) and (4.133) we have

$$\sum_{a \notin \hat{\mathcal{A}}^*(t)} N_a^{\text{opt}}(t) \hat{\Delta}_a(t) \leq \mathfrak{C}_{\{\hat{\theta}^*(t)\}}(\hat{\mu}(t)) \cdot f_\xi(t). \quad (4.134)$$

Combining Equations (4.129), (4.130) and (4.134) ends the proof.  $\square$

## 4.6.2 Upper bounds on the numbers of pulls

We recall that what we call optimal numbers of pulls at time step  $t \geq 1$  are the numbers of pulls  $(N_a^{\text{opt}}(t))_{a \in \mathcal{A}}$  solution of the optimisation problem.

**Lemma 15** (Upper bounds on current pulled sub-optimal arm). *For accuracy  $0 < \varepsilon < \varepsilon_0$ , for all sub-optimal  $a \notin \mathcal{A}^*$ , for all time step  $t \notin \mathcal{U}(\varepsilon)$  such that  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}} = a$ ,*

$$N_a(t) \leq N_a^{\text{opt}}(t) \leq \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a} \cdot f_\xi(t)$$

where

$$\lambda_t(\mu, \varepsilon) := \max_{a' \notin \mathcal{A}^*} \left( \frac{\Delta_{a'} + 2\varepsilon + \varepsilon_{a'}(t)}{\Delta_{a'} - 2\varepsilon - \varepsilon_{a'}(t)} \right)^2.$$

For accuracy  $0 < \varepsilon < \varepsilon_0$ , for all sub-optimal  $a \notin \mathcal{A}^*$ , for all time step  $t \notin \mathcal{U}(\varepsilon)$  such that  $a_{t+1} = a \neq \underline{a}_t$ ,

$$N_a(t) \leq \frac{f_\xi(t)}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)}.$$

For accuracy  $0 < \varepsilon < \varepsilon_0$ , for all sub-optimal  $a \notin \mathcal{A}^*$ ,

$$\forall t \geq 1, \quad N_a(t) \leq \max_{1 \leq s \leq t} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, s))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} \cdot f_\xi(t) + |\mathcal{U}(\varepsilon)| + 1.$$

For accuracy  $0 < \varepsilon < \varepsilon_0 \vee (1 - \mu^*)$ , for all sub-optimal  $a \notin \mathcal{A}^*$ , for all time step  $t \notin \mathcal{U}(\varepsilon)$  such that  $a_{t+1} = a \neq \underline{a}_t$ ,

$$N_a(t) \leq \frac{\Gamma \left[ f_\xi \left( \max_{1 \leq s \leq t} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, s))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} f_\xi(t) + |\mathcal{U}(\varepsilon)| + 1 \right) + 1 \right]}{(\gamma_t)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)}.$$

*Proof.* Let us consider a time step  $t \notin \mathcal{U}(\varepsilon)$  such that  $a_{t+1} = a$ .

Since  $t \notin \mathcal{U}(\varepsilon)$ , from Lemma 11 and Assumption 7 we have

$$\widehat{\mathcal{A}}^*(t) = \mathcal{A}^* = \{a^*\}. \quad (4.135)$$

In particular, we have

$$a_{t+1} = a \notin \widehat{\mathcal{A}}^*(t). \quad (4.136)$$

and Algorithm 9 implies

$$\bar{a}_t \notin \widehat{\mathcal{A}}^*(t). \quad (4.137)$$

Since  $0 < \varepsilon < (\mu^* - \mu_a)/3$ ,  $t \notin \mathcal{U}(\varepsilon)$  and  $a_{t+1} = a$ , Lemma 12, Lemma 11 and Equation (4.135) imply

$$\widehat{\mu}_a(t) \leq \mu_a + \varepsilon < \mu^* - \varepsilon \leq \widehat{\mu}^*(t) \leq \mu^* + \varepsilon. \quad (4.138)$$

Let us consider a time step  $t \notin \mathcal{U}(\varepsilon)$  such that  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}} = a$ . Then, Lemma 7 implies

$$N_a(t) \leq N_a^{\text{opt}}(t). \quad (4.139)$$

Since  $t \notin \mathcal{U}(\varepsilon)$ , we have  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t) = \mathcal{A}^*$  as shown in Equation (4.137) and  $a_{t+1} = \underline{a}_t = a_t^{\text{opt}}$ . Lemma 14 then implies

$$\forall a' \notin \widehat{\mathcal{A}}^*(t) = \mathcal{A}^*, \quad N_{a'}^{\text{opt}}(t) \leq \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_{a'}} \cdot \mathbf{f}_\xi(t). \quad (4.140)$$

By combining Equations (4.139) and (4.140), we get

$$N_a(t) \leq N_a^{\text{opt}}(t) \leq \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a} \cdot \mathbf{f}_\xi(t).$$

Let us consider a time step  $t \notin \mathcal{U}(\varepsilon)$  such that  $a_{t+1} = a \neq \underline{a}_t$ . From Equation (4.137), we have  $\bar{a}_t \notin \widehat{\mathcal{A}}^*(t)$ . Then Algorithm 9 implies

$$a_{t+1} \in \{\bar{a}_t, \hat{a}_t\}. \quad (4.141)$$

By combining previous Equation (4.141) and Lemma 7, we have

$$N_a(t) \leq \frac{\mathbf{f}_\xi(t)}{\text{kl}(\widehat{\mu}_a(t) | \widehat{\mu}^*(t))}. \quad (4.142)$$

From Equations (4.138) and (4.142), the classical monotonic properties of  $\text{kl}(\cdot | \cdot)$  imply

$$N_a(t) \leq \frac{\mathbf{f}_\xi(t)}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)}.$$

Let us consider a time step  $t \geq 1$ . We write the number of pulls  $N_a(t)$  as follows

$$\begin{aligned} N_a(t) &= 1 + \sum_{s \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{s+1}=a\}} \\ &= 1 + \sum_{s \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{s+1}=a, s \in \mathcal{U}(\varepsilon)\}} + \sum_{s \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{s+1}=\underline{a}_s = a_s^{\text{opt}}=a, s \notin \mathcal{U}(\varepsilon)\}} + \mathbb{I}_{\{a_{s+1}=a \neq \underline{a}_s, s \notin \mathcal{U}(\varepsilon)\}}. \end{aligned} \quad (4.143)$$

We first note that

$$\sum_{s \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{s+1}=a, s \in \mathcal{U}(\varepsilon)\}} \leq |\mathcal{U}(\varepsilon)|. \quad (4.144)$$

Then, we denote by

$$\tau_a(t) = \max \{s \in [|\mathcal{A}| + 1; t] : (a_{s+1} = \underline{a}_s = a_s^{\text{opt}} = a \text{ or } a_{s+1} = a \neq \underline{a}_s) \text{ and } s \notin \mathcal{U}(\varepsilon)\} \quad (4.145)$$

the last time step before time step  $t$  that does not belong to  $\mathcal{U}(\varepsilon)$  such that we pull arm  $a$ . Then, we have

$$N_a(\tau_a(t)) \leq \max_{1 \leq s \leq t} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, s))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} \times \mathbf{f}_\xi(t). \quad (4.146)$$

Furthermore, we note that

$$\begin{aligned} &\sum_{s \geq |\mathcal{A}|}^{t-1} \mathbb{I}_{\{a_{s+1}=\underline{a}_s = a_s^{\text{opt}}=a, s \notin \mathcal{U}(\varepsilon)\}} + \mathbb{I}_{\{a_{s+1}=a \neq \underline{a}_s, s \notin \mathcal{U}(\varepsilon)\}} \\ &\leq N_a(\tau_a(t)) \\ &\leq \max_{1 \leq s \leq t} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, s))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} \times \mathbf{f}_\xi(t). \end{aligned} \quad (4.147)$$

By combining Equations (4.143), (4.144) and (4.147), we get

$$\forall a \neq a^*, \forall t \geq 1, N_a(t) \leq 1 + |\mathcal{U}(\varepsilon)| + \max_{1 \leq s \leq t} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, s))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} \times f_\xi(t). \quad (4.148)$$

Let us consider a time step  $t \notin \mathcal{U}(\varepsilon)$  such that  $a_{t+1} = a \neq \underline{a}_t$ . From Equations (4.135) and (4.137), we have  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$ . Then Algorithm 9 implies

$$\bar{a}_t, a_t^{\text{opt}} \notin \hat{\mathcal{A}}^*(t) = \{a^*\}, \quad a_{t+1} \in \{\bar{a}_t, \dot{a}_t\}. \quad (4.149)$$

By combining previous Equation (4.149) and Lemma 7, we have

$$N_a(t) \leq \left( \mathbb{I}_{\{\bar{a}_t \notin \hat{\mathcal{A}}^*(t) \text{ and } \underline{a}_t \neq a_t^{\text{opt}}\}} \frac{\bar{I}_{\bar{a}_t}(t)}{I_{\dot{a}_t}(t)} \right) \frac{f_\xi(N_{a_t^{\text{opt}}}(t)) + 1}{(\gamma_t)^2 \hat{\mu}^*(t) (1 - \hat{\mu}^*(t)) \text{kl}(\hat{\mu}_a(t) | \hat{\mu}^*(t))}. \quad (4.150)$$

From Algorithm 9, we have

$$\mathbb{I}_{\{\bar{a}_t \notin \hat{\mathcal{A}}^*(t) \text{ and } \underline{a}_t \neq a_t^{\text{opt}}\}} \frac{\bar{I}_{\bar{a}_t}(t)}{I_{\dot{a}_t}(t)} \leq \Gamma. \quad (4.151)$$

From Equations (4.138) and (4.142) and the classical monotonic properties of  $\text{kl}(\cdot | \cdot)$ , we have

$$\hat{\mu}^*(t) (1 - \hat{\mu}^*(t)) \text{kl}(\hat{\mu}_a(t) | \hat{\mu}^*(t)) \geq (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon). \quad (4.152)$$

By combining Equations (4.148), (4.150), (4.151) and (4.152), we then get

$$N_a(t) \leq \frac{\Gamma \left[ f_\xi \left( 1 + |\mathcal{U}(\varepsilon)| + \max_{1 \leq s \leq t} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, s))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} f_\xi(t) \right) + 1 \right]}{(\gamma_t)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)}. \quad (4.153)$$

□

**Remark 14.** Proof of Lemma 15 highlights that for accuracy  $0 < \varepsilon < \varepsilon_0$ , for all time step  $t \notin \mathcal{U}(\varepsilon)$ , if  $a_{t+1} \notin \mathcal{A}^*$  then  $\bar{a}_t \notin \hat{\mathcal{A}}^*(t)$ .

### 4.6.3 Almost-sure upper bound on the cumulative regret

**Proposition 4** (Upper bound on the non-averaged cumulative regret). *For accuracy  $0 < \varepsilon < \varepsilon_0$ , for all time horizon  $T \geq |\mathcal{A}|$ ,*

$$\begin{aligned} \sum_{a \notin \mathcal{A}^*} N_a(T) \Delta_a &\leq \inf_{\tau \geq 1} \{ \Lambda_\tau(\mu, \varepsilon) \cdot C_\tau(\mu, \varepsilon) \cdot f_\xi(T) + \Sigma_\Delta \cdot \tau \} \\ &\quad + \mathbf{G}_T(\mu, \varepsilon) \cdot \mathbf{F}_\xi(\mu, \varepsilon, T) + \mathbf{G}_T(\mu, \varepsilon) \cdot f_\xi(|\mathcal{U}(\varepsilon)|) + |\mathcal{U}(\varepsilon)| \end{aligned}$$

where

$$\Lambda_\tau(\mu, \varepsilon) = \sup_{t \geq \tau} \lambda_t(\mu, \varepsilon) \quad C_\tau(\mu, \varepsilon) = \sup_{t \geq \tau} \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t))$$

$$\lambda_t(\mu, \varepsilon) = \max_{a' \notin \mathcal{A}^*} \left( \frac{\Delta_{a'} + 2\varepsilon + \varepsilon_{a'}(t)}{\Delta_{a'} - 2\varepsilon - \varepsilon_{a'}(t)} \right)^2$$

$$\begin{aligned}
F_\xi(\mu, \varepsilon, T) &= \max_{a \notin \mathcal{A}^*} f_\xi \left( \max_{t \geq 1} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} f_\xi(T) + 1 \right) + 1 \\
\mathbf{G}_T(\mu, \varepsilon) &= \sum_{a \notin \mathcal{A}^*} \frac{\Delta_a \Gamma}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \\
\Sigma_\Delta &= \sum_{a \in \mathcal{A}} \Delta_a.
\end{aligned}$$

*Proof.* Let us consider  $\tau \geq 1$ . For  $a \notin \mathcal{A}^*$ , let us consider

$$\bar{\tau}_a(\varepsilon, T) = \max \{t \in \llbracket 1, T-1 \rrbracket : t \notin \mathcal{U}(\varepsilon), a_{t+1} = \underline{a}_t = a_t^{\text{opt}} = a\} \quad (4.154)$$

$$\tau_a(\varepsilon, T) = \max \{t \in \llbracket 1, T-1 \rrbracket : t \notin \mathcal{U}(\varepsilon), a_{t+1} = a \neq \underline{a}_t\}. \quad (4.155)$$

Then, the number of pulls of sub-optimal arm  $a \notin \mathcal{A}^*$  satisfies

$$\begin{aligned}
N_a(T) &\leq 1 + \sum_{t=|A|}^{T-1} \mathbb{I}_{\{t \notin \mathcal{U}(\varepsilon), a_{t+1} = a \neq \underline{a}_t\}} + \sum_{t=|A|}^{T-1} \mathbb{I}_{\{t \notin \mathcal{U}(\varepsilon), a_{t+1} = a \neq \underline{a}_t\}} + \sum_{t \geq 1} \mathbb{I}_{\{t \in \mathcal{U}(\varepsilon), a_{t+1} = a\}} \\
&\leq N_a(\bar{\tau}_a(\varepsilon, T)) + N_a(\tau_a(\varepsilon, T)) + \sum_{t \geq 1} \mathbb{I}_{\{t \in \mathcal{U}(\varepsilon), a_{t+1} = a\}}.
\end{aligned} \quad (4.156)$$

In particular, previous Equation (4.156) implies

$$N_a(T) \leq \mathbb{I}_{\{\tau \leq \bar{\tau}_a(\varepsilon, T)\}} \cdot N_a(\bar{\tau}_a(\varepsilon, T)) + \tau + N_a(\tau_a(\varepsilon, T)) + \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1} = a, t \in \mathcal{U}(\varepsilon)\}}. \quad (4.157)$$

Since  $\bar{\tau}_a(\varepsilon, T) \notin \mathcal{U}(\varepsilon)$  and  $a_{\bar{\tau}_a(\varepsilon, T)+1} = \underline{a}_{\bar{\tau}_a(\varepsilon, T)} = a_{\bar{\tau}_a(\varepsilon, T)}^{\text{opt}} = a$ , Lemma 15 implies

$$N_a(\bar{\tau}_a(\varepsilon, T)) \leq N_a^{\text{opt}}(\bar{\tau}_a(\varepsilon, T)). \quad (4.158)$$

Since  $\bar{\tau}_a(\varepsilon, T) \notin \mathcal{U}(\varepsilon)$  and  $a_{\bar{\tau}_a(\varepsilon, T)+1} = a \notin \mathcal{A}^*$ , Remark 14 implies

$$\bar{\tau}_a(\varepsilon, T) \leq \bar{\tau}(\varepsilon, T) := \max \{t \in \llbracket 1, T-1 \rrbracket : t \notin \mathcal{U}(\varepsilon), a_{t+1} = \underline{a}_t = a_t^{\text{opt}} \notin \hat{\mathcal{A}}^*(t)\}. \quad (4.159)$$

By combining Equations (4.158) and (4.159) we get

$$\bar{\tau}_a(\varepsilon, T) \leq \bar{\tau}(\varepsilon, T) \quad N_a(\bar{\tau}_a(\varepsilon, T)) \leq N_a^{\text{opt}}(\bar{\tau}(\varepsilon, T)). \quad (4.160)$$

Since  $\bar{\tau}(\varepsilon, T) \notin \mathcal{U}(\varepsilon)$  and  $\bar{a}_{\bar{\tau}(\varepsilon, T)} \notin \hat{\mathcal{A}}^*(\bar{\tau}(\varepsilon, T))$ , from Lemma 14 we have

$$\mathbb{I}_{\{\tau \leq \bar{\tau}(\varepsilon, T) \leq T\}} \sum_{a \notin \mathcal{A}^*} N_a^{\text{opt}}(\bar{\tau}(\varepsilon, T)) \Delta_a \leq \Lambda_\tau(\mu, \varepsilon) \cdot C_\tau(\mu, \varepsilon) \cdot f_\xi(T). \quad (4.161)$$

Since  $\tau_a(\varepsilon, T) \notin \mathcal{U}(\varepsilon)$  and  $a_{\tau_a(\varepsilon, T)+1} = a \neq \underline{a}_{\tau_a(\varepsilon, T)}$ , Lemma 15 plus the monotony of  $(\gamma_t)_{t \geq 1}$ ,  $(f_\xi(t))_{t \geq 1}$  imply

$$\begin{aligned}
&N_a(\tau_a(\varepsilon, T)) \\
&\leq \frac{\Gamma \left[ f_\xi \left( \max_{t \geq 1} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} f_\xi(T) + |\mathcal{U}(\varepsilon)| + 1 \right) + 1 \right]}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \\
&\leq \frac{\Gamma \left[ f_\xi \left( \max_{t \geq 1} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathfrak{C}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} f_\xi(T) + 1 \right) + 1 \right]}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \\
&\quad + \frac{\Gamma}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} f_\xi(|\mathcal{U}(\varepsilon)|).
\end{aligned} \quad (4.162)$$

By combining Equations (4.157), (4.158), (4.160) and (4.162), we have for all sub-optimal arm  $a \notin \mathcal{A}^*$ ,

$$\begin{aligned}
& \Delta_a N_a(T) \\
\leq & \mathbb{I}_{\{\tau \leq \bar{\tau}(\varepsilon, T) \leq T\}} \Delta_a N_a^{\text{opt}}(\bar{\tau}(\varepsilon, T)) \\
& + \frac{\Delta_a \Gamma \left[ \mathbf{f}_\xi \left( \max_{t \geq 1} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathbf{e}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} \mathbf{f}_\xi(T) + 1 \right) + 1 \right]}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \\
& + \Delta_a \tau \\
& + \frac{\Delta_a \Gamma}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \mathbf{f}_\xi(|\mathcal{U}(\varepsilon)|) + \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a, t \in \mathcal{U}(\varepsilon)\}}.
\end{aligned} \tag{4.163}$$

By summing previous Equation (4.163) over the sub-optimal arms, we obtain

$$\begin{aligned}
& \sum_{a \notin \mathcal{A}^*} \Delta_a N_a(T) \\
\leq & \mathbb{I}_{\{\tau \leq \bar{\tau}(\varepsilon, T) \leq T\}} \sum_{a \notin \mathcal{A}^*} \Delta_a N_a^{\text{opt}}(\bar{\tau}(\varepsilon, T)) \\
& + \sum_{a \notin \mathcal{A}^*} \frac{\Delta_a \Gamma \left[ \mathbf{f}_\xi \left( \max_{t \geq 1} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathbf{e}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} \mathbf{f}_\xi(T) + 1 \right) + 1 \right]}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \\
& + \sum_{a \notin \mathcal{A}^*} \Delta_a \cdot \tau \\
& + \sum_{a \notin \mathcal{A}^*} \frac{\Delta_a \Gamma}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \mathbf{f}_\xi(|\mathcal{U}(\varepsilon)|) + \sum_{t \geq 1} \sum_{a \notin \mathcal{A}^*} \mathbb{I}_{\{a_{t+1}=a, t \in \mathcal{U}(\varepsilon)\}}.
\end{aligned} \tag{4.164}$$

Since

$$\sum_{t \geq 1} \sum_{a \notin \mathcal{A}^*} \mathbb{I}_{\{a_{t+1}=a, t \in \mathcal{U}(\varepsilon)\}} = \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1} \notin \mathcal{A}^*, t \in \mathcal{U}(\varepsilon)\}} \leq \sum_{t \geq 1} \mathbb{I}_{\{t \in \mathcal{U}(\varepsilon)\}} = |\mathcal{U}(\varepsilon)| \tag{4.165}$$

and

$$\begin{aligned}
\mathbf{F}_\xi(\mu, \varepsilon, T) &= \max_{a \notin \mathcal{A}^*} \mathbf{f}_\xi \left( \max_{t \geq 1} \left\{ \frac{\lambda_t(\mu, \varepsilon) \mathbf{e}_{\{\theta^*\}}(\mu(\varepsilon, t))}{\Delta_a}, \frac{1}{\text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \right\} \mathbf{f}_\xi(T) + 1 \right) + 1 \\
\mathbf{G}_T(\mu, \varepsilon) &= \sum_{a \notin \mathcal{A}^*} \frac{\Delta_a \Gamma}{(\gamma_T)^2 (\mu^* - \varepsilon) (1 - \mu^* - \varepsilon) \text{kl}(\mu_a + \varepsilon | \mu^* - \varepsilon)} \\
\Sigma_\Delta &= \sum_{a \notin \mathcal{A}^*} \Delta_a,
\end{aligned} \tag{4.166}$$

Equation (4.164) implies

$$\begin{aligned}
& \sum_{a \notin \mathcal{A}^*} \Delta_a N_a(T) \\
\leq & \mathbb{I}_{\{\tau \leq \bar{\tau}(\varepsilon, T) \leq T\}} \sum_{a \notin \mathcal{A}^*} \Delta_a N_a^{\text{opt}}(\bar{\tau}(\varepsilon, T)) \\
& + \mathbf{G}_T(\mu, \varepsilon) \cdot \mathbf{F}_\xi(\mu, \varepsilon, T) \\
& + \Sigma_\Delta \cdot \tau \\
& + \mathbf{G}_T(\mu, \varepsilon) \cdot \mathbf{f}_\xi(|\mathcal{U}(\varepsilon)|) + |\mathcal{U}(\varepsilon)|.
\end{aligned} \tag{4.167}$$

Combining Equations (4.161) and (4.167) ends the proof.  $\square$

## 4.6.4 Upper bound on the regret

In this section, we prove Theorem 3.

*Proof.* Theorem 3 is obtained as a corollary of Proposition 4. More precisely, the monotony of  $\lambda(\mu, \cdot)$  and  $\mathfrak{C}_{\{\theta^*\}}(\cdot) = \mathfrak{C}_\theta(\cdot)$  (see Lemma 16) simplify the expressions of  $\Lambda_\tau(\mu, \varepsilon)$  and  $C_\tau(\mu, \varepsilon)$ . The upper bounds  $\gamma_1 \leq 1/7$ ,  $\mathfrak{C}_\theta(\cdot) \leq \mathfrak{C}_0(\cdot)$  (see Lemma 16) and Pinsker's inequality simplify the terms  $F_\xi(\mu, \varepsilon, T)$  and  $G_T(\mu, \varepsilon)$ . Then, Theorem 3 is obtained by taking the expectation in both sides of the inequality. Indeed, since  $f_\xi(\cdot)$  is a concave function, we have  $\mathbb{E}_\nu[f_\xi(|\mathcal{U}(\varepsilon)|)] \leq f_\xi(\mathbb{E}_\nu[|\mathcal{U}(\varepsilon)|])$  and the constant  $C_{\xi,d,\varepsilon} = \mathbb{E}_\nu[|\mathcal{U}(\varepsilon)|]$ . From Lemma 8 and Equation (4.95), we have

$$C_{\xi,d,\varepsilon} = K_\xi^* + 2|\mathcal{A}|^2 K_{\xi,1} + \frac{(d+1)e^{d+1}}{2\varepsilon^2} \left( f_\xi \left( \frac{(d+1)e^{d+1}}{2\varepsilon_v^2} \right) \right)^2 + \sum_{n \geq 3} \frac{3|\mathcal{A}|^2}{n(\log(n))^\xi} + |\mathcal{A}|^2 \frac{e^{2\varepsilon^2}}{\varepsilon^2},$$

where  $K_\xi^*$  and  $K_{\xi,1}$  are explained in Lemma 8. □

We end this subsection by proving the following lemma that ensures a nice control of  $\mathfrak{C}_\theta(\mu(\cdot, \cdot))$  appearing in the first order term of the upper bound on the regret.

**Lemma 16** (Monotonic properties of  $\mathfrak{C}_\theta(\mu(\cdot, \cdot))$ ). *For  $0 < \varepsilon < \varepsilon_0$ ,  $(\mathfrak{C}_\theta(\mu(\varepsilon, t)))_{t \geq 1}$  and  $\mathfrak{C}_\theta(\mu(\cdot, \infty))$  are non-increasing and  $\lim_{\varepsilon \rightarrow 0} \lim_{t \rightarrow \infty} \mathfrak{C}_\theta(\mu(\varepsilon, t)) = \lim_{\varepsilon \rightarrow 0} \mathfrak{C}_\theta(\mu(\varepsilon, \infty)) = \mathfrak{C}_\theta(\mu)$ . Finally, we denote by*

$$\mathfrak{C}_0(\mu) := \min_{n \in \mathbb{R}_+^{\mathcal{A}}} \left\{ \sum_{a \in \mathcal{A}} n_a (\max(\mu) - \mu_a) \right. \\ \left. \text{s.t. } \forall a \notin \text{argmax}(\mu), \text{kl}(\mu_a | \max(\mu)) n_a \geq 1 \right\} = \sum_{a \notin \text{argmax}(\mu)} \frac{\max(\mu) - \mu_a}{\text{kl}(\mu_a | \max(\mu))}$$

the minimization problem  $\mathfrak{C}_{[-1,1]^{\mathcal{A}} \times \mathcal{A}}(\mu)$  when there is no structure, that is when  $\Theta = [-1, 1]^{\mathcal{A} \times \mathcal{A}}$ , for all  $\mu \in (0, 1)^{\mathcal{A}}$ . Then we have  $\mathfrak{C}_\theta(\cdot) \leq \mathfrak{C}_0(\cdot)$ .

*Proof.* We first note for all  $\mu \in (0, 1)^{\mathcal{A}}$ , by considering the change of variables  $n_a \leftarrow \Delta_a n_a$  if  $a \notin \text{argmax} \mu$  and  $n_a \leftarrow n_a$  if  $a \in \text{argmax} \mu$ , we have

$$\mathfrak{C}_\theta(\mu) := \min_{n \in \mathbb{R}_+^{\mathcal{A}}} \left\{ \sum_{a \in \mathcal{A}} n_a \text{ s.t. } \forall a \notin \text{argmax} \mu, \right. \\ \left. \min_{\substack{\theta \in \bar{\Theta}_a \\ \mu_{a'} \leq \max \mu - \theta_{a,a'}}} \sum_{a' \in \mathcal{A}} \frac{\text{kl}(\mu_{a'} | \max \mu - \theta_{a,a'})}{\max \mu - \mu_{a'}} n_{a'} \geq 1 \right\}. \quad (4.168)$$

For  $0 < \varepsilon' \leq \varepsilon < \varepsilon_0$  and  $t' \geq t \geq 1$ , we show that

$$\mathfrak{C}_\theta(\mu(\varepsilon', t')) \leq \mathfrak{C}_\theta(\mu(\varepsilon, t)). \quad (4.169)$$

Then, since  $\lim_{t \rightarrow \infty} \mu(\varepsilon, t) = \mu(\varepsilon, \infty)$  for  $0 < \varepsilon < \varepsilon_0$  and  $\lim_{\varepsilon \rightarrow 0} \mu(\varepsilon, \infty) = \mu$ , previous Equation (4.169) and Lemma 27 will imply

$$\lim_{\varepsilon \rightarrow 0} \lim_{t \rightarrow \infty} \mathfrak{C}_\theta(\mu(\varepsilon, t)) = \lim_{\varepsilon \rightarrow 0} \mathfrak{C}_\theta(\mu(\varepsilon, \infty)) = \mathfrak{C}_\theta(\mu). \quad (4.170)$$

Let  $0 < \varepsilon' \leq \varepsilon < \varepsilon_0$  and  $t' \geq t \geq 1$ . Let  $n \in \mathbb{R}_+^{\mathcal{A}}$  such that

$$\forall a \notin \text{argmax} \mu(\varepsilon, t), \min_{\theta \in \bar{\Theta}_a} \sum_{\substack{a' \in \mathcal{A} \\ \mu_{a'}(\varepsilon, t) \leq \max \mu(\varepsilon, t) - \theta_{a,a'}}} \frac{\text{kl}(\mu_{a'}(\varepsilon, t) | \max \mu(\varepsilon, t) - \theta_{a,a'})}{\max \mu(\varepsilon, t) - \mu_{a'}(\varepsilon, t)} n_{a'} \geq 1. \quad (4.171)$$

We first note that

$$\operatorname{argmax} \mu(\varepsilon, t) = \mu(\varepsilon', t') = \operatorname{argmax} \mu, \quad (4.172)$$

$$\max \mu(\varepsilon, t) = \mu^* - \varepsilon \leq \mu^* - \varepsilon' = \max \mu(\varepsilon', t'), \quad (4.173)$$

and

$$\forall a \notin \operatorname{argmax} \mu, \quad \mu_a(\varepsilon, t) = \mu_a + \varepsilon + \varepsilon(t) \geq \mu + \varepsilon' + \varepsilon(t') = \mu_a(\varepsilon', t'). \quad (4.174)$$

In particular, from Equations (4.172), (4.173) and (4.174), Lemma 25 implies

$$\forall a \notin \operatorname{argmax} \mu, \quad \forall \theta \in \overline{\Theta}_a, \quad \frac{\operatorname{kl}(\mu_{a'}(\varepsilon, t) | \max \mu(\varepsilon, t) - \theta_{a,a'})}{\max \mu(\varepsilon, t) - \mu_a(\varepsilon, t)} \leq \frac{\operatorname{kl}(\mu_{a'}(\varepsilon', t') | \max \mu(\varepsilon', t') - \theta_{a,a'})}{\max \mu(\varepsilon', t') - \mu_a(\varepsilon', t')}. \quad (4.175)$$

Combining all together Equations (4.171), (4.172), (4.173), (4.174), (4.175), we have

$$\forall a \notin \operatorname{argmax} \mu(\varepsilon', t'), \quad \min_{\theta \in \overline{\Theta}_a} \sum_{\substack{a' \in \mathcal{A} \\ \mu_{a'}(\varepsilon', t') \leq \max \mu(\varepsilon', t') - \theta_{a,a'}} \frac{\operatorname{kl}(\mu_{a'}(\varepsilon', t') | \max \mu(\varepsilon', t') - \theta_{a,a'})}{\max \mu(\varepsilon', t') - \mu_a(\varepsilon', t')} n_{a'} \geq 1. \quad (4.176)$$

Previous Equation 4.176 implies

$$\mathfrak{C}_\theta(\mu(\varepsilon', t')) \leq \sum_{a \in \mathcal{A}} n_a. \quad (4.177)$$

Since previous Equation (4.177) holds for all  $n \in \mathbb{R}_+^{\mathcal{A}}$  satisfying Equation (4.171), we have

$$\mathfrak{C}_\theta(\mu(\varepsilon', t')) \leq \mathfrak{C}_\theta(\mu(\varepsilon, t)). \quad (4.178)$$

□



## Chapter 5

# Concentration Inequalities for Structured Bandits

In this chapter, we state and prove Theorem 4, the main concentration result used in Chapter 4 for IMED-GS analysis. This key result enables to tighten some terms and reduce the burn-in phase as it enables the use of a parameter  $d < |\mathcal{A}| - 1$ . See also the numerical experiments in Section A.4 showing the practical benefit of this contribution.

**Theorem 4.** *For all set of Bernoulli distributions  $\nu = (\nu_a)_{a \in \mathcal{A}}$  with means  $(\mu_a)_{a \in \mathcal{A}} \subset (0, 1)^{\mathcal{A}}$ , for all  $\xi \geq 0$ , for all  $\underline{M} \geq 1$ , for all  $\mathcal{A}' \subset \mathcal{A}$ , for all  $\Phi \geq |\mathcal{A}'| + 1$ , for all  $\varepsilon > 0$ , for all  $n \geq n_{\xi, \mathcal{A}', \varepsilon}$ ,*

$$\begin{aligned} & \mathbb{P}_\nu \left( \bigcup_{t \geq 1} \left\{ \sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\hat{\mu}_a(t) | \mu_a - \varepsilon) \geq \Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_{\mathcal{A}'}(t))) \right\} \cap \Omega_{\mathcal{A}', \varepsilon, n}(t) \right) \\ & \leq \frac{e^{|\mathcal{A}'|+2}}{|\mathcal{A}'|^{|\mathcal{A}'|}} \log(n)^{|\mathcal{A}'|+1} [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}_{\nu, \xi, \mathcal{A}'}(n)]^{2|\mathcal{A}'|+1} e^{-[\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}_{\nu, \xi, \mathcal{A}'}(n)]}, \end{aligned}$$

where  $\Omega_{\mathcal{A}', \varepsilon, n}(t) = \{\underline{M} \leq N_{\mathcal{A}'}(t) \leq |\mathcal{A}'| \mathbf{f}_\xi(n) / 2\varepsilon^2\} \cap \{\forall a \in \mathcal{A}', \hat{\mu}_a(t) < \mu_a - \varepsilon\}$ ,

$$\underline{N}_{\nu, \xi, \mathcal{A}'}(n) = \underline{M} \vee \frac{\mathbf{f}_\xi(n)}{1 - \log(1 - \max_{a \in \mathcal{A}'} \mu_a)} \quad \text{and} \quad n_{\xi, \mathcal{A}', \varepsilon} = \left( e^2 \vee \frac{|\mathcal{A}'|}{2\varepsilon^2} \right) \left( \mathbf{f}_\xi \left( e^2 \vee \frac{|\mathcal{A}'|}{2\varepsilon^2} \right) \right)^2.$$

Although the statement looks a little intricate, it enables to obtain more precise control of the concentration terms compared to alternative tools present in the literature (e.g. in Magureanu et al. (2014)). Note the presence of the random term  $N_{\mathcal{A}'}(t)$  in event  $\Omega_{\mathcal{A}', \varepsilon, n}(t)$  in the left-hand side of the deviation inequality. Besides, it enables to show that setting  $\xi = 1$  is enough for the regret to be controlled.

*Proof.* The five main steps of the proof are the following. First, we establishing a lower bound on the aggregated number of pulls  $N_{\mathcal{A}'}(t)$  for  $t \geq 1$ . Second, we proceed to a change of measurement applying Lemma 24. Third, we use twice the peeling technique, for the aggregated number of pulls  $N_{\mathcal{A}'}(t)$  and for each ratio  $N_a(t)/N_{\mathcal{A}'}(t)$ ,  $a \in \mathcal{A}'$ . Fourth, we control each band of the peeling using Lemma 17 based on multivariate stochastic ordering, see Lemma 8 from Magureanu et al. (2014). Fifth, we apply the union bound.

Let us consider  $t \geq 1$  such that

$$\begin{aligned} & \forall a \in \mathcal{A}', \hat{\mu}_a(t) < \mu_a - \varepsilon \\ & \underline{M} \leq N_{\mathcal{A}'}(t) \leq |\mathcal{A}'| \mathbf{f}_\xi(n) / 2\varepsilon^2 \\ & \sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\hat{\mu}_a(t) | \mu_a - \varepsilon) \geq \Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_{\mathcal{A}'}(t))). \end{aligned} \tag{5.1}$$

### 1. Lower bound on $N_{\mathcal{A}'}(t)$ :

From Equation (5.1), we have  $\widehat{\mu}_a(t) < \mu_a - \varepsilon$  for all  $a \in \mathcal{A}'$ . Then, the monotonic properties of  $\text{kl}(\cdot | \cdot)$  imply

$$\forall a \in \mathcal{A}', \quad \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) \leq \text{kl}(0 | \mu_a) = -\log(1 - \mu_a) \quad (5.2)$$

and

$$\sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) \leq -\log(1 - \max_{a \in \mathcal{A}'} \mu_a) N_{\mathcal{A}'}(t). \quad (5.3)$$

From Equation (5.1), since  $N_{\mathcal{A}'}(t) \geq f_\xi(N_{\mathcal{A}'}(t))$ , Equation (5.1) and previous Equation (5.3) imply

$$N_{\mathcal{A}'}(t) \geq \underline{N} := \underline{M} \vee \frac{f_\xi(n)}{1 - \log(1 - \max_{a \in \mathcal{A}'} \mu_a) \vee e}. \quad (5.4)$$

### 2. Change of measurement :

From Equation (5.1), we have  $\widehat{\mu}_a(t) < \mu_a - \varepsilon < \mu_a$  for all  $a \in \mathcal{A}'$ . Then, Lemma 24 implies

$$\forall a \in \mathcal{A}', \quad \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) \leq \text{kl}(\widehat{\mu}_a(t) | \mu_a) - 2\varepsilon^2 \quad (5.5)$$

and

$$\sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) \leq \sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a) - 2\varepsilon^2 N_{\mathcal{A}'}(t). \quad (5.6)$$

From Lemma 28, previous Equation (5.6) implies

$$\sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\widehat{\mu}_a(t) | \mu_a) \geq \Phi \vee (f_\xi(n) - f_\xi(N_{\mathcal{A}'}(t))) + 2\varepsilon^2 N_{\mathcal{A}'}(t) \quad (5.7)$$

or similarly

$$\sum_{a \in \mathcal{A}'} \frac{N_a(t)}{N_{\mathcal{A}'}(t)} \text{kl}(\widehat{\mu}_a(t) | \mu_a) \geq \frac{\Phi \vee (f_\xi(n) - f_\xi(N_{\mathcal{A}'}(t))) + 2\varepsilon^2 N_{\mathcal{A}'}(t)}{N_{\mathcal{A}'}(t)}. \quad (5.8)$$

### 3. Peeling :

Let us consider

$$\rho = \sqrt{\frac{\Phi \vee f_\xi(n) + 2\varepsilon^2 \underline{N} - 1}{\Phi \vee f_\xi(n) + 2\varepsilon^2 \underline{N}}}. \quad (5.9)$$

Since  $\Phi \vee f_\xi(n) + 2\varepsilon^2 \underline{N} - 1 > \Phi - 1 \geq |\mathcal{A}'|$ , previous Equation (5.9) implies

$$\frac{|\mathcal{A}'|}{|\mathcal{A}'| + 1} \leq \rho^2 = 1 - \frac{1}{\Phi \vee f_\xi(n) + 2\varepsilon^2 \underline{N}} < 1. \quad (5.10)$$

Now we introduce decreasing geometric sequences  $(N_k)_{k \geq 0}$  and  $(x_\ell)_{\ell \geq 0}$  for the peeling :

$$\begin{cases} N_0 = \frac{|\mathcal{A}'| f_\xi(n)}{2\varepsilon^2} \\ N_{k+1} = \rho N_k, \quad k \geq 0 \end{cases} \quad \begin{cases} x_0 = 1 \\ x_{\ell+1} = \rho x_\ell, \quad \ell \geq 0. \end{cases} \quad (5.11)$$

By considering

$$K = \left\lceil \frac{\log(N_0/\underline{N})}{-\log(\rho)} \right\rceil \quad L = \left\lceil \frac{\log(N_0)}{-\log(\rho)} \right\rceil, \quad (5.12)$$

we have

$$N_{K+1} \leq \underline{N} \leq N_{\mathcal{A}'}(t) \quad x_{L+1} \leq \frac{1}{N_0} \leq \frac{N_a(t)}{N_{\mathcal{A}'}(t)}, \quad \forall a \in \mathcal{A}'. \quad (5.13)$$

Then, by considering the bands for  $k \in \llbracket 0, K \rrbracket$ , for  $(\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket$ ,

$$B(k, (\ell_a)_{a \in \mathcal{A}'}) = \left\{ t \geq 1 : \begin{array}{l} \forall a \in \mathcal{A}', \widehat{\mu}_a(t) < \mu_a \\ N_{k+1} \leq N_{\mathcal{A}'}(t) \leq N_k \\ \forall a \in \mathcal{A}', x_{\ell_{a+1}} \leq \frac{N_a(t)}{N_{\mathcal{A}'}(t)} \leq x_{\ell_a} \\ \sum_{a \in \mathcal{A}'} Z_a \geq \frac{\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_{\mathcal{A}'}(t))) + 2\varepsilon^2 N_{\mathcal{A}'}(t)}{N_{\mathcal{A}'}(t)} \end{array} \right\}, \quad (5.14)$$

we have

$$t \in \bigcup_{\substack{k \in \llbracket 0, K \rrbracket \\ (\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket}} B(k, (\ell_a)_{a \in \mathcal{A}'}), \quad (5.15)$$

where

$$\forall a \in \mathcal{A}', \quad Z_a = \frac{N_a(t)}{N_{\mathcal{A}'}(t)} \mathbf{kl}(\widehat{\mu}_a(t) | \mu_a). \quad (5.16)$$

Monotony of both function  $\mathbf{f}_\xi(\cdot)$  and sequence  $(N_k)_{k \geq 0}$  imply for  $k \in \llbracket 0, K \rrbracket$ , for  $(\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket$ ,

$$B(k, (\ell_a)_{a \in \mathcal{A}'}) \subset \overline{B}(k, (\ell_a)_{a \in \mathcal{A}'}) \cap \overline{S}(k, (\ell_a)_{a \in \mathcal{A}'}), \quad (5.17)$$

with

$$\overline{B}(k, (\ell_a)_{a \in \mathcal{A}'}) := \left\{ t \geq 1 : \begin{array}{l} \forall a \in \mathcal{A}', \widehat{\mu}_a(t) < \mu_a \\ N_{k+1} \leq N_{\mathcal{A}'}(t) \leq N_k \\ \forall a \in \mathcal{A}', x_{\ell_{a+1}} \leq \frac{N_a(t)}{N_{\mathcal{A}'}(t)} \leq x_{\ell_a} \end{array} \right\} \quad (5.18)$$

$$\overline{S}(k, (\ell_a)_{a \in \mathcal{A}'}) := \left\{ t \geq 1 : \sum_{a \in \mathcal{A}'} \mathbb{I}_{\{t \in \overline{B}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \frac{\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 N}{N_k} \right\},$$

where, according to Lemma 28,

$$\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0) \geq 0 \quad n \geq N_0, \quad (5.19)$$

since  $n \geq n_{\nu, \mathcal{A}', \varepsilon} = (e^2 \vee (|\mathcal{A}'|/2\varepsilon^2)) (\mathbf{f}_\xi(e^2 \vee (|\mathcal{A}'|/2\varepsilon^2)))^2$ . In particular, Equations (5.15) and (5.17)-(5.18) imply

$$t \in \bigcup_{\substack{k \in \llbracket 0, K \rrbracket \\ (\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket}} \overline{B}(k, (\ell_a)_{a \in \mathcal{A}'}) \cap \overline{S}(k, (\ell_a)_{a \in \mathcal{A}'}). \quad (5.20)$$

#### 4. Stochastic ordering :

The total probability rule implies for  $k \in \llbracket 0, K \rrbracket$ ,  $(\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket$ , for all  $\zeta \in \mathbb{R}_+^{\mathcal{A}'}$ ,

$$\begin{aligned}
& \mathbb{P}_\nu \left( \bigcap_{a \in \mathcal{A}'} \mathbb{I}_{\{t \in \bar{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \zeta_a \right) \\
&= \sum_{\substack{M \in \mathbb{N} \\ (n_a)_{a \in \mathcal{A}'} \subset \mathbb{N}}} \mathbb{P}_\nu \left( N_{\mathcal{A}'}(t) = M, \bigcap_{a \in \mathcal{A}'} N_a(t) = n_a, \mathbb{I}_{\{t \in \bar{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \zeta_a \right) \\
&= \sum_{\substack{M \in \mathbb{N} \\ (n_a)_{a \in \mathcal{A}'} \subset \mathbb{N}}} \prod_{a \in \mathcal{A}'} \mathbb{P}_\nu \left( N_{\mathcal{A}'}(t) = M, (N_{a'}(t))_{a' \in \mathcal{A}'} = (n_{a'})_{a' \in \mathcal{A}'}, \mathbb{I}_{\{t \in \bar{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \zeta_a \right) \\
&\leq \prod_{a \in \mathcal{A}'} \sum_{\substack{M \in \mathbb{N} \\ (n_a)_{a \in \mathcal{A}'} \subset \mathbb{N}}} \mathbb{P}_\nu \left( N_{\mathcal{A}'}(t) = M, (N_{a'}(t))_{a' \in \mathcal{A}'} = (n_{a'})_{a' \in \mathcal{A}'}, \mathbb{I}_{\{t \in \bar{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \zeta_a \right) \\
&= \prod_{a \in \mathcal{A}'} \mathbb{P}_\nu \left( \mathbb{I}_{\{t \in \bar{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \zeta_a \right). \tag{5.21}
\end{aligned}$$

By combining Previous Equation (5.21) and Lemma 32, we have for  $k \in \llbracket 0, K \rrbracket$ ,  $(\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket$ , for all  $\zeta \in \mathbb{R}_+^{\mathcal{A}'}$ ,

$$\begin{aligned}
\mathbb{P}_\nu \left( \bigcap_{a \in \mathcal{A}'} \mathbb{I}_{\{t \in \bar{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \zeta_a \right) &\leq \prod_{a \in \mathcal{A}'} \mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ \hat{\mu}_a(t) < \mu_a \\ N_a(t) \geq x_{\ell_a+1} N_{k+1}}} \text{kl}(\hat{\mu}_a(t) | \mu_a) \geq \zeta_a / x_{\ell_a} \right) \\
&\leq \prod_{a \in \mathcal{A}'} \exp \left( -\frac{x_{\ell_a+1}}{x_{\ell_a}} N_{k+1} \zeta_a \right) \\
&= \prod_{a \in \mathcal{A}'} \exp(-\rho N_{k+1} \zeta_a) \\
&= \exp \left( -\rho N_{k+1} \sum_{a \in \mathcal{A}'} \zeta_a \right). \tag{5.22}
\end{aligned}$$

Furthermore, for  $k \in \llbracket 0, K \rrbracket$ ,

$$\rho N_{k+1} \frac{\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N}}{N_k} = \rho^2 [\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N}]. \tag{5.23}$$

Using Equation (5.10) we have  $\rho^2 \geq |\mathcal{A}'| / |\mathcal{A}'| + 1$ , and previous Equation (5.23) implies

$$\rho N_{k+1} \frac{\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0) + 2\varepsilon^2 \underline{N}}{N_k} \geq \frac{|\mathcal{A}'|}{|\mathcal{A}'| + 1} \times \Phi \geq \frac{|\mathcal{A}'|}{|\mathcal{A}'| + 1} \times (|\mathcal{A}'| + 1) = |\mathcal{A}'|. \tag{5.24}$$

This Equation (5.24) is the key equation allowing the use of Lemma 17. From Lemma 17, Equations (5.22) and (5.24) imply for  $k \in \llbracket 0, K \rrbracket$ ,  $(\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket$ ,

$$\begin{aligned}
& \mathbb{P}_\nu \left( \sum_{a \in \mathcal{A}'} \mathbb{I}_{\{t \in \bar{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'})\}} Z_a \geq \frac{\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N}}{N_k} \right) \\
&\leq \frac{e^{|\mathcal{A}'|}}{|\mathcal{A}'|^{|\mathcal{A}'|}} [\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})]^{|\mathcal{A}'|} e^{-\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})}, \tag{5.25}
\end{aligned}$$

that is, according to Equation (5.17),

$$\begin{aligned}
& \mathbb{P}_\nu(t \in \overline{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'}) \cap \overline{\mathcal{S}}(k, (\ell_a)_{a \in \mathcal{A}'})) \\
& \leq \frac{e^{|\mathcal{A}'|}}{|\mathcal{A}'|^{|\mathcal{A}'|}} [\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})]^{|\mathcal{A}'|} e^{-\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})} \\
& \leq \frac{e^{|\mathcal{A}'|}}{|\mathcal{A}'|^{|\mathcal{A}'|}} [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|} e^{-\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})}.
\end{aligned} \tag{5.26}$$

### 5. Union bound :

The union bound writes

$$\begin{aligned}
& \mathbb{P}_\nu \left( t \in \bigcup_{\substack{k \in \llbracket 0, K \rrbracket \\ (\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket}} \overline{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'}) \cap \overline{\mathcal{S}}(k, (\ell_a)_{a \in \mathcal{A}'}) \right) \\
& \leq \sum_{\substack{k \in \llbracket 0, K \rrbracket \\ (\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket}} \mathbb{P}_\nu(t \in \overline{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'}) \cap \overline{\mathcal{S}}(k, (\ell_a)_{a \in \mathcal{A}'})) .
\end{aligned} \tag{5.27}$$

By combining both previous Equations (5.26) and (5.27), we have

$$\begin{aligned}
& \mathbb{P}_\nu \left( t \in \bigcup_{\substack{k \in \llbracket 0, K \rrbracket \\ (\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket}} \overline{\mathcal{B}}(k, (\ell_a)_{a \in \mathcal{A}'}) \cap \overline{\mathcal{S}}(k, (\ell_a)_{a \in \mathcal{A}'}) \right) \\
& \leq \sum_{\substack{k \in \llbracket 0, K \rrbracket \\ (\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket}} \frac{e^{|\mathcal{A}'|}}{|\mathcal{A}'|^{|\mathcal{A}'|}} [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|} e^{-\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})},
\end{aligned} \tag{5.28}$$

which implies, according to Equations (5.1) and (5.20),

$$\begin{aligned}
& \mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ \forall a \in \mathcal{A}', \widehat{\mu}_a(t) < \mu_a - \varepsilon \\ \underline{M} \leq N_{\mathcal{A}'}(t) \leq |\mathcal{A}'| \mathbf{f}_\xi(n) / 2\varepsilon^2}} \sum_{a \in \mathcal{A}'} N_a(t) \mathbf{kl}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) + \mathbf{f}_\xi(N_{\mathcal{A}'}(t)) \geq \mathbf{f}_\xi(n) \right) \\
& \leq \sum_{\substack{k \in \llbracket 0, K \rrbracket \\ (\ell_a)_{a \in \mathcal{A}'} \subset \llbracket 0, L \rrbracket}} \frac{e^{|\mathcal{A}'|}}{|\mathcal{A}'|^{|\mathcal{A}'|}} [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|} e^{-\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})} \\
& = \frac{e^{|\mathcal{A}'|}}{|\mathcal{A}'|^{|\mathcal{A}'|}} (K + 1) (L + 1)^{|\mathcal{A}'|} [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|} e^{-\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N})}.
\end{aligned} \tag{5.29}$$

We now simplify these terms in order to optimize the obtained upper bound. We note that Equations (5.10) and (5.19) plus Lemma (29) imply

$$K, L \leq \log(n) [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}] \tag{5.30}$$

and

$$\begin{aligned} \rho &< 1 \\ (K + 1)(L + 1)^{|\mathcal{A}'|} &\leq \log(n)^{|\mathcal{A}'|+1} [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|+1}. \end{aligned} \quad (5.31)$$

Furthermore, Equations (5.9) and (5.19) imply

$$\begin{aligned} &-\rho^2 (\Phi \vee (\mathbf{f}_\xi(n) - \mathbf{f}_\xi(N_0)) + 2\varepsilon^2 \underline{N}) \\ &\leq -\rho^2 (\Phi \vee \mathbf{f}_\xi(n) - \Phi \vee \mathbf{f}_\xi(N_0) + 2\varepsilon^2 \underline{N}) \\ &= -\frac{\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N} - 1}{\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}} (\Phi \vee \mathbf{f}_\xi(n) - \Phi \vee \mathbf{f}_\xi(N_0) + 2\varepsilon^2 \underline{N}) \\ &= -[\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N} - 1] + \frac{\Phi \vee \mathbf{f}_\xi(N_0)}{\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}} \\ &= -[\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}] + 1 + \frac{\Phi \vee \mathbf{f}_\xi(N_0)}{\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}} \\ &\leq -[\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}] + 2. \end{aligned} \quad (5.32)$$

Thus, by combining Equations (5.29), (5.31) and (5.32) we get

$$\begin{aligned} &\mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ \forall a \in \mathcal{A}', \hat{\mu}_a(t) < \mu_a - \varepsilon \\ \underline{M} \leq N_{\mathcal{A}'}(t) \leq |\mathcal{A}'| \mathbf{f}_\xi(n) / 2\varepsilon^2}} \sum_{a \in \mathcal{A}'} N_a(t) \text{kl}(\hat{\mu}_a(t) | \mu_a - \varepsilon) + \mathbf{f}_\xi(N_{\mathcal{A}'}(t)) \geq \mathbf{f}_\xi(n) \right) \\ &\leq \frac{e^{|\mathcal{A}'|+2}}{|\mathcal{A}'|^{|\mathcal{A}'|}} \log(n)^{|\mathcal{A}'|+1} [\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}]^{|\mathcal{A}'|+1} e^{-[\Phi \vee \mathbf{f}_\xi(n) + 2\varepsilon^2 \underline{N}]}, \end{aligned} \quad (5.33)$$

where  $\underline{N} = \underline{M} \vee \frac{\mathbf{f}_\xi(n)}{1 - \log(1 - \max_{a \in \mathcal{A}'} \mu_a)}$ . □

We remind in the following lemma a result from Magureanu et al. (2014) based on multivariate stochastic ordering.

**Lemma 17** (Stochastic ordering). *Let  $M > 0$ ,  $A \geq 2$ . Let  $Z \in \mathbb{R}^A$  be a random variable such that for all  $\zeta \in \mathbb{R}_+^A$ ,*

$$\mathbb{P}(Z \geq \zeta) \leq \exp\left(-M \sum_{a=1}^A \zeta_a\right).$$

Then for all  $z \geq A/M$ ,

$$\mathbb{P}\left(\sum_{a=1}^A Z_a \geq z\right) \leq \left(\frac{Mze}{A}\right)^A e^{-Mz}.$$

# Chapter 6

## Routine Bandits

### 6.1 The Routine Bandit Setting

A routine bandit problem is specified by a time horizon  $T \geq 1$  and a finite set of distributions  $\nu = (\nu_b)_{b \in \mathcal{B}}$  with means  $(\mu_{a,b})_{a \in \mathcal{A}, b \in \mathcal{B}}$ , where  $\mathcal{A}$  is a finite set of arms and  $\mathcal{B}$  is a finite set of bandit configurations. Each  $b \in \mathcal{B}$  can be seen as a classical multi-armed bandit problem defined by  $\nu_b = (\nu_{a,b})_{a \in \mathcal{A}}$ . At each period  $h \geq 1$  and for all time steps  $t \in \llbracket 1, T \rrbracket$ , the learner deals with a bandit  $b_\star^h \in \mathcal{B}$  and chooses an arm  $a_t^h \in \mathcal{A}$ , based only on the past. The learner then receives and observes a reward  $X_t^h \sim \nu_{a_t^h, b_\star^h}$ . The goal of the learner is to maximize the expected sum of rewards received over time (up to some unknown number of periods  $H \geq 1$ ). The distributions are unknown, which makes the problem non-trivial. The optimal strategy therefore consists in playing repeatedly on each period  $h$ , an optimal arm  $a_\star^h \in \operatorname{argmax}_{a \in \mathcal{A}} \mu_{a, b_\star^h}$ , which has mean  $\mu_\star^h = \mu_{a_\star^h, b_\star^h}$ . The goal of the learner is equivalent to minimizing the cumulative *regret* with respect to an optimal strategy:

$$R(\nu, H, T) = \mathbb{E}_\nu \left[ \sum_{h=1}^H \sum_{t=1}^T (\mu_\star^h - X_t^h) \right]. \quad (6.1)$$

**Related works** One of the closest setting to routine bandits is the sequential transfer scenario [Gheshlaghi Azar et al. \(2013\)](#), where the cardinality  $|\mathcal{B}|$  and quantities  $H$  and  $T$  are known ahead of time, and the instances in  $\mathcal{B}$  are either known perfectly or estimated with known confidence. Routine bandits also bear similarity with clustering bandits [Gentile et al. \(2014\)](#), a contextual bandit setting [Langford and Zhang \(2007\)](#) where contexts can be clustered into finite (unknown) clusters. While both settings are recurring bandit problems, routine bandits assume no information on users (including their number) but users are recurring for several iterations of interaction, while clustering bandits assume that each user is seen only once, but is characterized by features such that they can be associated with previously seen users. Finally, latent bandits [Maillard and Mannor \(2014\)](#) consider the less structured situation when the learner faces a possibly different user at every time.

**Assumptions and working conditions** The configuration  $\nu$ , the set of bandits  $\mathcal{B}$ , and the sequence of bandits  $(b_\star^h)_{h \geq 1}$  are *unknown* (in particular  $|\mathcal{B}|$  and the identity of user  $b_\star^h$  are unknown to the learner at time  $t$ ). The learner only knows that  $\nu \in \mathcal{D}$ , where  $\mathcal{D}$  is a given set of bandit configurations. In order to leverage information from the bandit instances encountered, we should consider that bandits reoccur. We denote by  $\beta_b^h = \sum_{h'=1}^h \mathbb{I}_{\{b_\star^{h'}=b\}}/h$  the frequency of bandit  $b \in \mathcal{B}$  at period  $h \geq 1$  and assume  $\beta_b^H > 0$ . The next two assumptions respectively allow for two bandit instances  $b$  and  $b'$  to be distinguishable from their means when  $b \neq b'$  and show consistency in their optimal strategy when  $b = b'$ .

**Assumption 10** (Separation). *Let us consider  $\gamma_\nu := \min_{b \neq b'} \min_{a \in \mathcal{A}} \{|\mu_{a,b} - \mu_{a,b'}|, 1\}$ . We assume  $\gamma_\nu > 0$ .*

**Assumption 11** (Unique optimal arm). *Each bandit  $b \in \mathcal{B}$  has a unique optimal arm  $a_b^*$ .*

Assumption 11 is standard. Finally, we consider normally-distributed rewards. Although most of our analysis (e.g., concentration) would extend to exponential families of dimension 1, Assumption 12 increases readability of the statements.

**Assumption 12** (Gaussian arms). *The set  $\mathcal{D}$  is the set of bandit configurations such that for all bandit  $b \in \mathcal{B}$ , for all arm  $a \in \mathcal{A}$ ,  $\nu_{a,b}$  is a one-dimensional Gaussian distribution with mean  $\mu_{a,b} \in \mathbb{R}$  and variance  $\sigma^2 = 1$ .*

For  $\nu \in \mathcal{D}$ , we define for an arm  $a \in \mathcal{A}$  and a bandit  $b \in \mathcal{B}$  their gap  $\Delta_{a,b} = \mu_b^* - \mu_{a,b}$  and their total number of pulls over  $H$  periods  $N_{a,b}(H, T) = \sum_{h=1}^H \sum_{t=1}^T \mathbb{I}_{\{a_t^h = a, b_t^h = b\}}$ . An arm is optimal for a bandit if their gap is equal to zero and sub-optimal if it is positive. Thanks to the chain rule, the regret rewrites as

$$R(\nu, H, T) = \sum_{b \in \mathcal{B}} \sum_{a \neq a_b^*} \mathbb{E}_\nu[N_{a,b}(H, T)] \Delta_{a,b}. \quad (6.2)$$

**Remark 15** (Fixed horizon time). *We assume the time horizon  $T$  to be the same for all periods  $h \in \llbracket 1, H \rrbracket$  out of clarity of exposure of the results and simplified definition of consistency (Definition 3). Considering a different time  $T_h$  for each  $h$  would indeed require a substantial rewriting of the statements (e.g. think of the regret lower bound), which we believe hinders readability and comparison to classical bandits.*

We conclude this section by adapting for completeness the known lower bound on the regret [Lai and Robbins \(1985\)](#); [Agrawal et al. \(1989\)](#); [Graves and Lai \(1997\)](#) for consistent strategies to the routine bandit setting. We defer the proof to Appendix C.1.

**Definition 3** (Consistent strategy). *A strategy is  $H$ -consistent on  $\mathcal{D}$  if for all configuration  $\nu \in \mathcal{D}$ , for all bandit  $b \in \mathcal{B}$ , for all sub-optimal arm  $a \neq a_b^*$ , for all  $\alpha > 0$ ,*

$$\lim_{T \rightarrow \infty} \mathbb{E}_\nu \left[ \frac{N_{a,b}(H, T)}{N_b(H, T)^\alpha} \right] = 0,$$

where  $N_b(H, T) = \beta_b^H HT$  is the number of time steps the learner has dealt with bandit  $b$ .

**Proposition 5** (Lower bounds on the regret). *Let us consider a consistent strategy. Then, for all configuration  $\nu \in \mathcal{D}$ , it must be that*

$$\liminf_{T \rightarrow \infty} \frac{R(\nu, H, T)}{\log(T)} \geq c_\nu^* := \sum_{b \in \mathcal{B}} \sum_{a \neq a_b^*} \frac{\Delta_{a,b}}{\text{KL}(\mu_{a,b} | \mu_b^*)},$$

where  $\text{KL}(\mu | \mu') = (\mu' - \mu)^2 / 2\sigma^2$  denotes the Kullback-Leibler divergence between one-dimensional Gaussian distributions with means  $\mu, \mu' \in \mathbb{R}$  and variance  $\sigma^2 = 1$ .

This lower bound differs (it is larger) from structured lower bound that can exclude some set of arms, as in [Agrawal et al. \(1989\)](#); [Maillard and Mannor \(2014\)](#) using prior knowledge on  $\mathcal{B}$ , which here is not available. On the other hand, we remark that the right hand side of the bound does not depend on  $H$ , which suggests that one at least asymptotically, one can learn from the recurring bandits. In the classical bandit setting, lower bounds on the regret [Lai and Robbins \(1985\)](#) have inspired the design of the well-known KLUCB [Garivier and Cappé \(2011\)](#) algorithm. In the next section, we build on this optimal strategy to propose a variant for the routine bandit.

## 6.2 The KLUCB–RB Strategy

Given the current period  $h$ , the general idea of this optimistic strategy consists in aggregating observations acquired in previous periods  $1 \dots h-1$  where bandit instances are tested to be the same as the current bandit  $b^h$ . To achieve this, KLUCB–RB relies both on concentration of observations gathered in previous periods and the consistency of the allocation strategy between different periods.



**Notations** The number of pulls, the sum of the rewards and the empirical mean of the rewards from the arm  $a$  in period  $h \geq 1$  at time  $t \geq 1$ , are respectively denoted by  $N_a^h(t) = \sum_{s=1}^t \mathbb{I}_{\{a_s^h=a\}}$ ,  $S_a^h(t) = \sum_{s=1}^t \mathbb{I}_{\{a_s^h=a\}} X_s^h$  and  $\hat{\mu}_a^h(t) = S_a^h(t)/N_a^h(t)$  if  $N_a^h(t) > 0$ , 0 otherwise.

**Strategy** For each period  $h \geq 1$  we compute an empirical best arm for bandit  $b_\star^h$  as the arm with maximum number of pulls in this period:  $\bar{a}_\star^h \in \operatorname{argmax}_{a \in \mathcal{A}} N_a^h(T)$ .<sup>1</sup> Similarly, in the current period  $h \geq 1$ , for each time step  $t \in \llbracket 1, T \rrbracket$ , we consider an arm with maximum number of pulls:  $\bar{a}_t^h \in \operatorname{argmax}_{a \in \mathcal{A}} N_a^h(t)$  (arbitrarily chosen). At each period  $h \in \llbracket 2, H \rrbracket$  each arm is pulled once. Then at each time step  $t \geq |\mathcal{A}| + 1$ , in order to possibly identify the current bandit  $b_\star^h$  with some bandits  $b_\star^k$  from a previous period  $k \in \llbracket 1, h-1 \rrbracket$ , we introduce for all arm  $a \in \mathcal{A}$ , the test statistics

$$Z_a^{k,h}(t) = \infty \cdot \mathbb{I}_{\{\bar{a}_t^h \neq \bar{a}_\star^k\}} + |\hat{\mu}_a^h(t) - \hat{\mu}_a^k(T)| - d(N_a^h(t), \delta^h(t)) - d(N_a^k(T), \delta^h(t)), \quad (6.3)$$

where the deviation for  $n \geq 1$  pulls with probability  $1 - \delta$ , for  $\delta > 0$ , and probability  $\delta^h(t)$  are, respectively,

$$d(n, \delta) = \sqrt{2 \left(1 + \frac{1}{n}\right) \frac{\log(\sqrt{n+1}/\delta)}{n}} \quad \delta^h(t) = \frac{1}{4|\mathcal{A}|} \times \frac{1}{h-1} \times \frac{1}{t(t+1)}.$$

The algorithm finally computes the test

$$\mathbf{T}^{k,h}(t) := \max_{a \in \mathcal{A}} Z_a^{k,h}(t) \leq 0. \quad (6.4)$$

After  $t$  rounds in current period  $h$ , the previous bandit  $b_\star^k$  is suspected of being the same as  $b_\star^h$  if the test  $\mathbf{T}^{k,h}(t)$  is true. From Eq. 6.3, we note that this requires the current mostly played arm to be the same as the arm that was mostly played in period  $k$ , which happens if there is consistency in the allocation strategy for both periods under Assumption 11. We then define aggregated numbers of pulls and averaged means: For all arm  $a \in \mathcal{A}$ , for all period  $h \geq 1$ , for all time step  $t \geq 1$ ,

$$\begin{aligned} \bar{N}_a^h(t) &:= N_a^h(t) + \sum_{k=1}^{h-1} \mathbb{I}_{\{\mathbf{T}^{k,h}(t)\}} N_a^k(T), & \bar{K}_t^h &:= \sum_{k=1}^{h-1} \mathbb{I}_{\{\mathbf{T}^{k,h}(t)\}}, \\ \bar{S}_a^h(t) &:= S_a^h(t) + \sum_{k=1}^{h-1} \mathbb{I}_{\{\mathbf{T}^{k,h}(t)\}} S_a^k(T), & \bar{\mu}_a^h(t) &= \bar{S}_a^h(t) / \bar{N}_a^h(t). \end{aligned}$$

and follow a KLUCB strategy by defining the index of arm  $a \in \mathcal{A}$  in period  $h \geq 1$  at time step  $t \geq 1$  as

$$u_a^h(t) = \min \left\{ U_a^h(t), \bar{U}_a^h(t) \right\}, \quad (6.5)$$

where

$$U_a^h(t) := \hat{\mu}_a^h(t) + \sqrt{\frac{2f(t)}{N_a^h(t)}}, \quad (6.6)$$

$$\bar{U}_a^h(t) := \bar{\mu}_a^h(t) + \sqrt{\frac{2f(\bar{K}_t^h T + t)}{\bar{N}_a^h(t)}}, \quad (6.7)$$

<sup>1</sup>ties are broken arbitrarily

with the function  $f$  being chosen, following Cappé et al. (2013) for classical bandits, as

$$f(x) := \log(x) + 3 \log \log(\max\{e, x\}), \forall x \geq 1.$$

One recognizes that Eq. 6.6 corresponds to the typical KLUCB upper bound for Gaussian distributions. The resulting KLUCB–RB strategy is summarized in Algorithm 10.

---

**Algorithm 10** KLUCB–RB

---

**Initialization** (period  $h=1$ ): follow a KLUCB strategy for bandit  $b_x^1$ .

**for** period  $h \geq 2$  **do**

    Pull each arm once

**for** time step  $t \in [|A|, T-1]$  **do**

        Compute for each previous period  $k \in [1, h-1]$  the test  $T^{k,h}(t) := \max_{a \in A} Z_a^{k,h}(t) \leq 0$

        Aggregate data from periods with positive test and compute for each arm  $a \in A$  the index  $u_a^h(t)$  according to equations (6.5)-(6.6)-(6.7).

        Pull an arm with maximum index  $a_{t+1}^h \in \operatorname{argmax}_{a \in A} u_a^h(t)$

**end for**

**end for**

---

**Theoretical guarantees** The next result shows that the number of sub-optimal pulls done by KLUCB–RB is upper-bounded in a near-optimal way.

**Theorem 5** (Upper bounds). *Let us consider a routine bandit problem specified by a set of Gaussian distributions  $\nu \in \mathcal{D}$  and a number of periods  $H \geq 1$ . Then under KLUCB–RB strategy, for all  $0 < \varepsilon < \varepsilon_\nu$ , for all bandit  $b \in \mathcal{B}$ , for all sub-optimal arm  $a \neq a_b^*$ ,*

$$\begin{aligned} \mathbb{E}_\nu[N_{a,b}(H, T)] &\leq \frac{f(\beta_b^H HT)}{\mathbf{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} \\ &\quad + \sum_{h=1}^H \mathbb{I}_{\{b^h = b\}} \left[ \tau_\nu^h + 4 |A| \left( \frac{1}{\varepsilon^2} + 1 \right) \left( 5 + \frac{8h f(hT)}{T \mathbf{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} \right) \right], \end{aligned}$$

where, for all period  $h \geq 2$ ,  $\tau_\nu^h := 2\varphi(8|A|[\varepsilon_\nu^{-2} + 65\gamma_\nu^{-2} \log(128|A|(4h)^{1/3}\gamma_\nu^{-2})])$ ,  $\varphi : x \geq 1 \mapsto x \log(x)$ ,  $\varepsilon_\nu = \min_{b \in \mathcal{B}} \min_{a \neq a_b^*} \Delta_{a,b}/2$  and  $\gamma_\nu = \min_{b \neq b'} \min_{a \in A} \{|\mu_{a,b} - \mu_{a,b'}|, 1\}$ .

This implies that the dependency on the time horizon  $T$  in these upper bounds is asymptotically optimal with regard to the lower bound on the regret given in Proposition 5. From Eq. 6.2, by considering the case when the time horizon  $T$  tends to infinity, we deduce that KLUCB–RB achieves asymptotic optimality.

**Corollary 4** (Asymptotic optimality). *With the same notations and under the assumptions as in Theorem 5, KLUCB–RB achieves*

$$\limsup_{T \rightarrow \infty} \frac{R(\nu, H, T)}{\log(T)} \leq c_\nu^*,$$

where  $c_\nu^*$  is defined as in Proposition 5.

For comparison, let us remark that under the strategy that runs a separate KLUCB type strategy for each period, the regret normalized by  $\log(T)$  asymptotically scales as  $H \sum_{b \in \mathcal{B}} \beta_b^H \sum_{a \neq a_b^*} \Delta_{a,b} / \mathbf{KL}(\mu_{a,b} | \mu_b^*)$ . KLUCB–RB strategy then performs better than this naive strategy by a factor of the order of  $H/|\mathcal{B}|$ . Also, up to our knowledge, this result is the first showing provably asymptotic optimal regret guarantee in a setting when an

agent attempts at transferring information from past to current bandits without contextual information. In the related but different settings considered in Gheshlaghi Azar et al. (2013); Gentile et al. (2014); Maillard and Mannor (2014), only logarithmic regret was shown, however asymptotic optimality was not proved for the considered strategies. Also, let us remind that  $|\mathcal{B}|$  does not need to be known ahead of time by the KLUCB-RB algorithm.

### 6.3 Sketch of Proof

This section contains a sketch of proof for Theorem 5. We refer to Appendix C.2 for more insights and detailed derivations. The first preoccupation is to ensure that KLUCB-RB is a consistent strategy. This is achieved by showing that KLUCB-RB aggregates observations that indeed come from the same bandits with high probability. In other words, we want to control the number of previously encountered bandits falsely identified as similar to the current one.

**Definition 4** (False positive). *At period  $h \geq 2$  and step  $t \geq 1$ , a previous period  $k \in \llbracket 1, h-1 \rrbracket$  is called a false positive if the test  $\mathbf{T}^{k,h}(t)$  is true while previous bandit  $b_*^k$  differs from current bandit  $b_*^h$ .*

Combining the triangle inequality and time-uniform Gaussian concentration inequalities (see e.g., Abbasi-Yadkori et al. (2011)), we prove necessary condition for having  $Z_a^{k,h}(t) \leq 0$  for some arm  $a \in \mathcal{A}$  at current period  $h$  and time step  $t$ , while having  $b_*^k \neq b_*^h$ .

**Lemma 18** (Condition for false positives). *If there exists a false positive at period  $h \geq 2$  and time step  $t > |\mathcal{A}|$ , then with probability  $1 - 1/t(t+1)$ , it must be that*

$$\min_{k \in \llbracket 1, h-1 \rrbracket : b_*^k \neq b_*^h} \min_{a \in \mathcal{A}} |\mu_{a,b_*^k} - \mu_{a,b_*^h}| \leq 4d \left( \frac{t}{|\mathcal{A}|}, \delta^h(t) \right).$$

The proof of this key result is provided in Appendix C.2.1. It relies on time-uniform concentration inequalities. We now introduce a few quantities.

Let us first consider at period  $h \geq 2$  the time step

$$t_\nu^h := \max \left\{ t \geq |\mathcal{A}| : \gamma_\nu \leq 4d \left( \frac{t}{|\mathcal{A}|}, \delta^h(t) \right) \right\} + 1, \quad (6.8)$$

beyond which there is no false positives with high probability. We define for all  $a \neq a_*^h$ , for all  $0 < \varepsilon < \varepsilon_\nu := \min_{b \in \mathcal{B}} \min_{a \neq a_*^h} \{\Delta_{a,b}, 1\}/2$  the subsets of times when there is a false positive

$$\mathcal{T}_a^h := \{t \geq t_\nu^h : a_{t+1}^h = a \text{ and } \mathcal{K}_+^h(t) \neq \mathcal{K}_*^h(t)\} \quad \mathcal{T}^h := \bigcup_{a \neq a_*^h} \mathcal{T}_a^h, \quad (6.9)$$

where we introduced for convenience the sets  $\mathcal{K}_+^h := \{k \in \llbracket 1, h-1 \rrbracket : \mathbf{T}^{k,h}(t) \text{ is true}\}$  and  $\mathcal{K}_*^h(t) := \{k \in \llbracket 1, h-1 \rrbracket : b_*^k = b_*^h\}$ . We also consider the times when the mean of the current pulled arm is poorly estimated or the best arm  $a_*^h$  is below its mean (either for the current period or by aggregation) and define

$$\mathcal{C}_{a,\varepsilon}^h := \left\{ t \geq 1 : a_{t+1}^h = a \text{ and } \left( |\widehat{\mu}_a^h(t) - \mu_a^h| > \varepsilon \text{ or } u_{a_*^h}^h(t) = U_{a_*^h}^h(t) < \mu_{a_*^h}^h \right) \right\}$$

$$\mathcal{C}_\varepsilon^h := \bigcup_{a \neq a_*^h} \mathcal{C}_{a,\varepsilon}^h \quad (6.10)$$

$$\overline{\mathcal{C}}_{a,\varepsilon}^h := \mathcal{T}_a^h \cup \left\{ t \geq t_\nu^h : t \notin \mathcal{T}^h, a_{t+1}^h = a \text{ and } \left( |\overline{\mu}_a^h(t) - \mu_a^h| > \varepsilon \text{ or } u_{a_*^h}^h(t) = \overline{U}_{a_*^h}^h(t) < \mu_{a_*^h}^h \right) \right\}$$

$$\bar{\mathcal{C}}_\varepsilon^h := \bigcup_{a \neq a_\star^h} \bar{\mathcal{C}}_{a,\varepsilon}^h. \quad (6.11)$$

The size of these (bad events) sets can be controlled by resorting to concentration arguments. The next lemma borrows elements of proof from [Combes and Proutiere \(2014b\)](#) for the estimation of the mean of current pulled arm and [Cappé et al. \(2013\)](#) for the effectiveness of the upper confidence bounds on the empirical means of optimal arms. We adapt these arguments to the routine-bandit setup, and provide additional details in the appendix.

**Lemma 19** (Bounded subsets of times). *For all period  $h \geq 2$ , for all arm  $a \in \mathcal{A}$ , for all  $0 < \varepsilon < \varepsilon_\nu$ ,*

$$\mathbb{E}_\nu [|\mathcal{T}^h|] \leq 1 \quad \mathbb{E}_\nu [|\mathcal{C}_{a,\varepsilon}^h|] \leq 4\varepsilon^{-2} + 2 \quad \mathbb{E}_\nu [|\bar{\mathcal{C}}_{a,\varepsilon}^h|] \leq 4\varepsilon^{-2} + 3.$$

By definition of the index (Eq. 6.7), we have

$$\begin{aligned} \forall t > |\mathcal{A}|, \quad N_a^h(t) \text{KL}(\hat{\mu}_a^h(t) | U_a^h(t)) &= f(t) \\ \bar{N}_a^h(t) \text{KL}(\bar{\mu}_a^h(t) | \bar{U}_a^h(t)) &= f(\bar{K}_t^h T + t). \end{aligned}$$

We then provide logarithmic upper bounds on the aggregated number of pulls  $\bar{N}_a^h(t)$  to deduce the consistency of KLUCB-RB strategy. The following non-trivial result combines standard techniques with the key mechanism of the algorithm.

**Lemma 20** (Consistency). *Under KLUCB-RB strategy for all period  $h \geq 2$ , for all  $0 < \varepsilon < \varepsilon_\nu$ , for all sub-optimal arm  $a \neq a_\star^h$ , for all  $t > |\mathcal{A}|$  such that  $a_{t+1}^h = a$ ,*

$$\text{if } t \notin \mathcal{C}_{a,\varepsilon}^h, \quad N_a^h(t) \leq \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)}, \quad \text{if } t \geq t_\nu^h \text{ and } t \notin \bar{\mathcal{C}}_{a,\varepsilon}^h, \quad \bar{N}_a^h(t) \leq \frac{f(\bar{K}_t^h T + t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)},$$

where  $\bar{K}_t^h := \min \left\{ \bar{K}_t^h, \beta_{b_\star^h}^{h-1}(h-1) \right\}$ . In particular this implies

$$\forall t \geq 1, \forall a \neq a_\star^h, \quad N_a^h(t) \leq \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)} + |\mathcal{C}_{a,\varepsilon}^h| + N_a^h(|\mathcal{A}| + 1),$$

where  $N_a^h(|\mathcal{A}| + 1) \leq 2$  and  $\mathbb{E}_\nu [|\mathcal{C}_{a,\varepsilon}^h|] \leq 4\varepsilon^{-2} + 2$ .

Thanks to Eq. 6.5 that involves the minimum of the aggregated index  $\bar{U}_a^h(t)$  on past episodes and (not aggregated) indexes  $U_a^h(t)$  for the current epoch, the proof proceeds by considering the appropriate sets of time, namely  $t \notin \mathcal{C}_{a,\varepsilon}^h$  or  $t \notin \bar{\mathcal{C}}_{a,\varepsilon}^h$  depending on the situation. In particular, we get for the considered  $a$  that the maximum index  $u_a^h(t)$  is either greater than  $u_{a_\star^h}^h(t) = U_{a_\star^h}^h(t)$  or  $u_{a_\star^h}^h(t) = \bar{U}_{a_\star^h}^h(t)$ , which in turns enable to have a control either on  $N_a^h(t)$  or  $\bar{N}_a^h(t)$ . In order to obtain the last statement, it essentially remains to consider the maximum time  $t' \in [|\mathcal{A}| + 1; t]$  such that  $a_{t'+1}^h = a$  and  $t' \notin \mathcal{C}_{a,\varepsilon}^h$ .

In order to be asymptotically optimal (in the sense of Corollary 4), the second preoccupation is to ensure with high probability that we aggregate all of the observations coming from current bandit  $b_\star^h$  when computing the indexes. From the definition of  $\mathcal{T}^h$  (Eq. 6.9) and Lemma 19, this amounts to ensure that the current most pulled arm and the most pulled arms of previous periods are the optimal arms of the corresponding periods with high probability. By using the consistency of KLUCB-RB, we prove necessary conditions for the most pulled arms being different from the optimal ones.

**Lemma 21** (Most pulled arms). *For all period  $h \geq 2$ , for all  $0 < \varepsilon < \varepsilon_\nu$ , for all  $t \geq t_\nu^h$  such that  $t \notin \mathcal{T}^h$  and  $\bar{a}_t^h \neq a_\star^h$ ,*

$$\frac{t + |\mathcal{K}_\star^h(t)|T}{2} - (f(t) + |\mathcal{K}_\star^h(t)|f(T)) \sum_{a \neq a_\star^h} \frac{1}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)} - (1 + |\mathcal{K}_\star^h(t)|)|\mathcal{A}| \leq \sum_{k \in \mathcal{K}_\star^h(t) \cup \{h\}} |\mathcal{C}_\varepsilon^k|.$$

Let us remind that  $\mathcal{K}_\star^h(t)$ , defined after Lemma 18, counts the previous phases before  $h$  facing the same bandit as the current one, and for which the most-played arm until then agree. Then, by combining Lemma 20 and Lemma 21 we obtain randomized upper bounds on the number of pulls of sub-optimal arms.

**Proposition 6** (Randomized upper bounds). *Under KLUCB-RB strategy, for all bandit  $b \in \mathcal{B}$ , for all sub-optimal arm  $a \neq a_b^\star$ , for all  $0 < \varepsilon < \varepsilon_\nu$ ,*

$$N_{a,b}(H, T) \leq \frac{f(\beta_b^H HT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^\star)} + \sum_{h=1}^H \mathbb{I}_{\{b_\star^h = b\}} \left[ T_{\nu, \varepsilon}^h + 4|\mathcal{C}_\varepsilon^h| + |\bar{\mathcal{C}}_\varepsilon^h| + \frac{f(hT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^\star)} \sum_{k=1}^h \frac{8|\mathcal{C}_\varepsilon^k|}{T} + \mathbb{I}_{\{T \in \mathcal{T}^k\}} \right],$$

where  $T_{\nu, \varepsilon}^h := \max \left\{ t \geq t_\nu^h : \frac{t}{4} - \sum_{a \neq a_\star^h} \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)} \leq |\mathcal{A}| \right\} + 1$  for  $h \geq 2$ , with  $t_\nu^h$  defined in Eq. (6.8).

We prove Theorem 5 by averaging the randomized upper bounds from Proposition 6.

## 6.4 Numerical Experiments

We now perform experiments to illustrate the performance of the proposed KLUCB-RB under different empirical conditions. We compare KLUCB-RB with a baseline strategy which consists in using a KLUCB that restarts from scratch at every new period, that is the default strategy when no information (features) is provided to share information across periods. We also include a comparison with the sequential transfer algorithm  $\text{tUCB}$  Gheshlaghi Azar et al. (2013) which constitutes interesting baseline to compare with, since it transfers the knowledge of past periods to minimize the regret in a very similar context. Through the periods  $h \in \llbracket 1, H \rrbracket$ ,  $\text{tUCB}$  incrementally estimates the mean vectors by the Robust Tensor Power method Anandkumar et al. (2013, 2014), then yielding a deviation of rate  $\mathcal{O}(1/\sqrt{h})$  over the empirical means. Thus, it needs to know in advance the total number of instances  $|\mathcal{B}|$ . Besides the RTP method requires the mean vectors to be linearly independent mutually, which forces the number of arms  $|\mathcal{A}|$  to be larger than  $|\mathcal{B}|$ , while KLUCB-RB can tackle this kind of distributions. The next comparisons between KLUCB-RB and  $\text{tUCB}$  will mainly illustrate the ability of the former to make large profits from the very first periods, while the later needs to get a sufficiently high confidence over the models estimates before beginning to use knowledge from the previous periods.

All experiments are repeated 100 times. Sequence  $(b^h)_{1 \leq h \leq H}$  is chosen randomly each time. All the different strategies are compared based on their cumulative regret (Eq. 6.1). Additional experiments are provided in Appendix C.3.

### 6.4.1 More Arms than Bandits: A Beneficial Case

We first investigate how Assumption 10 can be relaxed in practice. Indeed KLUCB-RB is designed such that only data from previous periods  $k < h$  for which the most pulled arm  $\bar{a}_\star^k$  is the same as the current most pulled arm  $\bar{a}_t^h$  may be aggregated. Consequently, let us define  $\gamma_\nu^\star := \min_{b \neq b'} \min_{a \in \mathcal{A}^\star} |\mu_{a,b} - \mu_{a,b'}|$  with  $\mathcal{A}^\star$  being the set of arms optimal on at least one instance  $b \in \mathcal{B}$ . Assuming that KLUCB-RB converges to the optimal action in a

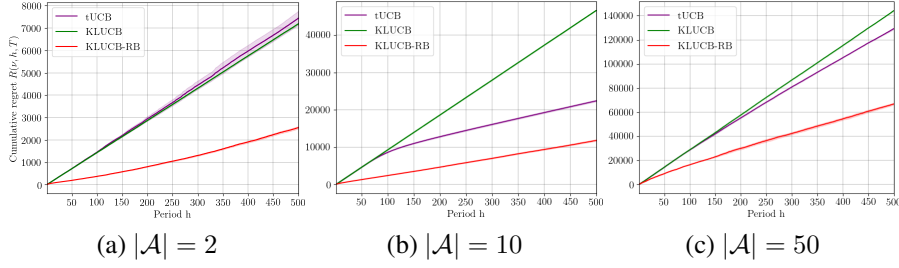


Figure 6.1: Cumulative regret of KLUCB, KLUCB-RB and  $t$ UCB along  $H = 500$  periods of  $T = 10^3$  rounds, for different action sets.

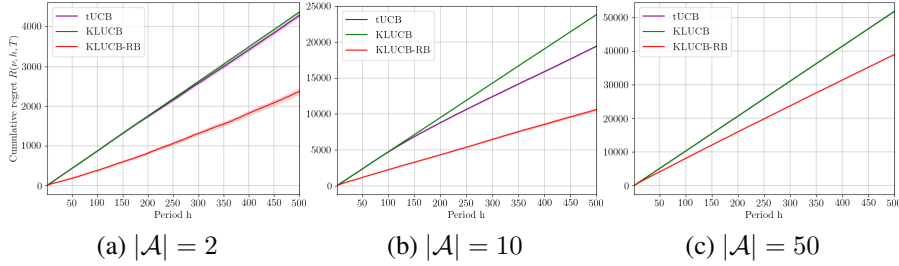


Figure 6.2: Cumulative regret of KLUCB, KLUCB-RB and  $t$ UCB along  $H = 500$  periods of  $T = 100$  rounds, for different action sets.

given period, it is natural in practice to relax Assumption 10 from  $\gamma_\nu > 0$  to  $\gamma_\nu^* > 0$ . Let us consider a routine two-bandit setting  $\mathcal{B} = \{b_1, b_2\}$  with actions  $\mathcal{A}$  such that

$$b_1 : (\mu_{1,b_1}, \mu_{2,b_1}) = \left(\frac{\Delta}{2}, -\frac{\Delta}{2}\right) \quad \text{and} \quad \forall a \geq 3, \mu_{a,b_1} = \mu \quad (6.12)$$

$$b_2 : (\mu_{1,b_2}, \mu_{2,b_2}) = \left(\frac{\Delta}{2} - \gamma, -\frac{\Delta}{2} + \gamma\right) \quad \text{and} \quad \forall a \geq 3, \mu_{a,b_2} = \mu, \quad (6.13)$$

with  $\mu = -\frac{\Delta}{2}$ , and  $\gamma = 0.85\Delta$ , and where  $\Delta = 10\sqrt{\frac{\log(HT)}{T}}$  is set to accommodate the convergence of KLUCB in the experiment. Note that Assumption 10 is not satisfied anymore since  $\gamma_\nu = 0$ , but that  $\gamma_\nu^* = \gamma$ . Fig. 6.1 shows the average cumulative regret with one standard deviation after  $H = 500$  periods of  $T = 10^3$  rounds on settings where  $|\mathcal{A}^*| = 2$  and  $|\mathcal{A}| \geq 2$ .

We observe that KLUCB-RB can largely benefit from relying on previous periods when the number of arms exceeds the number of optimal arms, which naturally happens when  $|\mathcal{A}| > |\mathcal{B}|$ . This can also happen for  $|\mathcal{A}| \leq |\mathcal{B}|$  if several bandits  $b \in \mathcal{B}$  share the same optimal arm. Besides, Fig. 6.2 shows a remake of the same experiment, that is  $\Delta = 10\sqrt{\frac{\log(H \times 10^3)}{10^3}}$ , where the number of rounds per period is decreased from  $10^3$  to  $T = 100$ . We can see that KLUCB-RB still yields good satisfying performances, although  $T$  is not large enough to enable a sure identification at each period of the current instance.

## 6.4.2 Increasing the Number of Bandit Instances

We now consider experiments where we switch among  $|\mathcal{B}| = 5$  four-armed bandits. This highlights the kind of settings which may cause more difficulties to KLUCB-RB in distinguishing the different instances: the lesser is the number of arms  $|\mathcal{A}|$  compared to the number of bandits  $|\mathcal{B}|$ , the harder it should be for KLUCB-RB to distinguish efficiently the different instances, in particular when the separation gaps are tight. Let us precise

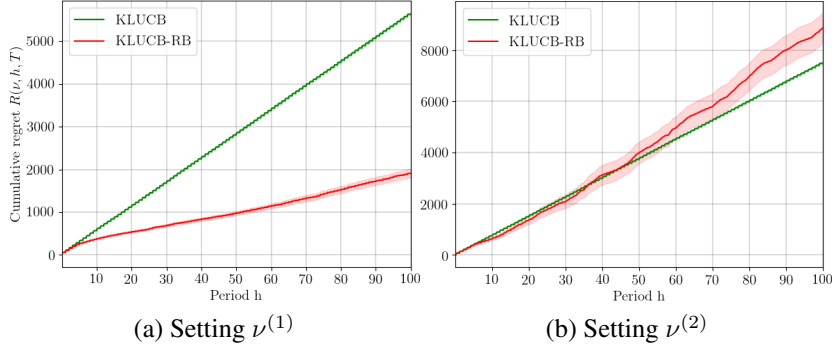


Figure 6.3: Cumulative regret of KLUCB and KLUCB-RB along  $H = 100$  periods of  $T = 5000$  rounds over three generated settings of  $|\mathcal{B}| = 5$  bandit instances with  $|\mathcal{A}| = 4$  arms per instance.

that  $\text{tUCB}$  cannot be tested on such settings, where the number of models  $|\mathcal{B}|$  exceeds the number of arms  $|\mathcal{A}|$ , since it requires that the mean vectors  $(\mu_{a,b})_{a \in \mathcal{A}}$  for all  $b$  in  $\mathcal{B}$  to be linearly independent.

Generating specific settings is far more complicated here than in cases where  $|\mathcal{B}| = 2$  because of the intrinsic dependency between regret gaps  $(\Delta_{a,b})_{a \in \mathcal{A}, b \in \mathcal{B}}$  and separation gaps  $(|\mu_{a,b} - \mu_{a,b'}|)_{a \in \mathcal{A}, b \neq b'}$ . Thus, distributions of bandits  $\nu \in \mathcal{D}$  used in the next experiments are generated randomly so that some conditions are satisfied (see Eq. 6.14, 6.15). Recall that  $\nu := (\nu_{b_1}, \dots, \nu_{b_{|\mathcal{B}|}})$  is the set of bandit configurations in the bandit set  $\mathcal{B}$ . We consider two different distributions  $\nu^{(1)}$  and  $\nu^{(2)}$ , resulting in associated sets of bandits  $\mathcal{B}_1$  and  $\mathcal{B}_2$ , satisfying the condition  $C(\nu)$  in order to ensure the convergence of algorithms at each period:

$$C(\nu) : \forall b \in \mathcal{B}, \quad 8\sqrt{\frac{\log(HT)}{T}} \leq \min_{a \neq a_b^*} \Delta_{a,b} \leq 12\sqrt{\frac{\log(HT)}{T}}. \quad (6.14)$$

Let  $\gamma(\alpha) := \alpha\sqrt{\frac{\log(HT)}{T}}$ . We generate two sets of bandits  $\mathcal{B}_1$  and  $\mathcal{B}_2$  such as to ensure that  $\nu^{(1)}$  and  $\nu^{(2)}$  satisfy

$$\gamma(12) \leq \gamma_{\nu^{(1)}}^* \leq \gamma(16) \quad \gamma(4) \leq \gamma_{\nu^{(2)}}^* \leq \gamma(8). \quad (6.15)$$

Fig. C.3 (Appendix C.3.3) shows the bandit instances in the two generated bandit sets.

All experiments are conducted under the fair frequency  $\beta = 1/|\mathcal{B}|$ . More precisely, once a period  $h \geq 1$  ends,  $b_*^{h+1}$  is sampled uniformly in  $\mathcal{B}$  and independently of the past sequence  $(b_*^k)_{1 \leq k \leq h}$ . Fig. 6.3 shows the average cumulative regret with one standard deviation after  $H = 100$  periods of  $T = 5000$  rounds for the two settings. We observe that the performance of KLUCB-RB is tied to the smallest sub-optimal gap for all bandit instances. Fig. 6.3a highlights that KLUCB-RB outperforms KLUCB if the minimal sub-optimal gap of each bandit is less than the characteristic smaller separation gap  $\gamma_\nu^*$ . This supports the observation from Sec. 6.4.1 that separation on optimal arms is sufficient. When arms are easier to separate than bandits, one might as well restart a classical KLUCB from scratch on each period (Fig. 6.3b). Note that situations where  $0 < \gamma_\nu \ll \min_{b \in \mathcal{B}} \min_{a \neq a_b^*} \Delta_{a,b}$  may not result in a catastrophic loss in learning performances if the arms in  $\mathcal{A}^*$  are *close enough* not to distort estimates computed on aggregated samples of from false positive models (see Appendix C.3).

### 6.4.3 Critical Settings

We saw previously that settings where bandit instances are difficult to distinguish may yield poor performance (see Section 6.4.2, Fig. 6.3b). Indeed, to determine if two estimated bandit models might result from the same bandit, both KLUCB-RB and  $\text{tUCB}$  rely on a compatibility over each arm, i.e. the intersection of confidence intervals. Therefore, it is generally harder to distinguish rollouts from many different distributions

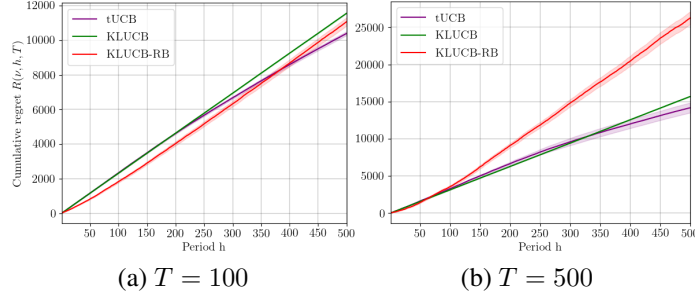


Figure 6.4: Cumulative regret of KLUCB, KLUCB-RB and  $t$ UCB along  $H = 500$  period for different numbers of rounds.

(that is the cardinal of  $|\mathcal{B}|$  is high) when  $|\mathcal{A}|$  is low and differences between arms are tight. To illustrate that, we consider an experiment on the setting described in Figure C.4 (Appendix C.3.3), composed of 4-armed bandits. We recall that  $t$ UCB requires in particular  $|\mathcal{A}| \geq |\mathcal{B}|$ . Thus we choose a set  $|\mathcal{B}|$  of cardinal 4 in order to include a comparison of our algorithm with  $t$ UCB.

Here we have  $|\mathcal{A}^*| = \{0, 1, 3\}$  and  $\gamma_\nu^* := \min_{b \neq b'} \min_{a \in \mathcal{A}^*} |\mu_{a,b} - \mu_{a,b'}| = 0.15$ , while the minimal regret gaps of each instances are  $(\min_{a \neq a'} \Delta_{a,b})_{b \in \mathcal{B}} = (0.74, 0.80, 0.81, 0.89)$ . Consequently, finding the optimal arm at each period independently is here far less difficult than separating the different instances. Such a setting is clearly unfavorable for KLUCB-RB and we expect KLUCB to perform better.

Fig. 6.4a and Fig. 6.4b the cumulative regret for the three strategies, along  $H = 500$  periods of  $T = 100$  and  $T = 500$  rounds respectively. As expected, KLUCB outperforms KLUCB-RB under this critical setting. On the other hand  $t$ UCB seems more robust and displays a cumulative regret trend that would be improving compared with KLUCB in the long run. One should still recall that  $t$ UCB requires knowing the cardinality of  $|\mathcal{B}|$ , while KLUCB-RB does not.

We may notice (Fig. 6.4a) that if the number of rounds  $T$  is sufficiently small, that is KLUCB does not have enough time to converge for each bandit, then KLUCB-RB does not perform significantly worse than KLUCB for the first periods. Then, as  $T$  rises (Fig. 6.4b), KLUCB begins to converge while KLUCB-RB still aggregate samples from confusing instances, which yields an explosion of the cumulative regret curve. We then expect for such setting that KLUCB will need far more longer periods ( $T \rightarrow \infty$ ) to reach a regime in which it will discard all false positive rollouts and takes advantage over KLUCB. On the contrary,  $t$ UCB takes advantage of the knowledge of  $|\mathcal{B}|$  and then waits to have enough confidence over the mean vectors of the 4 models to exploit them.

## 6.5 Conclusion

In this chapter we introduced the new routine bandits framework, for which we provided lower bounds on the regret (Proposition 5). This setting applies well to problems where, for example, customers anonymously return to interact with a system. These dynamics are known to be of interest to the community, as evidenced by the existing literature Gheshlaghi Azar et al. (2013); Gentile et al. (2014); Maillard and Mannor (2014). Routine bandits complement well these existing settings.

We then proposed the KLUCB-RB strategy (Alg. 10) to tackle the routine bandit setting by building on the seminal KLUCB algorithm for classical bandits. We proved upper bounds on the number of sub-optimal plays by KLUCB-RB (Theorem 5), which were used to prove asymptotic upper bounds on the regret (Corollary 4). This result shows the asymptotic optimality of the strategy and thanks to the proof technique that we considered, which is of independent interest, we further obtained finite-time regret guarantees with explicit quantities. We



indeed believe the proof technique may be useful to handle other structured setups beyond routine bandits. We finally provided extensive numerical experiments to highlight the situations where KLUCB-RB can efficiently leverage information from previously encountered bandit instances to improve over a classical KLUCB. More importantly, we highlighted the cost to pay for re-using observations from previous periods, and showed that easy tasks may be better tackled independently. This is akin to an agent that would behave badly by relying on a wrong inductive bias. Fortunately, there are many situations where one can leverage knowledge from bandit instances faced in the past. This would notably be the case if the agent has to select products to recommend from a large set ( $\mathcal{A}$ ) and it turns out that there exists a much smaller set of products ( $\mathcal{A}^*$ ) that is preferred by users (Sec. 6.4.1).

Our results notably show that transferring information from previously encountered bandits can be highly beneficial (e.g., see Fig. 6.1 and 6.3a). However, the lack of prior knowledge about previous instances (including the cardinality of the set of instances) introduces many challenges in transfer learning. For example, attempting to leverage knowledge from previous instances could result in negative transfer if bandits cannot be distinguished properly (e.g., see Fig. 6.4).

Therefore, reducing the cost incurred for separating bandit instances should constitute a relevant angle to tackle as future work. Another natural line of other future work could investigate extensions of KLUCB-RB to the recurring occurrence of other bandit instances, e.g., linear bandits, contextual bandits, and others.

# Chapter 7

## Conclusion

### 7.1 Summary of the presented contributions

First, we have revisited the setup of unimodal multi-armed bandits: We introduced a novel variant based on  $\text{IMED}$  algorithm. This algorithm does not separate exploration from exploitation rounds and is proven optimal for one-dimensional exponential family distributions. Remarkably,  $\text{IMED-UB}$  algorithm does not require any optimization procedure, which can be interesting for practitioners. We also provided a novel proof algorithm, in which we make explicit empirical lower and upper bounds, before tackling the handling of bad events by specific concentration tools, in particular Theorem 6 from Maillard (2018). This proof technique greatly simplifies and shortens the analysis of  $\text{IMED-UB}$ . Last, we provided numerical experiments that show the practical effectiveness of  $\text{IMED-UB}$ . Then, we extend our approach for unimodal bandits to what we called multimodal structure: We introduced  $\text{IMED-MB}$  algorithm that effectively explores local maximums by involving second order indexes (Equation 3.17). Interestingly, this second order exploration is made easy due to the consideration of  $\text{IMED}$  type indexes. Finally, we introduced the novel graph-structured bandit framework, for which we provided instance-dependent lower bounds on the regret (Proposition 3). This setting encompasses Unimodal and Lipschitz structures, which are known to be of interest to the community as evidenced by the existing literature (Saber et al. (2021a); Trinh et al. (2020); Combes and Proutiere (2014a); Magureanu et al. (2014)), and enables to have a unified treatment for such structures. We then proposed  $\text{IMED-GS}$  algorithm (Algorithm 9) to tackle the graph-structured bandit setting by building on popular  $\text{IMED}$  algorithm for unstructured bandits. We proved asymptotically optimal and fully explicit finite time guarantees on the regret, which are not very common in this framework. The analysis of the optimality of the proposed approach led us to state and prove a novel concentration inequality (Theorem 4) that can be useful for subsequent work. Finally, we show that novel  $\text{IMED-GS}$  algorithm has good performance with smaller samples (Figure 4.1).

### 7.2 Some personal satisfactions

The first problem I wanted to solve during my PhD was to propose an optimal and efficient algorithm for Aggregate of Bandits (Section A.1). This problem has been introduced to me by Odalric-Ambrym Maillard, my supervisor. I succeeded in solving this problem by considering the graph-structured structure, a more general setting, and by introducing  $\text{IMED-GS}$  algorithm. What I appreciated the most in solving this problem was to solving it by establishing strong links between  $\text{IMED}$  approaches and Theorem 6 from Maillard (2018), a theoretical paper of my supervisor on boundary crossing probabilities. Above all, I literally loved working with  $\text{IMED}$  algorithm from Honda and Takemura (2015). I appreciated promoting it at SCOOOL, our research team at Inria Lille, and teaching it as a teacher assistant at l’X or CentraleSupélec. These teaching opportunities granted me professional and personal enrichment, and have been possible thanks to Odalric-Ambrym Maillard,

to whom I am very grateful.

## 7.3 Future work

First of all, the content of Chapters 3, 4 and 5 about multimodal and graph-structured bandits could be exploited to follow up on future publications. Now, we briefly provide some possible guidelines for future work.

### 7.3.1 Algorithms with computational efficiency

When the number of arms  $|\mathcal{A}|$  is large, even for non-structured bandit problems, proposing algorithms that are computationally efficient is challenging. We illustrated in Figure 2 from [Saber et al. \(2020\)](#) how  $d$ -IMED-UB, algorithm based on IMED type indexes, can be efficient despite an increasing number of arms when assuming unimodal bandit structure. Further, the IMED indexes and its dichotomic exploration are easy to compute. That is why computation times under  $d$ -IMED-UB algorithm are quite reasonable. Thus,  $d$ -IMED-UB efficiency comes from IMED types indexes that allow second order exploration. We think that this mechanism of exploring the underlying structure based on second order IMED type indexes can be extended to structures other than the unimodal one.

### 7.3.2 IMED for Markov Decision Processes with known transition probabilities

An optimal algorithm for a given Markov Decision Process (MDP) follows optimal random cycle composed in particular of an optimal state-action pairs. Furthermore, there is two levels of sub-optimality for a sub-optimal state-action  $(s, a)$ :

- $s$  is a sub-optimal state but  $a$  is an optimal action in state  $s$  **or** state  $s$  is optimal but action  $a$  is sub-optimal
- $s$  is a sub-optimal state and  $a$  a sub-optimal action in state  $s$ .

Thus, a mechanism allowing second order exploration may be relevant to identify and deal with these two types of sub-optimality. We think our work that consists in proposing IMED approaches for bandit problems can be inspiring to provide algorithms for MDPs with efficient mechanism of exploring sub-optimal state-action pairs. However, unlike the case of bandits, the learner cannot choose an arbitrary state-action pair at each time step when dealing with a Markov Decision Process. To work around this difficulty, we may assume that the transition probabilities plus some bounds on the rewards are known. These additional assumptions would make it possible to better estimate how much does it cost (in terms of regret) to reach (before being able to explore) a specific state from the current state.

### 7.3.3 From routine bandit to non-stationary bandit problem

The routine bandit setting is a variant of the multi-armed bandit problem in which a learner faces every day one of  $\mathcal{B}$  many bandit instances. More specifically, at each period  $h \in \llbracket 1, H \rrbracket$ , the same bandit  $b_\star^h$  is considered during  $T > 1$  consecutive time steps, but the identity  $b_\star^h$  is unknown to the learner. Such a situation typically occurs in recommender systems when a learner may repeatedly serve the same user whose identity is unknown due to privacy issues. The numbers of time steps of episodes are supposed to be the same (equal to  $T$ ) only for convenience. Episodes with different numbers of time steps can be considered. For routine bandits, the end of each period is known. If we assume that the ends of periods are now unknown, then we recover a non-stationary bandit setting where we assume a finite (but unknown) possible changes of distribution. By combining bandit-identification tests with a KLUCB type algorithm, we introduced KLUCB for Routine Bandits

(KLUCB-RB) algorithm. A natural question to explore is: How to adapt KLUCB-RB when the end (and the beginning) of each period is now unknown ?

# Bibliography

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. (2011). Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320.
- Agrawal, R., Teneketzis, D., and Anantharam, V. (1989). Asymptotically efficient adaptive allocation schemes for controlled iid processes: Finite parameter space. *IEEE Transactions on Automatic Control*, 34(3).
- Anandkumar, A., Ge, R., Hsu, D., and Kakade, S. (2013). A tensor spectral approach to learning mixed membership community models. In *Conference on Learning Theory*, pages 867–881. PMLR.
- Anandkumar, A., Ge, R., Hsu, D. J., Kakade, S. M., and Telgarsky, M. (2014). Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15(1):2773–2832.
- Baransi, A., Maillard, O.-A., and Mannor, S. (2014). Sub-sampling for multi-armed bandits. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 115–131. Springer.
- Baudry, D., Kaufmann, E., and Maillard, O.-A. (2020). Sub-sampling for Efficient Non-Parametric Bandit Exploration. In *NeurIPS 2020*, Vancouver, Canada.
- Bubeck, S., Munos, R., Stoltz, G., and Szepesvári, C. (2008). Online optimization of X-armed bandits. In Koller, D., Schuurmans, D., Bengio, Y., and Bottou, L., editors, *Proceedings of the 22nd conference on advances in Neural Information Processing Systems*, NIPS '08, Vancouver, British Columbia, Canada. MIT Press.
- Burnetas, A. N. and Katehakis, M. N. (1997). Optimal adaptive policies for Markov decision processes. *Mathematics of Operations Research*, 22(1):222–255.
- Cappé, O., Garivier, A., Maillard, O.-A., Munos, R., and Stoltz, G. (2013). Kullback–Leibler upper confidence bounds for optimal sequential allocation. *Annals of Statistics*, 41(3):1516–1541.
- Cappé, O., Garivier, A., Maillard, O.-A., Munos, R., and Stoltz, G. (2013). Kullback–Leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516–1541.
- Combes, R., Magureanu, S., and Proutiere, A. (2017). Minimal exploration in structured stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 1763–1771.
- Combes, R. and Proutiere, A. (2014a). Unimodal bandits: Regret lower bounds and optimal algorithms. In *International Conference on Machine Learning*.
- Combes, R. and Proutiere, A. (2014b). Unimodal bandits: Regret lower bounds and optimal algorithms. In *Proceedings of the 31st International Conference on Machine Learning (ICML)*, pages 521–529.
- Cuvelier, T., Combes, R., and Gourdin, E. (2021a). Asymptotically optimal strategies for combinatorial semi-bandits in polynomial time. In *Algorithmic Learning Theory*, pages 505–528. PMLR.

- Cuvelier, T., Combes, R., and Gourdin, E. (2021b). Statistically efficient, polynomial-time algorithms for combinatorial semi-bandits. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 5(1):1–31.
- Degenne, R., Menard, P., Shang, X., and Valko, M. (2020a). Gamification of pure exploration for linear bandits. In III, H. D. and Singh, A., editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 2432–2442. PMLR.
- Degenne, R., Shao, H., and Koolen, W. (2020b). Structure adaptive algorithms for stochastic bandits. In III, H. D. and Singh, A., editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 2443–2452. PMLR.
- Durand, A., Maillard, O.-A., and Pineau, J. (2017). Streaming kernel regression with provably adaptive mean, variance, and regularization. *arXiv preprint arXiv:1708.00768*.
- Gao, X., Qi, H., Wen, X., Zheng, W., Lu, Z., and Hu, Z. (2019). Energy detection adjustment for fair co-existence of wi-fi and laa: A unimodal bandit approach. In *2019 IEEE 5th International Conference on Computer and Communications (ICCC)*, pages 1086–1091.
- Garivier, A. (2013). Informational confidence bounds for self-normalized averages and applications. In *Information Theory Workshop (ITW), 2013 IEEE*, pages 1–5. IEEE.
- Garivier, A. and Cappé, O. (2011). The KL-UCB algorithm for bounded stochastic bandits and beyond. In *Proceedings of the 24th Annual Conference on Learning Theory (COLT)*, pages 359–376.
- Garivier, A., Ménard, P., and Stoltz, G. (2016). Explore first, exploit next: The true shape of regret in bandit problems. *arXiv preprint arXiv:1602.07182*.
- Gentile, C., Li, S., and Zappella, G. (2014). Online clustering of bandits. In *International Conference on Machine Learning*, pages 757–765.
- Gheshlaghi Azar, M., Lazaric, A., and Brunskill, E. (2013). Sequential transfer in multi-armed bandit with finite set of models. In *Proc. NIPS*, pages 2220–2228. Curran Associates, Inc.
- Graves, T. L. and Lai, T. L. (1997). Asymptotically efficient adaptive choice of control laws in controlled markov chains. *SIAM journal on control and optimization*, 35(3):715–743.
- Honda, J. and Takemura, A. (2011). An asymptotically optimal policy for finite support models in the multi-armed bandit problem. *Machine Learning*, 85(3):361–391.
- Honda, J. and Takemura, A. (2015). Non-asymptotic analysis of a new bandit algorithm for semi-bounded rewards. *Machine Learning*, 16:3721–3756.
- Kaufmann, E., Korda, N., and Munos, R. (2012). Thompson sampling: An asymptotically optimal finite-time analysis. In *International Conference on Algorithmic Learning Theory (ALT)*, pages 199–213.
- Kleinberg, R. D., Niculescu-mizil, A., and Sharma, Y. (2008). Regret bounds for sleeping experts and bandits. In Servedio, R. A. and Zhang, T., editors, *Proceedings of the 21st annual Conference On Learning Theory*, volume 80 of *COLT '08*, pages 425–436, Helsinki, Finland. Omnipress.
- Korda, N., Kaufmann, E., and Munos, R. (2013). Thompson Sampling for 1-dimensional exponential family bandits. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1448–1456.

- Kunne, S., Maggi, L., Cohen, J., and Xu, X. (2020). Anytime backtrack unimodal bandits and applications to cloud computing. In *2020 IFIP Networking Conference (Networking)*, pages 82–90.
- Kveton, B., Szepesvari, C., Wen, Z., and Ashkan, A. (2015). Cascading bandits: Learning to rank in the cascade model. In *International Conference on Machine Learning*, pages 767–776. PMLR.
- Kveton, B., Zaheer, M., Szepesvari, C., Li, L., Ghavamzadeh, M., and Boutilier, C. (2020). Randomized exploration in generalized linear bandits. In Chiappa, S. and Calandra, R., editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 2066–2076. PMLR.
- Lai, T. L. (1987). Adaptive treatment allocation and the multi-armed bandit problem. *The Annals of Statistics*, pages 1091–1114.
- Lai, T. L. and Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22.
- Langford, J. and Zhang, T. (2007). The Epoch-Greedy Algorithm for Multi-armed Bandits with Side Information. In Platt, J. C., Koller, D., Singer, Y., Roweis, S. T., Platt, J. C., Koller, D., Singer, Y., and Roweis, S. T., editors, *NIPS*. MIT Press.
- Lattimore, T. and Szepesvari, C. (2017). The end of optimism? an asymptotic analysis of finite-armed linear bandits. In *Artificial Intelligence and Statistics*, pages 728–737.
- Lattimore, T. and Szepesvári, C. (2020). *Bandit algorithms*. Cambridge University Press.
- Lu, S., Wang, G., Hu, Y., and Zhang, L. (2019). Optimal algorithms for Lipschitz bandits with heavy-tailed rewards. In Chaudhuri, K. and Salakhutdinov, R., editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4154–4163. PMLR.
- Lu, T., Pál, D., and Pál, M. (2010). Contextual multi-armed bandits. In Teh, Y. W. and Titterton, M., editors, *Proceedings of the 13th international conference on Artificial Intelligence and Statistics*, volume 9, pages 485–492.
- Magureanu, S. (2018). *Efficient Online Learning under Bandit Feedback*. PhD thesis, KTH Royal Institute of Technology.
- Magureanu, S., Combes, R., and Proutière, A. (2014). Lipschitz bandits: Regret lower bounds and optimal algorithms. In *COLT 2014*.
- Maillard, O.-A. (2018). Boundary crossing probabilities for general exponential families. *Mathematical Methods of Statistics*, 27(1):1–31.
- Maillard, O.-A. and Mannor, S. (2014). Latent bandits. In *International Conference on Machine Learning (ICML)*.
- Peña, V. H., Lai, T. L., and Shao, Q.-M. (2008). *Self-normalized processes: Limit theory and Statistical Applications*. Springer Science & Business Media.
- Pesquerel, F., Saber, H., and Maillard, O.-A. (2021). Stochastic bandits with groups of similar arms. *International Conference on Neural Information Processing Systems (NeurIPS)*.
- Robbins, H. (1952). Some aspects of the sequential design of experiments. *Bulletin of the American Mathematics Society*, 58:527–535.

- Saber, H., Ménard, P., and Maillard, O.-A. (2020). Forced-exploration free strategies for unimodal bandits. *arXiv preprint arXiv:2006.16569*.
- Saber, H., Ménard, P., and Maillard, O.-A. (2021a). Indexed minimum empirical divergence for unimodal bandits. *International Conference on Neural Information Processing Systems (NeurIPS)*.
- Saber, H., Saci, L., Maillard, O.-A., and Durand, A. (2021b). Routine bandits: Minimizing regret on recurring problems. *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PPKD)*.
- Srinivas, N., Krause, A., Kakade, S., and Seeger, M. (2010). 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.
- Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294.
- Thompson, W. R. (1935). 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.
- Trinh, C., Kaufmann, E., Vernade, C., and Combes, R. (2020). Solving bernoulli rank-one bandits with unimodal thompson sampling. In *International Conference on Algorithmic Learning Theory*.
- Van Parys, B. and Golrezaeiand, N. (2020). Optimal learning for structured bandits. *Sloan School of Management, MIT*.
- Wang, T., Ye, W., Geng, D., and Rudin, C. (2020). Towards practical lipschitz bandits. *Proceedings of the 2020 ACM-IMS on Foundations of Data Science Conference*.
- Yu, J. Y. and Mannor, S. (2011). Unimodal bandits. In *ICML*, pages 41–48. Citeseer.



# Appendix A

## Graph-Structured Bandits: Complements

We reintroduce some relevant notations in the following paragraph.

**Notations.** Let  $\nu \in \mathcal{D}_\Theta$ . Let  $\mu^* = \max_{a \in \mathcal{A}} \mu_a$  be the optimal mean and  $\mathcal{A}^*(\nu) = \operatorname{argmax}_{a \in \mathcal{A}} \mu_a$  be the set of optimal arms 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 tower rule we can rewrite the regret as follows:

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

where  $N_a(t) = \sum_{s=1}^t \mathbb{1}_{\{a_s=a\}}$  is the number of pulls of arm  $a$  at time  $t$ . For arms  $a, a' \in \mathcal{A}$  and a relationship matrix  $\theta \in \Theta$  we define their relative gap  $\delta_{a,a'}(\theta) = \theta_{a,a'} - (\mu_a - \mu_{a'})$  and their algebraic gap  $d_{a,a'}(\theta) = \Delta_a - \delta_{a,a'}(\theta)$ . We note that under Assumption 4 for all  $\theta \in \Theta$ ,  $\delta_{a,a}(\theta) = 0$  and  $d_{a,a}(\theta) = \Delta_a \geq 0$ . Finally, we define  $\Theta^*(\nu) := \{\theta \in \Theta : \nu \in \mathcal{D}_\theta\}$  and  $\bar{\Theta}_a(\nu) = \{\theta \in \Theta^*(\nu) : \forall a' \in \mathcal{A}, \theta_{a,a'} \geq 0\}$  for  $a \in \mathcal{A}$ . When there is no possible confusion,  $\mathcal{A}^*(\nu)$ ,  $\Theta^*(\nu)$  and  $\bar{\Theta}_a(\nu)$  are simply denoted  $\mathcal{A}^*$ ,  $\Theta^*$  and  $\bar{\Theta}_a$ . We note that for all  $\theta \in \Theta^*$ , for all  $a, a' \in \mathcal{A}$ ,  $\delta_{a,a'}(\theta) \geq 0$ . For a subset of arms  $\mathcal{A}' \subset \mathcal{A}$ , we denote  $N_{\mathcal{A}'}(t) = \sum_{a' \in \mathcal{A}'} N_{a'}(t)$  the aggregated pulls from arms in  $\mathcal{A}' \subset \mathcal{A}$ . We remind that  $\Delta_{\min} = \min_{a \notin \mathcal{A}^*} \Delta_a$  and, with regard to Assumption 6 and 8, we

introduce for convenience the quantity  $\varepsilon_\nu = \frac{\delta_{\min}}{4} \wedge (1 - \mu^*)$ , where  $\delta_{\min} = \Delta_{\min}$  if  $\Theta$  is a singleton, and  $\delta_{\min} = \min_{a \neq a'} \{\Delta_{\min}, \theta_{a,a'}^* - (\mu_a - \mu_{a'})\}$  otherwise.

### A.1 Structures Unimodal, Lipschitz and Aggregate of Bandits

We assume for convenience that  $\mathcal{A} = \{1, \dots, |\mathcal{A}|\} \subset \mathbb{N}$ . Interestingly, the graph-structure encompasses the classical structures Lipschitz, Unimodal and Aggregate of bandits. We detail each case separately in the following.

**Lipschitz bandits.** The graph-structure can be specified to handle Lipschitz bandits<sup>1</sup> by assuming that  $\Theta = \{\theta\}$  is a singleton and there exists a positive constant  $k > 0$  such that

$$\forall a, a' \in \mathcal{A}, \quad \theta_{a,a'} = k |a - a'|. \quad (\text{A.2})$$

Indeed in such a case, the graph structure specializes to the following condition

$$\forall \nu \in \mathcal{D}_\Theta, \forall a, a' \in \mathcal{A}, \quad |\mu_a - \mu_{a'}| \leq k |a - a'|.$$

<sup>1</sup>The notion of Lipschitz bandit we refer to is the one used in [Magureanu et al. \(2014\)](#). We could have considered more general Lipschitz bandits.

**Unimodal bandits.** Likewise, a unimodal structure can be recovered by assuming

$$\Theta = \left\{ \theta \in \{0, 1\}^{\mathcal{A}^2} : \exists a^* \in \mathcal{A}, \forall a \neq a', \theta_{a,a'} = \begin{cases} 0 & \text{if } a' \in \llbracket a, a^* \rrbracket \text{ or } a' \in \llbracket a^*, a \rrbracket \\ 1 & \text{otherwise} \end{cases} \right\}. \quad (\text{A.3})$$

Indeed in such a case, the graph structure specializes into a classical unimodal structure

$$\forall \nu \in \mathcal{D}_\Theta, \exists a^* \in \mathcal{A}, \forall a < a' < a^* \text{ or } a > a' > a^*, \mu_a \leq \mu_{a'} \leq \mu_{a^*}.$$

**Aggregate of Bandits.** We now introduce yet another specialization of the graph structure, which shows the graph structure goes beyond the previous classical examples. Here, we assume  $\mathcal{A} = \mathcal{X} \times \mathcal{K}$  and view  $\nu = (\nu_a)_{a \in \mathcal{A}}$  as a set of bandits  $\nu = (\nu_x)_{x \in \mathcal{X}}$ , where for all  $x \in \mathcal{X}$ ,  $\nu_x = (\nu_{x,k})_{k \in \mathcal{K}}$  is a  $|\mathcal{K}|$ -multi-armed bandit. Each bandit  $\nu_x$ , for  $x \in \mathcal{X}$  can be seen as a customer segment and the  $|\mathcal{K}|$  arms as  $|\mathcal{K}|$  new customer offers. The goal of the learner is then to exploit the best twinings between customer segments  $\mathcal{X}$  and new customer offers  $\mathcal{K}$  while minimizing regret. We assume some similarities and dissimilarities between the customer segments in the appreciation of customer offers are provided. These similarities and dissimilarities are encoded in a matrix  $\omega = (\omega_{x,x'})_{x,x' \in \mathcal{X}} \subset [-1; 1]^{\mathcal{X}^2}$ , assumed to be known to the learner, such that for all customer segments  $x, x' \in \mathcal{X}$ , for all customer offer  $k \in \mathcal{K}$ ,

$$\mu_{x,k} - \mu_{x',k} \leq \omega_{x,x'}.$$

The matrix  $\omega$  measures the separation between the customer segments and satisfies a pseudo-metric property. That is, for all customer segments  $x, x', x'' \in \mathcal{X}$ ,  $\omega_{x,x} = 0$  and  $\omega_{x,x''} \leq \omega_{x,x'} + \omega_{x',x''}$ . Now, one can express this structure as a graph structure assuming that

$$\Theta = \{\theta\}, \quad \theta_{a,a'} = \omega_{x,x'}, \quad \forall a = (x, k), a' = (x', k') \in \mathcal{A} = \mathcal{X} \times \mathcal{K}. \quad (\text{A.4})$$

Indeed, this definition yields

$$\forall \nu \in \mathcal{D}_\Theta, \forall (x, k), (x', k') \in \mathcal{X} \times \mathcal{K}, \mu_{x,k} - \mu_{x',k'} \leq \omega_{x,x'}.$$

Let us note the Aggregate of bandits is reminiscent of contextual bandits, where similarities between means are dictated by similarities between contexts.

## A.2 Proof related to the regret lower bound

In this section, we provide a lower bound on the cumulative regret when assuming a graph structure. To this end, we follow the classical approach from [Lai and Robbins \(1985\)](#) that we apply to the case of a graph structure. In order to obtain non trivial lower bound we consider algorithms that are consistent (Definition 1).

For each sub-optimal arm  $a \notin \mathcal{A}^*(\nu)$ , we introduce the following distribution-dependent subset of relationship matrices

$$\Theta_a(\nu) := \{\theta \in \Theta^*(\nu) : \forall a^* \in \mathcal{A}^*(\nu), \theta_{a,a^*} > 0 \text{ and } \forall a' \notin \mathcal{A}^*(\nu), \theta_{a,a'} \geq 0\}, \quad (\text{A.5})$$

where  $\Theta^*(\nu) = \{\theta \in \Theta : \nu \in \mathcal{D}_\theta\}$ . Then, we derive from the notion of consistency asymptotic regret lower bounds by considering *most confusing configurations* for each relationship matrix in  $\Theta_a(\nu)$ . These lower bounds involve the set of *informative* sub-optimal arms

$$\mathcal{A}_a(\theta) := \{a' \in \mathcal{A} : d_{a,a'}(\theta) \geq 0\} = \{a' \in \mathcal{A} : \mu_{a'} \leq \mu^* - \theta_{a,a'}\}, \quad \theta \in \Theta_a(\nu). \quad (\text{A.6})$$

When “moving” sub-optimal arm  $a$  to make it optimal in a most confusing configuration,  $\mathcal{A}_a(\theta)$  represents the set of sub-optimal arms which must also be “moved” in order to ensure the “most confusing”<sup>2</sup> bandit for sub-optimal  $a$  belongs to the structure  $\mathcal{D}_\theta$ . By definition of  $\mathcal{A}_a(\theta)$ , for all  $a' \in \mathcal{A}$ , it holds that  $d_{a,a'}(\theta) < 0$  if  $a' \notin \mathcal{A}_a(\theta)$ . Since  $\theta_{a,a} = 0$  and  $\theta_{a,a^*} > 0$  for all  $a^* \in \mathcal{A}^*$ , we further get that  $a \in \mathcal{A}_a(\theta)$  and  $\mathcal{A}_a(\theta) \cap \mathcal{A}^* = \emptyset$ .

**Proposition 7** (Lower bounds on pulls). *Let us consider a consistent bandit algorithm. Then, for all configuration  $\nu \in \mathcal{D}_\Theta$ , for all sub-optimal arm  $a \notin \mathcal{A}^*(\nu)$ , for all relationship matrix  $\theta \in \Theta_a(\nu)$ , under Assumption 4 it must be that*

$$\forall 0 < \varepsilon < \varepsilon_a, \quad \liminf_{T \rightarrow \infty} \frac{1}{\log(T)} \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'} + \varepsilon) \geq 1,$$

where  $\varepsilon_a := \min \left\{ (-d_{a,a'}(\theta))_{a' \notin \mathcal{A}_a(\theta)}, (1 - \mu^*)/2 \right\} > 0$ .

The proof of this result uses a change-of-measure argument and follows classical proof techniques from the literature, see [Lai and Robbins \(1985\)](#); [Agrawal et al. \(1989\)](#); [Graves and Lai \(1997\)](#); [Cappé et al. \(2013\)](#). We detail it in the following sub-section for completeness.

## A.2.1 Proof of Proposition 7

Let us consider a sub-optimal arm  $a \notin \mathcal{A}^*$  and  $\theta \in \Theta_a(\nu)$ . From the definitions of  $\Theta_a$  (Eq. 4.2) and  $\varepsilon_a$  we respectively have

$$\forall a' \in \mathcal{A}, \quad \theta_{a,a'} \geq 0,$$

and

$$\forall a' \in \mathcal{A}_a(\theta), \quad \mu^* - \theta_{a,a'} + \varepsilon_a < 1 \quad \text{and} \quad \forall a' \notin \mathcal{A}_a(\theta), \quad \mu_{a'} \geq \mu^* - \theta_{a,a'} + \varepsilon_a. \quad (\text{A.7})$$

Let us then consider  $0 < \varepsilon < \varepsilon_a$  and the maximal confusing distributions  $\tilde{\nu} = (\tilde{\nu}_{a'})_{a' \in \mathcal{A}}$  for the sub-optimal arm  $a$  with means  $(\tilde{\mu}_{a'})_{a' \in \mathcal{A}}$  such that

$$\begin{aligned} \forall a' \in \mathcal{A}_a(\theta), \quad \tilde{\mu}_{a'} &= \mu^* - \theta_{a,a'} + \varepsilon \\ \forall a' \notin \mathcal{A}_a(\theta), \quad \tilde{\mu}_{a'} &= \mu_{a'}. \end{aligned} \quad (\text{A.8})$$

Then, since  $\theta_{a,a'} \geq 0$  for all  $a' \in \mathcal{A}_a(\theta)$ , we have for all  $a' \in \mathcal{A}$ ,

$$\tilde{\mu}_{a'} \leq \mu^* + \varepsilon,$$

with equality in previous inequality if and only if  $a' \in \mathcal{A}_a(\theta)$  and  $\theta_{a,a'} = 0$ . The set of optimal arms for  $\tilde{\nu}$  is then

$$\mathcal{A}^*(\tilde{\nu}) = \{a' \in \mathcal{A}_a(\theta) : \theta_{a,a'} = 0\} \ni a. \quad (\text{A.9})$$

Furthermore, from the definition of  $\mathcal{A}_a(\theta)$  (Eq. 4.3), we get

$$\forall a' \in \mathcal{A}_a(\theta), \quad \tilde{\mu}_{a'} > \mu_{a'}. \quad (\text{A.10})$$

Then, under Assumptions 4, Eq. A.7, A.8 and A.10 imply

$$\begin{aligned} \forall a', a'' \notin \mathcal{A}_a(\theta), \quad a' \neq a'', \quad \tilde{\mu}_{a'} - \tilde{\mu}_{a''} &= \mu_{a'} - \mu_{a''} \leq \theta_{a',a''} \\ \forall a' \notin \mathcal{A}_a(\theta), \quad \forall a'' \in \mathcal{A}_a(\theta), \quad \tilde{\mu}_{a'} - \tilde{\mu}_{a''} &< \mu_{a'} - \mu_{a''} \leq \theta_{a',a''} \\ \forall a', a'' \in \mathcal{A}_a(\theta), \quad a' \neq a'', \quad \tilde{\mu}_{a'} - \tilde{\mu}_{a''} &= \theta_{a,a''} - \theta_{a,a'} \leq \theta_{a',a''} \\ \forall a' \in \mathcal{A}_a(\theta), \quad \forall a'' \notin \mathcal{A}_a(\theta), \quad \tilde{\mu}_{a'} - \tilde{\mu}_{a''} &\leq \mu^* - \theta_{a,a'} + \varepsilon - (\mu^* - \theta_{a,a''} + \varepsilon) \leq \theta_{a',a''}. \end{aligned}$$

<sup>2</sup>These notions of “moving” and “most confusing” refer to the generic proof technique used to derive regret lower bounds. It involves a change-of-measure argument, from the initial configuration in which the arm is sub-optimal to another one chosen to make it optimal.

This implies

$$\tilde{\nu} \in \mathcal{D}_\theta \subset \mathcal{D}_\Theta. \quad (\text{A.11})$$

Let  $0 < c < 1$ . We will show that almost surely

$$\liminf_{T \rightarrow \infty} \frac{1}{\log(T)} \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) \text{kl}(\mu_{a'} | \tilde{\mu}_{a'}) \geq c.$$

We start with the following inequality

$$\begin{aligned} & \mathbb{P}_\nu \left( \liminf_{T \rightarrow \infty} \frac{1}{\log(T)} \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) \text{kl}(\mu_{a'} | \tilde{\mu}_{a'}) < c \right) \\ & \leq \liminf_{T \rightarrow \infty} \mathbb{P}_\nu \left( \frac{1}{\log(T)} \leq \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) \text{kl}(\mu_{a'} | \tilde{\mu}_{a'}) < c \right). \end{aligned}$$

Let us consider an horizon  $T \geq 1$  and let us introduce the event

$$\Omega_T = \left\{ \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) \text{kl}(\mu_{a'} | \tilde{\mu}_{a'}) < c \log(T) \right\}. \quad (\text{A.12})$$

We want to provide an upper bound on  $\mathbb{P}_\nu(\Omega_T)$  to ensure  $\lim_{T \rightarrow \infty} \mathbb{P}_\nu(\Omega_T) = 0$ . We start by taking advantage of the following lemma.

**Lemma 22** (Change of measure). *For every measurable event  $\Omega$  with respect to  $\nu$  and  $\tilde{\nu}$ ,*

$$\forall x \in \mathbb{R}, \quad \mathbb{P}_\nu(\Omega \cap \mathcal{C}_x) \leq \exp(x) \mathbb{P}_{\tilde{\nu}}(\Omega),$$

where  $\mathcal{C}_x = \left\{ \log\left(\frac{d\nu}{d\tilde{\nu}}(\psi)\right) \leq x \right\}$  and  $\psi = ((a_t), X_t)_{t=1..T}$  is the sequence of pulled arms and rewards.

Let  $\alpha \in (0, 1)$  and let us introduce the event

$$\mathcal{C}_{\alpha, T} = \left\{ \log\left(\frac{d\nu}{d\tilde{\nu}}(\psi)\right) \leq (1 - \alpha) \log(T) \right\}. \quad (\text{A.13})$$

Then we can decompose the probability  $\mathbb{P}_\nu(\Omega_T)$  as follows

$$\mathbb{P}_\nu(\Omega_T) = \mathbb{P}_\nu(\Omega_T \cap \mathcal{C}_{\alpha, T}) + \mathbb{P}_\nu(\Omega_T \cap \mathcal{C}_{\alpha, T}^c) \leq T^{1-\alpha} \mathbb{P}_{\tilde{\nu}}(\Omega_T) + \mathbb{P}_\nu(\Omega_T \cap \mathcal{C}_{\alpha, T}^c). \quad (\text{A.14})$$

Now, we control successively the terms  $T^{1-\alpha} \mathbb{P}_{\tilde{\nu}}(\Omega_T)$  and  $\mathbb{P}_\nu(\Omega_T \cap \mathcal{C}_{\alpha, T}^c)$  and show that they both tend to 0 as  $T$  tends to  $\infty$ .

We first provide an upper bound on  $\mathbb{I}_{\{\Omega_T\}}$  by noting that from Eq. A.9 we have

$$\begin{aligned} \Omega_T & \subset \left\{ \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) < \frac{c}{\min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \log(T) \right\} \\ & = \left\{ T < \frac{c}{\min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \log(T) + \sum_{a' \notin \mathcal{A}_a(\theta)} N_{a'}(T) \right\} \\ & \subset \left\{ T < \frac{c}{\min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \log(T) + \sum_{a' \notin \mathcal{A}^*(\tilde{\nu})} N_{a'}(T) \right\}. \end{aligned}$$

This implies

$$\mathbb{I}_{\{\Omega_T\}} \leq \frac{c}{\min_{a \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \frac{\log(T)}{T} + \sum_{a' \notin \mathcal{A}^*(\tilde{\nu})} \frac{N_{a'}(T)}{T}. \quad (\text{A.15})$$

Since we assume a consistent algorithm on  $\mathcal{D}_\Theta$  and  $\tilde{\nu} \in \mathcal{D}_\Theta$  (Eq. A.11), we know that

$$\forall a' \notin \mathcal{A}^*(\tilde{\nu}), \quad \mathbb{E}_{\tilde{\nu}} \left[ \frac{N_{a'}(T)}{T^\alpha} \right] = o(1), \quad (\text{A.16})$$

therefore from Eq. A.15 and A.16 we get

$$T^{1-\alpha} \mathbb{P}_{\tilde{\nu}}(\Omega_T) = o(1). \quad (\text{A.17})$$

Now we control the remaining term  $\mathbb{P}_{\nu}(\Omega_T \cap \mathcal{C}_{\alpha, T}^c)$ .

For each time  $t \in \llbracket 1, T \rrbracket$ , the reward  $X_t$  is sampled independently from the past and according to  $\nu_{a_t}$ . Hence the likelihood ratio rewrites

$$\frac{d\nu}{d\tilde{\nu}}(\psi) = \prod_{t=1}^T \frac{d\nu_{a_t}}{d\tilde{\nu}_{a_t}}(X_t), \quad \text{where } \frac{d\nu_a}{d\tilde{\nu}_a}(x) = \frac{\mu_a^x (1 - \mu_a)^{1-x}}{\tilde{\mu}_a^x (1 - \tilde{\mu}_a)^{1-x}}, \quad \forall a \in \mathcal{A}, \forall x \in \{0, 1\}. \quad (\text{A.18})$$

Thus, since for all  $a' \notin \mathcal{A}_a(\theta)$ ,  $\tilde{\mu}_a = \mu_a$ , the log-likelihood ratio is

$$\log \left( \frac{d\nu}{d\tilde{\nu}}(\psi) \right) = \sum_{a' \in \mathcal{A}_a(\theta)} \sum_{t=1}^T \mathbb{I}_{\{a_t = a'\}} \log \left( \frac{d\nu_{a'}}{d\tilde{\nu}_{a'}}(X_t) \right). \quad (\text{A.19})$$

Hence from Eq. A.12, A.13 and A.19 we can rewrite the set

$$\Omega \cap \mathcal{C}_{\alpha, T}^c = \left\{ \begin{array}{l} \sum_{a' \in \mathcal{A}_a(\theta)} \sum_{t=1}^T \mathbb{I}_{\{a_t = a'\}} \left[ \log \left( \frac{d\nu_{a'}}{d\tilde{\nu}_{a'}}(X_t) \right) - \text{kl}(\mu_{a'} | \tilde{\mu}_{a'}) \right] > (1 - \alpha - c) \log(T) \\ \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) \text{kl}(\mu_{a'} | \tilde{\mu}_{a'}) < c \log(T) \end{array} \right\}.$$

Let us introduce  $X_a^n = X_{\tau_a^n}$  where  $\tau_a^n = \min \{t \geq 1 : N_a(t) = n\}$  for all  $a \in \mathcal{A}$ . Note that the random variables  $\tau_a^n$  are predictable stopping times, since  $\{\tau_a^n = t\}$  is measurable with respect to the filtration generated by  $(a_1, X_1, \dots, a_{t-1}, X_{t-1})$ . For  $a' \in \mathcal{A}_a(\theta)$  and  $n \geq 1$ , let us consider

$$Z_{a'}^n = \frac{d\nu_{a'}}{d\tilde{\nu}_{a'}}(X_{a'}^n).$$

Then  $Z_{a'}^n$  is positive and bounded by  $B_{a'} = 1/\tilde{\mu}_{a'}(1 - \tilde{\mu}_{a'})$ , with mean  $\mathbb{E}_{\nu}[Z_{a'}^n] = \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})$ . Furthermore, the random variables  $Z_{a'}^n$ , for  $a' \in \mathcal{A}_a(\theta)$  and  $n \geq 1$ , are independent. Thus, it holds

$$\Omega_T \cap \mathcal{C}_{\alpha, T}^c \subset \left\{ \max_{m \in \mathcal{M}_T} \sum_{a' \in \mathcal{A}_a(\theta)} \sum_{n=1..m_{a'}} Z_{a'}^n - \mathbb{E}_{\nu}[Z_{a'}^n] > \left( \frac{1-\alpha}{c} - 1 \right) c \log(T) \right\}, \quad (\text{A.20})$$

where  $\mathcal{M}_T := \left\{ (m_{a'}) \in \llbracket 1, T \rrbracket^{\mathcal{A}_a(\theta)} : \sum_{a' \in \mathcal{A}_a(\theta)} m_{a'} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'}) < c \log(T) \right\}$ .

In the following, we control the asymptotic concentration of random variables  $(Z_{a'}^n)$  by applying Doob's maximal inequality. For  $a' \in \mathcal{A}_a(\theta)$  and  $\lambda > 0$ , let us introduce the super-martingale

$$(M_{a',m})_{m \geq 0} = \left( \exp \left( \lambda \sum_{n=1}^m (Z_{a'}^n - E[Z_{a'}^n]) - m \lambda^2 \frac{B_{a'}^2}{8} \right) \right)_{m \geq 0}. \quad (\text{A.21})$$

Then noting that

$$\forall (m_{a'}) \in \mathcal{M}_T, \frac{\sum_{a' \in \mathcal{A}_a(\theta)} \lambda^2 m_{a'} \frac{B_{a'}^2}{8}}{c \log(T)} < \frac{\sum_{a' \in \mathcal{A}_a(\theta)} \lambda^2 m_{a'} \frac{B_{a'}^2}{8}}{\sum_{a' \in \mathcal{A}_a(\theta)} m_{a'} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \leq \lambda^2 \frac{\max_{a' \in \mathcal{A}_a(\theta)} B_{a'}^2}{8 \min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})}, \quad (\text{A.22})$$

we obtain from Eq. A.20 and A.21 that

$$\begin{aligned} \Omega_T \cap \mathcal{C}_{\alpha,T}^c &\subset \left\{ \max_{(m_{a'}) \in \mathcal{M}_T} \prod_{a' \in \mathcal{A}_a(\theta)} M_{a',m_{a'}} > T \left[ \lambda \left( \frac{1-\alpha}{c} - 1 \right) - \lambda^2 \frac{\max_{a' \in \mathcal{A}_a(\theta)} B_{a'}^2}{8 \min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \right]^c \right\} \\ &\subset \left\{ \exists a' \in \mathcal{A}_a(\theta) : \max_{m \leq \bar{m}} M_{a',m} > T^\gamma \right\}, \end{aligned} \quad (\text{A.23})$$

where  $\bar{m} = \frac{c \log(T)}{\min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})}$  and  $\gamma = \left[ \lambda \left( \frac{1-\alpha}{c} - 1 \right) - \lambda^2 \frac{\max_{a' \in \mathcal{A}_a(\theta)} B_{a'}^2}{8 \min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \right] \frac{c}{|\mathcal{A}_a(\theta)|}$ .

In order to have  $\gamma > 0$  in (A.23), we impose:

- $0 < \alpha < 1 - c$  (this implies  $\frac{1-\alpha}{c} - 1 > 0$ )
- $\lambda \in \operatorname{argmax}_{\lambda' \geq 0} \left\{ \lambda' \left( \frac{1-\alpha}{c} - 1 \right) - \lambda'^2 \frac{\max_{a' \in \mathcal{A}_a(\theta)} B_{a'}^2}{8 \min_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \tilde{\mu}_{a'})} \right\} > 0$ .

Then from (A.23) we have

$$\begin{aligned} \mathbb{P}_\nu(\Omega_T \cap \mathcal{C}_{\alpha,T}^c) &\leq \sum_{a' \in \mathcal{A}_a(\theta)} \mathbb{P}_\nu \left( \max_{m \leq \bar{m}} M_{a',m} > T^\gamma \right) \quad (\text{Union bound}) \\ &\leq \sum_{a' \in \mathcal{A}_a(\theta)} \frac{\mathbb{E}_\nu[M_{a',0}]}{T^\gamma} \quad (\text{Doob's maximal inequality}) \\ &= \frac{|\mathcal{A}_a(\theta)|}{T^\gamma}. \end{aligned}$$

This implies the following control

$$\mathbb{P}_\nu(\Omega_T \cap E_T^c) = o(1). \quad (\text{A.24})$$

Finally, by combining (A.8), (A.12), (A.14), (A.17) and (A.24), we show

$$\forall 0 < \varepsilon < \varepsilon_a, \forall 0 < c < 1, \quad \liminf_{T \rightarrow \infty} \frac{1}{\log(T)} \sum_{a' \in \mathcal{A}_a(\theta)} N_{a'}(T) \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'} + \varepsilon) \geq c. \quad (\text{A.25})$$

The proof ends by doing  $c \rightarrow 1$ .

## A.2.2 Proof of Proposition 3

Let  $(T_k)_{k \in \mathbb{N}}$  be a sub-sequence such that

$$\liminf_{T \rightarrow \infty} \frac{R(T, \nu)}{\log(T)} = \lim_{k \rightarrow \infty} \frac{R(T_k, \nu)}{\log(T_k)}.$$

We assume that this limit is finite otherwise the result is straightforward. This implies in particular

$$\forall a \notin \mathcal{A}, \quad \limsup_{k \rightarrow \infty} \frac{\mathbb{E}_\nu [N_a(T_k)]}{\log(T_k)} < +\infty.$$

By Cantor's diagonal argument there exists an extraction of  $(T_k)_{k \in \mathbb{N}}$  denoted by  $(T'_k)_{k \in \mathbb{N}}$  such that for all  $a \notin \mathcal{A}^*$ , there exist  $N_a \in \mathbb{R}_+$  such that

$$\lim_{k' \rightarrow \infty} \frac{\mathbb{E}_\nu [N_a(T'_k)]}{\log(T'_k)} = N_a.$$

Hence we get

$$\liminf_{T \rightarrow \infty} \frac{R(T, \nu)}{\log(T)} = \sum_{a \notin \mathcal{A}^*} N_a \Delta_a.$$

But thanks to Proposition 7, under Assumptions 4, we have for all  $a \notin \mathcal{A}^*$  such that  $\Theta_a \neq \emptyset$ ,

$$\begin{aligned} & \inf_{\theta \in \Theta_a} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'}) N_a \\ &= \inf_{\theta \in \Theta_a} \min_{0 < \varepsilon < \varepsilon_a} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'} + \varepsilon) N_a \\ &= \inf_{\theta \in \Theta_a} \min_{0 < \varepsilon < \varepsilon_a} \lim_{k \rightarrow \infty} \mathbb{E}_\nu \left[ \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'} + \varepsilon) \frac{N_{a'}(T'_k)}{\log(T'_k)} \right] \\ &\geq \inf_{\theta \in \Theta_a} \min_{0 < \varepsilon < \varepsilon_a} \mathbb{E}_\nu \left[ \liminf_{k \rightarrow \infty} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'} + \varepsilon) \frac{N_{a'}(T'_k)}{\log(T'_k)} \right] \\ &\geq \inf_{\theta \in \Theta_a} \min_{0 < \varepsilon < \varepsilon_a} \mathbb{E}_\nu \left[ \liminf_{T \rightarrow \infty} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'} + \varepsilon) \frac{N_{a'}(T)}{\log(T)} \right] \geq 1. \end{aligned}$$

Therefore we obtain the lower bound

$$\begin{aligned} \liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log(T)} &\geq \mathfrak{C}_\Theta(\nu) := \inf_{n \in \mathbb{R}_+^{\mathcal{A}}} \sum_{a \notin \mathcal{A}^*} n_a \Delta_a \\ &\text{s.t. } \forall a \notin \mathcal{A}^*, \Theta_a \neq \emptyset, \inf_{\theta \in \Theta_a} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'}) n_{a'} \geq 1. \end{aligned}$$

To end the proof we show that for all  $n \in \mathbb{R}_+^{\mathcal{A}}$  for all sub-optimal arm  $a \notin \mathcal{A}^*$  such that  $\Theta_a \neq \emptyset$ ,

$$\inf_{\theta \in \Theta_a} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'}) n_{a'} = \min_{\theta \in \Theta_a} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'}) n_{a'}.$$

First, we rewrite

$$\inf_{\theta \in \Theta_a} \sum_{a' \in \mathcal{A}_a(\theta)} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'}) n_{a'} = \inf_{\theta \in \Theta_a} \sum_{a' \notin \mathcal{A}^*} n_{a'} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'}) \mathbb{I}_{\{\mu_{a'} \leq \mu^* - \theta_{a,a'}\}}.$$

Then, we simply note that the function

$$: \theta \in \bar{\Theta}_a \mapsto \sum_{a' \notin \mathcal{A}^*} n_{a'} \text{kl}(\mu_{a'} | \mu^* - \theta_{a,a'}) \mathbb{I}_{\{\mu_{a'} \leq \mu^* - \theta_{a,a'}\}}$$

is continuous.

### A.3 Technical results

In this section, some practical results, in particular about  $\text{kl}(\cdot | \cdot)$  and the continuity of minimization problem, are detailed.

**Lemma 23** (kl inequality 1). *For  $0 < p < q < 1$ ,*

$$\text{kl}(p|q) \leq \frac{(q-p)^2}{q(1-q)}.$$

*Proof.* Let us consider

$$f : x \in [p; q] \mapsto \text{kl}(p|x) - \frac{(x-p)^2}{x(1-x)}.$$

Then  $f$  admits a derivative in  $]p; q[$  and for  $x \in ]p; q[$ ,

$$\begin{aligned} f'(x) &= \frac{x-p}{x(1-x)} - \frac{2(x-p)}{x(1-x)} + \frac{(x-p)^2(1-x-x)}{x^2(1-x)^2} \\ &= -\frac{x-p}{x(1-x)} + \frac{(x-p)^2(1-2x)}{x^2(1-x)^2} \\ &= \frac{x-p}{x^2(1-x)^2} [-x(1-x) + (x-p)(1-2x)] \\ &= \frac{x-p}{x^2(1-x)^2} [-x^2 + 2px - p] \\ &= \frac{x-p}{x^2(1-x)^2} [-(x-p)^2 + p^2 - p] \\ &< 0. \end{aligned}$$

Thus,  $f$  is a decreasing function. Since  $f(p) = 0$ , we have  $f(q) \leq 0$ , which ends the proof.  $\square$

**Lemma 24** (kl inequality 2). *For  $0 < p < q < q' < 1$ ,*

$$\text{kl}(p|q') - \text{kl}(p|q) \geq 2(q-q')^2.$$

*Proof.* We note that

$$: x \in (0, q] \mapsto \text{kl}(x|q') - \text{kl}(x|q) = x \log\left(\frac{q}{q'}\right) + (1-x) \log\left(\frac{1-q}{1-q'}\right)$$

is an affine function with slope  $\log\left(\frac{q(1-q')}{q'(1-q)}\right) < 0$  since  $q < q'$ . This implies

$$\text{kl}(p|q') - \text{kl}(p|q) \geq \text{kl}(q|q') - \text{kl}(q|q) = \text{kl}(q|q').$$



Then Pinsker's inequality implies

$$\text{kl}(q|q') \geq 2(q - q')^2,$$

which ends the proof.  $\square$

**Lemma 25** (kl inequality 3). For  $z \geq 0$ , for  $0 < p' < p < q - z < q' - z < 1$ ,

$$\frac{\text{kl}(p|q-z)}{q-p} \leq \frac{\text{kl}(p'|q'-z)}{q'-p'}.$$

*Proof.* We show that functions

$$f : x \in [p', p] \mapsto \frac{\text{kl}(x|q-z)}{q-x} \tag{A.26}$$

$$g : y \in [q, q'] \mapsto \frac{\text{kl}(p'|y-z)}{y-p'} \tag{A.27}$$

are respectively decreasing and increasing functions. For  $x \in [p', p]$ ,

$$\begin{aligned} f'(x) &= \frac{1}{q-x} \times \frac{\partial \text{kl}}{\partial p}(x|q-z) + \frac{1}{(q-x)^2} \times \text{kl}(x|q-z) \\ &= \frac{1}{(q-x)^2} \times \left[ \text{kl}(x|q-z) - \left( -(q-z) \frac{\partial \text{kl}}{\partial p}(x|q-z) \right) \right]. \end{aligned} \tag{A.28}$$

Since  $x < p < q - z$ , from Equation (A.28) and Lemma 26 we get

$$f'(x) = -\frac{\text{kl}(q-z|x)}{(q-x)^2} \leq 0. \tag{A.29}$$

For  $y \in [q, q']$ ,

$$\begin{aligned} g'(y) &= \frac{1}{y-p'} \times \frac{\partial \text{kl}}{\partial q}(p'|y-z) - \frac{1}{(y-p')^2} \times \text{kl}(p'|y-z) \\ &= \frac{1}{y-p'} \times \frac{y-z-p'}{(y-z)(1-y+z)} - \frac{1}{(y-p')^2} \times \text{kl}(p'|y-z) \\ &= \frac{(y-z)(y-z-p') - (y-z)(1-y+z)\text{kl}(p'|y-z)}{(y-p')^2(y-z)(1-y+z)}. \end{aligned} \tag{A.30}$$

Since  $p' < q - z < y - z$ , from Equation (A.30) and Lemma 23 we get

$$\begin{aligned} g'(y) &\geq \frac{(y-z)(y-z-p') - (y-z-p')^2}{(y-p')^2(y-z)(1-y+z)} \\ &= \frac{p'(y-z-p')}{(y-p')^2(y-z)(1-y+z)} \\ &\geq 0. \end{aligned} \tag{A.31}$$

$\square$

**Lemma 26** (kl inequality 4). For  $0 < p < q < 1$ ,

$$-(q-p) \frac{\partial \text{kl}}{\partial p}(p|q) = \text{kl}(p|q) + \text{kl}(q|p).$$

*Proof.* We have

$$\text{kl}(p|q) = p \log(p) - p \log(q) + (1-p) \log(1-p) - (1-p) \log(1-q) \quad (\text{A.32})$$

$$\text{kl}(q|p) = q \log(q) - q \log(p) + (1-q) \log(1-q) - (1-q) \log(1-p), \quad (\text{A.33})$$

which implies

$$\begin{aligned} \frac{\partial \text{kl}}{\partial p}(p|q) &= \log(p) + p \times \frac{1}{p} - \log(q) - \log(1-p) + (1-p) \times \frac{(-1)}{1-p} + \log(1-q) \\ &= \log(p) + 1 - \log(q) - 1 - \log(1-p) + \log(1-q) \\ &= \log(p) - \log(q) - \log(1-p) + \log(1-q) \end{aligned} \quad (\text{A.34})$$

and

$$\begin{aligned} \text{kl}(p|q) + \text{kl}(q|p) &= (p-q) \log(p) + (q-p) \log(q) + (q-p) \log(1-p) + (p-q) \log(1-q) \\ &= -(q-p) [\log(p) - \log(q) - \log(1-p) + \log(1-q)] \\ &= -(q-p) \frac{\partial \text{kl}}{\partial p}(p|q). \end{aligned} \quad (\text{A.35})$$

□

**Lemma 27** (Continuity of minimization problem).  $:\nu' \in \mathcal{D}_\Theta \mapsto c_\Theta(\nu') \in \mathbb{R}_+$  is well defined and continuous.

Please refer to Section B.2 from [Combes et al. \(2017\)](#) for a proof of Lemma 27.

**Lemma 28** (Real analysis 1). For  $x \geq e$ ,

$$\forall y \geq x (\mathbf{f}_\xi(x))^2, \quad \frac{y}{\mathbf{f}_\xi(y)} \geq x.$$

**Lemma 29** (Real analysis 2). For  $x > 1$ ,

$$\frac{1}{-\log(1-1/x)} \leq x.$$

## A.4 Additional experiments

In this section we introduce PO-IMED-GS algorithm, an inspired but simplified version of IMED-GS, that may be appealing to the practitioner. Then, we provide additional experiments in which the regret is averaged over many random structured configurations for each of the three considered structures (Unimodal, Lipschitz and Aggregate of bandits).

### A.4.1 PO-IMED-GS algorithm

We introduce in this subsection PO-IMED-GS algorithm, for Pareto-Optimal IMED-GS algorithm. PO-IMED-GS is inspired from IMED-GS and simply consists in pulling the arm with minimal graph-structured IMED-type index at each time step. Pulling this arm intuitively ensures the constrains in the optimization problem of the regret lower bound are asymptotically satisfied (Pareto optimality). In particular, PO-IMED-GS algorithm does not solve any optimization problem Please refer to (4.9) for the definition of these graph-structured indexes. Hence this algorithm is simpler to implement, interpret and also has lower computational complexity,

although does not a priori enjoy as refined theoretical guarantees as IMED-GS. PO-IMED-GS algorithm is summarized in Algorithm 11.

---

**Algorithm 11** PO-IMED-GS

---

- 1: **Input:** Structure  $\Theta$ ,  $\xi$ .
  - 2: Pull each arm once
  - 3: **for**  $t = |\mathcal{A}| \dots T - 1$  **do**
  - 4:   Pull arbitrarily  $a_{t+1} \in \operatorname{argmin}_{a \in \mathcal{A}} \bar{I}_a(t)$ , see (4.9).
  - 5: **end for**
- 

## A.4.2 Regrets Averaged on Random Structured Configurations

For each structure of Figure A.1, we considered  $|\mathcal{A}| = 18$  arms, a time horizon of  $T = 3000$ , and results averaged over 100 randomly generated structured configurations with 10 independent runs for each obtained structured configuration (that is 1000 independent runs in total). We compare each time IMED-GS to IMED algorithm for unstructured bandits and PO-IMED-GS algorithm. We also compare IMED-GS to OSSB algorithm for generic structured bandits. For Unimodal structure we add specific comparison with OSUB and UTS that are specialized to this structure, and for Lipschitz structure we add numerical comparison with CKL-UCB. We further report the IMED-GS run with setting  $d = |\mathcal{A}| - 1$ . Here the parameter  $\Gamma$  is set to  $|\mathcal{A}|^{1.5}$  and  $\xi = 1$ . The sequence  $(\gamma_t)_{t \geq 1}$  is set to  $(\gamma_1 \log(t)^{-0.25})_{t \geq 1}$ . The Lipschitz constant is set to  $k = 0.03$  and applies for all the random configurations sampled for the experiment regarding the Lipschitz structure. For the Aggregate of bandits, the set of arms is decomposed as  $\mathcal{A} = \mathcal{X} \times \mathcal{K}$ , with  $\mathcal{X} = \llbracket 1, 2 \rrbracket$  and  $\mathcal{K} = \llbracket 1, 9 \rrbracket$ . This means we assume an aggregate of 2 bandits with 9 arms each. The relationship matrix is set to  $\omega_{x,x'} = 0.05 |x - x'|$  and applies for all the random configurations sampled for the experiment regarding the Aggregate of bandits.

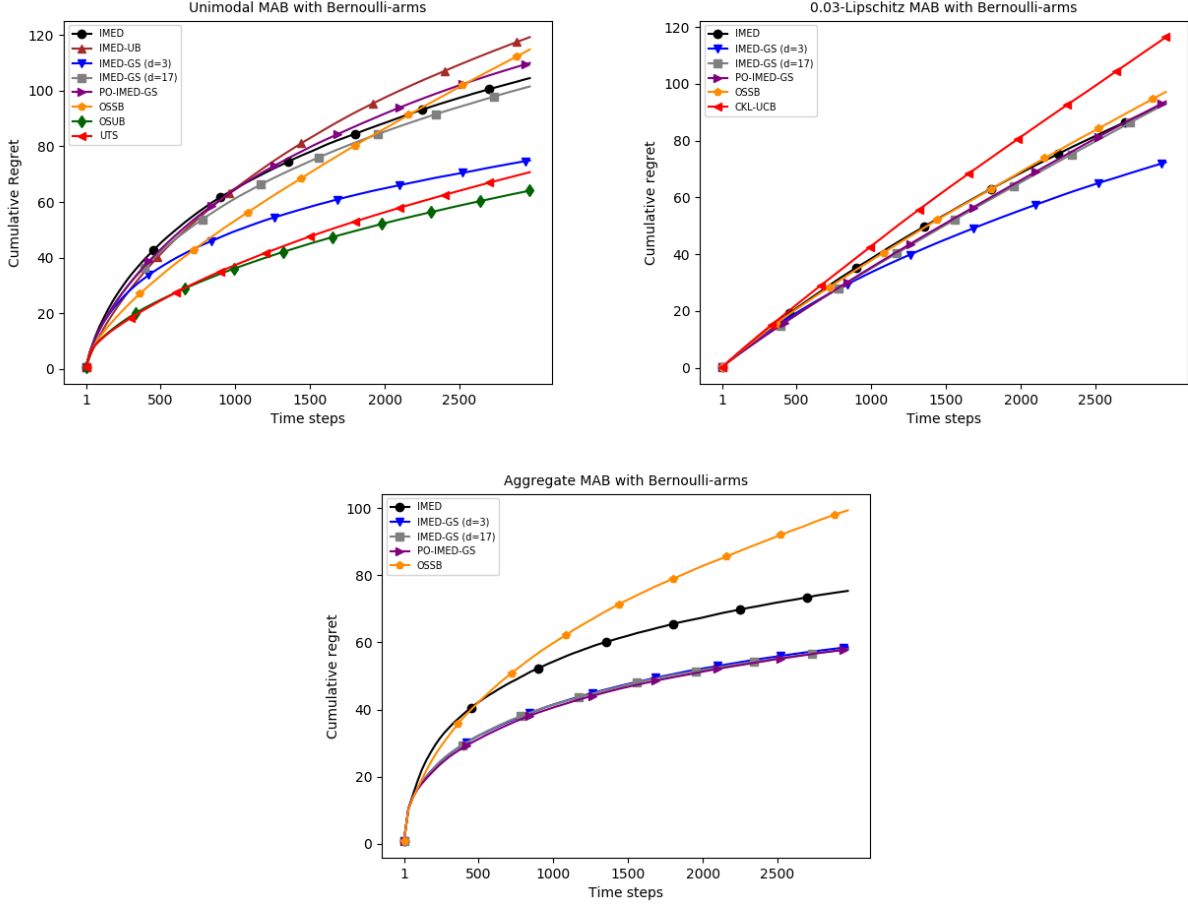


Figure A.1: Comparison of IMED-GS to other algorithms averaged over 100 randomly generated structured bandit instances.

**Discussion.** In the experiments, IMED-GS performs well for each considered structure and competes with specialized state-of-the-art algorithms for the Unimodal and Lipschitz structures. Furthermore, it seems that there is little benefit to use IMED-GS with no burning phase (when the parameter  $d = |\mathcal{A}| - 1$ ) over PO-IMED-GS. This seems to confirm that reaching Pareto-Optimality provides good results in practice. Note that there is no burn-in phase in PO-IMED-GS algorithm, which makes relevant the comparison with IMED-GS with parameter  $d$  set to  $|\mathcal{A}| - 1$ . Finally, the theoretical benefit consisting in introducing the parameter  $d$  (less than  $|\mathcal{A}| - 1$ ) in order to reduce the constant  $C_{\xi, d, \varepsilon} = O(\varepsilon^{-2})$  in the upper bound on the regret (Theorem 3) implies practical benefit as it is highlighted in Figure A.1. Indeed, it appears that IMED-GS with parameter  $d$  set to 3 outperforms or competes with IMED-GS with parameter  $d$  set to  $|\mathcal{A}| - 1$ . This is a practical illustration of the benefit of having introduced the refined concentration inequality (Theorem 4). We recall that with the parameter  $d$  we replace the current informative sets of arms  $\widehat{\mathcal{A}}_a(t)$  for  $a \in \mathcal{A}$  with the sets  $\widehat{\mathcal{A}}_a^{(d)}(t) \subset \widehat{\mathcal{A}}_a(t)$  consisting of  $\{a\}$  plus the  $d$ -th most pulled arms from  $\widehat{\mathcal{A}}_a(t) \setminus \{a\}$ . It is justified by the fact that in the beginning, no structure can reasonably be exploited due to the poor estimates.

# Appendix B

## Generic Tools

### B.1 Non-reliable current means

In this section, we define and study relevant subsets of time steps for which the current mean of a specific arm is not reliable. These subsets of times appear in `IMED-UB` finite-time analysis. Similar subsets appear in finite-time analysis of `IMED-GS` (Section 4.5.3). Note that the definitions and the stated properties of these subsets of time steps are independent from the considered algorithms.

For all arms  $a, a' \in \mathcal{A}$  and for all accuracy  $\varepsilon > 0$ , let  $\mathcal{E}_{a,a'}^+(\varepsilon)$  be the set of times where the current mean of arm  $a$   $\varepsilon$ -deviates from above while arm  $a$  has more pulls than the current pulled arm  $a'$ ,

$$\mathcal{E}_{a,a'}^+(\varepsilon) := \{t \in \llbracket 1, T-1 \rrbracket : a_{t+1} = a', N_{a'}(t) \leq N_a(t), \widehat{\mu}_a(t) \geq \mu_a + \varepsilon\}. \quad (\text{B.1})$$

We similarly define

$$\mathcal{E}_{a,a'}^-(\varepsilon) := \{t \in \llbracket 1, T-1 \rrbracket : a_{t+1} = a', N_{a'}(t) \leq N_a(t), \widehat{\mu}_a(t) \leq \mu_a - \varepsilon\}. \quad (\text{B.2})$$

We also define

$$\mathcal{E}_{a,a'}(\varepsilon) = \mathcal{E}_{a,a'}^+(\varepsilon) \cup \mathcal{E}_{a,a'}^-(\varepsilon). \quad (\text{B.3})$$

**Definition 5** (KL-log deviation). *For  $\varepsilon > 0$ , the couple of arms  $(a, a') \in \mathcal{A}^2$  shows  $\varepsilon^-$ -KL-log deviation at time step  $t \geq 1$  if the following conditions are satisfied*

- (1)  $a_{t+1} = a'$
- (2)  $\widehat{\mu}_a(t) \leq \mu_a - \varepsilon$
- (3)  $\log(N_{a'}(t)) \leq N_a(t) \mathbf{KL}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) + \log(N_a(t))$ .

For all couple of arms  $(a, a') \in \mathcal{A}^2$  and for all accuracy  $\varepsilon > 0$ , let  $\mathcal{K}_{a,a'}^-(\varepsilon)$  be the set of times where couple of arms  $(a, a')$  shows  $\varepsilon^-$ -KL-log deviation, that is

$$\mathcal{K}_{a,a'}^-(\varepsilon) := \left\{ t \in \llbracket 1, T-1 \rrbracket : \begin{array}{l} (1) \ a_{t+1} = a' \\ (2) \ \widehat{\mu}_a(t) \leq \mu_a - \varepsilon \\ (3) \ \log(N_{a'}(t)) \leq N_a(t) \mathbf{KL}(\widehat{\mu}_a(t) | \mu_a - \varepsilon) + \log(N_a(t)) \end{array} \right\}. \quad (\text{B.4})$$

We note that

$$\mathcal{E}_{a,a'}^-(\varepsilon) \subset \mathcal{K}_{a,a'}^-(\varepsilon).$$

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

**Lemma 30** (Bounded subsets of times). For  $\varepsilon > 0$ , for  $(a, a') \in \mathcal{A}^2$ ,

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^+(\varepsilon)|], \mathbb{E}_\nu [|\mathcal{E}_{a,a'}^-(\varepsilon)|] \leq \frac{2\sigma_\varepsilon^2 e^{\varepsilon^2/2\sigma_\varepsilon^2}}{\varepsilon^2}$$

$$\mathbb{E}_\nu [|\mathcal{K}_{a,a'}^-(\varepsilon) \setminus \mathcal{E}_{a,a'}^-(\varepsilon)|] \leq 1 + c_\varepsilon^{-1} + 2C_\varepsilon \sqrt{\log(c_\varepsilon T)},$$

where  $\sigma_\varepsilon^2 = \max_{a \in \mathcal{A}} \left\{ \mathbb{V}_{X \sim p(\mu')} (X) : \mu' \in [\mu_a - \varepsilon, \mu_a] \right\}$ ,  $c_\varepsilon, C_\varepsilon > 0$  are the constants involved in Theorem 6.

*Proof.* We start by proving  $\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^-(\varepsilon)|] \leq \frac{2\sigma_\varepsilon^2 e^{\varepsilon^2/2\sigma_\varepsilon^2}}{\varepsilon^2}$ . The proof that  $\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^+(\varepsilon)|] \leq \frac{2\sigma_\varepsilon^2 e^{\varepsilon^2/2\sigma_\varepsilon^2}}{\varepsilon^2}$  is similar.

We write

$$|\mathcal{E}_{a,a'}^-(\varepsilon)| = \sum_{t=1}^{T-1} \mathbb{I}_{\{a_{t+1}=a', N_{a'}(t) \leq N_a(t), \mu_a - \hat{\mu}_a(t) \geq \varepsilon\}}. \quad (\text{B.5})$$

Considering the stopped stopping times  $\tau_n = \inf \{t \geq 1, N_{a'}(t) = n\}$  we will rewrite the sum of indicators and use Lemma 32.

$$\begin{aligned} |\mathcal{E}_{a,a'}^-(\varepsilon)| &\leq \sum_{t \geq 1} \mathbb{I}_{\{a_{t+1}=a', N_{a'}(t) \leq N_a(t), \mu_a - \hat{\mu}_a(t) \geq \varepsilon\}} \\ &\leq \sum_{n \geq 1} \mathbb{I}_{\{n-1 \leq N_a(\tau_n-1), \mu_a - \hat{\mu}_a(\tau_n-1) \geq \varepsilon\}} \\ &\leq 1 + \sum_{n \geq 2} \mathbb{I}_{\{n-1 \leq N_a(\tau_n-1), \mu_a - \hat{\mu}_a(\tau_n-1) \geq \varepsilon\}}. \end{aligned} \quad (\text{B.6})$$

Taking the expectation of Equation (B.6), by optional skipping it comes

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^-(\varepsilon)|] \leq 1 + \sum_{n \geq 1} \mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ N_a(t) \geq n}} \hat{\mu}_a(t) \leq \mu_a - \varepsilon \right). \quad (\text{B.7})$$

From Lemma 32, previous Equation (B.7) implies

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^-(\varepsilon)|] \leq 1 + \sum_{n \geq 1} \exp(-n \text{KL}(\mu_a - \varepsilon | \mu_a)). \quad (\text{B.8})$$

From Lemma 31, previous Equation (B.8) implies

$$\mathbb{E}_\nu [|\mathcal{E}_{a,a'}^-(\varepsilon)|] \leq \sum_{n \geq 0} \exp(-n \varepsilon^2 / 2\sigma_\varepsilon^2) = \frac{1}{1 - e^{-\varepsilon^2/2\sigma_\varepsilon^2}}, \quad (\text{B.9})$$

where  $\sigma_\varepsilon^2 = \max_{a \in \mathcal{A}} \left\{ \mathbb{V}_{X \sim p(\mu')} (X) : \mu' \in [\mu_a - \varepsilon, \mu_a] \right\}$ . Finally we note that

$$\frac{1}{1 - e^{-\varepsilon^2/2\sigma_\varepsilon^2}} = \frac{e^{\varepsilon^2/2\sigma_\varepsilon^2}}{e^{\varepsilon^2/2\sigma_\varepsilon^2} - 1} \leq \frac{2\sigma_\varepsilon^2 e^{\varepsilon^2/2\sigma_\varepsilon^2}}{\varepsilon^2}.$$

We now show that  $\mathbb{E}_\nu [|\mathcal{K}_{a,a'}^-(\varepsilon) \setminus \mathcal{E}_{a,a'}^-(\varepsilon)|] \leq 1 + c_\varepsilon^{-1} + 2C_\varepsilon \sqrt{\log(c_\varepsilon T)}$ .

We write

$$\begin{aligned} & |\mathcal{K}_{a,a'}^-(\varepsilon) \setminus \mathcal{E}_{a,a'}^-(\varepsilon)| \\ &= \sum_{t=1}^{T-1} \mathbb{I}_{\{a_{t+1}=a', 1 \leq N_a(t) < N_{a'}(t), \hat{\mu}_a(t) \leq \mu_a - \varepsilon, \log(N_{a'}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon) + \log(N_a(t))\}}. \end{aligned} \quad (\text{B.10})$$

Considering the stopped stopping times  $\tau_n = \inf \{t \geq 1, N_{a'}(t) = n\}$  we shall rewrite the sum given by  $\sum_{t \in [1, T-1]} \mathbb{I}_{\{a_{t+1}=a', 1 \leq N_a(t) < N_{a'}(t), \hat{\mu}_a(t) \leq \mu_a - \varepsilon, \log(N_{a'}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon) + \log(N_a(t))\}}$  and use boundary crossing probabilities for one-dimensional exponential family distributions.

$$\begin{aligned} & |\mathcal{K}_{a,a'}^-(\varepsilon) \setminus \mathcal{E}_{a,a'}^-(\varepsilon)| \\ & \leq \sum_{t=1}^{T-1} \mathbb{I}_{\{a_{t+1}=a', 1 \leq N_a(t) < N_{a'}(t), \hat{\mu}_a(t) \leq \mu_a - \varepsilon, \log(N_{a'}(t)) \leq N_a(t) \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon) + \log(N_a(t))\}} \\ & = \sum_{t=1}^{T-1} \sum_{n=1}^{T-1} \mathbb{I}_{\{\tau_{n+1}=t+1\}} \mathbb{I}_{\{1 \leq N_a(\tau_{n+1}-1) < n, \hat{\mu}_a(\tau_{n+1}-1) \leq \mu_a - \varepsilon\}} \times \\ & \quad \mathbb{I}_{\{\log(n) \leq N_a(\tau_{n+1}-1) \text{KL}(\hat{\mu}_a(\tau_{n+1}-1)|\mu_a - \varepsilon) + \log(N_a(\tau_{n+1}-1))\}} \\ & = \sum_{n=1}^{T-1} \mathbb{I}_{\{1 \leq N_a(\tau_{n+1}-1) < n, \hat{\mu}_a(\tau_{n+1}-1) \leq \mu_a - \varepsilon\}} \times \\ & \quad \mathbb{I}_{\{\log(n) \leq N_a(\tau_{n+1}-1) \text{KL}(\hat{\mu}_a(\tau_{n+1}-1)|\mu_a - \varepsilon) + \log(N_a(\tau_{n+1}-1))\}} \sum_{t=1}^{T-1} \mathbb{I}_{\{\tau_{n+1}=t+1\}} \\ & \leq \sum_{n=1}^{T-1} \mathbb{I}_{\{1 \leq N_a(\tau_{n+1}-1) < n, \hat{\mu}_a(\tau_{n+1}-1) \leq \mu_a - \varepsilon, \log(n) \leq N_a(\tau_{n+1}-1) \text{KL}(\hat{\mu}_a(\tau_{n+1}-1)|\mu_a - \varepsilon) + \log(N_a(\tau_{n+1}-1))\}} \\ & = \sum_{n=2}^{T-1} \mathbb{I}_{\{1 \leq N_a(\tau_{n+1}-1) < n, \hat{\mu}_a(\tau_{n+1}-1) \leq \mu_a - \varepsilon, \log(n) \leq N_a(\tau_{n+1}-1) \text{KL}(\hat{\mu}_a(\tau_{n+1}-1)|\mu_a - \varepsilon) + \log(N_a(\tau_{n+1}-1))\}}. \end{aligned} \quad (\text{B.11})$$

From Equation (B.11), we get

$$\begin{aligned} & |\mathcal{K}_{a,a'}^-(\varepsilon) \setminus \mathcal{E}_{a,a'}^-(\varepsilon)| \\ & \leq \sum_{n=2}^{T-1} \mathbb{I}_{\{1 \leq N_a(\tau_{n+1}-1) < n, N_a(\tau_{n+1}-1) \text{KL}(\hat{\mu}_a(\tau_{n+1}-1)|\mu_a - \varepsilon) \geq \log(n/N_a(\tau_{n+1}-1)), \hat{\mu}_a(\tau_{n+1}-1) < \mu_a - \varepsilon\}}. \end{aligned} \quad (\text{B.12})$$

Taking the expectation of Equation (B.12), it comes

$$\begin{aligned} & \mathbb{E}_\nu [|\mathcal{K}_{a,a'}^-(\varepsilon) \setminus \mathcal{E}_{a,a'}^-(\varepsilon)|] \\ & \leq \sum_{n=2}^{T-1} \mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ \hat{\mu}_a(t) < \mu_a - \varepsilon \\ 1 \leq N_a(t) \leq n}} N_a(t) \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon) \geq \log(n/N_a(t)) \right). \end{aligned} \quad (\text{B.13})$$

From Theorem 6, previous Equation (B.13) implies

$$\begin{aligned}
& \mathbb{E}_\nu \left[ \left| \mathcal{K}_{a,a'}^-(\varepsilon) \setminus \mathcal{E}_{a,a'}^-(\varepsilon) \right| \right] & (B.14) \\
& \leq 1 + c_\varepsilon^{-1} + C_\varepsilon \sum_{n \geq 1+c_\varepsilon^{-1}}^{T-1} \frac{c_\varepsilon}{c_\varepsilon n \sqrt{\log(c_\varepsilon n)}} \\
& \leq 1 + c_\varepsilon^{-1} + C_\varepsilon \int_{c_\varepsilon^{-1}}^T \frac{c_\varepsilon dx}{c_\varepsilon x \sqrt{\log(c_\varepsilon x)}} \\
& = 1 + c_\varepsilon^{-1} + 2C_\varepsilon \sqrt{\log(c_\varepsilon T)}. & (B.15)
\end{aligned}$$

□

## B.2 Concentration of measure

In this section, Pinsker's inequality for one-dimensional exponential family distributions is reminded. Please refer to Lemma 3 from Cappé et al. (2013) for more insights. We also state two concentration results from Maillard (2018).

**Lemma 31** (Pinsker's inequality). *For  $\mu < \mu'$ , it holds that*

$$\text{KL}(\mu|\mu') \geq \frac{(\mu' - \mu)^2}{2\sigma^2},$$

where  $\sigma^2 = \max \left\{ \mathbb{V}_{X \sim p(\mu'')} (X) : \mu'' \in [\mu, \mu'] \right\}$ .

**Lemma 32** (Time-uniform concentration). *For all arm  $a \in \mathcal{A}$ , for  $x < \mu_a$ ,  $m \geq 1$ , we have*

$$\mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ N_a(t) \geq m}} \hat{\mu}_a(t) < x \right) \leq \exp(-m \text{KL}(x|\mu_a)).$$

**Theorem 6** (Boundary crossing probabilities). *For all arm  $a \in \mathcal{A}$ , for all  $\varepsilon > 0$ , for all  $n \geq 1$ , we have*

$$\mathbb{P}_\nu \left( \bigcup_{\substack{t \geq 1 \\ \hat{\mu}_a(t) < \mu_a - \varepsilon \\ 1 \leq N_a(t) \leq n}} N_a(t) \text{KL}(\hat{\mu}_a(t)|\mu_a - \varepsilon) \geq \log(n/N_a(t)) \right) \leq \frac{C_\varepsilon}{n \sqrt{\log(c_\varepsilon n)}},$$

where  $c_\varepsilon, C_\varepsilon > 0$  are explained in Maillard (2018).



# Appendix C

## Routine Bandits : Proof and Additional Experiments

### C.1 Proof of Proposition 5

Let us denote by  $\mathcal{S}$  the routine bandit setting and by  $\mathcal{S}_0$  the setting resulting from the routine bandit setting and the *additional assumption* that now the sequence of bandits  $(b_\star^h)_{h \in \llbracket 1, H \rrbracket}$  is known to the learner. Then, since a consistent strategy for  $\mathcal{S}$  is also consistent for  $\mathcal{S}_0$  (in the sense of Definition 3), we deduce Proposition 5 from Lemma 33.

**Lemma 33** (Lower bounds on the regret for  $\mathcal{S}_0$ ). *Let us consider a consistent strategy for the setting  $\mathcal{S}_0$ . Then, for all configuration  $\nu \in \mathcal{D}$ , it must be that*

$$\liminf_{T \rightarrow \infty} \frac{R(\nu, H, T)}{\log(T)} \geq c_\nu^\star := \sum_{b \in \mathcal{B}} \sum_{a \neq a_b^\star} \frac{\Delta_{a,b}}{\text{KL}(\mu_{a,b} | \mu_b^\star)}.$$

*Proof.* Since the sequence of bandits  $(b_\star^h)_{h \in \llbracket 1, H \rrbracket}$  is known to the learner and since there is no shared information between the bandits at first glance, the setting  $\mathcal{S}_0$  amounts to consider each of the  $|\mathcal{B}|$  bandits  $(\nu_b)_{b \in \mathcal{B}}$  as a separate problem, where  $\nu_b := (\nu_{a,b})_{a \in \mathcal{A}}$  for  $b \in \mathcal{B}$ . Then, from the known lower bound on the regret for the classical multi-armed bandit problem [Lai and Robbins \(1985\)](#), we get under the assumption of consistency for all bandit  $b \in \mathcal{B}$ ,

$$\liminf_{T \rightarrow \infty} \frac{1}{\log(N_b(H, T))} \sum_{a \neq a_b^\star} N_{a,b}(T) \Delta_{a,b} \geq \sum_{a \neq a_b^\star} \frac{\Delta_{a,b}}{\text{KL}(\mu_{a,b} | \mu_b^\star)}, \text{ where } N_b(H, T) = \beta_b^H HT.$$

From previous inequalities and Eq. 6.2, we conclude that

$$\begin{aligned} \liminf_{T \rightarrow \infty} \frac{R(\nu, H, T)}{\log(T)} &\geq \sum_{b \in \mathcal{B}} \liminf_{T \rightarrow \infty} \frac{\log(\beta_b^H HT)}{\log(T)} \liminf_{T \rightarrow \infty} \frac{1}{\log(\beta_b^H HT)} \sum_{a \neq a_b^\star} N_{a,b}(T) \Delta_{a,b} \\ &\geq \sum_{b \in \mathcal{B}} \sum_{a \neq a_b^\star} \frac{\Delta_{a,b}}{\text{KL}(\mu_{a,b} | \mu_b^\star)}, \end{aligned}$$

by Fatou's Lemma and since we have  $\liminf_n u_n v_n \geq \liminf_n u_n \liminf_n v_n$  for all positive real-valued sequences  $u, v$ .  $\square$

## C.2 Proof of Theorem 5

From Proposition 6, we have the following inequality

$$N_{a,b}(H, T) \leq \frac{f(\beta_b^H HT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} + \sum_{h=1}^H \mathbb{I}_{\{b_h^* = b\}} \left[ T_{\nu, \varepsilon}^h + 4|\mathcal{C}_\varepsilon^h| + |\bar{\mathcal{C}}_\varepsilon^h| + \frac{f(hT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} \sum_{k=1}^h \frac{8|\mathcal{C}_\varepsilon^k|}{T} + \mathbb{I}_{\{T \in \mathcal{T}^k\}} \right], \quad (\text{C.1})$$

where for all  $h \geq 1$ ,  $\mathbb{P}_\nu(T \in \mathcal{T}^h) \leq 1/T(T+1)$ ,  $\mathbb{E}_\nu[|\mathcal{C}_\varepsilon^h|]$ ,  $\mathbb{E}_\nu[|\bar{\mathcal{C}}_\varepsilon^h|] \leq 4|\mathcal{A}|\varepsilon^{-2} + 3$  according to Lemma 38 and  $T_{\nu, \varepsilon}^h \leq \tau_\nu^h$  according to Lemma 34 stated below.

By taking the expectation in Eq. C.1 then it comes

$$\begin{aligned} & \mathbb{E}_\nu[N_{a,b}(H, T)] \\ & \leq \frac{f(\beta_b^H HT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} \\ & \quad + \sum_{h=1}^H \mathbb{I}_{\{b_h^* = b\}} \left[ \tau_\nu^h + 5 \times (4|\mathcal{A}|\varepsilon^{-2} + 3) + \frac{f(hT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} \sum_{k=1}^h \frac{32|\mathcal{A}|\varepsilon^{-2} + 24}{T} + 1/T(T+1) \right]. \end{aligned}$$

We conclude the proof of Theorem 5 by using the two following inequalities

$$\begin{aligned} 4|\mathcal{A}|\varepsilon^{-2} + 3 & \leq 4|\mathcal{A}|(\varepsilon^{-2} + 1) \\ \frac{32|\mathcal{A}|\varepsilon^{-2} + 24}{T} + \frac{1}{T(T+1)} & \leq \frac{32|\mathcal{A}|(\varepsilon^{-2} + 1)}{T}. \end{aligned}$$

**Lemma 34** (Upper bound on  $T_{\nu, \varepsilon}^h$ ). *With the same notations as Proposition 6, for all  $0 < \varepsilon < \varepsilon_\nu$ ,*

$$T_{\nu, \varepsilon}^h \leq \tau_\nu^h := 2\varphi(8|\mathcal{A}|[\varepsilon_\nu^{-2} + 65\gamma_\nu^{-2} \log(128|\mathcal{A}|(4h)^{1/3}\gamma_\nu^{-2})]),$$

where  $\varphi: x \geq 1 \mapsto x \log(x)$ .

*Proof.* We first show that

$$t^h < 130|\mathcal{A}|\gamma_\nu^{-2} \log(128(4h)^{1/3}|\mathcal{A}|\gamma_\nu^{-2}). \quad (\text{C.2})$$

Let us consider  $t \geq 3|\mathcal{A}|$ . We have

$$d\left(\frac{t}{|\mathcal{A}|}, \delta^h(t)\right) = \sqrt{2\left(1 + \frac{|\mathcal{A}|}{t}\right) \frac{\log\left(4|\mathcal{A}|^3(h-1)\sqrt{t/|\mathcal{A}|+1}(t/|\mathcal{A}|)(t/|\mathcal{A}|+1/|\mathcal{A}|)\right)}{t/|\mathcal{A}|}}.$$

Since  $1/|\mathcal{A}| < 1$  and  $t/|\mathcal{A}| \geq 3$ , we have

$$\sqrt{t/|\mathcal{A}|+1}(t/|\mathcal{A}|)(t/|\mathcal{A}|+1/|\mathcal{A}|) \leq \sqrt{t/|\mathcal{A}|+1}(t/|\mathcal{A}|)(t/|\mathcal{A}|+1) \leq (t/|\mathcal{A}|)^3.$$

Then, since  $1+|\mathcal{A}|/t < 1+1/3$  and  $h-1 \leq h$ , we get

$$d\left(\frac{t}{|\mathcal{A}|}, \delta^h(t)\right) \leq \sqrt{\frac{8|\mathcal{A}|(4h)^{1/3}}{\Phi((4h)^{1/3}t)}}$$

where  $\Phi : x \geq 3 \mapsto x / \log(x) \geq \Phi(3)$ .  $\Phi(\cdot)$  is a one-to-one function and  $\forall y \geq \Phi(3)$ ,  $\Phi^{-1}(y) \leq y \log(y) + 2 \log(y)$ . Thus we have

$$\begin{aligned} \gamma_\nu \leq 4d \left( \frac{t}{|\mathcal{A}|}, \delta^h(t) \right) &\Rightarrow t \leq (4h)^{-1/3} \Phi^{-1}(128(4h)^{1/3} |\mathcal{A}| \gamma_\nu^{-2}) \\ &\leq 128 |\mathcal{A}| \gamma_\nu^{-2} \log(128(4h)^{1/3} |\mathcal{A}| \gamma_\nu^{-2}) \\ &\quad + 2(4h)^{-1/3} \log(128(4h)^{1/3} |\mathcal{A}| \gamma_\nu^{-2}) . \end{aligned}$$

In particular, we get the following implication

$$\gamma_\nu \leq 4d \left( \frac{t}{|\mathcal{A}|}, \delta^h(t) \right) \Rightarrow t < 130 |\mathcal{A}| \gamma_\nu^{-2} \log(128(4h)^{1/3} |\mathcal{A}| \gamma_\nu^{-2}) - 1$$

and  $t_\nu^h < 130 |\mathcal{A}| \gamma_\nu^{-2} \log(128(4h)^{1/3} |\mathcal{A}| \gamma_\nu^{-2})$ .

Furthermore, we have

$$\sum_{a \neq a_\star^h} \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)} < 2 |\mathcal{A}| \varepsilon_\nu^{-2} f(t) . \quad (\text{C.3})$$

By combining Eq. C.2 and Eq. C.3, from the definition of  $T_{\nu,\varepsilon}^h$  (see Proposition 6) we get

$$T_{\nu,\varepsilon}^h \leq \max \{ t \geq 10 : t - 8 |\mathcal{A}| (\varepsilon_\nu^{-2} f(t) - 65 \gamma_\nu^{-2} \log(128(4h)^{1/3} |\mathcal{A}| \gamma_\nu^{-2})) \leq 0 \} . \quad (\text{C.4})$$

We finally prove Lemma 34 by applying Lemma 35 with  $c = 8 |\mathcal{A}| \varepsilon_\nu^{-2}$  and  $c' = 130 |\mathcal{A}| \gamma_\nu^{-2} \log(128(4h)^{1/3} |\mathcal{A}| \gamma_\nu^{-2})$ .  $\square$

**Lemma 35.** *For all  $c, c' > 10$ , it holds*

$$\max \{ t \geq 10 : t - cf(t) - c' \leq 0 \} \leq 2\varphi(c + c') ,$$

where  $\varphi : x \geq 1 \mapsto x \log(x)$ .

*Proof.* It can be shown that  $(t - cf(t) - c')_{t \geq 2\varphi(c+c')}$  is non-decreasing by standard derivative analysis and that  $2\varphi(c+c') - cf(2\varphi(c+c')) - c \geq 0$ .  $\square$

In the following we prove the results stated in Section 6.3.

### C.2.1 Proof of Lemma 18

In this subsection we control the number previously encountered bandits falsely identified as different to the current one (see Definition 6) in addition to false positives and prove Lemma 36, an extension of Lemma 18.

**Definition 6** (False negative). *At period  $h \geq 2$  and step  $t \geq 1$ , a previous period  $k \in \llbracket 1, h-1 \rrbracket$  is called a false negative if the test  $\mathbf{T}^{k,h}(t)$  is false while previous bandit  $b_\star^k$  corresponds to current bandit  $b_\star^h$ .*

We prove necessary conditions for having false positives or false negatives.

**Lemma 36** (Condition for false positives/negatives). *At period  $h \geq 2$  and time step  $t > |\mathcal{A}|$ , for all period  $k \in \llbracket 1, h-1 \rrbracket$ , with probability  $1 - 1/(h-1)t(t+1)$ ,*

$$\begin{aligned} k \text{ is a false positive} &\implies b_\star^k \neq b_\star^h \quad \text{and} \quad \min_{a \in \mathcal{A}} |\mu_{a,b_\star^k} - \mu_{a,b_\star^h}| \leq 4d \left( \frac{t}{|\mathcal{A}|}, \delta^h(t) \right) \\ k \text{ is a false negative} &\iff b_\star^k = b_\star^h \quad \text{and} \quad \bar{a}_t^k \neq \bar{a}_t^h . \end{aligned}$$

*Proof.* From Lemma 42, with probability  $1 - 4|\mathcal{A}|\delta^h(t) = 1 - 1/(h-1)t(t+1)$ , it holds,

$$\forall a \in \mathcal{A}, \quad |\widehat{\mu}_a^h(t) - \mu_a^h| \leq d(N_a^h(t), \delta^h(t)) \quad \text{and} \quad |\widehat{\mu}_a^k(T) - \mu_a^k| \leq d(N_a^k(T), \delta^h(t)). \quad (\text{C.5})$$

False negative: Here we assume that  $b_*^k = b_*^h$ . By the triangle inequality, this implies

$$\forall a \in \mathcal{A}, \quad |\widehat{\mu}_a^h(t) - \widehat{\mu}_a^k(T)| = |(\widehat{\mu}_a^h(t) - \mu_a^h) - (\widehat{\mu}_a^k(T) - \mu_a^k)| \leq |\widehat{\mu}_a^h(t) - \mu_a^h| + |\widehat{\mu}_a^k(T) - \mu_a^k|. \quad (\text{C.6})$$

By combining Eq.C.5 and C.6, with probability  $1 - 1/(h-1)t(t+1)$ , we have

$$\forall a \in \mathcal{A}, \quad |\widehat{\mu}_a^h(t) - \widehat{\mu}_a^k(T)| - d(N_a^h(t), \delta^h(t)) - d(N_a^k(T), \delta^h(t)) \leq 0.$$

Then, from the definitions of the random variables  $(Z_a^{k,h}(t))_{a \in \mathcal{A}}$  (Eq. 6.3) and the test  $\mathbf{T}^{k,h}(t)$  (Eq. 6.4), this implies with probability  $1 - 1/(h-1)t(t+1)$ ,

$$\max_{a \in \mathcal{A}} Z_a^{k,h}(t) \leq \infty \cdot \mathbb{I}_{\{\bar{a}_t^h \neq \bar{a}_*^k\}}, \quad \mathbf{T}^{k,h}(t) = (\bar{a}_t^h = \bar{a}_*^k).$$

Thus, with probability  $1 - 1/(h-1)t(t+1)$ , period  $k$  is a false negative if, and only if,  $\bar{a}_t^h \neq \bar{a}_*^k$ .

False positive: Here we assume that period  $k$  is a false positive. In particular, we have  $b_*^k \neq b_*^h$ . By the triangle inequality, this implies

$$\forall a \in \mathcal{A}, \quad |\widehat{\mu}_a^h(t) - \widehat{\mu}_a^k(T)| \geq |\mu_a^h - \mu_a^k| - |\widehat{\mu}_a^h(t) - \mu_a^h| - |\widehat{\mu}_a^k(T) - \mu_a^k|. \quad (\text{C.7})$$

By combining Eq.C.5 and C.7, with probability  $1 - 1/(h-1)t(t+1)$ , we have

$$\forall a \in \mathcal{A}, \quad Z_a^{k,h}(t) \geq \infty \cdot \mathbb{I}_{\{\bar{a}_t^h \neq \bar{a}_*^k\}} + \min_{a \in \mathcal{A}} |\mu_a^h - \mu_a^k| - 2d(N_a^h(t), \delta^h(t)) - 2d(N_a^k(T), \delta^h(t)). \quad (\text{C.8})$$

Since period is assumed to be a false positive, we have  $\max_{a \in \mathcal{A}} Z_a^{k,h}(t) \leq 0$  and Eq. C.8 implies that, with probability  $1 - 1/(h-1)t(t+1)$ ,

$$\bar{a}_t^h = \bar{a}_*^k, \quad \min_{a \in \mathcal{A}} |\mu_a^h - \mu_a^k| \leq 2d\left(N_{\bar{a}_t^h}^h(t), \delta^h(t)\right) + 2d\left(N_{\bar{a}_*^k}^k(T), \delta^h(t)\right). \quad (\text{C.9})$$

Since  $N_{\bar{a}_t^h}^h(t) \geq t/|\mathcal{A}|$ ,  $N_{\bar{a}_*^k}^k(T) \geq T/|\mathcal{A}|$  ( $\bar{a}_t^h$  and  $\bar{a}_*^k$  are most pulled arms) and  $\delta^h(T) \leq \delta^h(t)$ , the monotonic properties of  $d(\cdot, \cdot)$  and Eq. C.9, imply that, with probability  $1 - 1/(h-1)t(t+1)$ ,

$$\min_{a \in \mathcal{A}} |\mu_a^h - \mu_a^k| \leq 2d\left(\frac{t}{|\mathcal{A}|}, \delta^h(t)\right) + 2d\left(\frac{T}{|\mathcal{A}|}, \delta^h(T)\right).$$

We conclude the proof of Lemma 36 by using Lemma 37 stated below. □

**Lemma 37** (Monotonic properties of  $d(\cdot, \cdot)$ ). *For all period  $h \geq 2$ ,  $(d(t/|\mathcal{A}|, \delta^h(t)))_{t > |\mathcal{A}|}$  is non-increasing.*

*Proof.* For all time step  $t > |\mathcal{A}|$ , a direct calculation gives

$$d\left(\frac{t}{|\mathcal{A}|}, \delta^h(t)\right) = \sqrt{2|\mathcal{A}|\left(1 + \frac{|\mathcal{A}|}{t}\right)\left(\frac{1}{2} \frac{\log(t/|\mathcal{A}|+1)}{t} + \frac{\log(4|\mathcal{A}|(h-1))}{t} + \frac{\log(t+1)}{t} + \frac{\log(t)}{t}\right)}.$$

Then, in order to prove Lemma 37, it is sufficient to note that  $(\log(t/|\mathcal{A}|+1)/t)_{t \geq 1}$ ,  $(\log(t+1)/t)_{t \geq 2}$  and  $(\log(t)/t)_{t \geq 3}$  are non-increasing. □

## C.2.2 Proof of Lemma 19

Let us consider the subsets of times when the mean of the current pulled arm is poorly estimated

$$\begin{aligned}\mathcal{E}_{a,\varepsilon}^h &:= \{t > |\mathcal{A}| : a_{t+1}^h = a \text{ and } |\widehat{\mu}_a^h(t) - \mu_a^h| > \varepsilon\} & \mathcal{E}_\varepsilon^h &:= \bigcup_{a \neq a_\star^h} \mathcal{E}_{a,\varepsilon}^h \\ \overline{\mathcal{E}}_{a,\varepsilon}^h &:= \{t \geq t_\nu^h : t \notin \mathcal{T}^h, a_{t+1}^h = a \text{ and } |\overline{\mu}_a^h(t) - \mu_a^h| > \varepsilon\} & \overline{\mathcal{E}}_\varepsilon^h &:= \bigcup_{a \neq a_\star^h} \overline{\mathcal{E}}_{a,\varepsilon}^h\end{aligned}$$

and the subsets of times when the best arm  $a_\star^h$  is below its mean

$$\begin{aligned}\mathcal{U}_a^h &:= \{t > |\mathcal{A}| : a_{t+1}^h = a \text{ and } u_{a_\star^h}^h(t) = U_{a_\star^h}^h(t) < \mu_\star^h\} & \mathcal{U}^h &:= \bigcup_{a \neq a_\star^h} \mathcal{U}_a^h \\ \overline{\mathcal{U}}_a^h &:= \{t \geq t_\nu^h : t \notin \mathcal{T}^h, a_{t+1}^h = a \text{ and } u_{a_\star^h}^h(t) = \overline{U}_{a_\star^h}^h(t) < \mu_\star^h\} & \overline{\mathcal{U}}^h &:= \bigcup_{a \neq a_\star^h} \overline{\mathcal{U}}_a^h.\end{aligned}$$

Then we have

$$\begin{aligned}\mathcal{C}_{a,\varepsilon}^h &= \mathcal{T}_a^h \cup \mathcal{E}_{a,\varepsilon}^h \cup \mathcal{U}_a^h & \mathcal{C}_\varepsilon^h &= \mathcal{T}^h \cup \mathcal{E}_\varepsilon^h \cup \mathcal{U}^h \\ \overline{\mathcal{C}}_{a,\varepsilon}^h &= \mathcal{T}_a^h \cup \overline{\mathcal{E}}_{a,\varepsilon}^h \cup \overline{\mathcal{U}}_a^h & \overline{\mathcal{C}}_\varepsilon^h &= \mathcal{T}^h \cup \overline{\mathcal{E}}_\varepsilon^h \cup \overline{\mathcal{U}}^h\end{aligned}$$

and deduce Lemma 19 from the extended Lemma 38.

**Lemma 38** (Bounded subsets of times). *For all period  $h \geq 2$ , for all arm  $a \in \mathcal{A}$ , for all  $0 < \varepsilon < \varepsilon_\nu$ ,*

$$\forall t \in \llbracket 1, T \rrbracket, \mathbb{P}_\nu(t \in \mathcal{T}^h) \leq \frac{1}{t(t+1)}, \quad \mathbb{E}_\nu[|\mathcal{E}_{a,\varepsilon}^h|], \mathbb{E}_\nu[|\overline{\mathcal{E}}_{a,\varepsilon}^h|] \leq 4\varepsilon^{-2}, \quad \mathbb{E}_\nu[|\mathcal{U}^h|], \mathbb{E}_\nu[|\overline{\mathcal{U}}^h|] \leq 2.$$

*This implies*

$$\mathbb{E}_\nu[|\mathcal{T}^h|] \leq 1, \quad \mathbb{E}_\nu[|\mathcal{C}_{a,\varepsilon}^h|], \mathbb{E}_\nu[|\overline{\mathcal{C}}_{a,\varepsilon}^h|] \leq 4\varepsilon^{-2} + 3, \quad \mathbb{E}_\nu[|\mathcal{C}_\varepsilon^h|], \mathbb{E}_\nu[|\overline{\mathcal{C}}_\varepsilon^h|] \leq 4|\mathcal{A}|\varepsilon^{-2} + 3.$$

*Proof.* Subset  $\mathcal{T}^h$ : From Lemma 38 and the definition of  $t_\nu^h$  (see Eq. 6.8), for all  $t \geq t_\nu^h$ , with probability  $1 - 1/t(t+1)$ , there is no false positive and if a previous period  $k \in \llbracket 1, h-1 \rrbracket$  is a false negative then  $b_\star^k = b_\star^h$  and  $\overline{a}_\star^k \neq \overline{a}_\star^h$  (the most pulled arms are different). From the definition of  $\mathcal{T}^h$  (see Eq. 6.9) this implies that for all  $t \geq t_\nu^h$ , with probability  $1 - 1/t(t+1)$ ,  $t \notin \mathcal{T}^h$ . That is  $\forall t \geq t_\nu^h, \mathbb{P}_\nu(t \in \mathcal{T}^h) \leq 1/t(t+1)$ . Since on the other hand, we have

$$|\mathcal{T}^h| = \sum_{t=t_\nu^h}^T \mathbb{I}_{\{t \in \mathcal{T}^h\}},$$

by taking expectation on both sides, it comes

$$\mathbb{E}_\nu[|\mathcal{T}^h|] = \sum_{t=t_\nu^h}^T \mathbb{P}_\nu(t \in \mathcal{T}^h) \leq \sum_{t=t_\nu^h}^T \frac{1}{t(t+1)} \leq 1.$$

We note that for  $1 \leq t < t_\nu^h$ , it holds that  $t \notin \mathcal{T}^h$  and  $\mathbb{P}_\nu(t \in \mathcal{T}^h) = 0 \leq 1/t(t+1)$ .

Subset  $\mathcal{E}_{a,\varepsilon}^h$ :

Since we have

$$|\mathcal{E}_{a,\varepsilon}^h| = \sum_{t>|\mathcal{A}|}^T \mathbb{I}_{\{a_{t+1}^h = a, |\widehat{\mu}_a^h(t) - \mu_a^h| > \varepsilon\}},$$

by taking the expectation on both sides, it comes

$$\mathbb{E}_\nu [|\mathcal{E}_{a,\varepsilon}^h|] \leq \sum_{t=1}^T \mathbb{P}_\nu(a_{t+1}^h = a, |\widehat{\mu}_a^h(t) - \mu_a^h| > \varepsilon). \quad (\text{C.10})$$

Then, by combining Eq. C.10 and Lemma 41, we prove  $\mathbb{E}_\nu [|\mathcal{E}_{a,\varepsilon}^h|] \leq 4\varepsilon^{-2}$ .

Subset  $\overline{\mathcal{E}}_{a,\varepsilon}^h$ : From the definitions of  $t_\nu^h$  and  $\mathcal{T}^h$  (Eq. 6.8 and 6.9), we get the following inclusion

$$\{t \geq t_\nu^h : t \notin \mathcal{T}^h, a_{t+1}^h = a, |\overline{\mu}_a^h(t) - \mu_a^h| > \varepsilon\} \subset \{t \geq t_\nu^h : a_{t+1}^h = a, |\widehat{\mu}_a^{\mathcal{K}_*^h(t),h}(t) - \mu_a^h| > \varepsilon\}, \quad (\text{C.11})$$

where  $\mathcal{K}_*^h(t) := \{k \in \llbracket 1, h-1 \rrbracket : b_*^k = b_*^h \text{ and } \overline{a}_*^k = \overline{a}_*^h\}$  and  $N_a^{\mathcal{K},h}(t) = \sum_{k \in \mathcal{K}} N_a^k(T) + N_a^h(t)$ ,  $S_a^{\mathcal{K},h}(t) = \sum_{k \in \mathcal{K}} S_a^k(T) + S_a^h(t)$ ,  $\widehat{\mu}_a^{\mathcal{K},h}(t) = S_a^{\mathcal{K},h}(t) / N_a^{\mathcal{K},h}(t)$ ,  $\forall \mathcal{K} \subset \mathcal{K}^h := \{k \in \llbracket 1, h-1 \rrbracket : b_*^k = b_*^h\}$ .

Thus, by defining  $\mathcal{K}_t := \mathcal{K}_*^h(t)$  if  $t \geq t_\nu^h$  and  $t \notin \mathcal{T}^h$ ,  $\emptyset$  otherwise, Eq. C.11 implies

$$\forall t \geq t_\nu^h, \mathbb{P}_\nu(t \notin \mathcal{T}^h, a_{t+1}^h = a, |\overline{\mu}_a^h(t) - \mu_a^h| > \varepsilon) \leq \mathbb{P}_\nu(a_{t+1}^h = a, |\widehat{\mu}_a^{\mathcal{K}_t,h}(t) - \mu_a^h| > \varepsilon). \quad (\text{C.12})$$

Since we have

$$|\overline{\mathcal{E}}_{a,\varepsilon}^h| = \sum_{t=t_\nu^h}^T \mathbb{I}_{\{t \notin \mathcal{T}^h, a_{t+1}^h = a, |\overline{\mu}_a^h(t) - \mu_a^h| > \varepsilon\}},$$

by taking the expectation on both sides and using inequalities from Eq. C.12, it comes

$$\begin{aligned} \mathbb{E}_\nu [|\overline{\mathcal{E}}_{a,\varepsilon}^h|] &= \sum_{t=t_\nu^h}^T \mathbb{P}_\nu(t \notin \mathcal{T}^h, a_{t+1}^h = a, |\overline{\mu}_a^h(t) - \mu_a^h| > \varepsilon) \\ &\leq \sum_{t=t_\nu^h}^T \mathbb{P}_\nu(a_{t+1}^h = a, |\widehat{\mu}_a^{\mathcal{K}_t,h}(t) - \mu_a^h| > \varepsilon). \end{aligned} \quad (\text{C.13})$$

Then, by combining Eq. C.13 and Lemma 41, we prove  $\mathbb{E}_\nu [|\overline{\mathcal{E}}_{a,\varepsilon}^h|] \leq 4\varepsilon^{-2}$ .

Subset  $\mathcal{U}^h$ :

By definition of the index (Eq. 6.6), we have

$$\forall t > |\mathcal{A}|, N_{a_*^h}^h(t) \text{KL}(\widehat{\mu}_{a_*^h}^h(t) | U_{a_*^h}^h(t)) = f(t). \quad (\text{C.14})$$

Since  $\widehat{\mu}_{a_*^h}^h(t) \leq U_{a_*^h}^h(t)$  for all  $t > |\mathcal{A}|$ , from the monotony of  $\text{KL}(x|\cdot)$  on  $[x, +\infty)$ , it comes

$$\forall t > |\mathcal{A}| \text{ such that } U_{a_*^h}^h(t) \leq \mu_{a_*^h}^h, \quad \text{KL}(\widehat{\mu}_{a_*^h}^h(t) | \mu_{a_*^h}^h) \geq \text{KL}(\widehat{\mu}_{a_*^h}^h(t) | U_{a_*^h}^h(t)). \quad (\text{C.15})$$

From Eq. C.14 and C.15 we deduce that

$$\mathcal{U}^h \subset \left\{ t > |\mathcal{A}| : N_{a_*^h}^h(t) \text{KL}(\widehat{\mu}_{a_*^h}^h(t) | \mu_{a_*^h, b_*^h}^h) \geq f(t) \right\}. \quad (\text{C.16})$$

From Eq. C.16 plus the union bound, it comes

$$|\mathcal{U}^h| \leq \sum_{t>|\mathcal{A}|}^T \mathbb{I}_{\left\{ N_{a_*^h}^h(t) \text{KL}(\widehat{\mu}_{a_*^h}^h(t) | \mu_{a_*^h, b_*^h}^h) \geq f(t) \right\}}. \quad (\text{C.17})$$

By taking the expectation on both sides in previous inequality (Eq. C.17), we have

$$\mathbb{E}_\nu [|\mathcal{U}^h|] \leq \sum_{t > |\mathcal{A}|}^T \mathbb{P}_\nu \left( N_{a_*^h}^h(t) \text{KL} \left( \widehat{\mu}_{a_*^h}^h(t) \middle| \mu_{a_*^h, b_*^h} \right) \geq f(t) \right). \quad (\text{C.18})$$

Combining Eq. C.18 and Lemma 42, it comes

$$\mathbb{E}_\nu [|\mathcal{U}^h|] \leq \sum_{t > |\mathcal{A}|} t^{-1} \log(t)^{-2}.$$

This implies  $\mathbb{E}_\nu [|\mathcal{U}^h|] \leq 2$  since it can be shown that

$$\sum_{t > |\mathcal{A}|} t^{-1} \log(t)^{-2} \leq \int_{t \geq |\mathcal{A}|}^{\infty} t^{-1} \log(t)^{-2} dt = \frac{1}{\log(|\mathcal{A}|)} \leq \frac{1}{\log(2)} \leq 2.$$

Subset  $\overline{\mathcal{U}}^h$ : From the definition of subset  $\mathcal{T}^h$  (see Eq. 6.9), we have

$$\{t \geq t_\nu^h : t \notin \mathcal{T}^h\} \subset \left\{ t \geq t_\nu^h : \overline{N}_{a_*^h}^h(t) = N_{a_*^h}^{\mathcal{K}_*^h(t), h}(t), \overline{\mu}_{a_*^h}^h(t) = \widehat{\mu}_{a_*^h}^{\mathcal{K}_*^h(t), h}(t), \overline{K}_t^h = |\mathcal{K}_*^h(t)| \right\}, \quad (\text{C.19})$$

where  $\mathcal{K}_*^h(t) := \{k \in \llbracket 1, h-1 \rrbracket : b_*^k = b_*^h \text{ and } \overline{a}_*^k = \overline{a}_*^h\} \subset \mathcal{K}^h$ ,  $\mathcal{K}^h := \{k \in \llbracket 1, h-1 \rrbracket : b_*^k = b_*^h\}$ ,  $N_a^{\mathcal{K}, h}(t) = \sum_{k \in \mathcal{K}} N_a^k(T) + N_a^h(t)$ ,  $S_a^{\mathcal{K}, h}(t) = \sum_{k \in \mathcal{K}} S_a^k(T) + S_a^h(t)$  and  $\widehat{\mu}_a^{\mathcal{K}, h}(t) = S_a^{\mathcal{K}, h}(t) / N_a^{\mathcal{K}, h}(t)$  for all  $\mathcal{K} \subset \mathcal{K}^h$  and  $a \in \mathcal{A}$ .

By definition of the index (Eq. 6.7), we have

$$\forall t \geq t_\nu^h, \quad \overline{N}_{a_*^h}^h(t) \text{KL} \left( \overline{\mu}_{a_*^h}^h(t) \middle| \overline{U}_{a_*^h}^h(t) \right) = f \left( \overline{K}_t^h T + t \right). \quad (\text{C.20})$$

Since  $\overline{\mu}_{a_*^h}^h(t) \leq \overline{U}_{a_*^h}^h(t)$  for all  $t \geq t_\nu^h$ , from the monotony of  $\text{KL}(x|\cdot)$  on  $[x, +\infty)$ , it comes

$$\forall t \geq t_\nu^h \text{ such that } \overline{U}_{a_*^h}^h(t) \leq \mu_{a_*^h}^h, \quad \text{KL} \left( \overline{\mu}_{a_*^h}^h(t) \middle| \mu_{a_*^h}^h \right) \geq \text{KL} \left( \overline{\mu}_{a_*^h}^h \middle| \overline{U}_{a_*^h}^h(t) \right). \quad (\text{C.21})$$

By defining  $\mathcal{K}_t := \mathcal{K}_*^h(t)$  if  $t \geq t_\nu^h$  and  $t \notin \mathcal{T}^h$ ,  $\emptyset$  otherwise, from Eq. C.19, C.20 and C.21 we deduce that

$$\overline{\mathcal{U}}^h \subset \left\{ t \geq t_\nu^h : N_{a_*^h}^{\mathcal{K}_t, h}(t) \text{KL} \left( \widehat{\mu}_{a_*^h}^{\mathcal{K}_t, h}(t) \middle| \mu_{a_*^h, b_*^h} \right) \geq f(|\mathcal{K}_t| T + t) \right\}. \quad (\text{C.22})$$

Since we have

$$\begin{aligned} & \left\{ t \geq t_\nu^h : N_{a_*^h}^{\mathcal{K}_t, h}(t) \text{KL} \left( \widehat{\mu}_{a_*^h}^{\mathcal{K}_t, h}(t) \middle| \mu_{a_*^h, b_*^h} \right) \geq f(|\mathcal{K}_t| T + t) \right\} \\ &= \bigcup_{K=0}^{h-1} \left\{ t \geq t_\nu^h : |\mathcal{K}_t| = K, N_{a_*^h}^{\mathcal{K}_t, h}(t) \text{KL} \left( \widehat{\mu}_{a_*^h}^{\mathcal{K}_t, h}(t) \middle| \mu_{a_*^h, b_*^h} \right) \geq f(KT + t) \right\} \end{aligned}$$

by using the inclusion from Eq. C.22 plus the union bound, it comes

$$\left| \overline{\mathcal{U}}^h \right| \leq \sum_{K=0}^{h-1} \sum_{t=t_\nu^h}^T \mathbb{I} \left\{ |\mathcal{K}_t| = K, N_{a_*^h}^{\mathcal{K}_t, h}(t) \text{KL} \left( \widehat{\mu}_{a_*^h}^{\mathcal{K}_t, h}(t) \middle| \mu_{a_*^h, b_*^h} \right) \geq f(KT + t) \right\}. \quad (\text{C.23})$$

By taking the expectation on both sides in previous inequality (Eq. C.23), we have

$$\mathbb{E}_\nu \left[ \left| \overline{\mathcal{U}}^h \right| \right] \leq \sum_{K=0}^{h-1} \sum_{t=t_\nu^h}^T \mathbb{P}_\nu \left( |\mathcal{K}_t| = K, N_{a_*^h}^{\mathcal{K}_t, h}(t) \text{KL} \left( \widehat{\mu}_{a_*^h}^{\mathcal{K}_t, h}(t) \middle| \mu_{a_*^h, b_*^h} \right) \geq f(KT + t) \right). \quad (\text{C.24})$$

Combining Eq. C.24 and Lemma 42, it comes

$$\begin{aligned}\mathbb{E}_\nu \left[ \left| \bar{\mathcal{U}}^h \right| \right] &\leq \sum_{K=0}^{h-1} \sum_{t=t_\nu^h}^T (KT+t)^{-1} \log(KT+t)^{-2} \\ &\leq \sum_{t \geq t_\nu^h} t^{-1} \log(t)^{-2}.\end{aligned}$$

This implies  $\mathbb{E}_\nu \left[ \left| \bar{\mathcal{U}}^h \right| \right] \leq 2$  since it can be shown that

$$\sum_{t \geq t_\nu^h} t^{-1} \log(t)^{-2} \leq \int_{t=t_\nu^h-1}^{\infty} t^{-1} \log(t)^{-2} dt = \frac{1}{\log(t_\nu^h-1)} \leq \frac{1}{\log(2)} \leq 2.$$

Subsets  $\mathcal{C}_{a,\varepsilon}^h, \mathcal{C}_\varepsilon^h, \bar{\mathcal{C}}_{a,\varepsilon}^h$  and  $\bar{\mathcal{C}}_\varepsilon^h$ : We conclude the proof of Lemma 38 by taking the expectation on both sides in the following inequalities and by using the bounds on subsets  $\mathcal{T}^h, \mathcal{U}^h, \bar{\mathcal{U}}^h, \mathcal{E}_{a,\varepsilon}^h$  and  $\bar{\mathcal{E}}_{a,\varepsilon}^h$ .

$$\begin{aligned}|\mathcal{C}_{a,\varepsilon}^h| &\leq |\mathcal{U}^h| + |\mathcal{E}_{a,\varepsilon}^h| \\ |\mathcal{C}_\varepsilon^h| &\leq |\mathcal{U}^h| + \sum_{a \neq a_\star^h} |\mathcal{E}_{a,\varepsilon}^h| \\ |\bar{\mathcal{C}}_{a,\varepsilon}^h| &\leq |\mathcal{T}^h| + |\bar{\mathcal{U}}^h| + |\bar{\mathcal{E}}_{a,\varepsilon}^h| \\ |\bar{\mathcal{C}}_\varepsilon^h| &\leq |\mathcal{T}^h| + |\bar{\mathcal{U}}^h| + \sum_{a \neq a_\star^h} |\bar{\mathcal{E}}_{a,\varepsilon}^h|.\end{aligned}$$

□

### C.2.3 Proof of Lemma 20

From Lemma 38, we have the bound  $\mathbb{E}_\nu \left[ |\mathcal{C}_{a,\varepsilon}^h| \right] \leq 4\varepsilon^{-2} + 2$  and  $\mathbb{E}_\nu \left[ |\bar{\mathcal{C}}_{a,\varepsilon}^h| \right] \leq 4\varepsilon^{-2} + 3$ .

Let us consider  $t > |\mathcal{A}|$  such that  $t \notin \mathcal{C}_{a,\varepsilon}^h = \mathcal{E}_{a,\varepsilon}^h \cup \mathcal{U}_a^h$  and  $a_{t+1}^h = a$ . By definition of the index (Eq. 6.6), we have

$$N_a^h(t) \text{KL}(\hat{\mu}_a^h(t) | U_a^h(t)) = f(t). \quad (\text{C.25})$$

Since  $a_{t+1}^h = a$ , it follows from the KLUCB-RB strategy that

$$u_{a_\star^h}^h(t) \leq u_a^h(t) \leq U_a^h(t). \quad (\text{C.26})$$

Since  $a_{t+1}^h = a$ , we have  $t \notin \mathcal{U}^h$  and

$$\mu_\star^h \leq u_{a_\star^h}^h(t) = U_{a_\star^h}^h(t). \quad (\text{C.27})$$

Since  $\varepsilon < \varepsilon_\nu$  and since  $a$  is a sub-optimal arm, we have

$$\mu_a^h + \varepsilon < \mu_\star^h. \quad (\text{C.28})$$

Since  $a_{t+1}^h = a$ , we have  $t \notin \mathcal{E}_{a,\varepsilon}^h$  and

$$\hat{\mu}_a^h(t) \leq \mu_a^h + \varepsilon. \quad (\text{C.29})$$

Then Eq. C.26, C.27, C.28 and C.29 imply

$$\hat{\mu}_a^h(t) \leq \mu_a^h + \varepsilon < \mu_\star^h \leq U_a^h(t). \quad (\text{C.30})$$



Combining Eq. C.25 and Eq. C.30, it holds

$$\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h) \leq \text{KL}(\hat{\mu}_a^h(t) | U_a^h(t)) \quad \text{and} \quad N_a^h(t) \text{KL}(\mu_a^h + \varepsilon | \mu_\star^h) \leq f(t).$$

Let us consider  $t \geq t_\nu^h$  such that  $t \notin \bar{\mathcal{C}}_{a,\varepsilon}^h = \mathcal{T} \cup \bar{\mathcal{E}}_{a,\varepsilon}^h \cup \bar{\mathcal{U}}_a^h$  and  $a_{t+1}^h = a$ . By definition of the index (Eq. 6.7), we have

$$\bar{N}_a^h(t) \text{KL}(\bar{\mu}_a^h(t) | \bar{U}_a^h(t)) = f(\bar{K}_t^h T + t). \quad (\text{C.31})$$

Since  $a_{t+1}^h = a$ , it follows from the KLUCB-RB strategy that

$$u_{a_\star^h}^h(t) \leq u_a^h(t) \leq \bar{U}_a^h(t). \quad (\text{C.32})$$

Since  $a_{t+1}^h = a$ , we have  $t \notin \mathcal{T}^h \cup \bar{\mathcal{U}}^h$  and

$$\mu_\star^h \leq u_{a_\star^h}^h(t) = \bar{U}_{a_\star^h}^h(t). \quad (\text{C.33})$$

Since  $\varepsilon < \varepsilon_\nu$  and since  $a$  is a sub-optimal arm, we have

$$\mu_a^h + \varepsilon < \mu_\star^h. \quad (\text{C.34})$$

Since  $a_{t+1}^h = a$ , we have  $t \notin \mathcal{T}^h \cup \bar{\mathcal{E}}_{a,\varepsilon}^h$  and

$$\bar{\mu}_a^h(t) \leq \mu_a^h + \varepsilon. \quad (\text{C.35})$$

Then Eq. C.32, C.33, C.34 and C.35 imply

$$\bar{\mu}_a^h(t) \leq \mu_a^h + \varepsilon < \mu_\star^h \leq \bar{U}_a^h(t). \quad (\text{C.36})$$

Combining Eq. C.31 and Eq. C.36, it holds

$$\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h) \leq \text{KL}(\bar{\mu}_a^h(t) | \bar{U}_a^h(t)) \quad \text{and} \quad \bar{N}_a^h(t) \text{KL}(\mu_a^h + \varepsilon | \mu_\star^h) \leq f(\bar{K}_t^h T + t).$$

In order to conclude the proof it remains to show that  $\bar{K}_t^h \leq \beta_{b_\star^h}^{h-1}(h-1)$ . Since  $a_{t+1}^h = a$ , we have  $t \notin \mathcal{T}^h$  and we deduce from the definition of  $\mathcal{T}^h$  (see Eq. 6.9) that

$$\bar{K}_t^h = |\mathcal{K}_\star^h(t)| \leq |\{k \in \llbracket 1, h-1 \rrbracket : b_\star^k = b_\star^h\}| = \beta_{b_\star^h}^{h-1}(h-1),$$

where  $\mathcal{K}_\star^h(t) := \{k \in \llbracket 1, h-1 \rrbracket : b_\star^k = b_\star^h \text{ and } \bar{a}_\star^k = \bar{a}_t^h\}$ .

Finally, we prove the last statement of Lemma 20. For all sub-optimal arm  $a \in \mathcal{A}$ , for all period  $h \geq 1$ , for all time step  $t > |\mathcal{A}|$ , we denote by

$$\tau_a^h(t) = \max \{t' \in \llbracket |\mathcal{A}| + 1; t \rrbracket : a_{t'+1}^h = a \quad \text{and} \quad t' \notin \mathcal{C}_{a,\varepsilon}^h\} \quad (\text{C.37})$$

the last time step before time step  $t$  that does not belong to  $\mathcal{C}_{a,\varepsilon}^h$  such that we pull arm  $a$  in period  $h$ . In particular, we have

$$N_a^h(\tau_a^h(t)) \leq \frac{f(\tau_a^h(t))}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)} \leq \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)}. \quad (\text{C.38})$$

Then, from Eq. C.37 and Eq. C.38 we have

$$\begin{aligned}
N_a^h(t) &= N_a^h(|\mathcal{A}|+1) + \sum_{t' > |\mathcal{A}|}^{t-1} \mathbb{I}\{a_{t'+1}^h = a\} \\
&= N_a^h(|\mathcal{A}|+1) + \sum_{t' > |\mathcal{A}|}^{t-1} \mathbb{I}\{a_{t'+1}^h = a, t' \in \mathcal{C}_{a,\varepsilon}^h\} + \sum_{t' > |\mathcal{A}|}^{t-1} \mathbb{I}\{a_{t'+1}^h = a, t' \notin \mathcal{C}_{a,\varepsilon}^h\} \\
&\leq N_a^h(|\mathcal{A}|+1) + |\mathcal{C}_{a,\varepsilon}^h| + N_a^h(\tau_a^h(t)) \\
&\leq N_a^h(|\mathcal{A}|+1) + |\mathcal{C}_{a,\varepsilon}^h| + \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)}.
\end{aligned}$$

## C.2.4 Proof of Lemma 21

Let us consider  $t \geq t_\nu^h$  such that  $t \notin \mathcal{T}^h$  and  $\bar{a}_t^h \neq a_\star^h$ . Since  $t \notin \mathcal{T}^h$  (see Eq. 6.9),

$$\forall a \in \mathcal{A}, \quad \bar{N}_a^h(t) = N_a^h(t) + \sum_{k \in \mathcal{K}_\star^h(t)} N_a^k(T), \quad (\text{C.39})$$

where  $\mathcal{K}_\star^h(t) := \{k \in \llbracket 1, h-1 \rrbracket : \bar{a}_\star^k = \bar{a}_t^h\}$ . Since for all  $k \in \mathcal{K}_\star^h$ ,  $\bar{a}_\star^k \in \text{argmax}_{a \in \mathcal{A}} N_a^k(T)$ , from Eq. C.39 we deduce that  $\bar{a}_t^h \in \text{argmax}_{a \in \mathcal{A}} \bar{N}_a^h(t)$ . Since  $\bar{a}_t^h \neq a_\star^h$ , this implies

$$\bar{N}_{a_\star^h}^h(t) \leq \bar{N}_{\bar{a}_t^h}^h(t) \quad \text{and} \quad \bar{N}_{\bar{a}_t^h}^h(t) \leq \sum_{a \neq a_\star^h} \bar{N}_a^h(t). \quad (\text{C.40})$$

Furthermore, since  $t \notin \mathcal{T}^h$  (see Eq. 6.9), we have

$$\bar{K}_t^h := |\mathcal{K}_+^h(t)| = |\mathcal{K}_\star^h(t)|. \quad (\text{C.41})$$

Then it comes

$$\bar{N}_{a_\star^h}^h(t) = |\mathcal{K}_\star^h(t)| T + t - \sum_{a \neq a_\star^h} \bar{N}_a^h(t). \quad (\text{C.42})$$

Then Eq. C.39, C.40 and C.42 imply

$$\frac{|\mathcal{K}_\star^h(t)| T}{2} + \frac{t}{2} \leq \sum_{a \neq a_\star^h} N_a^h(t) + \sum_{k \in \mathcal{K}_\star^h(t)} N_a^k(T). \quad (\text{C.43})$$

For  $a \neq a_\star^h$  and for  $k \in \mathcal{K}_\star^h(t)$ , the arm  $a$  is sub-optimal for bandit  $b_\star^k = b_\star^h$ . Thus, from Lemma 20, we have

$$\begin{aligned}
\forall a \neq a_\star^h, \forall k \in \mathcal{K}_\star^h(t), \quad N_a^h(t) &\leq \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)} + |\mathcal{C}_{a,\varepsilon}^h| + N_a^h(|\mathcal{A}|+1) \\
N_a^k(T) &\leq \frac{f(T)}{\text{KL}(\mu_a^k + \varepsilon | \mu_\star^k)} + |\mathcal{C}_{a,\varepsilon}^k| + N_a^k(|\mathcal{A}|+1).
\end{aligned} \quad (\text{C.44})$$

Then, by combining Eq. C.43 and Eq. C.44, we get

$$\begin{aligned}
&\frac{t}{2} - \left( \sum_{a \neq a_\star^h} \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon | \mu_\star^h)} + N_a^h(|\mathcal{A}|+1) \right) \\
&+ \frac{|\mathcal{K}_\star^h(t)| T}{2} - \left( \sum_{k \in \mathcal{K}_\star^h(t)} \sum_{a \neq a_\star^h} \frac{f(T)}{\text{KL}(\mu_a^k + \varepsilon | \mu_\star^k)} + N_a^k(|\mathcal{A}|+1) \right) \\
&\leq \sum_{k \in \mathcal{K}_\star^h(t) \cup \{h\}} \sum_{a \neq a_\star^h} |\mathcal{C}_{a,\varepsilon}^k|
\end{aligned} \quad (\text{C.45})$$

We finally prove Lemma 21 from Eq. C.45 and the following inequalities

$$\begin{aligned} \forall k \in \mathcal{K}_*^h(t) \cup \{h\}, \quad a_*^k &= a_*^h, \\ \forall k \in \mathcal{K}_*^h(t) \cup \{h\}, \quad \sum_{a \neq a_*^h} N_a^k(|\mathcal{A}|+1) &= \sum_{a \neq a_*^k} N_a^k(|\mathcal{A}|+1) \leq |\mathcal{A}|, \\ \forall k \in \mathcal{K}_*^h(t) \cup \{h\}, \quad \sum_{a \neq a_*^h} |\mathcal{C}_{a,\varepsilon}^k| &= \sum_{a \neq a_*^k} |\mathcal{C}_{a,\varepsilon}^k| = |\mathcal{C}_\varepsilon^k|. \end{aligned}$$

## C.2.5 Proof of Proposition 6

We first deduce Lemma 39 from Lemma 21.

**Lemma 39** (Conditions for misidentifying the best arms). *For all period  $h \geq 1$ , for all  $0 < \varepsilon < \varepsilon_\nu$ , for all  $t \geq T_{\nu,\varepsilon}^h$ ,*

$$(t \notin \mathcal{T}^h \text{ and } \bar{a}_t^h \neq a_*^h) \iff (t < 4|\mathcal{C}_\varepsilon^h| \text{ or } \exists k \in \llbracket 1, h \rrbracket, T < 8|\mathcal{C}_\varepsilon^k|).$$

This implies

$$(T \notin \mathcal{T}^h \text{ and } \bar{a}_*^h \neq a_*^h) \iff \exists k \in \llbracket 1, h \rrbracket, T < 8|\mathcal{C}_\varepsilon^k|.$$

We respectively refer to Proposition 6, Eq. 6.9 and Eq. 6.11 for the definitions of  $T_{\nu,\varepsilon}^h$ ,  $\mathcal{T}^h$  and  $\mathcal{C}_\varepsilon^h$ .

The proof of Lemma 39 is deferred to the Section C.2.5. We prove Proposition 6 in the following.

Let us introduce the subset  $\mathcal{P}$  of pairs period-time when there is false positives or false negatives, or when the mean of the current pulled arm is underestimated, or when the index of the best arm is below its mean, or when the most pulled arms are different from the best arms. More formally,

$$\mathcal{P} := \left\{ (h, t) \in \llbracket 1, h \rrbracket \times \llbracket 1, T \rrbracket : \left( t \geq T_{\nu,\varepsilon}^h, t \in \bar{\mathcal{C}}_\varepsilon^h \cup \mathcal{M}_\varepsilon^h \right) \text{ or } \left( \exists k \in \llbracket 1, h-1 \rrbracket, T \in \mathcal{T}^k \cup \mathcal{M}_\varepsilon^k \right) \right\}, \quad (\text{C.46})$$

where  $\mathcal{M}_\varepsilon^h := \{t \geq T_{\nu,\varepsilon}^h : t \notin \mathcal{T}^h \text{ and } \bar{a}_t^h \neq a_*^h\}$ , for all period  $h \geq 1$ .

Then, for a bandit  $b \in \mathcal{B}$  and a sub-optimal arm  $a \in \mathcal{A}$ , from Lemma 20 we have

$$N_{a,b}(H, T) \leq \frac{f(\beta_b^H HT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} + \sum_{h=1}^H \sum_{t=0}^{T-1} \mathbb{I}_{\{b_*^h=b, a_{t+1}^h=a, t < T_{\nu,\varepsilon}^h \text{ or } (h,t) \in \mathcal{P}\}}. \quad (\text{C.47})$$

From the definitions of  $\mathcal{P}$  (Eq. C.46),  $T_{\nu,\varepsilon}^h$  (Proposition 6) and  $\bar{\mathcal{C}}_\varepsilon^h$  (Eq. 6.11) for  $h \geq 1$ , this implies

$$\begin{aligned} N_{a,b}(H, T) &\leq \frac{f(\beta_b^H HT)}{\text{KL}(\mu_{a,b} + \varepsilon | \mu_b^*)} \\ &+ \sum_{h=1}^H \sum_{t=0}^{T-1} \mathbb{I}_{\{b_*^h=b, a_{t+1}^h=a, t < T_{\nu,\varepsilon}^h\}} \\ &+ \sum_{h=1}^H \sum_{t=0}^{T-1} \mathbb{I}_{\{b_*^h=b, a_{t+1}^h=a, t \in \bar{\mathcal{C}}_\varepsilon^h\}} \\ &+ \sum_{h=1}^H \sum_{t=0}^{T-1} \mathbb{I}_{\{b_*^h=b, a_{t+1}^h=a, t \notin \bar{\mathcal{C}}_{a,\varepsilon}^h\}} \mathbb{I}_{\{\exists k \in \llbracket 1, h-1 \rrbracket, T \in \mathcal{T}^k\}} \\ &+ \sum_{h=1}^H \sum_{t=0}^{T-1} \mathbb{I}_{\{b_*^h=b, a_{t+1}^h=a, t \notin \bar{\mathcal{C}}_{a,\varepsilon}^h\}} \mathbb{I}_{\{t \in \mathcal{M}_\varepsilon^h \text{ or } \exists k \in \llbracket 1, h-1 \rrbracket, T \in \mathcal{M}_\varepsilon^k\}}. \end{aligned} \quad (\text{C.48})$$

Furthermore, from Lemma 39 we have for all period  $h \geq 1$ ,

$$\mathbb{I}_{\{t \in \mathcal{M}_\varepsilon^h \text{ or } \exists k \in [1, h-1], T \in \mathcal{M}_\varepsilon^k\}} \leq \mathbb{I}_{\{t < 4|\mathcal{C}_\varepsilon^h|\}} + \sum_{k=1}^h \mathbb{I}_{\{T < 8|\mathcal{C}_\varepsilon^k|\}}. \quad (\text{C.49})$$

By combining Eq.C.48 and Eq.C.49, we get

$$N_{a,b}(H, T) \leq \frac{f(\beta_b^H HT)}{\text{KL}(\mu_{a,b} + \varepsilon|\mu_b^*)} + \sum_{h=1}^H \mathbb{I}_{\{b_\star^h = b\}} \left[ T_{\nu, \varepsilon}^h + 4|\mathcal{C}_\varepsilon^h| + \left| \bar{\mathcal{C}}_\varepsilon^h \right| + \left( \sum_{t=0}^{T-1} \mathbb{I}_{\{b_\star^h = b, a_{t+1}^h = a, t \notin \bar{\mathcal{C}}_{a, \varepsilon}^h\}} \right) \left( \sum_{k=1}^h \mathbb{I}_{\{T < 8|\mathcal{C}_\varepsilon^k|\}} + \mathbb{I}_{\{T \in \mathcal{T}^k\}} \right) \right]. \quad (\text{C.50})$$

Since the arm  $a$  is sub-optimal for the bandit  $b$ , the consistency (Lemma 20) implies

$$\forall h \geq 1, \quad \sum_{t=0}^{T-1} \mathbb{I}_{\{b_\star^h = b, a_{t+1}^h = a, t \notin \bar{\mathcal{C}}_{a, \varepsilon}^h\}} \leq \frac{f(hT)}{\text{KL}(\mu_{a,b} + \varepsilon|\mu_b^*)}. \quad (\text{C.51})$$

In addition, the following Markov's type inequalities are satisfied

$$\forall k \geq 1, \quad \mathbb{I}_{\{T < 8|\mathcal{C}_\varepsilon^k|\}} \leq \frac{8|\mathcal{C}_\varepsilon^k|}{T}. \quad (\text{C.52})$$

By combining Eq. C.50, C.51 and C.52, we prove Proposition 6, that is

$$N_{a,b}(H, T) \leq \frac{f(\beta_b^H HT)}{\text{KL}(\mu_{a,b} + \varepsilon|\mu_b^*)} + \sum_{h=1}^H \mathbb{I}_{\{b_\star^h = b\}} \left[ T_{\nu, \varepsilon}^h + 4|\mathcal{C}_\varepsilon^h| + \left| \bar{\mathcal{C}}_\varepsilon^h \right| + \frac{f(hT)}{\text{KL}(\mu_{a,b} + \varepsilon|\mu_b^*)} \sum_{k=1}^h \frac{8|\mathcal{C}_\varepsilon^k|}{T} + \mathbb{I}_{\{T \in \mathcal{T}^k\}} \right].$$

### Proof of Lemma 39

Let us consider a period  $h \geq 1$ ,  $0 < \varepsilon < \varepsilon_\nu$ , and a time step all  $t \geq T_{\nu, \varepsilon}^h$  such that  $t \notin \mathcal{T}^h$  and  $\bar{a}_t^h \neq a_\star^h$ . Then, since  $T_{\nu, \varepsilon}^h \geq t_\nu^h$ , from Lemma 21 we have

$$\frac{t + |\mathcal{K}_\star^h(t)|T}{2} - (1 + |\mathcal{K}_\star^h(t)|)|\mathcal{A}| - (f(t) + |\mathcal{K}_\star^h(t)|f(T)) \sum_{a \neq a_\star^h} \frac{1}{\text{KL}(\mu_a^h + \varepsilon|\mu_\star^h)} \leq \sum_{k \in \mathcal{K}_\star^h(t) \cup \{h\}} |\mathcal{C}_\varepsilon^k|. \quad (\text{C.53})$$

Furthermore, by definition of  $T_{\nu, \varepsilon}^h$ , since  $t \geq T_{\nu, \varepsilon}^h$ , we have

$$\begin{aligned} \frac{t}{2} - \sum_{a \neq a_\star^h} \frac{f(t)}{\text{KL}(\mu_a^h + \varepsilon|\mu_\star^h)} - |\mathcal{A}| &> \frac{t}{4} \\ \frac{T}{2} - \sum_{a \neq a_\star^h} \frac{f(T)}{\text{KL}(\mu_a^h + \varepsilon|\mu_\star^h)} - |\mathcal{A}| &> \frac{T}{4}. \end{aligned} \quad (\text{C.54})$$

By respectively combining Eq. C.53 and Eq. C.54, we thus deduce

$$\begin{aligned} |\mathcal{K}_\star^h(t)| = 0 &\Rightarrow t \leq 4|\mathcal{C}_\varepsilon^h| \\ |\mathcal{K}_\star^h(t)| \geq 1 &\Rightarrow \exists k \in [1, h], T \leq 8|\mathcal{C}_\varepsilon^k| \end{aligned}$$

which implies Lemma 39.

## C.2.6 Tools from Concentration of Measure

This subsection gathers useful concentration lemmas that do not depend on the considered strategy.

**Notations** For all period  $h \geq 2$ , for all time step  $t > |\mathcal{A}|$ , for each (possible random) subset of past periods  $\mathcal{K} \subset \mathcal{K}^h := \{k \in \llbracket 1, h-1 \rrbracket : b_*^k = b_*^h\}$ , for all arm  $a \in \mathcal{A}$ , we define  $N_a^{\mathcal{K},h}(t) := \sum_{k \in \mathcal{K}} N_a^k(T) + N_a^h(t)$ ,  $S_a^{\mathcal{K},h}(t) := \sum_{k \in \mathcal{K}} S_a^k(T) + S_a^h(t)$  and  $\widehat{\mu}_a^{\mathcal{K},h}(t) := S_a^{\mathcal{K},h}(t)/N_a^{\mathcal{K},h}(t)$ .

In particular, for  $\mathcal{K} = \mathcal{K}_*^h(t) := \{k \in \llbracket 1, h-1 \rrbracket : b_*^k = b_*^h \text{ and } \bar{a}_*^k = \bar{a}_*^h\}$ , we have  $N_a^{\mathcal{K}_*^h(t),h}(t) = \bar{N}_a^h(t)$  and  $\widehat{\mu}_a^{\mathcal{K}_*^h(t),h}(t) = \bar{\mu}_a^h(t)$  when  $t \geq t_\nu^h$  and  $t \notin \mathcal{T}^h$  (see Eq. 6.8 and 6.9).

Uniform bounds based on the Laplace method (method of mixtures for sub-Gaussian random variables, see Peña et al. (2008)) are given in Lemma 40.

**Lemma 40** (Uniform sub-Gaussian concentration). *For all period  $h \geq 2$ , for all time step  $t > |\mathcal{A}|$ , for all arm  $a \in \mathcal{A}$ , for all  $\delta \in (0, 1)$ , it holds*

$$\begin{aligned} \mathbb{P}_\nu(\widehat{\mu}_a^h(t) - \mu_{a,b_*^h} \geq d(N_a^h(t), \delta)) &\leq \delta \\ \mathbb{P}_\nu(\mu_{a,b_*^h} - \widehat{\mu}_a^h(t) \geq d(N_a^h(t), \delta)) &\leq \delta, \end{aligned}$$

where  $d(n, \delta) = \sqrt{2(1+1/n) \log(\sqrt{n+1}/\delta)}/n$ , for all  $n \geq 1$ .

Lemma 41 reformulates Lemma B.1 from Combes and Proutiere (2014b).

**Lemma 41** (Concentration inequalities). *For all period  $h \geq 2$ , for all arm  $a \in \mathcal{A}$ , for all  $\varepsilon \in (0, 1/2)$ , and all possibly random subset of periods  $\mathcal{K}_t$  such that the random variable  $N_a^{\mathcal{K}_t,h}(t)$  is a random stopping time, it holds*

$$\sum_{t \geq 1} \mathbb{P}_\nu(a_{t+1}^h = a, |\widehat{\mu}_a^{\mathcal{K}_t,h}(t) - \mu_a^h| \geq \varepsilon) \leq 4\varepsilon^{-2}.$$

Lemma 42 reformulates Theorem 1 from Garivier (2013).

**Lemma 42** (Self-normalized inequalities). *For all period  $h \geq 2$ , for all time step  $t > |\mathcal{A}|$ , for all arm  $a \in \mathcal{A}$ , for all  $K \in \llbracket 0, h-1 \rrbracket$ , for all  $\delta > 0$  and all possibly random subset of periods  $\mathcal{K}_t$  such that the random variable  $N_a^{\mathcal{K}_t,h}(t)$  is a random stopping time, it holds*

$$\mathbb{P}_\nu(|\mathcal{K}_t| = K, N_a^{\mathcal{K}_t,h}(t) \text{KL}(\widehat{\mu}_a^{\mathcal{K}_t,h}(t) | \mu_{a,b_*^h}) \geq \delta) \leq 2e^{\lceil \delta \log(KT+t) \rceil} \exp(-\delta).$$

In particular, this implies for  $\delta = f(KT+t)$ ,

$$\mathbb{P}_\nu(|\mathcal{K}_t| = K, N_a^{\mathcal{K}_t,h}(t) \text{KL}(\widehat{\mu}_a^{\mathcal{K}_t,h}(t) | \mu_{a,b_*^h}) \geq f(KT+t)) \leq (KT+t)^{-1} \log(KT+t)^{-2}.$$

## C.3 Additional Experiments: Ideal Cases for which Bandits are Close Enough on the Subset of Optimal Arms

This section provides additional experiments where we investigate some favorable distributions  $\nu$  where it is hard to separate the different bandits from each other. All experiments are repeated 100 times.

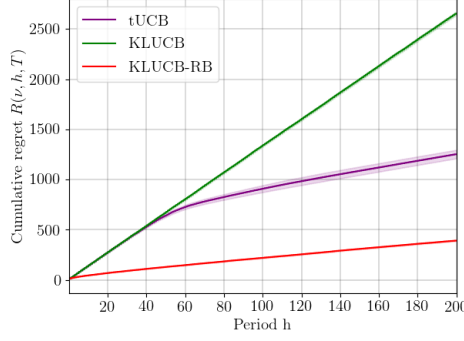


Figure C.1: Cumulative regret of KLUCB, KLUCB-RB and tUCB along  $H = 200$  periods of  $T = 10^3$  rounds for the bandit set  $\mathcal{B} = \{b\}$ .

### C.3.1 A Single Instance

Let us first make a remark in the trivial limit case of a unique bandit, that is to say  $\mathcal{B} = \{b\}$ . In such cases, playing KLUCB-RB is obviously equivalent to playing for  $T_{total} := HT$  rounds a KLUCB strategy on the bandit instance  $b$ , with an additional term  $(H - 1) \sum_{a \in \mathcal{A}} \Delta_{a,b}$  in the final cumulative regret due to the initialization at each period. Figure C.1 highlights this fact for the two-armed bandit  $b$  defined in Eq. C.55, over  $H = 200$  periods of  $T = 10^3$  rounds.

$$b : (\mu_{1,b}, \mu_{2,b}) = \left(-\frac{\Delta}{2}, \frac{\Delta}{2}\right) \quad \text{where} \quad \Delta = 10\sqrt{\frac{\log(HT)}{T}}. \quad (\text{C.55})$$

Although the case  $|\mathcal{B}| = 1$  is not an interesting one since there is no switches between different bandits instances, it enables to understand what happens when  $|\mathcal{B}| > 1$  and bandits are similar, that is  $\max_{a \in \mathcal{A}^*} \max_{b \neq b'} |\mu_{a,b} - \mu_{a,b'}|$  approaches 0. Besides, it highlights the need for tUCB to see a sufficient number of periods before exploiting the estimated models of the bandits.

### C.3.2 Similarity of Different Instances on the Optimal Subset $\mathcal{A}^*$

Let us consider routines over two bandits  $b_1$  and  $b_2$  composed of two arms such that  $(\mu_{1,b_2}, \mu_{2,b_2}) = (\mu_{1,b_1} + \gamma, \mu_{2,b_1} - \gamma)$  and  $a_{b_1}^* = 2$ . If  $\gamma > \Delta_{2,b_1}/2$ , arms arrangements are different in both instances and these cases are studied in subsection 6.4.1. Otherwise we have  $a_{b_1}^* = a_{b_2}^* = 2$  if ever  $0 < \gamma < \Delta_{2,b_1}/2$ . Although separation of instances is particularly hard in such cases, samples aggregation from false positive periods does not perturb the empirical means arrangement, and thus yields great performances for KLUCB-RB. To explain how this kind of distribution generalizes to settings composed of arbitrary numbers of bandits and arms, we present in Fig. C.2b a distribution  $\nu$  such that  $|\mathcal{B}| = 5$  and  $|\mathcal{A}| = 4$ . In this setting, we have  $\mathcal{A}^* = \{1, 4\}$ . Considering distributions of bandits restricted to  $\mathcal{A}^*$ ,  $\mathcal{B}$  naturally decomposes into 3 clusters  $\mathcal{C}^{(1)} := \{b_1, b_4\}$ ,  $\mathcal{C}^{(2)} := \{b_2, b_3\}$  and  $\mathcal{C}^{(3)} := \{b_5\}$  so that

$$\forall i \in \{1, 2, 3\}, \forall b, b' \in \mathcal{C}^{(i)}, \forall a \in \mathcal{A}^*, |\mu_{a,b} - \mu_{a,b'}| < \frac{1}{2} \min_{y \in \mathcal{C}^{(i)}} \min_{x \in \mathcal{A}, x \neq a_y^*} \Delta_{x,y} \quad (\text{C.56})$$

which entails in particular

$$\forall i \in \{1, 2, 3\}, \exists a^{(i)} \in \mathcal{A}^*, \forall b \in \mathcal{C}^{(i)}, a_b^* = a^{(i)},$$

and

$$\forall i \in \{1, 2, 3\}, \forall b \in \mathcal{C}^{(i)}, \forall b' \notin \mathcal{C}^{(i)}, |\mu_{a^{(i)},b} - \mu_{a^{(i)},b'}| > \min_{x \in \mathcal{C}^{(i)}} \min_{a \neq a_x^*} \Delta_{a,x}. \quad (\text{C.57})$$

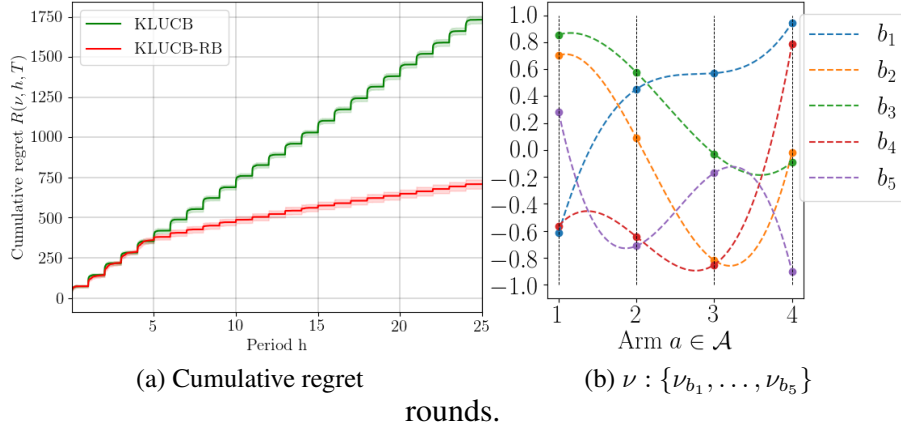


Figure C.2: KLUCB-RB and KLUCB performances on a clustered distribution according to  $\mathcal{A}^*$ , along  $H = 25$  periods of  $2 \times 10^4$  rounds.

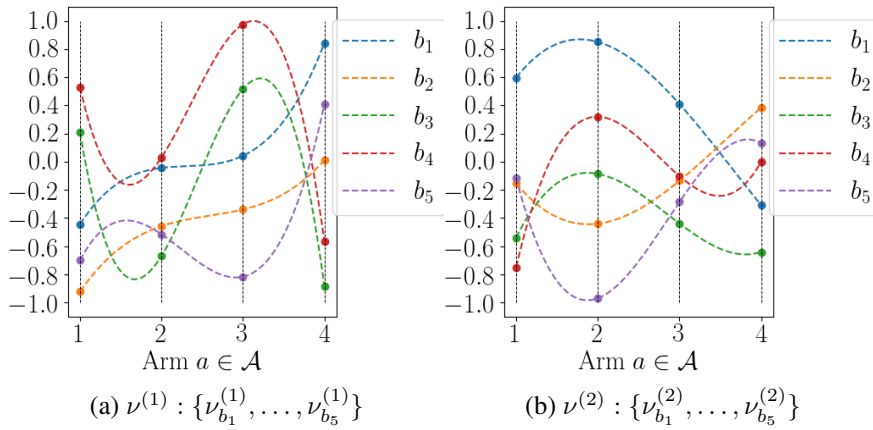


Figure C.3: Distribution  $\nu$  for each bandit in sets  $\mathcal{B}_1$  and  $\mathcal{B}_2$ .

On the one hand, Eq. C.56 sums up that different bandits from a same cluster are hard to distinguish, in comparison with the difficulty of learning each instance independently. On the other hand, Eq. C.57 implies that clusters are easy to separate from each other. Besides Eq. C.56 also implies that the permutation of  $\mathcal{A}$  sorting arms according to an increasing order is the same for all instances from a same cluster. Thus KLUCB-RB is expected to perform well for this kind of arms distributions. Fig. C.2a shows the cumulative regret curves with one standard deviation obtained on this setting, along  $H = 25$  periods of  $T = 2 \times 10^4$  rounds. As expected, it highlights that a positive cluster effect causes an improvement in regret minimization. In practice, KLUCB-RB naturally clusterizes the previously seen periods while the current period index  $h$  increases. More specifically, noting  $\mathcal{C}(h)$  the cluster containing  $b_\star^h$ , KLUCB-RB makes all the different bandits from  $\mathcal{C}(h)$  share their samples with  $b_\star^h$  for a large amount of rounds, which enables to boost the minimization of regret across period  $h$ .

### C.3.3 Complement of Sections 6.4.2 and 6.4.3

Figure C.3 shows the generated settings used in experiments of Section 6.4.2, and Figure C.4 the setting used in Section 6.4.3. More specifically, each sub-figure displays the expected reward for each of the four arms, in each of the bandits, for the three considered bandit sets.

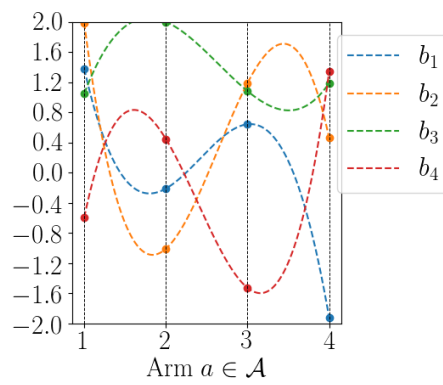


Figure C.4: Distribution  $\nu$  for each bandit in the critical setting.



# Appendix D

## Bandits with Groups of Similar Arms

This work is the subject of [Pesquerel et al. \(2021\)](#) and has been initiated and mainly driven by Fabien Pesquerel. The proofs are not reproduced in this chapter. We refer to [Pesquerel et al. \(2021\)](#) and its Appendix for more details.

### D.1 Introduction

Motivated by various practical reasons, one may want to restrict to a subset  $\mathcal{D}_{\text{sim}} \subset \mathcal{D}$  of allowed bandit configurations instead of the full set  $\mathcal{D}$ . In this chapter study a variant of the multi-armed bandit problem in which the reward function,  $\mu : a \in \mathcal{A} \rightarrow \mu_a$ , is assumed to satisfy a cluster-like structural property. A bandit configuration  $\nu$  is said to satisfy the **q-equivalence property** if for every arm  $a \in \mathcal{A}$ , there are at least  $q - 1$  distinct arms having the same expected value:

$$\forall a \in \mathcal{A}, \quad |\{a' \in \mathcal{A} : \mu_{a'} = \mu_a\}| \geq q.$$

Assuming the set of arms  $\mathcal{A}$  and base distributions  $\mathcal{D}$  is known to the learner, we denote by  $\mathcal{D}_{\text{sim}}(q)$  the set of bandit configurations having the q-equivalence property.

**Definition 7** (Arm equivalence and equivalence class). *Given a bandit configuration  $\nu$ , two arms  $a, a' \in \mathcal{A}$  are said to be equivalent if their associated distributions have the same expected values:*

$$a \sim a' \Leftrightarrow \mu_a = \mu_{a'}$$

*An equivalence class  $c$  in  $\nu$  is a maximal subset of arms in  $\mathcal{A}$  having the same mean, i.e., for all arm  $a, a'$  in  $c$ ,  $\mu_a = \mu_{a'}$  and for all arm  $a \in c$  and  $a' \in \mathcal{A} \setminus c$ ,  $\mu_a \neq \mu_{a'}$ .*

This situation typically appears in practical situations when each arm can be described with a list of categorical attributes, and the (unknown) mean reward function only depends on a subset of them, the others being redundant. In this case,  $q$  is naturally linked to the number of attributes considered redundant (or useless descriptors), and the number of categories of each attribute. Precisely,  $q = \prod_{i \in \mathcal{R}} c_i$  where  $\mathcal{R}$  is the set of redundant attributes and  $c_i$  the number of categories for attribute  $i$ . The learner may know that there exists such a structure while not knowing a closed form formula mapping the list of categorical attributes to the significant subset. In this case,  $q$  might be a lower bound on the sizes of the class since the set  $\mathcal{R}$  might not be the largest possible one or because the number of redundant attributes depends on the number of relevant attributes. In all cases, the smallest possible number of redundant attributes can be naturally linked to  $q$ . We hereafter consider the learner only knows  $q$  but would like to exploit the prior knowledge of this structure in a bandit problem.

**Goal.** For the structure  $\mathcal{D}_{\text{sim}}(q)$ , as we show in Theorem 7 below, the term  $\mathfrak{C}_{\mathcal{D}_{\text{sim}}(q)}(\mu)$  unfortunately makes appear in general a combinatorial optimization problem. This makes resorting to OSSB or any strategy targeting exact asymptotic optimality a daunting task for the practitioner. In this chapter, our goal is to provide a computationally efficient strategy adapted to the structure  $\mathcal{D}_{\text{sim}}(q)$ , that is able to reach optimality up to controlled error term.

**Outline and contributions** The rest of this chapter is organized as follows. In section D.2, we derive a lower bound on the regret for the structured set of bandit configurations  $\mathcal{D}_{\text{sim}}(q)$ . This bound makes appear two components, one that we call *non-combinatorial* as optimizing it can be done efficiently, and a second term that we term *combinatorial* as it involves solving a combinatorial problem. Interestingly, using in Lemma 43 and Theorem 9 that the contribution of the combinatorial part of the lower bound can be controlled. Owing to this key insight, we introduce in section D.3, IMED-EC, an adaptation of the IMED strategy from Honda and Takemura (2015) to the structured set  $\mathcal{D}_{\text{sim}}(q)$ . One advantage of IMED over a KLUCB alternative is its reduced complexity, which translates to the equivalence class setup. At each time step, the complexity of computing the next arm to be pulled by IMED-EC is no more than the one of sorting a list of  $|\mathcal{A}|$  elements once the IMED indexes have been computed is at most the one of sorting a list of  $|\mathcal{A}|$  elements, which is only  $\log |\mathcal{A}|$  times larger than looking for the minimal IMED index. In Section D.4, we prove that IMED-EC achieves a controlled asymptotic regret that matches the non-combinatorial part of the lower bound and is at most (less than) 2 times from the optimal regret bound. Last, we illustrate the benefit of the IMED-EC over its unstructured version in section D.5, where it shows a substantial improvement. Our experiments also highlights the robustness of the algorithm to a misspecified parameter  $q$ , which is a desirable feature for the practitioner.

## D.2 A regret lower bound with combinatorial and non-combinatorial parts

In this section, we derive a lower bound on the number of pulls of suboptimal arms that involves a combinatorial optimization problem. Using that lower bound, we derive a simple algorithm, IMED-EC, that does not involve any optimization problem. While not being asymptotically optimal, we will show in the next section that our algorithm have an upper bound on its regret that is no more than a fraction of the unstructured regret.

**Definition 8** (Confusing instance). *Given a bandit configuration  $\nu \in \mathcal{D}_{\text{sim}}(q)$ , a real number  $\lambda$  and a subset  $c_q \subseteq \mathcal{A}$  of  $q$  equivalent arms in  $\nu$ , we denote by  $\mathcal{D}_{\text{sim}}(q, \nu, c_q, \lambda)$  the set of all bandit configurations having the same set of arms as  $\nu$  and such that for all  $\nu' \in \mathcal{D}_{\text{sim}}(q, \nu, c_q, \lambda)$ ,  $\nu' \in \mathcal{D}_{\text{sim}}(q)$  and for every arm  $a$  in  $c_q$ ,  $\mu'_a \geq \lambda$ . When  $\lambda > \mu^*$ , and  $c_q$  is a subset of a suboptimal class, a bandit configuration in  $\mathcal{D}_{\text{sim}}(q, \nu, c_q, \lambda)$  is called a **confusing instance** of  $\nu$ . Similarly to the notation introduced above, we will use the notation  $\mathcal{D}_{\text{sim}}(q, \mu, c_q, \lambda)$  to specify the set of means of bandit configurations in  $\mathcal{D}_{\text{sim}}(q, \nu, c_q, \lambda)$ .*

The aim of an asymptotic lower bound on the number of pulls of a suboptimal arm is to mathematically understand the minimal asymptotic amount of exploration an algorithm should perform.

**Theorem 7** (Asymptotic lower bound). *Let  $q \in \mathbb{N} - \{0\}$  be a positive integer and  $\nu \in \mathcal{D}_{\text{sim}}(q)$  be a bandit configuration having the  $q$ -equivalence property. Let  $c \subset \mathcal{A}$  be a suboptimal equivalence class in  $\nu$ . Assuming uniform consistency, for all suboptimal arm  $a$ ,*

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

*we have the following asymptotic bandit dependent lower bound on the number of pulls of arms in  $c$ :*

$$\liminf_{T \rightarrow \infty} \frac{\min_{c_q \subseteq c} \sum_{a \in c_q} \mathbb{E}_\nu [N_a(T)] \text{KL}(\mu_a | \mu^*) + \inf_{\mu' \in \mathcal{D}_{\text{sim}}(q, \mu, c_q, \lambda)} \sum_{a \notin c_q} \mathbb{E}_\nu [N_a(T)] \text{KL}(\mu_a | \mu'_a)}{\log T} \geq 1 \quad (\text{D.1})$$

where  $c_q$  is any subset of  $c$  having  $q$  distinct arms within it.

While this lower bound involves a combinatorial optimization term, one can distinguish between two regimes depending on the size of the suboptimal class. The *combinatorial regime* and the *non-combinatorial regime*.

**Non-combinatorial regime** For a suboptimal class  $c$ , if  $|c| = q$  or  $|c| \geq 2q$ , then the lower bound reduces to

$$\liminf_{T \rightarrow \infty} \frac{\min_{c_q \subseteq c} \sum_{a \in c_q} \mathbb{E}_\nu[N_a(T)] \text{KL}(\mu_a | \mu^*)}{\log T} \geq 1.$$

We call this the *non-combinatorial regime* because the minimum over all  $q$ -partitions of  $c$  is in fact the sum of the  $q$  smallest elements of  $\{\mathbb{E}_\nu[N_a(T)] \text{KL}(\mu_a | \lambda)\}_{a \in c}$ . The search amongst all the  $q$ -partitions of  $c$  amount to a research of the  $q$  smallest elements which is not more complex than sorting a list of  $|c|$  elements.

**Lemma 43.** *Let  $\nu \in \mathcal{D}_{\text{sim}}(q)$  be a bandit configuration having the  $q$ -equivalence property. Let  $c$  be a suboptimal class in the non-combinatorial regime, then, under assumption 1 and 2,*

$$\liminf_{T \rightarrow \infty} \frac{\sum_{a \in c} \mathbb{E}_\nu[N_a(T)]}{\log T} \geq \frac{|c|}{q} \frac{1}{\text{KL}(\mu_a | \lambda)} \quad (\text{D.2})$$

While we do not have information about individual number of time an arm in a class has been sampled, lemma 43 roughly tells us than on average, the lower bound on the minimal amount of exploration of an arm in a suboptimal class has been divided by  $q$ .

**Lemma 44.** *If all suboptimal classes are in the non-combinatorial regime, the regret may be asymptotically lower bounded by*

$$\liminf_{T \rightarrow \infty} \frac{R(\nu, T)}{\log T} \geq \frac{1}{q} \sum_{a \in \mathcal{A} \setminus \mathcal{A}^*} \frac{\Delta_a}{\text{KL}(\mu_a | \lambda)}. \quad (\text{D.3})$$

Lemma 44 informs us that in the non-combinatorial regime, the classical lower bound on the regret given by equation (1.3) has been divided by  $q$ .

**Combinatorial regime** For a suboptimal class  $c$  to be in the combinatorial regime, we need  $q < |c| < 2q$ . In that case, the lower bound (D.1) involves a combinatorial optimization problem. The difficulty arising from the term

$$\inf_{\mu' \in \mathcal{D}_{\text{sim}}(q, \mu, c_q, \lambda)} \sum_{a \notin c_q} \mathbb{E}_\nu[N_a(T)] \text{KL}(\mu_a | \mu'_a),$$

is two fold. First, while we could have thought that summing on the reminder  $c \setminus c_q$  would be enough, the summand has to be on  $a \notin c_q$  as a whole. Indeed, the residual  $c \setminus c_q$  may be of size  $q - 1$  meaning that it might cost less to move an arm from another class to the residual in order to complete it rather than moving all the reminder. Second, while we could have thought that moving elements from one class of  $\nu$  to another might be enough, the *infimum* has to be taken on  $\mathcal{D}_{\text{sim}}(q, \mu, c_q, \lambda)$ . Indeed, the residual  $c \setminus c_q$  may be of size  $q - 1$  and the *nearest* class might be of size exactly  $q$ . In this case, it may cost less to move all the  $2q - 1$  distributions in between the two classes and create a new one rather than merging one of the two with the other.

**Lemma 45.** *Let  $\nu \in \mathcal{D}_{\text{sim}}(q)$  be a bandit configuration having the  $q$ -equivalence property and  $c$  be a suboptimal class in the combinatorial regime. Then, we have*

$$\liminf_{T \rightarrow \infty} \frac{\sum_{a \in c} \mathbb{E}_\nu[N_a(T)]}{\log T} \geq \frac{\log(T)}{2q} \sum_{a \in c} \frac{1}{\text{KL}(\mu_a | \lambda)}. \quad (\text{D.4})$$

Those equations can be compared to the equation (D.2) from the non-combinatorial regime. We emphasize the fact that the lower bounds given by equation (D.4) are not the *largest* possible lower bound and hence do not provide as much information about the algorithmically achievable regret as the largest one given by equation (D.1). However, together with a regret upper bound on the algorithm IMED-EC, those quantities will help us control the asymptotic discrepancy between IMED-EC's regret and the asymptotic lower bound given by Theorem 7.

### D.3 Information Minimization for bandits with equivalence class

The algorithm we present, IMED-EC, depends on the (*weak*) indexes introduced in the IMED paper by Honda and Takemura (2015). At each time step  $t$ , for each arm  $a \in \mathcal{A}$ , we can compute its IMED index as

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

where  $\hat{\mu}^*(t) = \max_{a \in \mathcal{A}} \hat{\mu}_a(t)$  and for each arm  $a \in \mathcal{A}$ ,  $\hat{\mu}_a(t)$  is the empirical mean of arm  $a$  computed with samples from this arm collected up to time  $t$ ,  $\hat{\mu}_a(t) = \frac{1}{N_a(t)} \sum_{s=1}^t X_s \mathbb{I}_{\{a_s=a\}}$ . Let  $\nu \in \mathcal{D}_{\text{sim}}(q)$  be a bandit configuration having the  $q$ -equivalence property. We denote by  $\mathcal{A}^*(t) = \arg \max_{a \in \mathcal{A}} \hat{\mu}_a(t)$  the set of empirical optimal arms at time  $t$ . We will denote by  $\mathcal{A}_q(t)$  the set of arms having the  $q$  smallest IMED indexes (breaking ties randomly so that this set has size  $q$ ). We will also consider the two following quantities for each time  $t$ :

$$I^*(t) = \min_{a \in \mathcal{A}^*(t)} I_a(t) = \min_{a \in \mathcal{A}^*(t)} \log N_a(t)$$

$$I(t) = \min_{\substack{\mathcal{A}' \subset \mathcal{A} \\ |\mathcal{A}'|=q}} \sum_{a' \in \mathcal{A}'} I_{a'}(t) = \sum_{a' \in \mathcal{A}_q(t)} I_{a'}(t)$$

$I(t)$  can be computed efficiently by summing the  $q$  smallest elements of the list of IMED indexes. The complexity of computing  $I(t)$  once the IMED indexes have been computed is at most the one of sorting a list of  $|\mathcal{A}|$  elements,  $\mathcal{O}(|\mathcal{A}| \log |\mathcal{A}|)$ , which is only  $\log |\mathcal{A}|$  times larger than looking for the minimal IMED index. Using selection algorithms, we may even achieve a better mean time complexity. The IMED-EC algorithm is presented in Algorithm 12.

---

#### Algorithm 12 IMED for Equivalent classes

---

```

Pull each arm once
for  $t = |\mathcal{A}| \dots T - 1$  do
  if  $I^*(t) \leq I(t)$  then
    Pull  $a_{t+1} \in \arg \min_{a \in \mathcal{A}^*(t)} N_a(t)$  (chosen arbitrarily)
  else
    Pull  $a_{t+1} \in \arg \min_{a \notin \mathcal{A}^*(t)} I_a(t)$  (chosen arbitrarily)
  end if
end for

```

---

While the original problem involves combinatorial quantities, those are not involved in the IMED-EC algorithm. From a time complexity viewpoint, this makes this algorithm on par with other popular algorithms such as UCB, KLUCB, and IMED algorithm. On the contrary, the general structure algorithm OSSB involves solving a combinatorial optimization problem at each time step which makes it numerically inefficient. We are not aware of a general relaxation method for this algorithm that we could compare IMED-EC with. It is interesting to note that in the case where  $q = 1$ , the IMED-EC algorithms coincide with the IMED algorithm.

**Intuition** For an arm  $a$ ,  $N_a(t)\text{KL}(\hat{\mu}_a(t)|\hat{\mu}^*(t))$  may be interpreted as the opposite of a *log-likelihood of optimality* of that arm.  $\log N_a(t)$  is linked to the log-frequency of play of that arm, the frequency of play of an arm being interpreted as the probability of pulling that arm is a sequence of length  $t$ . The IMED algorithm thus can be intuitively understood as an algorithm matching an empirical log-probability with a log-frequency of play. In our setting, there is at least  $q$  elements in each group. It therefore makes sense to test for the optimality of a group rather single elements. Since all arms are independent, it makes sense to sum the *log-likelihood of optimality* on all the  $q$ -partitions of the set of arms. Since we have the intuition that this first part is the logarithm of a product of probability, we may compare it to the product of the frequencies. Therefore, we get that important quantities are the sum of IMED indexes for each  $q$  partition of the arms, seen as a comparison between the optimality of this group of  $q$  elements and the associated frequency of play of that group. The minimal IMED index is the one whose frequency of play is the lowest compared to its *likelihood of optimality*, similarly for the sum of IMED indexes.

## D.4 Regret analysis

In this section, we now detail the main bound on the regret of IMED-EC .

**Theorem 8** (Upper bound on the number of pulls). *Under the IMED-EC algorithms, the number of pulls of a suboptimal arm  $a$  is upper bounded by:*

$$\mathbb{E}_\nu[N_a(T)] \leq \frac{\log T}{q \text{KL}(\mu_a|\mu^*)} (1 + \alpha(\varepsilon)) + f(\varepsilon) \quad (\text{D.5})$$

where  $0 < \varepsilon < \frac{1}{3} \min_{a \in \mathcal{A} \setminus \mathcal{A}^*} (\mu^* - \mu_a) \alpha$  and  $\alpha$  and  $f$  tends to 0 as  $\varepsilon$  tends to 0.

**Corollary 5.** *Under the IMED-EC algorithms, the number of pulls of a suboptimal arm  $a$  is upper bounded by:*

$$\min_{c_q \subseteq c} \sum_{a \in c_q} \mathbb{E}_\nu[N_a(T)] \text{KL}(\mu_a|\mu^*) \leq (1 + \alpha(\varepsilon)) \log T + g(\varepsilon) \quad (\text{D.6})$$

where  $0 < \varepsilon < \frac{1}{3} \min_{a \in \mathcal{A} \setminus \mathcal{A}^*} (\mu^* - \mu_a) \alpha$  and  $\alpha$  and  $g$  tends to 0 as  $\varepsilon$  tends to 0.

**Theorem 9** (Asymptotic upper bound on the number of pulls). *Under the IMED-EC algorithms, the number of pulls of a suboptimal arm  $a$  is asymptotically upper bounded by:*

$$\liminf_{t \rightarrow +\infty} \frac{\mathbb{E}_\nu[N_a(T)]}{\log T} \leq \frac{1}{q \text{KL}(\mu_a|\mu^*)} \quad (\text{D.7})$$

**Discussion** This upper bound shows that in particular, the number of pulls of a suboptimal class,  $\sum_{a \in c} \mathbb{E}_\nu[N_a(T)]$  is asymptotically no more than  $\frac{|c|}{q \text{KL}(\mu_a|\mu^*)} \log T$ . This hence matches the lower bound in the *non-combinatorial regime*. In the *combinatorial regime*, along with equation (D.4), this regret upper bound shows that

$$\frac{|c|}{q \text{KL}(\mu_a|\mu^*)} \geq \liminf_{T \rightarrow \infty} \sum_{a \in c} \frac{\mathbb{E}_\nu[N_a(T)]}{\log T} \geq \frac{1}{2} \cdot \frac{|c|}{q \text{KL}(\mu_a|\mu^*)}$$

proving that the regret of the proposed IMED-EC does not differ from the optimal lower bound by a factor more than 2. This is striking result.

Full proof of Theorem 9 and Theorem 8 are provided in Appendix C from [Pesquerel et al. \(2021\)](#).

## D.5 Experiments

In this section, we support our theoretical analysis by conducting three sets of experiments. The Python code used to perform those experiments is available on Github. We support our empirical evidences using plots of cumulative regrets. In this section, all the experiments are conducted using gaussian distributions whose means are between 0 and 1 and of unit standard deviation. Those graphs are representative of all the experiments that we conducted.

**Balanced class, perfect knowledge** In this set of experiments, see Figure D.1, we focus on the bandit configurations in which all equivalence classes have the same cardinality and assume that we know the number of elements per class. This setting is interesting for two reasons. First, one can compute the theoretical lower-bound without solving a combinatorial optimization problem. Second, the theoretical analysis shows that  $\text{IMED-EC}$  is asymptotically optimal in this case. This setting will thus allow us to numerically grasp what happens in the most structured case. We compare  $\text{IMED-EC}$  to unspecialized bandit algorithm, UCB, IMED and  $\text{KLUCB}$ . To make the comparison fairer we also compare  $\text{IMED-EC}$  to  $\text{OSSB}$ , an algorithm specialized in structured bandit. Since  $\text{OSSB}$  has to solve a combinatorial optimization problem at each time step, we cannot carry experiments on large set of arms while comparing  $\text{IMED-EC}$  to it. In this particular setting,

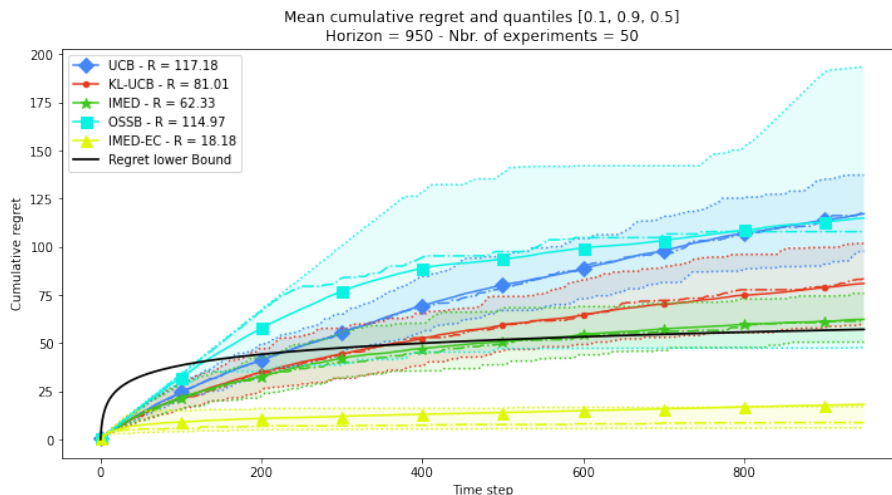


Figure D.1: 3 classes, 3 distributions per class - set of means = (0.1, 0.3, 0.6, 0.8)

we see that while  $\text{OSSB}$  and  $\text{IMED-EC}$  are provably asymptotically optimal,  $\text{IMED-EC}$  numerically performs better in finite time horizon. We recall that it is furthermore numerically more efficient since it does not involve any combinatorial optimization. Without too much surprises,  $\text{IMED-EC}$  also outperforms unspecialized algorithm.

**Unperfect knowledge** In the experiment plotted Figure D.2, we leverage the knowledge hypothesis and assume that we only know a lower bound on the number of elements per class while the classes are still balanced. We compare  $\text{IMED-EC}$  to unspecialized bandit algorithm, IMED and  $\text{KLUCB}$ . We drop  $\text{OSSB}$  from our test bed due to the computational burden of solving a combinatorial optimization problem at each time step. We can see that the finite time cumulative regret of  $\text{IMED-EC}$  indeed is much smaller than the regret of the unspecialized algorithms.

**Influence of the parameter  $q$**  Here we show the numerical robustness of  $\text{IMED-EC}$  with respect to the lower bound parameter  $q$  on the number of elements per classes. On the same bandit problem, we compare

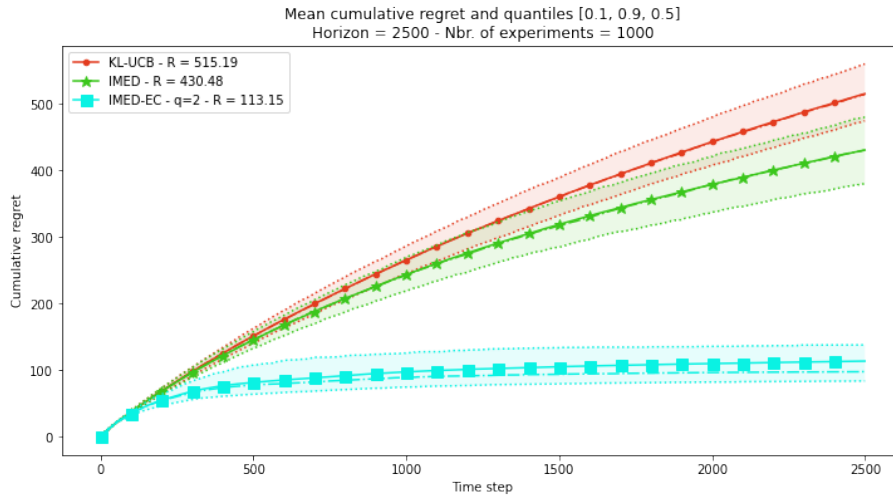


Figure D.2: 7 classes, 8 distributions per class - set of means =  $(0.1, 0.3, 0.4, 0.5, 0.6, 0.75, 0.9)$

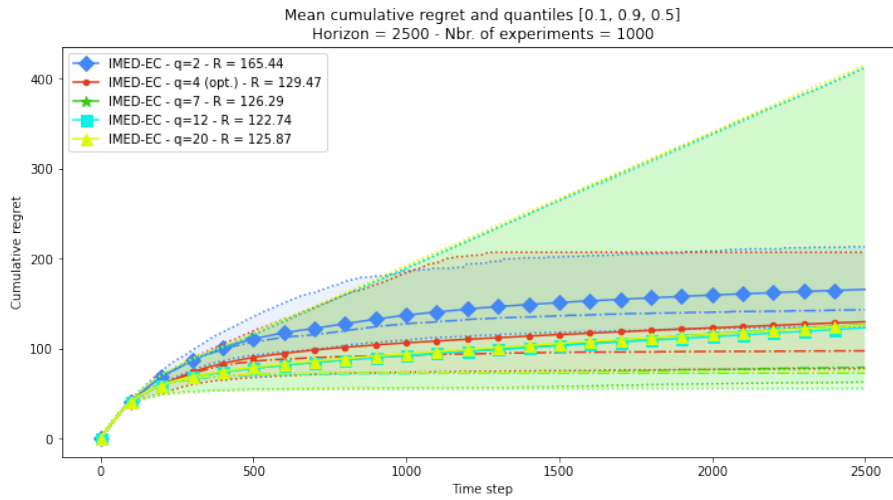


Figure D.3: 7 classes, unbalanced - set of means =  $(0.1, 0.3, 0.4, 0.5, 0.6, 0.75, 0.9)$

different instances of IMED-EC where different values of  $q$  are used. In the legend, *opt.* stands for optimal and corresponds to the largest valid lower bound on the number of elements per class, *i.e.* the minimal number of elements in a class. The experiments Figure D.3 is performed on a bandit problem with 7 classes and an uneven number of distributions per class. The smallest class has 4 elements and the largest 23. While  $q$  increases up to the minimum cardinality of a class, we see that the performances of IMED-EC increases. It is rather remarkable that once we go beyond that theoretical threshold, the performances of IMED-EC do not deteriorate. We even found it difficult to find setting to deteriorate them at all. While the expected regret does not seem to deteriorate, we sometimes see that the tails of the regret widen as it can be seen on the plot Figure D.3 for  $q = 7$  and  $q = 20$  since the 0.9 quantile curves are so large for those values of  $q$ . We interpret part of this robustness to the fact that the relaxation induced in IMED-EC makes the algorithm over explore compared to what the true lower bound suggests. Increasing  $q$  reduces the exploration and therefore may improve the performances of the algorithm. However, this robustness is observed even in the case where the classes are balanced. This interpretation thus does not explain everything about the numerical robustness of IMED-EC. This type of experiment does not take more than roughly 10 to 15 minutes on a notebook run in Google Colab depending on the number of arms, the horizon and the number of runs. This supports the numerical efficiency of the relaxation made in IMED-EC.

## D.6 Conclusion

In this chapter, we introduced  $\text{IMED-EC}$ , a numerically efficient algorithm to solve a structured bandit problem for which we derived a lower bound involving a combinatorial optimization problem. While not being asymptotically optimal, we proved that the asymptotic regret of  $\text{IMED-EC}$  is always smaller than the unstructured one and that we can control the discrepancy with respect to the structured regret lower bound by a factor of at most 2.