



**HAL**  
open science

# KL-UCB-switch: optimal regret bounds for stochastic bandits from both a distribution-dependent and a distribution-free viewpoints

Aurélien Garivier, Hédi Hadiji, Pierre Ménard, Gilles Stoltz

► **To cite this version:**

Aurélien Garivier, Hédi Hadiji, Pierre Ménard, Gilles Stoltz. KL-UCB-switch: optimal regret bounds for stochastic bandits from both a distribution-dependent and a distribution-free viewpoints. 2019. hal-01785705v2

**HAL Id: hal-01785705**

**<https://hal.science/hal-01785705v2>**

Preprint submitted on 5 Nov 2019 (v2), last revised 28 Jun 2022 (v3)

**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.

# KL-UCB-switch: optimal regret bounds for stochastic bandits from both a distribution-dependent and a distribution-free viewpoints

Aurélien Garivier\* — Hédi Hadiji† — Pierre Ménard‡ — Gilles Stoltz†,§

This draft: November 5, 2019

First draft: May 14, 2018 (arXiv:1805.05071v1)

---

## Abstract

In the context of  $K$ -armed stochastic bandits with distribution only assumed to be supported by  $[0, 1]$ , we introduce the first algorithm, called KL-UCB-switch, that enjoys *simultaneously* a distribution-free regret bound of optimal order  $\sqrt{KT}$  and a distribution-dependent regret bound of optimal order as well, that is, matching the  $\kappa \ln T$  lower bound by Lai and Robbins [1985] and Burnetas and Katehakis [1996]. This self-contained contribution simultaneously presents state-of-the-art techniques for regret minimization in bandit models, and an elementary construction of non-asymptotic confidence bounds based on the empirical likelihood method for bounded distributions.

Keywords:  $K$ -armed stochastic bandits, distribution-dependent regret bounds, distribution-free regret bounds

## 1. Introduction and brief literature review

Great progress has been made, over the last decades, in the understanding of the stochastic  $K$ -armed bandit problem. In this simplistic and yet paradigmatic sequential decision model, an agent can at each step  $t \in \mathbb{N}^*$  sample one out of  $K$  independent sources of randomness and receive the corresponding outcome as a reward. The most investigated challenge is to minimize the regret, which is defined as the difference between the cumulated rewards obtained by the agent and by an oracle knowing in hindsight the distribution with largest expectation.

After Thompson's seminal paper (Thompson, 1933) and Gittins' Bayesian approach in the 1960s, Lai and his co-authors wrote in the 1980s a series of articles laying the foundations of a frequentist analysis of bandit strategies based on confidence regions. Lai and Robbins [1985] provided a general asymptotic lower bound, for parametric bandit models: for any reasonable strategy, the regret after  $T$  steps grows at least as  $\kappa \ln(T)$ , where  $\kappa$  is an informational complexity measure of the problem. In the 1990s, Agrawal [1995] and Burnetas and Katehakis [1996] analyzed the UCB algorithm (see also the later analysis by Auer et al., 2002a), a simple procedure where at step  $t$  the arm with highest upper

---

\*Univ. Lyon, ENS de Lyon, UMPA UMR 5669, LIP UMR 5668, [aurelien.garivier@ens-lyon.fr](mailto:aurelien.garivier@ens-lyon.fr)

†Laboratoire de mathématiques d'Orsay, Université Paris-Sud, CNRS, Université Paris-Saclay, Orsay, France; [hedi.hadiji@math.u-psud.fr](mailto:hedi.hadiji@math.u-psud.fr) and [gilles.stoltz@math.u-psud.fr](mailto:gilles.stoltz@math.u-psud.fr)

‡Inria Lille Nord Europe; [pierre.menard@inria.fr](mailto:pierre.menard@inria.fr)

§HEC Paris, Jouy-en-Josas, France; [stoltz@hec.fr](mailto:stoltz@hec.fr)

confidence bound is chosen. The same authors also extended the lower bound by Lai and Robbins to non-parametric models.

In the early 2000s, the much noticed contributions of Auer et al. [2002a] and Auer et al. [2002b] promoted three important ideas.

1. First, a bandit strategy should not address only specific statistical models, but general and non-parametric families of probability distributions, e.g., bounded distributions.
2. Second, the regret analysis should not only be asymptotic, but should provide finite-time bounds.
3. Third, a good bandit strategy should be competitive with respect to two concurrent notions of optimality: distribution-dependent optimality (it should reach the asymptotic lower bound of Lai and Robbins and have a regret not much larger than  $\kappa \ln(T)$ ) and distribution-free optimality (the maximal regret over all considered probability distributions should be of the optimal order  $\sqrt{KT}$ ).

These efforts were pursued by further works in those three directions. Maillard et al. [2011] and Garivier and Cappé [2011] simultaneously proved that the distribution-dependent lower bound could be reached with exactly the right multiplicative constant in simple settings (for example, for binary rewards) and provided finite-time bounds to do so. They were followed by similar results for other index policies like BayesUCB (Kaufmann et al., 2012) or Thompson sampling (Korda et al., 2013).

Initiated by Honda and Takemura for the IMED algorithm (see Honda and Takemura, 2015 and references to earlier works of the authors therein) and followed by Cappé et al. [2013] for the KL-UCB algorithm, the use of the *empirical likelihood method* for the construction of the upper confidence bounds was proved to be optimal as far as distribution-dependent bounds are concerned. The analysis for IMED was led for all (semi-)bounded distributions, while the analysis for KL-UCB was only successfully achieved in some classes of distributions (e.g., bounded distributions with finite supports). A contribution in passing of the present article is to also provide optimal distribution-dependent bounds for KL-UCB for families of bounded distributions.

On the other hand, classical UCB strategies were proved not to enjoy distribution-free optimal regret bounds. A modified strategy named MOSS was proposed by Audibert and Bubeck [2009] to address this issue: minimax (distribution-free) optimality was proved, but distribution-dependent optimality was then not considered. It took a few more years before Ménard and Garivier [2017] and Lattimore [2016] proved that, in simple parametric settings, a strategy can enjoy, at the same time, regret bounds that are optimal both from a distribution-dependent and a distribution-free viewpoints.

**Main contributions.** In this work, we generalize the latter bi-optimality result to the non-parametric class of distributions with bounded support, say,  $[0, 1]$ . Namely, we propose the KL-UCB-switch algorithm, a bandit strategy belonging to the family of upper-confidence-bounds strategies. We prove that it is simultaneously optimal from a distribution-free viewpoint (Theorem 1) and from a distribution-dependent viewpoint in the considered class of distributions (Theorem 2).

We go one step further by providing, as Honda and Takemura [2015] already achieved for IMED, a second-order term of the optimal order  $-\ln(\ln(T))$  in the distribution-dependent bound (Theorem 3). This explains from a theoretical viewpoint why simulations consistently show strategies having a regret smaller than the main term of the lower bound of Lai and Robbins [1985]. Note that, to the best of our knowledge, IMED is not proved to enjoy an optimal distribution-free regret bound; only a distribution-dependent regret analysis was provided for it. And according to the numerical experiments (see Section 3) IMED indeed does not seem to be optimal from a distribution-free viewpoint.

Beyond these results, we took special care of the clarity and simplicity of all the proofs, and all our bounds are finite time, with closed-form expressions. In particular, we provide for the first time an elementary analysis of performance of the KL-UCB algorithm on the class of all distributions over a bounded interval. The study of KL-UCB in Cappé et al. [2013] indeed remained somewhat intricate and limited to finitely supported distributions. Furthermore, our simplified analysis allowed us to derive similar optimality results for the anytime version of this new algorithm, with little if no additional effort (see Theorems 4 and 5).

**Organization of the paper.** Section 2 presents the main contributions of this article: a description of the KL-UCB-switch algorithm, the precise statement of the aforementioned theorems, and corresponding results for an anytime version of the KL-UCB-switch algorithm. Section 3 discusses some numerical experiments comparing the performance of an empirically tuned version of the KL-UCB-switch algorithm to competitors like IMED or KL-UCB. The focus is not only set on the growth of the regret with time, but also on its dependency with respect to the number  $K$  of arms. Section 4 contains the statements and the proofs of several results that were already known before, but for which we sometimes propose a simpler derivation. All technical results needed in this article are stated and proved from scratch (e.g., on the  $\mathcal{K}_{\text{inf}}$  quantity that is central to the analysis of IMED and KL-UCB, and on the analysis of the performance of MOSS), though sometimes in appendix, which makes our paper fully self-contained. These known results are used as building blocks in Section 5 and 6, where the main results of this article are proved: Section 5 is devoted to distribution-free bounds, while Section 6 focuses on distribution-dependent bounds. An appendix provides the proofs of the classical material presented in Section 4, whenever these proofs did not fit in a few lines: anytime analysis of the MOSS strategy (Appendix A) and proofs of the regularity and deviation results on the  $\mathcal{K}_{\text{inf}}$  quantity mentioned above (Appendix B), which might be of independent interest. It also features the proof of a sophisticated distribution-dependent regret bound in the case of a known  $T$ : a regret bound with an optimal second order term (Appendix C).

## 2. Setting and statement of the main results

We consider the simplest case of a bounded stochastic bandit problem, with finitely many arms indexed by  $a \in \{1, \dots, K\}$  and with rewards in  $[0, 1]$ . We denote by  $\mathcal{P}[0, 1]$  the set of probability distributions over  $[0, 1]$ : each arm  $a$  is associated with an unknown probability distribution  $\nu_a \in \mathcal{P}[0, 1]$ . We call  $\underline{\nu} = (\nu_1, \dots, \nu_K)$  a bandit problem over  $[0, 1]$ . At each round  $t \geq 1$ , the player pulls the arm  $A_t$  and gets a real-valued reward  $Y_t$  drawn independently at random according to the distribution  $\nu_{A_t}$ . This reward is the only piece of information available to the player.

A typical measure of the performance of a strategy is given by its *regret*. To recall its definition, we denote by  $\mathbb{E}(\nu_a) = \mu_a$  the expected reward of arm  $a$  and by  $\Delta_a$  its gap to an optimal arm:

$$\mu^* = \max_{a=1, \dots, K} \mu_a \quad \text{and} \quad \Delta_a = \mu^* - \mu_a.$$

Arms  $a$  such that  $\Delta_a > 0$  are called sub-optimal arms. The expected regret of a strategy equals

$$R_T = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T Y_t \right] = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T \mu_{A_t} \right] = \sum_{a=1}^K \Delta_a \mathbb{E}[N_a(T)] \quad \text{where} \quad N_a(T) = \sum_{t=1}^T \mathbb{1}_{\{A_t=a\}}.$$

The first equality above follows from the tower rule. To control the expected regret, it is thus sufficient to control the  $\mathbb{E}[N_a(T)]$  quantities for sub-optimal arms  $a$ .

**Reminder of the existing lower bounds.** The distribution-free lower bound of Auer et al. [2002b] states that for all strategies, for all  $T \geq 1$  and all  $K \geq 2$ ,

$$\sup_{\underline{\nu}} R_T \geq \frac{1}{20} \min \left\{ \sqrt{KT}, T \right\}, \tag{1}$$

where the supremum is taken over all bandit problems  $\underline{\nu}$  over  $[0, 1]$ . Hence, a strategy is called optimal from a distribution-free viewpoint if there exists a numerical constant  $C$  such that for all  $K \geq 2$ , for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , for all  $T \geq 1$ , the regret is bounded by  $R_T \leq C\sqrt{KT}$ .

We denote by  $\mathcal{P}[0, 1]$  the set of all distributions over  $[0, 1]$ . The key quantity in stating distribution-dependent lower bounds is based on KL, the Kullback-Leibler divergence between two probability distributions. We recall its definition: consider two probability distributions  $\nu, \nu'$  over  $[0, 1]$ . We write

$\nu \ll \nu'$  when  $\nu$  is absolutely continuous with respect to  $\nu'$ , and denote by  $d\nu/d\nu'$  the density (the Radon-Nikodym derivative) of  $\nu$  with respect to  $\nu'$ . Then,

$$\text{KL}(\nu, \nu') = \begin{cases} \int_{[0,1]} \ln\left(\frac{d\nu}{d\nu'}\right) d\nu & \text{if } \nu \ll \nu'; \\ +\infty & \text{otherwise.} \end{cases}$$

Now, the key information-theoretic quantity for stochastic bandit problems is given by an infimum of Kullback-Leibler divergences: for  $\nu_a \in \mathcal{P}[0, 1]$  and  $x \in [0, 1]$ ,

$$\mathcal{K}_{\text{inf}}(\nu_a, x) = \inf \left\{ \text{KL}(\nu_a, \nu'_a) : \nu'_a \in \mathcal{P}[0, 1] \text{ and } \mathbb{E}(\nu'_a) > x \right\}$$

where  $\mathbb{E}(\nu'_a)$  denotes the expectation of the distribution  $\nu'_a$  and where by convention, the infimum of the empty set equals  $+\infty$ . Because of this convention, we may equivalently define  $\mathcal{K}_{\text{inf}}$  as

$$\mathcal{K}_{\text{inf}}(\nu_a, x) = \inf \left\{ \text{KL}(\nu_a, \nu'_a) : \nu'_a \in \mathcal{P}[0, 1] \text{ with } \nu_a \ll \nu'_a \text{ and } \mathbb{E}(\nu'_a) > x \right\}. \quad (2)$$

As essentially proved by Lai and Robbins [1985] and Burnetas and Katehakis [1996]—see also Garivier et al., 2018—, for any “reasonable” strategy, for any bandit problem  $\underline{\nu}$  over  $[0, 1]$ , for any sub-optimal arm  $a$ ,

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq \frac{1}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)}. \quad (3)$$

A strategy is called optimal from a distribution-dependent viewpoint if the reverse inequality holds with a lim sup instead of a lim inf, for any bandit problem  $\underline{\nu}$  over  $[0, 1]$  and for any sub-optimal arm  $a$ .

By a “reasonable” strategy above, we mean a strategy that is uniformly fast convergent on  $\mathcal{P}[0, 1]$ , that is, such that for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , for all sub-optimal arms  $a$ ,

$$\forall \alpha > 0, \quad \mathbb{E}[N_a(T)] = o(T^\alpha);$$

there exist such strategies, for instance, the UCB strategy already mentioned above. For uniformly super-fast convergent strategies, that is, strategies for which there actually exists a constant  $C$  such for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , for all sub-optimal arms  $a$ ,

$$\frac{\mathbb{E}[N_a(T)]}{\ln T} \leq \frac{C}{\Delta_a^2}$$

(again, UCB is such a strategy), the lower bound above can be strengthened into: for any bandit problem  $\underline{\nu}$  over  $[0, 1]$ , for any sub-optimal arm  $a$ ,

$$\mathbb{E}[N_a(T)] \geq \frac{\ln T}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} - \Omega(\ln(\ln T)), \quad (4)$$

see Garivier et al. [2018, Section 4]. This order of magnitude  $-\ln(\ln T)$  for the second-order term in the regret bound is optimal, as follows from the upper bound exhibited by Honda and Takemura [2015, Theorem 5].

## 2.1. The KL-UCB-switch algorithm

---

**Algorithm 1** Generic index policy

---

**Inputs:** index functions  $U_a$

**Initialization:** Play each arm  $a = 1, \dots, K$  once and compute the  $U_a(K)$

**for**  $t = K, \dots, T - 1$  **do**

Pull an arm  $A_{t+1} \in \arg \max_{a=1, \dots, K} U_a(t)$

Get a reward  $Y_{t+1}$  drawn independently at random according to  $\nu_{A_{t+1}}$

**end for**

---

For any index policy as described above, we have  $N_a(t) \geq 1$  for all arms  $a$  and  $t \geq K$  and may thus define, respectively, the empirical distribution of the rewards associated with arm  $a$  up to round  $t$  included and their empirical mean:

$$\hat{\nu}_a(t) = \frac{1}{N_a(t)} \sum_{s=1}^t \delta_{Y_s} \mathbb{1}_{\{A_s=a\}} \quad \text{and} \quad \hat{\mu}_a(t) = \mathbb{E}[\hat{\nu}_a(t)] = \frac{1}{N_a(t)} \sum_{s=1}^t Y_s \mathbb{1}_{\{A_s=a\}},$$

where  $\delta_y$  denotes the Dirac point-mass distribution at  $y \in [0, 1]$ .

The MOSS algorithm (see Audibert and Bubeck 2009) uses the index functions

$$U_a^{\text{M}}(t) \stackrel{\text{def}}{=} \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \ln_+ \left( \frac{T}{KN_a(t)} \right)}, \quad (5)$$

where  $\ln_+$  denotes the non-negative part of the natural logarithm,  $\ln_+ = \max\{\ln, 0\}$ .

We also consider a slight variation of the KL-UCB algorithm (see Cappé et al. 2013), which we call KL-UCB<sup>+</sup> and which relies on the index functions

$$U_a^{\text{KL}}(t) \stackrel{\text{def}}{=} \sup \left\{ \mu \in [0, 1] \mid \mathcal{K}_{\text{inf}}(\hat{\nu}_a(t), \mu) \leq \frac{1}{N_a(t)} \ln_+ \left( \frac{T}{KN_a(t)} \right) \right\}. \quad (6)$$

We introduce a new algorithm KL-UCB-switch. The novelty here is that this algorithm switches from the KL-UCB-type index to the MOSS index once it has pulled an arm more than  $f(T, K)$  times. The purpose is to capture the good properties of both algorithms. In the sequel we will take  $f(T, K) = \lfloor (T/K)^{1/5} \rfloor$ . More precisely, we define the index functions

$$U_a(t) = \begin{cases} U_a^{\text{KL}}(t) & \text{if } N_a(t) \leq f(T, K), \\ U_a^{\text{M}}(t) & \text{if } N_a(t) > f(T, K). \end{cases}$$

The reasons for the choice of a threshold  $f(T, K) = \lfloor (T/K)^{1/5} \rfloor$  will become clear in the proof of Theorem 1. Note that asymptotically KL-UCB-switch should behave like KL-UCB-type algorithm, as for large  $T$  we expect the number of pulls of a sub-optimal arm to be of order  $N_a(t) \sim \ln(T)$  and optimal arms to have been played linearly many times, entailing  $U_a^{\text{M}}(t) \approx U_a^{\text{KL}}(t) \approx \hat{\mu}_a(t)$ .

Since we are considering distributions over  $[0, 1]$ , the data-processing inequality for Kullback-Leibler divergences ensures (see, e.g., Garivier et al., 2018, Lemma 1) that for all  $\nu \in \mathcal{P}[0, 1]$  and all  $\mu \in (\mathbb{E}(\nu), 1)$ ,

$$\mathcal{K}_{\text{inf}}(\nu, \mu) \geq \inf_{\nu': \mathbb{E}(\nu') > \mu} \text{KL}(\text{Ber}(\mathbb{E}(\nu)), \text{Ber}(\mathbb{E}(\nu'))) = \text{KL}(\text{Ber}(\mathbb{E}(\nu)), \text{Ber}(\mu)),$$

where  $\text{Ber}(p)$  denotes the Bernoulli distribution with parameter  $p$ . Therefore, by Pinsker's inequality for Bernoulli distributions,

$$\mathcal{K}_{\text{inf}}(\nu, \mu) \geq 2(\mathbb{E}(\nu) - \mu)^2, \quad \text{thus} \quad U_a^{\text{KL}}(t) \leq U_a^{\text{M}}(t) \quad (7)$$

for all arms  $a$  and all rounds  $t \geq K$ . In particular, this actually shows that KL-UCB-switch interpolates between KL-UCB and MOSS,

$$U_a^{\text{KL}}(t) \leq U_a(t) \leq U_a^{\text{M}}(t). \quad (8)$$

## 2.2. Optimal distribution-dependent and distribution-free regret bounds (known horizon $T$ )

We first consider a fixed and beforehand-known value of  $T$ . The proofs of the two theorems below are provided in Sections 5 and 6, respectively.

**Theorem 1** (Distribution-free bound). *Given  $T \geq 1$ , the regret of the KL-UCB-switch algorithm, tuned with the knowledge of  $T$  and the switch function  $f(T, K) = \lfloor (T/K)^{1/5} \rfloor$ , is uniformly bounded over all bandit problems  $\underline{\nu}$  over  $[0, 1]$  by*

$$R_T \leq (K - 1) + 23\sqrt{KT}.$$

KL-UCB-switch thus enjoys a distribution-free regret bound of optimal order  $\sqrt{KT}$ , see (1). The MOSS strategy by Audibert and Bubeck [2009] already enjoyed this optimal distribution-free regret bound but its construction (relying on a sub-Gaussian assumption) prevents it from being optimal from a distribution-dependent viewpoint.

**Theorem 2** (Distribution-dependent bound). *Given  $T \geq 1$ , the KL-UCB-switch algorithm, tuned with the knowledge of  $T$  and the switch function  $f(T, K) = \lfloor (T/K)^{1/5} \rfloor$ , ensures that for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , for all sub-optimal arms  $a$ ,*

$$\mathbb{E}[N_a(T)] \leq \frac{\ln T}{\mathcal{K}_{\inf}(\nu_a, \mu^*)} + \mathcal{O}_T((\ln T)^{2/3}),$$

where a finite-time, closed-form expression of the  $\mathcal{O}_T((\ln T)^{2/3})$  term is given by (39) for the choice  $\delta = (\ln T)^{-1/3}$ .

By considering the exact same algorithm but by following a more sophisticated proof we may in fact get a stronger result, whose (extremely technical) proof is deferred to Appendix C.

**Theorem 3** (Distribution-dependent bound with a second-order term). *We actually have, when  $\mu^* \in (0, 1)$  and  $T \geq K/(1 - \mu^*)$ ,*

$$\mathbb{E}[N_a(T)] \leq \frac{\ln T - \ln \ln T}{\mathcal{K}_{\inf}(\nu_a, \mu^*)} + \mathcal{O}_T(1),$$

where a finite-time, closed-form expression of the  $\mathcal{O}_T(1)$  term is provided in (57).

KL-UCB-switch thus enjoys a distribution-dependent regret bounds of optimal orders, see (3) and (4). This optimal order was already reached by the IMED strategy by Honda and Takemura [2015] on the model  $\mathcal{P}[0, 1]$ . The KL-UCB algorithm studied, e.g., by Cappé et al. [2013], only enjoyed optimal regret bounds for more limited models; for instance, for distributions over  $[0, 1]$  with finite support. In the analysis of KL-UCB-switch we actually provide in passing an analysis of KL-UCB for the model  $\mathcal{P}[0, 1]$  of all distributions over  $[0, 1]$ .

### 2.3. Adaptation to the horizon $T$ (an anytime version of KL-UCB-switch)

A standard doubling trick fails to provide a meta-strategy that would not require the knowledge of  $T$  and have optimal  $\mathcal{O}(\sqrt{KT})$  and  $(1 + o(1))(\ln T)/\mathcal{K}_{\inf}(\nu_a, \mu^*)$  bounds. Indeed, there are first, two different rates,  $\sqrt{T}$  and  $\ln T$ , to accommodate simultaneously and each would require different regime lengths, e.g.,  $2^r$  and  $2^{2r}$ , respectively, and second, any doubling trick on the distribution-dependent bound would result in an additional multiplicative constant in front of the  $1/\mathcal{K}_{\inf}(\nu_a, \mu^*)$  factor. This is why a dedicated anytime version of our algorithm is needed.

For technical reasons, it was useful in our proof to perform some additional exploration, which deteriorates the second-order terms in the regret bound. Indeed, we define the augmented exploration function (which is non-decreasing) by

$$\varphi(x) = \ln_+(x(1 + \ln_+^2 x)) \tag{9}$$

and the associated index functions by

$$U_a^{\text{KL-A}}(t) \stackrel{\text{def}}{=} \sup \left\{ \mu \in [0, 1] \mid \mathcal{K}_{\text{inf}}(\widehat{\nu}_a(t), \mu) \leq \frac{1}{N_a(t)} \varphi\left(\frac{t}{KN_a(t)}\right) \right\} \quad (10)$$

$$U_a^{\text{M-A}}(t) \stackrel{\text{def}}{=} \widehat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \varphi\left(\frac{t}{KN_a(t)}\right)}. \quad (11)$$

A careful comparison of (10) and (11) to (5) and (6) shows that  $U_a^{\text{KL-A}}(t) \leq U_a^{\text{KL}}(t)$  and

$$U_a^{\text{M-A}}(t) \leq U_a^{\text{M},\varphi}(t) \stackrel{\text{def}}{=} \widehat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \varphi\left(\frac{T}{KN_a(t)}\right)} \quad (12)$$

when all these quantities are based on the same past (i.e., when they are defined for the same algorithm).

The -A in the superscripts stands for ‘‘augmented’’ or for ‘‘anytime’’ as this augmented exploration gives rise to the anytime version of KL-UCB-switch, which simply relies on the index

$$U_a^{\text{A}}(t) = \begin{cases} U_a^{\text{KL-A}}(t) & \text{if } N_a(t) \leq f(t, K) \\ U_a^{\text{M-A}}(t) & \text{if } N_a(t) > f(t, K) \end{cases} \quad (13)$$

where  $f(T, K) = \lfloor (t/K)^{1/5} \rfloor$ . Note that the thresholds  $f(t, K)$  when the switches occur from the sub-index  $U_a^{\text{KL-A}}(t)$  to the other sub-index  $U_a^{\text{M-A}}(t)$  now vary with  $t$  (and we cannot exclude that a switch back may occur).

For this anytime version of KL-UCB-switch, the same ranking of (sub-)indexes holds as the one (8) for our first version of KL-UCB-switch relying on the horizon  $T$ :

$$U_a^{\text{KL-A}}(t) \leq U_a^{\text{A}}(t) \leq U_a^{\text{M-A}}(t). \quad (14)$$

The performance guarantees are indicated in the next two theorems, whose proofs may be found in Sections 5 and 6, respectively. The distribution-free analysis is essentially the same as in the case of a known horizon, although the additional exploration required an adaptation of most of the calculations. Note also that the simulations detailed below suggest that all anytime variants of the KL-UCB algorithms (KL-UCB-switch included) behave better without the additional exploration required, i.e., with  $\ln_+$  as the exploration function.

**Theorem 4** (Anytime distribution-free bound). *The regret of the anytime version of KL-UCB-switch algorithm above, tuned with the switch function  $f(t, K) = \lfloor (t/K)^{1/5} \rfloor$ , is uniformly bounded over all bandit problems  $\underline{\nu}$  over  $[0, 1]$  as follows: for all  $T \geq 1$ ,*

$$R_T \leq (K - 1) + 44\sqrt{KT}.$$

**Theorem 5** (Anytime distribution-dependent bound). *The anytime version of KL-UCB-switch algorithm above, tuned with the switch function  $f(t, K) = \lfloor (t/K)^{1/5} \rfloor$ , ensures that for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , for all sub-optimal arms  $a$ , for all  $T \geq 1$ ,*

$$\mathbb{E}[N_a(T)] \leq \frac{\ln T}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} + \mathcal{O}_T((\ln T)^{6/7})$$

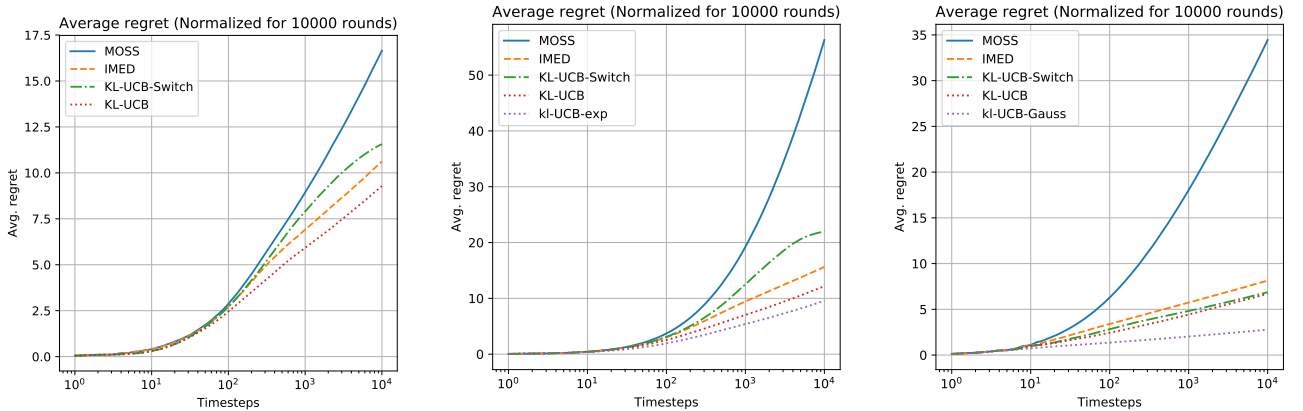
where a finite-time, closed-form expression of the  $\mathcal{O}_T((\ln T)^{6/7})$  term is given by Equation (32) for the choice  $\delta = (\ln T)^{-1/7}$ .



### 3. Numerical experiments

We provide here some numerical experiments comparing the different algorithms we refer to in this work. The KL-UCB-switch, KL-UCB, and MOSS algorithms are used in their anytime versions as described in Section 2.1 and Section 2.3. However, we stick to the natural exploration function  $\ln_+(t/(KN_a(t)))$ , i.e., without extra-exploration. For KL-UCB-switch we actually consider a slightly delayed switch function, different from the one in our theoretical analysis:  $f(t, K) = \lfloor t/K \rfloor^{8/9}$ , which generally exhibits a good empirical performance. While our choice  $f(t, K) = \lfloor t/K \rfloor^{1/5}$  appeared to be a good choice for minimizing the theoretical upper bounds, many other choices (such as the one considered in the experiments below) would also have been possible, at the cost of larger constants in one of the two regret bounds.

**Distribution-dependent bounds.** We compare in Figure 1 the distribution-dependent behaviors of the algorithms. For the two scenarios with truncated exponential or Gaussian rewards we also consider the appropriate version of the kl-UCB algorithm for one-parameter exponential family (see Cappé et al., 2013), with the same exploration function as for the other algorithms; we call these algorithms kl-UCB-exp or kl-UCB-Gauss, respectively. The parameters of the middle and right scenarios were chosen in a way that, even with the truncation, the kl-UCB algorithms have a significantly better performance than the other algorithms. This is the case because they are able to exploit the shape of the underlying distributions. Note that the kl-UCB-Gauss algorithm reduces to the MOSS algorithm with the constant  $2\sigma^2$  instead of  $1/2$ . As expected, the regret of KL-UCB-switch lies between the one of MOSS and the one of KL-UCB.



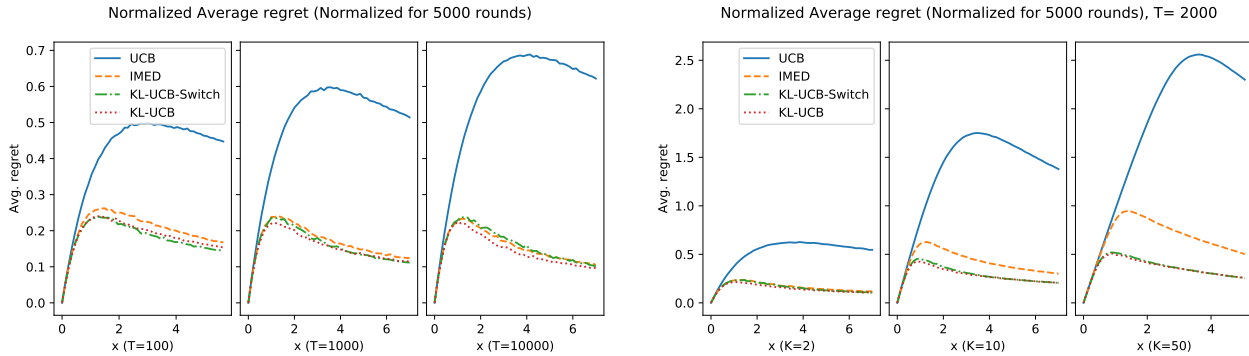
**Figure 1:** Regrets approximated over 10,000 runs, shown on a logarithmic scale; distributions of the arms consist of:

*Left:* Bernoulli distributions with parameters (0.9, 0.8)

*Middle:* Exponential distributions with expectations (0.15, 0.12, 0.10, 0.05), truncated on  $[0, 1]$ ,

*Right:* Gaussian distributions with means (0.7, 0.5, 0.3, 0.2) and same standard deviation  $\sigma = 0.1$ , truncated on  $[0, 1]$

**Distribution-free bounds.** Here we also consider the UCB algorithm of Auer et al. [2002a] with the exploration function  $\ln(t)$ . We plot the behavior of the normalized regret,  $R_T/\sqrt{KT}$ , either as a function of  $T$  (Figure 2 left) or of  $K$  (Figure 2 right). This quantity should remain bounded as  $T$  or  $K$  increases. KL-UCB-switch and KL-UCB have a normalized regret that does not depend too much on  $T$  and  $K$  (KL-UCB may perhaps satisfy a distribution-free bound of the optimal order, but we were unable to prove this fact). The regrets of UCB and IMED seem to suffer from a sub-optimal dependence in  $K$ .



**Figure 2:** Expected regret  $R_T/\sqrt{KT}$ , approximated over 5,000 runs

*Left:* as a function of  $x$ , for a Bernoulli bandit problem with parameters  $(0.8, 0.8 - x\sqrt{K/T})$  and for time horizons  $T \in \{100, 1000, 10000\}$

*Right:* as a function of  $x$ , for a Bernoulli bandit problem with parameters  $(0.8, 0.8 - x\sqrt{K/T}, \dots, 0.8 - x\sqrt{K/T})$  and  $K$  arms, where  $K \in \{2, 10, 50\}$

## 4. Results (more or less) extracted from the literature

We gather in this section results that are all known and published elsewhere (or almost). For the sake of self-completeness we provide a proof of each of them (sometimes this proof is shorter or simpler than the known proofs, and we then comment on this fact). *Readers familiar with the material described here are urged to move to the next section.*

### 4.1. Optional skipping—how to go from global times $t$ to local times $n$

The trick detailed here is standard in the bandit literature, see, e.g., its application in Auer et al. [2002a]. It is sometimes called optional skipping, and sometimes, optional sampling; we pick the first terminology, following what seems to be the preferred terminology in probability theory<sup>1</sup>. In any case, the original reference is Theorem 5.2 of Doob [1953, Chapter III, p. 145]; one can also check Chow and Teicher [1988, Section 5.3] for a more recent reference.

Doob’s optional skipping enables the rewriting of various quantities like  $U_a(t)$ ,  $\hat{\mu}_a(t)$ , etc., that are indexed by the global time  $t$ , into versions indexed by the local number of times  $N_a(t) = n$  that the specific arm considered has been pulled so far. The corresponding quantities will be denoted by  $U_{a,n}$ ,  $\hat{\mu}_{a,n}$ , etc.

The reindexation is possible as soon as the considered algorithm pulls each arm infinitely often; it is the case for all algorithms considered in this article (exploration never stops even if it becomes rare after a certain time).

We denote by  $\mathcal{F}_0 = \{\emptyset, \Omega\}$  the trivial  $\sigma$ -algebra and by  $\mathcal{F}_t$  the  $\sigma$ -algebra generated by  $A_1, Y_1, \dots, A_t, Y_t$ , when  $t \geq 1$ . We fix an arm  $a$ . For each  $n \geq 1$ , we denote by

$$\tau_{a,n} = \min\{t \geq 1 : N_a(t) = n\}$$

the round at which arm  $a$  was pulled for the  $n$ -th time. Now, Doob’s optional skipping ensures that the random variables  $X_{a,n} = Y_{\tau_{a,n}}$  are independent and identically distributed according to  $\nu_a$ .

We can then define, for instance, for  $n \geq 1$ ,

$$\hat{\mu}_{a,n} = \frac{1}{n} \sum_{k=1}^n X_{a,k}$$

and have the equality  $\hat{\mu}_a(t) = \hat{\mu}_{a,N_a(t)}$  for  $t \geq K$ . Here is an example of how to use this rewriting.

**Example 1** (Controlling an empirical average). Recall that  $N_a(t) \geq 1$  for  $t \geq K$  and  $N_a(t) \leq t - K + 1$  as each arm was pulled once in the first rounds. Given a subset  $\mathcal{E} \subseteq [0, 1]$ , we get the inclusion

$$\{\hat{\mu}_a(t) \in \mathcal{E}\} = \bigcup_{n=1}^{t-K+1} \{\hat{\mu}_a(t) \in \mathcal{E} \text{ and } N_a(t) = n\} = \bigcup_{n=1}^{t-K+1} \{\hat{\mu}_{a,n} \in \mathcal{E} \text{ and } N_a(t) = n\}$$

so that, by a union bound,

$$\mathbb{P}[\hat{\mu}_a(t) \in \mathcal{E}] \leq \sum_{n=1}^{t-K+1} \mathbb{P}[\hat{\mu}_{a,n} \in \mathcal{E} \text{ and } N_a(t) = n] \leq \sum_{n=1}^{t-K+1} \mathbb{P}[\hat{\mu}_{a,n} \in \mathcal{E}].$$

The last sum above only deals with independent and identically distributed random variables; we took care of all dependency issues that are so present in bandit problems. The price to pay, however, is that we bounded one probability by a sum of probabilities.

<sup>1</sup>The abstract of a recent article by Simons et al. [2002] reads: “A general set of distribution-free conditions is described under which an i.i.d. sequence of random variables is preserved under optional skipping. This work is motivated by theorems of J.L. Doob (1936) and Z. Ignatov (1977), unifying and extending aspects of both.”

Actually, a more careful use of optional skipping would be

$$\mathbb{P}[\widehat{\mu}_a(t) \in \mathcal{E}] \leq \mathbb{P}\left[\bigcup_{n=1}^{t-K+1} \{\widehat{\mu}_{a,n} \in \mathcal{E}\}\right] = \mathbb{P}\left[\exists n \in \{1, \dots, t-K+1\} : \widehat{\mu}_{a,n} \in \mathcal{E}\right].$$

## 4.2. Maximal version of Hoeffding's inequality

The maximal version of Hoeffding's inequality (Proposition 6) is a standard result from Hoeffding [1963]. It was already used in the original analysis of MOSS (Audibert and Bubeck, 2009). For our slightly simplified analysis of MOSS (see Section 4.3), we will rather rely on Corollary 7, a consequence of Proposition 6 obtained by integrating it.

**Proposition 6.** *Let  $X_1, \dots, X_n$  be a sequence of i.i.d. random variables bounded in  $[0, 1]$  and let  $\widehat{\mu}_n$  denote their empirical mean. Then for all  $u \geq 0$  and for all  $N \geq 1$ :*

$$\mathbb{P}\left[\max_{n \geq N} (\widehat{\mu}_n - \mu) \geq u\right] \leq e^{-2Nu^2}. \quad (15)$$

**Corollary 7.** *Under the same assumptions, for all  $\varepsilon \geq 0$ ,*

$$\mathbb{E}\left[\left(\max_{n \geq N} (\mu - \widehat{\mu}_n - \varepsilon)\right)^+\right] \leq \sqrt{\frac{\pi}{8}} \sqrt{\frac{1}{N}} e^{-2N\varepsilon^2}. \quad (16)$$

Of course, Proposition 6 and Corollary 7 hold by symmetry with  $\mu - \widehat{\mu}_n$  instead of  $\widehat{\mu}_n - \mu$ .

**Proof:** By the Fubini-Tonelli theorem, an integration of the maximal deviation inequality (15) yields

$$\begin{aligned} \mathbb{E}\left[\left(\max_{n \geq N} (\mu - \widehat{\mu}_n - \varepsilon)\right)^+\right] &= \int_0^{+\infty} \mathbb{P}\left[\max_{n \geq N} (\widehat{\mu}_n - \mu - \varepsilon) \geq u\right] du \\ &\leq \int_0^{+\infty} e^{-2N(u+\varepsilon)^2} du \leq e^{-2N\varepsilon^2} \int_0^{+\infty} e^{-2Nu^2} du = \sqrt{\frac{\pi}{8}} \sqrt{\frac{1}{N}} e^{-2N\varepsilon^2}. \quad \square \end{aligned}$$

## 4.3. Distribution-free bound for the MOSS algorithm

Such a distribution-free bound was already provided in the literature, both for a known horizon  $T$  (see Audibert and Bubeck, 2009) and for an anytime version (see Degenne and Perchet, 2016). We only provide a slightly shorter and more focused proof of these results based on Corollary 7 and indicate an intermediate result—see (17)—that will be useful for us in the analysis of our new KL-UCB-switch algorithm. We do not claim any improvement on the results themselves, just a clarification of the existing proofs.

Our proof is slightly shorter and more focused for two reasons. First, in the two references mentioned, the peeling trick was used on the probabilities of deviations (see Proposition 6) and had to be performed separately and differently for each deviation  $u$ ; then, these probabilities were integrated to obtain a control on the needed expectations. In contrast, we perform the peeling trick directly on the expectations at hand, and we do so by applying it only once, based on Corollary 7 and at fixed times depending solely on  $T$ . Second, unlike the two mentioned references, we do not attempt to simultaneously build a distribution-free and some type of distribution-dependent bound. This raised technical difficulties because of the correlations between the choices of the arms and the observed rewards. The idea of our approach is to focus solely on the distribution-free regime, for which we notice that some crude bounding neglecting the correlations suffice (i.e., our analysis deals with all sub-optimal arms in the same way, independently of how often they are played).

For a known horizon  $T$ , we denote by  $A_{t+1}^M$  the arm played by the index strategy maximizing, at each step  $t + 1$  with  $t \geq K$ , the quantities (5):

$$U_a^M(t) \stackrel{\text{def}}{=} \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \ln_+ \left( \frac{T}{KN_a(t)} \right)}.$$

The superscripts M in  $A_{t+1}^M$  and  $U_a^M(t)$  stand for MOSS. We do so not to mix it with the arm  $A_{t+1}$  played by the KL-UCB-switch strategy (no superscript), but of course, once an arm  $a$  was sufficiently pulled, we have  $A_{t+1} = A_{t+1}^M$  by definition of the KL-UCB-switch strategy.

Appendix A provides the proof of the following regret bound. We denote by  $a^*$  an optimal arm, i.e., an arm such that  $\mu_a = \mu^*$ .

**Proposition 8.** *For a known horizon  $T \geq 1$ , for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , MOSS achieves a regret bound smaller than  $R_T \leq (K - 1) + 17\sqrt{KT}$ . More precisely, with the notation of optional skipping (Section 4.1), we have the inequalities*

$$\begin{aligned} R_T &= T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T \mu_{A_t^M} \right] \\ &\leq (K - 1) + \underbrace{\sum_{t=K+1}^T \mathbb{E} \left[ (\mu^* - U_{a^*}^M(t-1))^+ \right]}_{\leq 13\sqrt{KT}} + \underbrace{\sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E} \left[ \left( \hat{\mu}_{a,n} + \sqrt{\frac{\ln_+(T/(Kn))}{2n}} - \mu_a - \sqrt{\frac{K}{T}} \right)^+ \right]}_{\leq 4\sqrt{KT}} \end{aligned} \quad (17)$$

**Remark 1.** The proof (see Remark 4) actually reveals that for a known horizon  $T \geq 1$ , for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , and for all strategies (not only MOSS), the following bound holds:

$$\sum_{t=K+1}^T \mathbb{E} \left[ (\mu^* - U_{a^*}^M(t-1))^+ \right] \leq 13\sqrt{KT}.$$

We will re-use this fact to state a similar remark below (Remark 2), which will be useful for Part 2 of the proof lying in Section 5.

Our proof in Appendix A reveals that designing an adaptive version of MOSS comes at no effort. For this adaptive version we will also want to possibly explore more. We will do so by considering an augmented exploration function  $\varphi$ , that is, a function  $\varphi \geq \ln_+$  as in (9). We therefore define MOSS-anytime (M-A) as relying on the indexes defined in (11), which we copy here:

$$U_a^{M-A}(t) \stackrel{\text{def}}{=} \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \varphi \left( \frac{t}{KN_a(t)} \right)}.$$

We denote by  $A_{t+1}^{M-A}$  the arm picked as  $\arg \max_{a=1, \dots, K} U_a^{M-A}(t)$ .

**Proposition 9.** *For all horizons  $T \geq 1$ , for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , MOSS-anytime achieves a regret bound smaller than  $R_T \leq (K - 1) + c\sqrt{KT}$  where  $c = 30$  for  $\varphi = \ln_+$  and  $c = 33$  for the augmented exploration function  $\varphi(x) = \ln_+(x(1 + \ln_+^2 x))$  defined in (9). More precisely, with the notation of optional skipping (Section 4.1), we have the inequalities*

$$\begin{aligned} R_T &= T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T \mu_{A_t^{M-A}} \right] \\ &\leq (K - 1) + \underbrace{\sum_{t=K+1}^T \mathbb{E} \left[ (\mu^* - U_{a^*}^{M-A}(t-1))^+ \right]}_{\leq 26\sqrt{KT}} + \underbrace{\sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E} \left[ \left( \hat{\mu}_{a,n} + \sqrt{\frac{\varphi(T/(Kn))}{2n}} - \mu_a - \sqrt{\frac{K}{T}} \right)^+ \right]}_{\leq 4\sqrt{KT} \text{ for } \varphi = \ln_+ \text{ and } 7\sqrt{KT} \text{ for } \varphi(x) = \ln_+(x(1 + \ln_+^2 x))} \end{aligned} \quad (18)$$

**Remark 2.** Similarly to above, the proof (see Remark 4) actually reveals that for a known horizon  $T \geq 1$ , for all bandit problems  $\nu$  over  $[0, 1]$ , and for all strategies (not only MOSS-anytime), the following bound holds:

$$\sum_{t=K+1}^T \mathbb{E} \left[ (\mu^* - U_{a^*}^{\text{M-A}}(t-1))^+ \right] \leq 26\sqrt{KT}.$$

This remark will be useful for Part 2 of the proof lying in Section 5.

#### 4.4. Regularity and deviation/concentration results on $\mathcal{K}_{\text{inf}}$

We start with a quantification of the (left-)regularity of  $\mathcal{K}_{\text{inf}}$  and then provide a deviation and a concentration result on  $\mathcal{K}_{\text{inf}}$ .

##### 4.4.1. Regularity of $\mathcal{K}_{\text{inf}}$

The lower left-semi-continuity (19) first appeared as Lemma 7 in Honda and Takemura [2015], see also Garivier et al. [2018, Lemma 3] for a later but simpler proof. The upper left-semi-continuity (20) relies on the same arguments as (7), namely, the data-processing inequality for Kullback-Leibler divergences and Pinsker's inequality. These two inequalities are proved in detail in Appendix B; the proposed proofs are slightly simpler or lead to sharper bounds than in the mentioned references.

**Lemma 10** (regularity of  $\mathcal{K}_{\text{inf}}$ ). *For all  $\nu \in \mathcal{P}[0, 1]$  and all  $\mu \in (0, 1)$ ,*

$$\forall \varepsilon \in [0, \mu], \quad \mathcal{K}_{\text{inf}}(\nu, \mu) \leq \mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon) + \frac{\varepsilon}{1 - \mu}, \quad (19)$$

and

$$\forall \varepsilon \in [0, \mu - \mathbb{E}(\nu)], \quad \mathcal{K}_{\text{inf}}(\nu, \mu) \geq \mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon) + 2\varepsilon^2. \quad (20)$$

We draw two consequences from Lemma 10: the left-continuity of  $\mathcal{K}_{\text{inf}}$  and a useful inclusion in terms of level sets.

**Corollary 11.** *For all  $\nu \in \mathcal{P}[0, 1]$ , the function  $\mathcal{K}_{\text{inf}}(\nu, \cdot) : \mu \in (0, 1) \mapsto \mathcal{K}_{\text{inf}}(\nu, \mu)$  is left-continuous. In particular, on the one hand,  $\mathcal{K}_{\text{inf}}(\nu, \mathbb{E}(\nu)) = 0$  whenever  $\mathbb{E}(\nu) \in (0, 1)$ , and on the other hand, for all  $\nu \in \mathcal{P}[0, 1]$  and  $\mu \in (0, 1)$ ,*

$$\mathcal{K}_{\text{inf}}(\nu, \mu) = \inf \left\{ \text{KL}(\nu, \nu') : \nu' \in \mathcal{P}[0, 1] \text{ and } \mathbb{E}(\nu') \geq \mu \right\}.$$

**Proof:** The left-continuity follows from a sandwich argument via the upper bound (19) and the lower bound  $\mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon) \leq \mathcal{K}_{\text{inf}}(\nu, \mu)$  that holds for all  $\varepsilon \in [0, \mu]$  by the very definition of  $\mathcal{K}_{\text{inf}}$ . The fact that  $\mathcal{K}_{\text{inf}}(\nu, \mathbb{E}(\nu) - \varepsilon) = 0$  for all  $\varepsilon \in (0, \mathbb{E}(\nu)]$  thus entails, in particular, that  $\mathcal{K}_{\text{inf}}(\nu, \mathbb{E}(\nu)) = 0$ .  $\square$

**Corollary 12.** *For all  $\nu \in \mathcal{P}[0, 1]$ , all  $\mu \in (0, 1)$ , all  $u > 0$ , and all  $\varepsilon > 0$ ,*

$$\{\mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon) > u\} \subseteq \{\mathcal{K}_{\text{inf}}(\nu, \mu) > u + 2\varepsilon^2\}.$$

**Proof:** We apply (20) and merely need to explain why the condition  $\varepsilon \in [0, \mu - \mathbb{E}(\nu)]$  therein is satisfied. Indeed,  $\mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon) > u > 0$  indicates in particular that  $\mu - \varepsilon > \mathbb{E}(\nu)$ , or put differently,  $\varepsilon < \mu - \mathbb{E}(\nu)$ .  $\square$

#### 4.4.2. Deviation results on $\mathcal{K}_{\text{inf}}$

We provide two deviation results on  $\mathcal{K}_{\text{inf}}$ : first, in terms of probabilities of deviations and next, in terms of expected deviations.

The first deviation inequality was essentially provided by Cappé et al. [2013, Lemma 6]. For the sake of completeness, we recall its proof in Section B.

**Proposition 13** (deviation result on  $\mathcal{K}_{\text{inf}}$ ). *Let  $\hat{\nu}_n$  denote the empirical distribution associated with a sequence of  $n \geq 1$  i.i.d. random variables with distribution  $\nu$  over  $[0, 1]$  with  $\mathbb{E}(\nu) \in (0, 1)$ . Then, for all  $u \geq 0$ ,*

$$\mathbb{P}\left[\mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mathbb{E}(\nu)) \geq u\right] \leq e(2n+1)e^{-nu}.$$

A useful corollary in terms of expected deviations can now be stated.

**Corollary 14** (integrated deviations for  $\mathcal{K}_{\text{inf}}$ ). *Under the same assumptions, for all  $\varepsilon > 0$ , the index*

$$U_{\varepsilon, n} = \sup\left\{\mu \in [0, 1] \mid \mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mu) \leq \varepsilon\right\}$$

*satisfies*

$$\mathbb{E}\left[(\mathbb{E}(\nu) - U_{\varepsilon, n})^+\right] \leq (2n+1)e^{-n\varepsilon} \sqrt{\frac{\pi}{n}}$$

**Proof:** By the Fubini-Tonelli theorem, just as in the proof of Corollary 7 (for the first two equalities), and subsequently using the definition of  $U_{\varepsilon, n}$  as a supremum (for the third equality, together with the left-continuity of  $\mathcal{K}_{\text{inf}}$  deriving from Lemma 10), we have

$$\begin{aligned} \mathbb{E}\left[(\mathbb{E}(\nu) - U_{\varepsilon, n})^+\right] &= \int_0^{+\infty} \mathbb{P}\left[\mathbb{E}(\nu) - U_{\varepsilon, n} > u\right] du = \int_0^{+\infty} \mathbb{P}\left[U_{\varepsilon, n} < \mathbb{E}(\nu) - u\right] du \\ &= \int_0^{+\infty} \mathbb{P}\left[\mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mathbb{E}(\nu) - u) > \varepsilon\right] du. \end{aligned}$$

Now, Corollary 12 (for the first inequality) and the deviation inequality of Proposition 13 (for the second inequality) indicate that for all  $u > 0$ ,

$$\mathbb{P}\left[\mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mathbb{E}(\nu) - u) > \varepsilon\right] \leq \mathbb{P}\left[\mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mathbb{E}(\nu)) > \varepsilon + 2u^2\right] \leq e(2n+1)e^{-n(\varepsilon+2u^2)}.$$

Combining all elements, we get

$$\mathbb{E}\left[(\mathbb{E}(\nu) - U_{\varepsilon, n})^+\right] \leq e(2n+1)e^{-n\varepsilon} \int_0^{+\infty} e^{-2nu^2} du = e(2n+1)e^{-n\varepsilon} \frac{1}{2} \sqrt{\frac{\pi}{2n}}.$$

from which the stated bound follows, as  $e/(2\sqrt{2}) \leq 1$ .  $\square$

#### 4.4.3. Concentration result on $\mathcal{K}_{\text{inf}}$

The next proposition is similar in spirit to Honda and Takemura [2015, Proposition 11] but is better suited to our needs. We prove it in Appendix B.

**Proposition 15** (concentration result on  $\mathcal{K}_{\text{inf}}$ ). *With the same notation and assumptions as in the previous proposition, consider a real number  $\mu \in (\mathbb{E}(\nu), 1)$  and define*

$$\gamma = \frac{1}{\sqrt{1-\mu}} \left( 16e^{-2} + \ln^2\left(\frac{1}{1-\mu}\right) \right). \quad (21)$$

*Then for all  $x < \mathcal{K}_{\text{inf}}(\nu, \mu)$ ,*

$$\mathbb{P}\left[\mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mu) \leq x\right] \leq \begin{cases} \exp(-n\gamma/8) \leq \exp(-n/4) & \text{if } x \leq \mathcal{K}_{\text{inf}}(\nu, \mu) - \gamma/2 \\ \exp\left(-n(\mathcal{K}_{\text{inf}}(\nu, \mu) - x)^2/(2\gamma)\right) & \text{if } x > \mathcal{K}_{\text{inf}}(\nu, \mu) - \gamma/2 \end{cases}.$$

## 5. Proofs of the distribution-free bounds: Theorems 1 and 4

The two proofs are extremely similar; we show, for instance, Theorem 4 and explain how to adapt the proof for Theorem 1. The first steps of the proof(s) use the exact same arguments as in the proofs of the performance bounds of MOSS (Propositions 8 and 9, see Appendix A) in the exact same order. We explain below why we had to copy them and had to resort to the intermediary bounds for MOSS stated in the indicated propositions.

We recall that we denote by  $a^*$  an optimal arm, i.e., an arm such that  $\mu_a = \mu^*$ . We first apply a trick introduced by Bubeck and Liu [2013]: by definition of the index policy, for  $t \geq K$ ,

$$U_{a^*}^A(t) \leq \max_{a=1,\dots,K} U_a^A(t) = U_{A_{t+1}}^A(t)$$

so that the regret of KL-UCB-switch is bounded by

$$R_T = \sum_{t=1}^T \mathbb{E}[\mu^* - \mu_{A_t}] \leq (K-1) + \sum_{t=K+1}^T \mathbb{E}[\mu^* - U_{a^*}^A(t-1)] + \sum_{t=K+1}^T \mathbb{E}[U_{A_t}^A(t-1) - \mu_{A_t}]. \quad (22)$$

*Part 1:* We first deal with the second sum in (22) and successively use  $x \leq \delta + (x - \delta)^+$  for all  $x$  and  $\delta$  for the first inequality; the fact that  $U_a^A(t) \leq U_a^{M-A}(t) \leq U_a^{M,\varphi}(t)$  by (12) and (14), for the second inequality; and optional skipping (Section 4.1) for the third inequality, keeping in mind that pairs  $(a, n)$  such  $A_t = a$  and  $N_a(t-1) = n$  correspond to at most one round  $t \in \{K+1, \dots, T\}$ :

$$\begin{aligned} \sum_{t=K+1}^T \mathbb{E}[U_{A_t}^A(t-1) - \mu_{A_t}] &\leq \sqrt{KT} + \sum_{t=K+1}^T \mathbb{E} \left[ \left( U_{A_t}^A(t-1) - \mu_{A_t} - \sqrt{\frac{K}{T}} \right)^+ \right] \\ &\leq \sqrt{KT} + \sum_{t=K+1}^T \mathbb{E} \left[ \left( U_{A_t}^{M,\varphi}(t-1) - \mu_{A_t} - \sqrt{\frac{K}{T}} \right)^+ \right] \end{aligned} \quad (23)$$

$$\leq \sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E} \left[ \left( U_{a,n}^{M,\varphi} - \mu_a - \sqrt{\frac{K}{T}} \right)^+ \right] \quad (24)$$

where we recall that

$$U_{a,n}^{M,\varphi} = \hat{\mu}_{a,n} + \sqrt{\frac{1}{2n} \varphi\left(\frac{T}{Kn}\right)}.$$

We now apply one of the bounds of Proposition 9 to further bound the sum at hand by

$$\sum_{t=K+1}^T \mathbb{E}[U_{A_t}^A(t-1) - \mu_{A_t}] \leq \sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E} \left[ \left( U_{a,n}^{M,\varphi} - \mu_a - \sqrt{\frac{K}{T}} \right)^+ \right] \leq 7\sqrt{KT}.$$

**Remark 3.** We may now explain why we copied the beginning of the proof of Proposition 9 and why we cannot just say that the ranking  $U_a^A(t) \leq U_a^{M-A}(t)$  entails that the regret of the anytime version of KL-UCB-switch is bounded by the regret of the anytime version of MOSS. Indeed, it is difficult to relate

$$\sum_{t=K+1}^T \mathbb{E}[U_{A_t}^{M-A}(t-1) - \mu_{A_t}] \quad \text{and} \quad \sum_{t=K+1}^T \mathbb{E}[U_{A_t}^{M-A}(t-1) - \mu_{A_t}^{M-A}]$$

as the two series of arms  $A_t$  (picked by KL-UCB-switch) and  $A_t^{M-A}$  (picked by the adaptive version of MOSS) cannot be related. Hence, it is difficult to directly bound quantities like (23). However, the proof of the performance bound of MOSS relies on optional skipping and considers, in some sense, all possible values  $a$  for the arms picked: it controls the quantity (24), which appears as a regret bound that is achieved by all index policies with indexes smaller than the ones of the anytime version of MOSS.



*Part 2:* We now deal with the first sum in (22). We take positive parts, get back to the definition (13) of  $U_{a^*}^A(t-1)$ , and add some extra non-negative terms:

$$\begin{aligned} & \sum_{t=K+1}^T \mathbb{E}[\mu^* - U_{a^*}^A(t-1)] \leq \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*}^A(t-1))^+\right] \\ & \leq \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*}^{\text{KL-A}}(t-1))^+ \mathbf{1}_{\{N_{a^*}(t-1) \leq f(t-1, K)\}}\right] + \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*}^{\text{M-A}}(t-1))^+ \underbrace{\mathbf{1}_{\{N_{a^*}(t-1) > f(t-1, K)\}}}_{\leq 1}\right] \\ & \leq \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*}^{\text{KL-A}}(t-1))^+ \mathbf{1}_{\{N_{a^*}(t-1) \leq f(t-1, K)\}}\right] + \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*}^{\text{M-A}}(t-1))^+\right]. \end{aligned}$$

Now, the bound (18) of Proposition 9, together with the Remark 2, indicates that

$$\sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*}^{\text{M-A}}(t-1))^+\right] \leq 26\sqrt{KT}.$$

Note that Remark 2 exactly explains that for the sum above we do not bump into the issues raised in Remark 3 for the other sum in (22).

*Part 3: Integrated deviations in terms of  $\mathcal{K}_{\text{inf}}$  divergence.* We showed so far that the distribution-free regret bound of the anytime version of KL-UCB-switch was given by the (intermediary) regret bound (18) of Proposition 9, which is smaller than  $(K-1) + 33\sqrt{KT}$ , plus

$$\begin{aligned} \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*}^{\text{KL-A}}(t-1))^+ \mathbf{1}_{\{N_{a^*}(t-1) \leq f(t-1, K)\}}\right] &= \sum_{t=K}^{T-1} \mathbb{E}\left[(\mu^* - U_{a^*}^{\text{KL-A}}(t))^+ \mathbf{1}_{\{N_{a^*}(t) \leq f(t, K)\}}\right] \\ &\leq \sum_{t=K}^{T-1} \sum_{n=1}^{f(t, K)} \mathbb{E}\left[(\mu^* - U_{a^*, t, n}^{\text{KL-A}})^+\right] \end{aligned} \quad (25)$$

where we applied optional skipping (Section 4.1) and where we denoted

$$U_{a^*, t, n}^{\text{KL-A}} = \sup\left\{\mu \in [0, 1] \mid \mathcal{K}_{\text{inf}}(\hat{v}_{a^*, n}, \mu) \leq \frac{1}{n} \varphi\left(\frac{t}{Kn}\right)\right\} \quad (26)$$

the counterpart of the quantity  $U_{a^*}^{\text{KL-A}}(t)$  defined in (10). Here, the additional subscript  $t$  in  $U_{a^*, t, n}^{\text{KL-A}}$  refers to the denominator of  $t/(Kn)$  in the  $\varphi(t/(Kn))$  term.

Now, Corollary 14 exactly indicates that for each given  $t$  and all  $n \geq 1$ ,

$$\mathbb{E}\left[(\mu^* - U_{a^*, t, n}^{\text{KL-A}})^+\right] \leq (2n+1) \sqrt{\frac{\pi}{n}} \exp\left(-\varphi\left(\frac{t}{Kn}\right)\right).$$

The  $t$  considered are such that  $t \geq K$  and thus,  $f(t, K) \leq (t/K)^{1/5} \leq t/K$ . Therefore, the considered  $n$  are such that  $1 \leq n \leq f(t, K)$  and thus,  $t/(Kn) \geq 1$ . Given that  $\varphi \geq \ln_+$ , we proved

$$\mathbb{E}\left[(\mu^* - U_{a^*, t, n}^{\text{KL-A}})^+\right] \leq (2n+1) \sqrt{\frac{\pi}{n}} \frac{Kn}{t} = \frac{K\sqrt{\pi}}{t} (2n+1)\sqrt{n}.$$

We sum this bound over  $n \in \{1, \dots, f(t/K)\}$ , using again that  $f(t, K) \leq (t/K)^{1/5}$ :

$$\sum_{n=1}^{f(t, K)} \mathbb{E}\left[(\mu^* - U_{a^*, t, n}^{\text{KL-A}})^+\right] \leq \frac{K\sqrt{\pi}}{t} \sum_{n=1}^{f(t, K)} \underbrace{(2n+1)\sqrt{n}}_{\leq 3f(t, K)^{3/2}} \leq \frac{3K\sqrt{\pi}}{t} \underbrace{f(t, K)^{5/2}}_{\leq (t/K)^{1/2}} \leq 3\sqrt{\pi} \sqrt{\frac{K}{t}}.$$

We substitute this inequality into (25):

$$\begin{aligned}
 \sum_{t=K+1}^T \mathbb{E} \left[ (\mu^* - U_{a^*}^{\text{KL-A}}(t-1))^+ \mathbf{1}_{\{N_{a^*}(t-1) \leq f(t-1, K)\}} \right] &\leq \sum_{t=K}^{T-1} \sum_{n=1}^{f(t, K)} \mathbb{E} \left[ (\mu^* - U_{a^*, t, n}^{\text{KL-A}})^+ \right] \\
 &\leq 3\sqrt{\pi} \underbrace{\sum_{t=K}^{T-1} \sqrt{\frac{K}{t}}}_{\leq 2\sqrt{KT}, \text{ see (40)}} \leq 6\sqrt{\pi} \sqrt{KT} \leq 11\sqrt{KT}.
 \end{aligned}$$

The final regret bound is obtained as the sum of this  $11\sqrt{KT}$  bound plus the  $(K-1) + 33\sqrt{KT}$  bound obtained above. This concludes the proof of Theorem 4.

*Part 4: Adaptations needed for Theorem 1*, i.e., to analyze the version of KL-UCB-switch relying on the knowledge of the horizon  $T$ . Parts 1 and 2 of the proof remain essentially unchanged, up to the (intermediary) regret bound to be applied now: (17) of Proposition 8, which is smaller than  $(K-1) + 17\sqrt{KT}$ . The additional regret bound, accounting, as we did in Part 3, for the use of KL-UCB-indexes for small  $T$ , is no larger than

$$\begin{aligned}
 \sum_{t=K}^{T-1} \sum_{n=1}^{f(T, K)} (2n+1) \sqrt{\frac{\pi}{n}} \exp\left(-\ln_+ \left(\frac{T}{Kn}\right)\right) &= \sum_{t=K}^{T-1} \sum_{n=1}^{f(T, K)} (2n+1) \sqrt{\frac{\pi}{n}} \frac{Kn}{T} = K\sqrt{\pi} \sum_{n=1}^{f(T, K)} \underbrace{(2n+1)\sqrt{n}}_{\leq 3f(T, K)^{3/2}} \\
 &\leq 3\sqrt{\pi} K f(T, K)^{5/2} \leq 3\sqrt{\pi} K \sqrt{\frac{T}{K}} \leq 6\sqrt{KT}.
 \end{aligned}$$

This yields the claimed  $(K-1) + 23\sqrt{KT}$  bound.

## 6. Proofs of the distribution-dependent bounds: Theorems 2 and 5

The proofs below can be adapted (simplified) to provide an elementary analysis of performance of the KL-UCB algorithm on the class of all distributions over a bounded interval, by keeping only Parts 1 and 2 of the proofs below. The study of KL-UCB in Cappé et al. [2013] remained somewhat intricate and limited to finitely supported distributions.

We provide first an anytime analysis, i.e., the proof of Theorem 5, and then explain the simplifications in the analysis (and improvements in the second-order terms in the regret bound) arising when the horizon  $T$  is known, i.e., as far as the proof of Theorem 2 is concerned.

### 6.1. Proof of Theorem 5

The proof starts as in Cappé et al. [2013]. We fix a sub-optimal arm  $a$ . Given  $\delta \in (0, \mu^*)$  sufficiently small (to be determined by the analysis), we first decompose  $\mathbb{E}[N_a(T)]$  as

$$\begin{aligned} \mathbb{E}[N_a(T)] &= 1 + \sum_{t=K}^{T-1} \mathbb{P}[A_{t+1} = a] \\ &= 1 + \sum_{t=K}^{T-1} \mathbb{P}[U_a^\Delta(t) < \mu^* - \delta \text{ and } A_{t+1} = a] + \sum_{t=K}^{T-1} \mathbb{P}[U_a^\Delta(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a]. \end{aligned}$$

We then use that by definition of the index policy,  $A_{t+1} = a$  only if  $U_a^\Delta(t) \geq U_{a^*}^\Delta(t)$ , where we recall that  $a^*$  denotes an optimal arm (i.e., an arm such that  $\mu_a = \mu^*$ ). We also use  $U_{a^*}^\Delta(t) \geq U_{a^*}^{\text{KL-}\Delta}(t)$ , which was stated in (14). We get

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq 1 + \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}^\Delta(t) < \mu^* - \delta \text{ and } A_{t+1} = a] + \sum_{t=K}^{T-1} \mathbb{P}[U_a^\Delta(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a] \\ &\leq 1 + \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}^{\text{KL-}\Delta}(t) < \mu^* - \delta] + \sum_{t=K}^{T-1} \mathbb{P}[U_a^\Delta(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a]. \end{aligned}$$

Finally, by the definition (13) of  $U_a^\Delta(t)$ , we proved so far

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq 1 + \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}^{\text{KL-}\Delta}(t) < \mu^* - \delta] \\ &\quad + \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{KL-}\Delta}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(t, K)] \\ &\quad + \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{M-}\Delta}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(t, K)]. \end{aligned} \quad (27)$$

We now deal with each of the three sums above.

*Part 1:* We first deal with the first sum in (27) and to that end, fix some  $t \in \{K, \dots, T-1\}$ . By the definition (10) of  $U_{a^*}^{\text{KL-}\Delta}(t)$  as a supremum,

$$\mathbb{P}[U_{a^*}^{\text{KL-}\Delta}(t) < \mu^* - \delta] \leq \mathbb{P}\left[\mathcal{K}_{\text{inf}}(\widehat{v}_{a^*}(t), \mu^* - \delta) > \frac{1}{N_{a^*}(t)} \varphi\left(\frac{t}{KN_{a^*}(t)}\right)\right].$$

By optional skipping (see Section 4.1), applied with some care,

$$\begin{aligned} \mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*}(t), \mu^* - \delta) > \frac{1}{N_{a^*}(t)} \varphi \left( \frac{t}{KN_{a^*}(t)} \right) \right] \\ \leq \mathbb{P} \left[ \exists n \in \{1, \dots, t - K + 1\} : \mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,n}, \mu^* - \delta) > \frac{1}{n} \varphi \left( \frac{t}{Kn} \right) \right]. \end{aligned}$$

Now, for  $n \geq \lfloor t/K \rfloor + 1$  and given the definition (9) of  $\varphi$ , we have  $\varphi(t/(Kn)) = 0$ . By definition,  $\mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,n}, \mu^* - \delta) > 0$  requires in particular that the expectation  $\widehat{\mu}_{a^*,n}$  of  $\widehat{\nu}_{a^*,n}$  be smaller than  $\mu^* - \delta$ . This fact, together with a union bound, implies

$$\begin{aligned} \mathbb{P} \left[ \exists n \in \{1, \dots, t - K + 1\} : \mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,n}, \mu^* - \delta) > \frac{1}{n} \varphi \left( \frac{t}{Kn} \right) \right] \\ \leq \mathbb{P} \left[ \exists n \geq \lfloor t/K \rfloor + 1 : \widehat{\mu}_{a^*,n} \leq \mu^* - \delta \right] + \sum_{n=1}^{\lfloor t/K \rfloor} \mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,n}, \mu^* - \delta) > \frac{1}{n} \varphi \left( \frac{t}{Kn} \right) \right]. \end{aligned}$$

Hoeffding's maximal inequality (Proposition 6) upper bounds the first term by  $\exp(-2\delta^2 t/K)$ , while Corollary 12 and Proposition 13 provide the upper bound

$$\mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,n}, \mu^* - \delta) > \frac{1}{n} \varphi \left( \frac{t}{Kn} \right) \right] \leq e(2n + 1) \exp \left( -n \left( 2\delta^2 + \varphi(t/(Kn))/n \right) \right).$$

Collecting all inequalities, we showed so far that

$$\mathbb{P} [U_{a^*}^{\text{KL-A}}(t) < \mu^* - \delta] \leq \exp(-2\delta^2 t/K) + \sum_{n=1}^{\lfloor t/K \rfloor} e(2n + 1) \exp \left( -2n\delta^2 - \varphi(t/(Kn)) \right).$$

Summing over  $t \in \{K, \dots, T - 1\}$ , using the formula for geometric series, on the one hand, and performing some straightforward (and uninteresting) calculation detailed below in Lemma 16 on the other hand, we finally bound the first sum in (27) by

$$\begin{aligned} \sum_{t=K}^{T-1} \mathbb{P} [U_{a^*}^{\text{KL-A}}(t) < \mu^* - \delta] &\leq \sum_{t=K}^{T-1} \exp(-2\delta^2 t/K) + \sum_{t=K}^{T-1} \sum_{n=1}^{\lfloor t/K \rfloor} e(2n + 1) \exp \left( -2n\delta^2 - \varphi(t/(Kn)) \right) \\ &\leq \frac{1}{1 - e^{-2\delta^2/K}} + \frac{e(3 + 8K)}{(1 - e^{-2\delta^2})^3}. \end{aligned}$$

This concludes the first part of this proof.

*Part 2: We then deal with the second sum in (27). We introduce*

$$\widetilde{U}_a^{\text{KL-A}}(t) \stackrel{\text{def}}{=} \sup \left\{ \mu \in [0, 1] \mid \mathcal{K}_{\text{inf}}(\widehat{\nu}_a(t), \mu) \leq \frac{1}{N_a(t)} \varphi \left( \frac{T}{KN_a(t)} \right) \right\}$$

that only differs from the original index  $U_a^{\text{KL-A}}(t)$  defined in (10) by the replacement of  $t/(Kn)$  by  $T/(Kn)$  as the argument of  $\varphi$ . Therefore, we have  $\widetilde{U}_a^{\text{KL-A}}(t) \geq U_a^{\text{KL-A}}(t)$ . Replacing also  $f(t, K)$  by the larger quantity  $f(T, K)$ , the second sum in (27) is therefore bounded by

$$\begin{aligned} &\sum_{t=K}^{T-1} \mathbb{P} [U_a^{\text{KL-A}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(t, K)] \\ &\leq \sum_{t=K}^{T-1} \mathbb{P} [\widetilde{U}_a^{\text{KL-A}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(T, K)] \\ &\leq \sum_{n=1}^{f(T, K)} \sum_{t=K}^{T-1} \mathbb{P} [\widetilde{U}_a^{\text{KL-A}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) = n]. \end{aligned} \tag{28}$$

Optional skipping (see Section 4.1) indicates that for each value of  $n$ ,

$$\begin{aligned} \sum_{t=K}^{T-1} \mathbb{P} \left[ \tilde{U}_a^{\text{KL-A}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) = n \right] \\ = \sum_{t=K}^{T-1} \mathbb{P} \left[ U_{a^*, T, n}^{\text{KL-A}} \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) = n \right] \end{aligned}$$

where  $U_{a^*, T, n}^{\text{KL-A}}$  was defined in (26). We now observe that the events  $\{A_{t+1} = a \text{ and } N_a(t) = n\}$  are disjoint as  $t$  varies in  $\{K, \dots, T-1\}$ . Therefore,

$$\sum_{t=K}^{T-1} \mathbb{P} \left[ U_{a^*, T, n}^{\text{KL-A}} \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) = n \right] \leq \mathbb{P} \left[ U_{a^*, T, n}^{\text{KL-A}} \geq \mu^* - \delta \right].$$

All in all, we proved so far that

$$\sum_{t=K}^{T-1} \mathbb{P} \left[ U_a^{\text{KL-A}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(t, K) \right] \leq \sum_{n=1}^{f(T, K)} \mathbb{P} \left[ U_{a^*, T, n}^{\text{KL-A}} \geq \mu^* - \delta \right]. \quad (29)$$

Now, note that the supremum in (26) is taken over a closed interval, as  $\mathcal{K}_{\text{inf}}$  is non-decreasing in its second argument (by its definition as an infimum) and as  $\mathcal{K}_{\text{inf}}$  is left-continuous (Corollary 11). This supremum is therefore a maximum. Hence, by distinguishing the cases where  $U_{a^*, T, n}^{\text{KL-A}} = \mu^* - \delta$  and  $U_{a^*, T, n}^{\text{KL-A}} > \mu^* - \delta$ , we have the equality of events

$$\left\{ U_{a^*, T, n}^{\text{KL-A}} \geq \mu^* - \delta \right\} = \left\{ \mathcal{K}_{\text{inf}}(\hat{\nu}_{a, n}, \mu^* - \delta) \leq \frac{1}{n} \varphi \left( \frac{T}{Kn} \right) \right\}.$$

We assume that  $\delta \in (0, \mu^*)$  is sufficiently small for

$$\delta < \frac{1 - \mu^*}{2} \mathcal{K}_{\text{inf}}(\nu_a, \mu^*)$$

to hold and introduce

$$n_1 = \left\lceil \frac{\varphi(T/K)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - 2\delta/(1 - \mu^*)} \right\rceil \geq 1.$$

For  $n \geq n_1$ , by definition of  $n_1$ ,

$$\frac{1}{n} \varphi \left( \frac{T}{Kn} \right) \leq \underbrace{\frac{\varphi(T/(Kn))}{\varphi(T/K)}}_{\leq 1} \left( \mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \frac{2\delta}{1 - \mu^*} \right) \leq \mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \frac{2\delta}{1 - \mu^*}$$

while by the regularity property (19), we have  $\mathcal{K}_{\text{inf}}(\hat{\nu}_{a, n}, \mu^* - \delta) \geq \mathcal{K}_{\text{inf}}(\hat{\nu}_{a, n}, \mu^*) - \frac{\delta}{1 - \mu^*}$ . We therefore proved that for  $n \geq n_1$ ,

$$\begin{aligned} \mathbb{P} \left[ U_{a^*, T, n}^{\text{KL-A}} \geq \mu^* - \delta \right] &= \mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\hat{\nu}_{a, n}, \mu^* - \delta) \leq \frac{1}{n} \varphi \left( \frac{T}{Kn} \right) \right] \\ &\leq \mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\hat{\nu}_{a, n}, \mu^*) \leq \mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \frac{\delta}{1 - \mu^*} \right]. \end{aligned}$$

Therefore we may resort to the concentration inequality on  $\mathcal{K}_{\text{inf}}$  stated as Proposition 15. We set  $x = \mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \delta/(1 - \mu^*)$  and simply sum the bounds obtained in the two regimes considered therein:

$$\mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\hat{\nu}_{a, n}, \mu^* - \delta) \leq \mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \frac{\delta}{1 - \mu^*} \right] \leq e^{-n/4} + \exp \left( - \frac{n\delta^2}{2\gamma_*(1 - \mu^*)^2} \right)$$

where  $\gamma_\star$  was defined in (21). For  $n \leq n_1 - 1$ , we bound the probability at hand by 1. Combining all these arguments together yields

$$\begin{aligned} \sum_{n=1}^{f(T,K)} \mathbb{P}\left[U_{a^\star, T, n}^{\text{KL-A}} \geq \mu^\star - \delta\right] &\leq n_1 - 1 + \sum_{n=n_1}^{f(T,K)} e^{-n/4} + \sum_{n=n_1}^{f(T,K)} \exp\left(-\frac{n\delta^2}{2\gamma_\star(1-\mu^\star)^2}\right) \\ &\leq \frac{\varphi(T/K)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^\star) - 2\delta/(1-\mu^\star)} + \underbrace{\frac{1}{1-e^{-1/4}}}_{\leq 5} + \underbrace{\frac{1}{1-e^{-\delta^2/(2\gamma_\star(1-\mu^\star)^2)}}}_{=\mathcal{O}(1/\delta^2)} \end{aligned}$$

where the second inequality follows from the formula for geometric series and from the definition of  $n_1$ .

*Part 3:* We then deal with the third sum in (27). This sum involves the indexes  $U_a^{\text{M-A}}(t)$  only when  $N_a(t) > f(t, K)$ , that is, when  $N_a(t) \geq f(t, K) + 1$ , where  $f(t, K) = \lfloor (t/K)^{1/5} \rfloor$ . Under the latter condition, the indexes are actually bounded by

$$U_a^{\text{M-A}}(t) \stackrel{\text{def}}{=} \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \varphi\left(\frac{t}{KN_a(t)}\right)} \leq \hat{\mu}_a(t) + \underbrace{\sqrt{\frac{1}{2(t/K)^{1/5}} \varphi((t/K)^{4/5})}}_{\rightarrow 0 \text{ as } t \rightarrow \infty}.$$

We denote by  $T_0(\Delta_a, K)$  the smallest time  $T_0$  such that for all  $t \geq T_0$ ,

$$\sqrt{\frac{1}{2(t/K)^{1/5}} \varphi((t/K)^{4/5})} \leq \frac{\Delta_a}{4}. \quad (30)$$

This time  $T_0$  only depends on  $K$  and  $\Delta_a$ ; a closed-form upper bound on its value could be easily provided. With this definition, we already have that the sum of interest may be bounded by

$$\begin{aligned} &\sum_{t=K}^{T-1} \mathbb{P}\left[U_a^{\text{M-A}}(t) \geq \mu^\star - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(t, K)\right] \\ &\leq T_0(\Delta_a, K) + \sum_{t=T_0(\Delta_a, K)}^{T-1} \mathbb{P}\left[\hat{\mu}_a(t) + \Delta_a/4 \geq \mu^\star - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(t, K)\right] \\ &\leq T_0(\Delta_a, K) + \sum_{t=T_0(\Delta_a, K)}^{T-1} \mathbb{P}\left[\hat{\mu}_a(t) \geq \mu_a + \Delta_a/2 \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(t, K)\right] \end{aligned}$$

where for the second inequality, we assumed that  $\delta \in (0, \mu^\star)$  is sufficiently small for

$$\delta < \frac{\Delta_a}{4}$$

to hold. Optional skipping (see Section 4.1), using that the events  $\{A_{t+1} = a \text{ and } N_a(t) = n\}$  are disjoint as  $t$  varies, as already done between (28) and (29), provides the upper bound

$$\begin{aligned} &\sum_{t=T_0(\Delta_a, K)}^{T-1} \mathbb{P}\left[\hat{\mu}_a(t) \geq \mu_a + \Delta_a/2 \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(t, K)\right] \\ &\leq \sum_{n \geq 1} \mathbb{P}\left[\hat{\mu}_{a,n} \geq \mu_a + \Delta_a/2\right] \leq \sum_{n \geq 1} e^{-n\Delta_a^2/2} = \frac{1}{1-e^{-\Delta_a^2/2}} \end{aligned}$$

where the second inequality is due to Hoeffding's inequality (in its non-maximal version, see Proposition 6). A summary of the bound thus provided in this part is:

$$\sum_{t=K}^{T-1} \mathbb{P}\left[U_a^{\text{M-A}}(t) \geq \mu^\star - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(t, K)\right] \leq T_0(\Delta_a, K) + \frac{1}{1-e^{-\Delta_a^2/2}} = \mathcal{O}(1)$$

where  $T_0(\Delta_a, K)$  was defined in (30).

*Part 4: Conclusion of the proof of Theorem 5.* Collecting all previous bounds and conditions, we proved that when  $\delta \in (0, \mu^*)$  is sufficiently small for

$$\delta < \min \left\{ \frac{1 - \mu^*}{2} \mathcal{K}_{\text{inf}}(\nu_a, \mu^*), \frac{\Delta_a}{4} \right\} \quad (31)$$

to hold, then

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq \frac{\varphi(T/K)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - 2\delta/(1 - \mu^*)} \\ &+ \underbrace{\frac{e(3 + 8K)}{(1 - e^{-2\delta^2})^3}}_{=\mathcal{O}(1/\delta^6)} + \underbrace{\frac{1}{1 - e^{-2\delta^2/K}} + \frac{1}{1 - e^{-\delta^2/(2\gamma_*(1 - \mu^*)^2)}}}_{=\mathcal{O}(1/\delta^2)} + \underbrace{T_0(\Delta_a, K) + \frac{1}{1 - e^{-\Delta_a^2/2}} + 6}_{=\mathcal{O}(1)} \end{aligned} \quad (32)$$

where

$$\frac{\varphi(T/K)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - 2\delta/(1 - \mu^*)} = \frac{\ln T + \ln \ln T + \mathcal{O}(1)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - 2\delta/(1 - \mu^*)} = \frac{\ln T + \ln \ln T}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} + \mathcal{O}(\delta \ln T).$$

The leading term in this regret bound is  $\ln T / \mathcal{K}_{\text{inf}}(\nu_a, \mu^*)$ , while the order of magnitude of the smaller-order terms is given by

$$\delta \ln T + \frac{1}{\delta^6} = \mathcal{O}((\ln T)^{6/7})$$

for  $\delta$  of the order of  $(\ln T)^{-1/7}$ . When  $T$  is sufficiently large, this value of  $\delta$  is smaller than the required threshold (31).

It only remains to state and prove Lemma 16 (used at the very end of the first part of the proof above).

**Lemma 16.** *We have the bound*

$$\sum_{t=K}^{T-1} \sum_{n=1}^{\lfloor t/K \rfloor} e(2n+1) \exp\left(-2n\delta^2 - \varphi(t/(Kn))\right) \leq \frac{e(3+8K)}{(1 - e^{-2\delta^2})^3}.$$

**Proof:** The double sum can be rewritten, by permuting the order of summations, as

$$\begin{aligned} \sum_{t=K}^{T-1} \sum_{n=1}^{\lfloor t/K \rfloor} e(2n+1) \exp\left(-2n\delta^2 - \varphi(t/(Kn))\right) &= \sum_{n=1}^{\lfloor T/K \rfloor} \sum_{t=Kn}^{T-1} e(2n+1) \exp\left(-2n\delta^2 - \varphi(t/(Kn))\right) \\ &= \sum_{n=1}^{\lfloor T/K \rfloor} e(2n+1) \exp(-2n\delta^2) \sum_{t=Kn}^{T-1} \exp\left(-\varphi(t/(Kn))\right). \end{aligned}$$

We first fix  $n \geq 1$  and use that  $t \mapsto \exp(-\varphi(t/(Kn)))$  is non-increasing to get

$$\sum_{t=Kn}^{T-1} \exp\left(-\varphi(t/(Kn))\right) \leq 1 + \int_{Kn}^{T-1} \exp\left(-\varphi(t/(Kn))\right) dt = 1 + Kn \int_1^{(T-1)/(Kn)} \exp(-\varphi(u)) du$$

where we operated the change of variable  $u = t/(Kn)$ . Now, by the change of variable  $v = \ln(u)$ ,

$$\begin{aligned} \int_1^{(T-1)/(Kn)} \exp(-\varphi(u)) du &\leq \int_1^{+\infty} \exp(-\varphi(u)) du = \int_1^{+\infty} \frac{1}{u(1 + \ln^2(u))} du \\ &= \int_0^{+\infty} \frac{1}{1 + v^2} dv = [\arctan]_0^{+\infty} = \frac{\pi}{2}. \end{aligned}$$

All in all, we proved so far that

$$\begin{aligned} \sum_{t=K}^{T-1} \sum_{n=1}^{\lfloor t/K \rfloor} e(2n+1) \exp\left(-2n\delta^2 - \varphi(t/(Kn))\right) &\leq \sum_{n=1}^{\lfloor T/K \rfloor} e(2n+1)(1+Kn\pi/2) \exp(-2n\delta^2) \\ &\leq \sum_{n=1}^{+\infty} e(1+(2+K\pi/2)n+K\pi n^2) \exp(-2n\delta^2). \end{aligned}$$

To conclude our calculation, we use that by differentiation of series, for all  $\theta > 0$ ,

$$\begin{aligned} \sum_{m=0}^{+\infty} e^{-m\theta} &= \frac{1}{1-e^{-\theta}} \\ -\sum_{m=1}^{+\infty} m e^{-m\theta} &= \frac{-e^{-\theta}}{(1-e^{-\theta})^2} \quad \text{thus} \quad \sum_{m=1}^{+\infty} m e^{-m\theta} \leq \frac{1}{(1-e^{-\theta})^2} \end{aligned} \quad (33)$$

$$\sum_{m=1}^{+\infty} m^2 e^{-m\theta} = \frac{e^{-\theta}(1+e^{-\theta})}{(1-e^{-\theta})^3} \leq \frac{2}{(1-e^{-\theta})^3}. \quad (34)$$

Hence, taking  $\theta = 2\delta^2$ ,

$$\sum_{n=1}^{+\infty} e(1+(2+K\pi/2)n+K\pi n^2) \exp(-2n\delta^2) \leq \frac{e}{1-e^{-2\delta^2}} + \frac{e(2+K\pi/2)}{(1-e^{-2\delta^2})^2} + \frac{2eK\pi}{(1-e^{-2\delta^2})^3} \leq \frac{e(3+8K)}{(1-e^{-2\delta^2})^3}$$

which concludes the proof of this lemma.  $\square$

## 6.2. Proof of Theorem 2

We adapt (simplify) the proof of Theorem 5, by replacing the thresholds  $f(t, K)$  by  $f(T, K)$ , by taking  $\varphi = \ln_+$ , etc. To that end, we start with a similar decomposition,

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq 1 + \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}^{\text{KL}}(t) < \mu^* - \delta] \\ &\quad + \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}^{\text{KL}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(T, K)] \\ &\quad + \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{M}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(T, K)]. \end{aligned} \quad (35)$$

The first sum is bounded using exactly the same arguments as in the proof of Theorem 5 (optional skipping, Hoeffding's maximal inequality, Corollary 12 and Proposition 13): for all  $t \in \{K, \dots, T-1\}$ ,

$$\begin{aligned} &\mathbb{P}[U_{a^*}^{\text{KL}}(t) < \mu^* - \delta] \\ &\leq \mathbb{P}[\exists n \geq \lfloor T/K \rfloor + 1 : \hat{\mu}_{a^*, n} \leq \mu^* - \delta] + \sum_{n=1}^{\lfloor T/K \rfloor} \mathbb{P}\left[\mathcal{K}_{\text{inf}}(\hat{\nu}_{a^*, n}, \mu^* - \delta) > \frac{1}{n} \ln\left(\frac{T}{Kn}\right)\right] \\ &\leq \exp(-2\delta^2 T/K) + e(2n+1) \exp\left(-n(2\delta^2 + \ln(T/(Kn))/n)\right) \\ &= \exp(-2\delta^2 T/K) + \frac{eK}{T} (2n^2 + n) \exp(-2n\delta^2). \end{aligned}$$



Summing over  $t$  and substituting the bounds (33)–(34), we proved

$$\begin{aligned} \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}^{\text{KL}}(t) < \mu^* - \delta] &\leq T \exp(-2\delta^2 T/K) + \frac{eK}{T} \left( \frac{4}{(1 - e^{-2\delta^2})^3} + \frac{1}{(1 - e^{-2\delta^2})^2} \right) \\ &\leq T \exp(-2\delta^2 T/K) + \frac{5eK}{T(1 - e^{-2\delta^2})^3} \end{aligned}$$

For the second sum in (35), we note that the initial manipulations in Part 2 of the proof of Theorem 5 are unnecessary in the case of Theorem 2; we may directly start at (28) and the rest of the arguments used and calculation performed then hold word for word, under the same condition  $\delta < (1 - \mu^*)\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)/2$ . We get, with the same notation,

$$\begin{aligned} \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{KL}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(T, K)] \\ \leq \frac{\ln(T/K)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - 2\delta/(1 - \mu^*)} + 5 + \frac{1}{1 - e^{-\delta^2/(2\gamma_*(1 - \mu^*)^2)}}. \quad (36) \end{aligned}$$

The third sum in (35) involves the indexes  $U_a^{\text{M}}(t)$  only under the condition  $N_a(t) > f(T, K)$ , in which case  $N_a(t) \geq (T/K)^{1/5}$  and

$$U_a^{\text{M}}(t) \stackrel{\text{def}}{=} \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \ln_+ \left( \frac{T}{KN_a(t)} \right)} \leq \hat{\mu}_a(t) + \sqrt{\frac{1}{2(T/K)^{1/5}} \ln_+ ((T/K)^{4/5})}.$$

We mimic the proof scheme of Part 3 of the proof of Theorem 5 and start by assuming that  $T$  is sufficiently large for

$$\sqrt{\frac{1}{2(T/K)^{1/5}} \ln_+ ((T/K)^{4/5})} \leq \frac{\Delta_a}{4} \quad (37)$$

to hold. Under the same condition  $\delta < \Delta_a/4$ , we get, by a careful application of optional skipping using that the events  $\{A_{t+1} = a \text{ and } N_a(t) = n\}$  are disjoint as  $t$  varies and by Hoeffding's inequality,

$$\begin{aligned} \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{M}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) > f(T, K)] \\ \leq \sum_{n=f(T,K)+1}^{T-1} \mathbb{P}[\hat{\mu}_{a,n} \geq \mu_a + \Delta_a/2] \leq \sum_{n \geq f(T,K)+1} e^{-n\Delta_a^2/2} \leq \frac{1}{1 - e^{-\Delta_a^2/2}}. \quad (38) \end{aligned}$$

Collecting all bounds, we proved that whenever  $T$  is sufficiently large for (37) to hold and whenever  $\delta$  is sufficiently small for (31) to hold,

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq \frac{\ln(T/K)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - 2\delta/(1 - \mu^*)} \\ &\quad + \underbrace{\frac{1}{1 - e^{-\delta^2/(2\gamma_*(1 - \mu^*)^2)}}}_{=\mathcal{O}(1/\delta^2)} + \underbrace{\frac{1}{1 - e^{-\Delta_a^2/2}}}_{=\mathcal{O}(1)} + 6 + T \exp(-2\delta^2 T/K) + \underbrace{\frac{5eK}{T(1 - e^{-2\delta^2})^3}}_{=\mathcal{O}(1/(T\delta^6))}. \quad (39) \end{aligned}$$

The leading term in this regret bound is  $\ln T/\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)$ , while the order of magnitude of the smaller-order terms is given by

$$\delta \ln T + \frac{1}{\delta^2} + T \exp(-2\delta^2 T/K) + \frac{1}{T\delta^6} = \mathcal{O}((\ln T)^{2/3})$$

for  $\delta$  of the order of  $(\ln T)^{-1/3}$ . When  $T$  is sufficiently large, this value of  $\delta$  is smaller than the required threshold (31).

## Acknowledgements

This work was partially supported by the CIMI (Centre International de Mathématiques et d’Informatique) Excellence program. The authors acknowledge the support of the French Agence Nationale de la Recherche (ANR), under grants ANR-13-BS01-0005 (project SPADRO) and ANR-13-CORD-0020 (project ALICIA). Aurélien Garivier acknowledges the support of the Project IDEXLYON of the University of Lyon, in the framework of the Programme Investissements d’Avenir (ANR-16-IDEX-0005).

## References

- R. Agrawal. Sample mean based index policies with  $o(\ln n)$  regret for the multi-armed bandit problem. *Advances in Applied Probability*, 27(4):1054–1078, 1995.
- J.-Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In *Proceedings of the 22nd Annual Conference on Learning Theory, COLT’09*, pages 217–226, 2009.
- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002a.
- P. Auer, N. Cesa-Bianchi, Y. Freund, and R.E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002b.
- S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities. A nonasymptotic theory of independence*. Oxford University Press, 2013.
- S. Bubeck and C.-Y. Liu. Prior-free and prior-dependent regret bounds for Thompson sampling. arXiv:1304.5758, 2013.
- A.N. Burnetas and M.N. Katehakis. Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*, 17(2):122–142, 1996.
- O. Cappé, A. Garivier, O.-A. Maillard, R. Munos, and G. Stoltz. Kullback–Leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516–1541, 2013.
- Y. Chow and H. Teicher. *Probability Theory*. Springer, 1988.
- R. Degenne and V. Perchet. Anytime optimal algorithms in stochastic multi-armed bandits. In *Proceedings of the 2016 International Conference on Machine Learning, ICML’16*, pages 1587–1595, 2016.
- J.L. Doob. *Stochastic Processes*. Wiley Publications in Statistics. John Wiley & Sons, 1953.
- A. Garivier and O. Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *Proceedings of the 24th Annual Conference on Learning Theory, COLT’11*, 2011.
- A. Garivier, P. Ménard, and G. Stoltz. Explore first, exploit next: The true shape of regret in bandit problems. *Mathematics of Operations Research*, 2018. To appear; meanwhile, see arXiv preprint arXiv:1602.07182.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- J. Honda and A. Takemura. Non-asymptotic analysis of a new bandit algorithm for semi-bounded rewards. *Journal of Machine Learning Research*, 16:3721–3756, 2015.

- A. Hoorfar and M. Hassani. Inequalities on the Lambert  $W$  function and hyperpower function. *Journal of Inequalities in Pure and Applied Mathematics*, 9(2):Article 51, 2008.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On bayesian upper confidence bounds for bandit problems. In *Proceedings of the 2012 International Conference on Artificial Intelligence and Statistics*, AISTats’12, pages 592–600, 2012.
- N. Korda, E. Kaufmann, and R. Munos. Thompson sampling for 1–dimensional exponential family bandits. In *Proceedings of the 26th International Conference on Neural Information Processing Systems*, NIPS’13, pages 1448–1456, 2013.
- T.L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- T. Lattimore. Regret analysis of the anytime optimally confident UCB algorithm. arXiv:1603.08661, 2016.
- O.-A. Maillard, R. Munos, and G. Stoltz. Finite-time analysis of multi-armed bandits problems with Kullback-Leibler divergences. In *Proceedings of the 24th annual Conference on Learning Theory*, COLT’11, 2011.
- P. Ménard and A. Garivier. A minimax and asymptotically optimal algorithm for stochastic bandits. In *Proceedings of the 2017 Algorithmic Learning Theory Conference*, ALT’17, 2017.
- G. Simons, L. Yang, and Y.-C. Yao. Doob, Ignatov and optional skipping. *Annals of Probability*, 30(4):1933–1958, 2002.
- W.R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25:285–294, 1933.

# Appendix for the article

“KL-UCB-switch: optimal regret bounds for stochastic bandits  
from both a distribution-dependent and a distribution-free viewpoints”

by Aurélien Garivier, Hédi Hadiji, Pierre Ménard, Gilles Stoltz

---

## Content and structure

- A. A simplified proof of the regret bounds for MOSS(-anytime)
- B. Proofs of the regularity and deviation/concentration results on  $\mathcal{K}_{\text{inf}}$ 
  - 1. Proof of the regularity lemma (Lemma 10)
  - 2. A useful tool: a variational formula for  $\mathcal{K}_{\text{inf}}$  (statement)
  - 3. Proof of the deviation result (Proposition 13)
  - 4. Proof of the concentration result (Proposition 15)
- C. Proof of Theorem 3 (with the  $-\ln \ln T$  term in the regret bound)
- D. Proof of the variational formula (Lemma 18)

## A. A simplified proof of the regret bounds for MOSS(-anytime)

This section provides the proofs of Propositions 8 and 9. To emphasize the similarity of the analyses in the anytime and non-anytime cases, we present both of them in a unified fashion. The indexes used only differ by the replacement of  $T$  by  $t$  in the logarithmic exploration term in case  $T$  is unknown, see (5) and (11), which we both state with a generic exploration function  $\varphi$ . Indeed, compare

$$U_a^M(t) = \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \varphi\left(\frac{T}{KN_a(t)}\right)} \quad \text{and} \quad U_a^{M-A}(t) = \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \varphi\left(\frac{t}{KN_a(t)}\right)}.$$

We will denote by

$$U_{a,\tau}^{\text{GM}}(t) = \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \varphi\left(\frac{\tau}{KN_a(t)}\right)}$$

the index of the generic MOSS (GM) strategy, so that  $U_a^M(t) = U_{a,T}^{\text{GM}}(t)$  and  $U_a^{M-A}(t) = U_{a,t}^{\text{GM}}(t)$ . This GM strategy considers a sequence  $(\tau_K, \dots, \tau_{T-1})$  of integers, either  $\tau_t \equiv T$  for MOSS or  $\tau_t = t$  for MOSS-anytime, and picks at each step  $t+1$  with  $t \geq K$ , an arm  $A_{t+1}^{\text{GM}}$  with maximal index  $U_{a,\tau_t}^{\text{GM}}(t)$ . For a given  $t$ , we denote by  $U_{a,\tau_t,n}^{\text{GM}}$  the quantities corresponding to  $U_{a,\tau_t}^{\text{GM}}(t)$  by optional skipping (see Section 4.1).

We provide below an analysis for increasing exploration functions  $\varphi : (0, +\infty) \rightarrow [0, +\infty)$  such that  $\varphi$  vanishes on  $(0, 1]$  and  $\varphi \geq \ln_+$ , properties that are all satisfied for the two exploration functions stated in Proposition 9. The general result is stated as the next proposition.

**Proposition 17.** *For all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , for all  $T \geq 1$  and all sequences  $(\tau_K, \dots, \tau_{T-1})$  bounded by  $T$ , the regret of the generic MOSS strategy described above, with an increasing exploration function  $\varphi \geq \ln_+$  vanishing on  $(0, 1]$ , is smaller than*

$$R_T \leq (K-1) + \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*,\tau_{t-1}}^{\text{GM}}(t-1))^+\right] + \sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E}\left[(U_{a,T,n}^{\text{GM}} - \mu_a - \sqrt{K/T})^+\right]$$

where

$$U_{a,T,n}^{\text{GM}} = \hat{\mu}_{a,n} + \sqrt{\frac{1}{2n} \varphi\left(\frac{T}{Kn}\right)}.$$

In addition,

$$\sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*,\tau_{t-1}}^{\text{GM}}(t-1))^+\right] \leq \underbrace{20\sqrt{\frac{\pi}{8}}}_{\leq 12.6} \sum_{t=K}^{T-1} \sqrt{\frac{K}{\tau_t}}$$

and

$$\sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E}\left[(U_{a,T,n}^{\text{GM}} - \mu_a - \sqrt{K/T})^+\right] \leq \sqrt{KT} \left(1 + \frac{\pi}{4} + \frac{1}{\sqrt{2}} \int_1^{+\infty} u^{-3/2} \sqrt{\varphi(u)} du\right).$$

The bounds of Propositions 8 and 9, including the intermediary bounds (17) and (18), follow from this general result, up to the following straightforward calculation. On the one hand, in the known horizon case  $\sum 1/\sqrt{\tau_t} \leq T/\sqrt{T} = \sqrt{T}$ , whereas in the anytime case,

$$\sum_{t=K}^{T-1} 1/\sqrt{\tau_t} = \sum_{t=K}^{T-1} 1/\sqrt{t} \leq \int_0^T \frac{1}{\sqrt{u}} du = 2\sqrt{T}. \quad (40)$$

On the other hand, by the change of variable  $u = e^{v^2}$ ,

$$\int_1^{+\infty} u^{-3/2} \sqrt{\ln(u)} du = 2 \int_0^{+\infty} v^2 e^{-v^2/2} dv = \sqrt{2\pi}$$

and, using well-known inequalities like  $\sqrt{x+x'} \leq \sqrt{x} + \sqrt{x'}$  and  $\ln(1+x) \leq x$  for  $x, x' \geq 0$ ,

$$\begin{aligned} \int_1^{+\infty} \sqrt{u^{-3} \ln(u(1+\ln^2(u)))} du &\leq \int_1^{+\infty} \sqrt{u^{-3} \ln(u)} du + \int_1^{+\infty} \sqrt{u^{-3} \ln(1+\ln^2(u))} du \\ &\leq \int_1^{+\infty} \sqrt{u^{-3} \ln(u)} du + \int_1^{+\infty} \sqrt{u^{-3} \ln^2(u)} du \\ &= 2 \int_0^{+\infty} v^2 e^{-v^2/2} dv + 2 \int_0^{+\infty} v^3 e^{-v^2/2} dv = \sqrt{2\pi} + 4. \end{aligned}$$

The constant 17 of Proposition 8 is obtained as an upper bound on the sum of  $12.6 \leq 13$  and  $1 + \pi/4 + \sqrt{\pi} \leq 3.6 \leq 4$ . The constants 30 and 33 of Proposition 9 are respectively obtained as upper bounds on the sum of  $2 \times 12.6 \leq 26$  and  $1 + \pi/4 + \sqrt{\pi} \leq 4$ , and on the sum of  $2 \times 12.6 \leq 26$  and  $1 + \pi/4 + \sqrt{\pi} + 4/\sqrt{2} \leq 6.4 \leq 7$ .

**Proof:** The beginning of this proof is completely similar to the beginning of the proof provided in Section 5.

The first step is standard, see Bubeck and Liu [2013]. By definition of the index policy, for  $t \geq K$ ,

$$U_{a^*, \tau_t}^{\text{GM}}(t) \leq \max_{a=1, \dots, K} U_{a, \tau_t}^{\text{GM}}(t) = U_{A_{t+1}^{\text{GM}}, \tau_t}^{\text{GM}}(t)$$

so that the regret of the strategy is smaller than

$$R_T = \sum_{t=1}^T \mathbb{E}[\mu^* - \mu_{A_t^{\text{GM}}}] \leq (K-1) + \sum_{t=K+1}^T \mathbb{E}[\mu^* - U_{a^*, \tau_{t-1}}^{\text{GM}}(t-1)] + \sum_{t=K+1}^T \mathbb{E}[U_{A_t^{\text{GM}}, \tau_{t-1}}^{\text{GM}}(t-1) - \mu_{A_t^{\text{GM}}}] . \quad (41)$$

The term  $K-1$  above accounts for the initial  $K$  rounds, when each arm is played once.

*A preliminary transformation of the right-hand side of (41).* We successively use the fact that the index  $U_{a, \tau}^{\text{GM}}(t-1)$  increases with  $\tau$  since  $\varphi$  is increasing (for the first inequality below),  $x \leq \delta + (x - \delta)^+$  for all  $x$  and  $\delta$  (for the second inequality), and optional skipping (Section 4.1, for the third inequality), keeping in mind that pairs  $(a, n)$  such  $A_t^{\text{GM}} = a$  and  $N_a(t-1) = n$  correspond to at most one round  $t \in \{K+1, \dots, T\}$ :

$$\begin{aligned} \sum_{t=K+1}^T \mathbb{E}[U_{A_t^{\text{GM}}, \tau_{t-1}}^{\text{GM}}(t-1) - \mu_{A_t^{\text{GM}}}] &\leq \sum_{t=K+1}^T \mathbb{E}[U_{A_t^{\text{GM}}, T}^{\text{GM}}(t-1) - \mu_{A_t^{\text{GM}}}] \\ &\leq \sqrt{KT} + \sum_{t=K+1}^T \mathbb{E}\left[\left(U_{A_t^{\text{GM}}, T}^{\text{GM}}(t-1) - \mu_{A_t^{\text{GM}}} - \sqrt{\frac{K}{T}}\right)^+\right] \\ &\leq \sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E}\left[\left(U_{a, T, n}^{\text{GM}} - \mu_a - \sqrt{\frac{K}{T}}\right)^+\right]. \end{aligned}$$

While the last two inequalities may seem very crude, it turns out they are sharp enough to obtain the claimed distribution-free bounds. Moreover, they get rid of the bothersome dependencies among the arms that are contained in the choice of the arms  $A_t^{\text{GM}}$ . Therefore, we have shown that the right-hand side of (41) is bounded by

$$\begin{aligned} &(K-1) + \sum_{t=K+1}^T \mathbb{E}[\mu^* - U_{a^*, \tau_{t-1}}^{\text{GM}}(t-1)] + \sum_{t=K+1}^T \mathbb{E}[U_{A_t^{\text{GM}}, \tau_{t-1}}^{\text{GM}}(t-1) - \mu_{A_t^{\text{GM}}}] \\ &\leq (K-1) + \sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*, \tau_{t-1}}^{\text{GM}}(t-1))^+\right] + \sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E}\left[(U_{a, T, n}^{\text{GM}} - \mu_a - \sqrt{K/T})^+\right]. \quad (42) \end{aligned}$$

This inequality actually holds for all choices of sequences  $(\tau_t)_{K \leq t \leq T-1}$  with  $\tau_t \leq T$ . The first sum in the right-hand side of (42) depends on the specific value of  $(\tau_t)_{K \leq t \leq T-1}$ , and thus, on the specific MOSS algorithm considered, but the second sum only depends on  $T$ .

*Control of the left deviations of the best arm*, that is, of the first sum in (41) and (42). For each given round  $t \in \{K, \dots, T-1\}$ , we decompose

$$\mathbb{E} \left[ (\mu^* - U_{a^*, \tau_t}^{\text{GM}}(t))^+ \right] = \mathbb{E} \left[ (\mu^* - U_{a^*, \tau_t}^{\text{GM}}(t))^+ \mathbf{1}_{\{N_{a^*}(t) < \tau_t/K\}} \right] + \mathbb{E} \left[ (\mu^* - U_{a^*, \tau_t}^{\text{GM}}(t))^+ \mathbf{1}_{\{N_{a^*}(t) \geq \tau_t/K\}} \right].$$

The two pieces are handled differently. The second one is dealt with first by using  $U_{a^*, \tau_t}^{\text{GM}}(t) \geq \hat{\mu}_{a^*}(t)$ , which actually holds with equality given  $N_{a^*}(t) \geq \tau_t/K$ , and second, by optional skipping (Section 4.1) and by the integrated version of Hoeffding's inequality (Corollary 7):

$$\begin{aligned} \mathbb{E} \left[ (\mu^* - U_{a^*, \tau_t}^{\text{GM}}(t))^+ \mathbf{1}_{\{N_{a^*}(t) \geq \tau_t/K\}} \right] &\leq \mathbb{E} \left[ (\mu^* - \hat{\mu}_{a^*}(t))^+ \mathbf{1}_{\{N_{a^*}(t) \geq \tau_t/K\}} \right] \\ &\leq \mathbb{E} \left[ \max_{n \geq \tau_t/K} (\mu^* - \hat{\mu}_{a^*, n})^+ \right] \leq \sqrt{\frac{\pi}{8}} \sqrt{\frac{K}{\tau_t}}. \end{aligned} \quad (43)$$

When the arm has not been pulled often enough, we resort to a “peeling trick”. We consider a real number  $\beta > 1$  and further decompose the event  $\{N_{a^*}(t) < \tau_t/K\}$  along the geometric grid  $x_\ell = \beta^{-\ell} \tau_t/K$ , where  $\ell = 0, 1, 2, \dots$  (the endpoints  $x_\ell$  are not necessarily integers, and some intervals  $[x_{\ell+1}, x_\ell]$  may contain no integer, but none of these facts is an issue):

$$\begin{aligned} \mathbb{E} \left[ (\mu^* - U_{a^*, \tau_t}^{\text{GM}}(t))^+ \mathbf{1}_{\{N_{a^*}(t) < \tau_t/K\}} \right] &= \sum_{\ell=0}^{+\infty} \mathbb{E} \left[ (\mu^* - U_{a^*, \tau_t}^{\text{GM}}(t))^+ \mathbf{1}_{\{x_{\ell+1} \leq N_{a^*}(t) < x_\ell\}} \right] \\ &\leq \sum_{\ell=0}^{+\infty} \mathbb{E} \left[ \max_{x_{\ell+1} \leq n < x_\ell} (\mu^* - U_{a^*, \tau_t, n}^{\text{GM}})^+ \right] \end{aligned}$$

where in the second inequality, we applied optional skipping (Section 4.1) once again. Now for any  $\ell$ , the summand can be controlled as follows, first, by  $\varphi \geq \ln_+ = \ln$  on  $[1, +\infty)$ , second, by using  $n < x_\ell$  and third, by Corollary 7:

$$\begin{aligned} \mathbb{E} \left[ \max_{x_{\ell+1} \leq n < x_\ell} (\mu^* - U_{a^*, \tau_t, n}^{\text{GM}})^+ \right] &= \mathbb{E} \left[ \max_{x_{\ell+1} \leq n < x_\ell} \left( \mu^* - \hat{\mu}_{a^*, n} - \sqrt{\frac{1}{2n} \varphi\left(\frac{\tau_t}{Kn}\right)} \right)^+ \right] \\ &\leq \mathbb{E} \left[ \max_{x_{\ell+1} \leq n < x_\ell} \left( \mu^* - \hat{\mu}_{a^*, n} - \sqrt{\frac{1}{2n} \ln\left(\frac{\tau_t}{Kn}\right)} \right)^+ \right] \\ &\leq \mathbb{E} \left[ \max_{x_{\ell+1} \leq n < x_\ell} \left( \mu^* - \hat{\mu}_{a^*, n} - \sqrt{\frac{1}{2x_\ell} \ln\left(\frac{\tau_t}{Kx_\ell}\right)} \right)^+ \right] \\ &\leq \sqrt{\frac{\pi}{8}} \sqrt{\frac{1}{x_{\ell+1}}} \exp\left(-\frac{x_{\ell+1}}{x_\ell} \ln\left(\frac{\tau_t}{Kx_\ell}\right)\right) \\ &= \sqrt{\frac{\pi}{8}} \sqrt{\frac{1}{x_{\ell+1}}} (\beta^{-\ell})^{1/\beta} = \sqrt{\frac{\pi}{8}} \sqrt{\frac{K}{\tau_t}} \beta^{1/2 + \ell(1/2 - 1/\beta)}. \end{aligned}$$

The above series is summable whenever  $\beta \in (1, 2)$ . For instance we may choose  $\beta = 3/2$ , for which

$$\sum_{\ell=0}^{+\infty} \left(\frac{3}{2}\right)^{1/2 + \ell(1/2 - 2/3)} = \sqrt{\frac{3}{2}} \sum_{\ell=0}^{+\infty} \alpha^\ell = \frac{1}{1 - \alpha} \sqrt{\frac{3}{2}} \leq 19 \quad \text{where} \quad \alpha = \left(\frac{3}{2}\right)^{(1/2 - 2/3)} \in (0, 1)$$

Therefore we have shown that

$$\mathbb{E} \left[ (\mu^* - U_{a^*, \tau_t}^{\text{GM}}(t))^+ \mathbf{1}_{\{N_{a^*}(t) < \tau_t/K\}} \right] \leq 19 \sqrt{\frac{\pi}{8}} \sqrt{\frac{K}{\tau_t}}. \quad (44)$$

Combining this bound with (43) and summing over  $t$ , we proved that the first sum in (42) is bounded as

$$\sum_{t=K+1}^T \mathbb{E} \left[ (\mu^* - U_{a^*, \tau_{t-1}}^{\text{GM}}(t-1))^+ \right] \leq 20 \sqrt{\frac{\pi}{8}} \sum_{t=K}^{T-1} \sqrt{\frac{K}{\tau_t}} \quad (45)$$

**Remark 4.** The proof technique reveals that the bound (45) obtained in this step of the proof actually holds even if the arms are pulled according to a strategy that is not a generic MOSS strategy. This is because we never used which specific arms  $A_t^{\text{GM}}$  were pulled: we only distinguished according to how many times  $a^*$  was pulled and resorted to optional skipping.

*Control of the right deviations of all arms*, that is, of the second sum in (42). As  $(x+y)^+ \leq x^+ + y^+$  for all real numbers  $x, y$ , and as  $\varphi$  vanishes on  $(0, 1]$ , we have, for all  $a$  and  $n \geq 1$ ,

$$\begin{aligned} (U_{a,T,n}^{\text{GM}} - \mu_a - \sqrt{K/T})^+ &\leq (\hat{\mu}_{a,n} - \mu_a - \sqrt{K/T})^+ + \sqrt{\frac{1}{2n} \varphi\left(\frac{T}{Kn}\right)} \\ &= (\hat{\mu}_{a,n} - \mu_a - \sqrt{K/T})^+ + \begin{cases} 0 & \text{if } n \geq T/K \\ \sqrt{\frac{1}{2n} \varphi\left(\frac{T}{Kn}\right)} & \text{if } n < T/K. \end{cases} \end{aligned}$$

Therefore, for each arm  $a$ ,

$$\sum_{n=1}^T \mathbb{E} \left[ (U_{a,T,n}^{\text{GM}} - \mu_a - \sqrt{K/T})^+ \right] \leq \sum_{n=1}^T \mathbb{E} \left[ (\hat{\mu}_{a,n} - \mu_a - \sqrt{K/T})^+ \right] + \sum_{n=1}^{\lfloor T/K \rfloor} \sqrt{\frac{1}{2n} \varphi\left(\frac{T}{Kn}\right)}. \quad (46)$$

We are left with two pieces to deal with separately. For the first sum in (46), we exploit the integrated version of Hoeffding's inequality (Corollary 7),

$$\begin{aligned} \sum_{n=1}^T \mathbb{E} \left[ (\hat{\mu}_{a,n} - \mu_a - \sqrt{K/T})^+ \right] &\leq \sqrt{\frac{\pi}{8}} \sum_{n=1}^T \sqrt{\frac{1}{n}} e^{-2n(\sqrt{K/T})^2} \leq \sqrt{\frac{\pi}{8}} \int_0^T \sqrt{\frac{1}{x}} e^{-2xK/T} dx \\ &= \sqrt{\frac{\pi}{8}} \sqrt{\frac{T}{2K}} \int_0^{+\infty} \frac{e^{-u}}{\sqrt{u}} du = \frac{\pi}{4} \sqrt{\frac{T}{K}}, \end{aligned} \quad (47)$$

where we used the equalities  $\int_0^{+\infty} (e^{-u}/\sqrt{u}) du = 2 \int_0^{+\infty} e^{-v^2} dv = \sqrt{\pi}$ .

For the second sum in (46), we also resort to a sum–integral comparison (which exploits the fact that  $\varphi$  is increasing) and perform the change of variable  $u = T/(Kx)$ :

$$\sum_{n=1}^{\lfloor T/K \rfloor} \sqrt{\frac{1}{2n} \varphi\left(\frac{T}{Kn}\right)} \leq \int_0^{T/K} \sqrt{\frac{1}{2x} \varphi\left(\frac{T}{Kx}\right)} dx = \sqrt{\frac{T}{2K}} \int_1^{+\infty} u^{-3/2} \sqrt{\varphi(u)} du.$$

*Conclusion.* Getting back to (41) and (42) and collecting all the bounds above, we showed the desired



bounds,

$$\begin{aligned}
 R_T &\leq (K-1) + \underbrace{\sum_{t=K+1}^T \mathbb{E}\left[(\mu^* - U_{a^*, \tau_{t-1}}^{\text{GM}}(t-1))^+\right]}_{\leq} + \underbrace{\sum_{t=K+1}^T \mathbb{E}\left[U_{A_t^{\text{GM}}, \tau_{t-1}}^{\text{GM}}(t-1) - \mu_{A_t^{\text{GM}}}\right]}_{\leq} \\
 &\leq (K-1) + 20\sqrt{\frac{\pi}{8}} \sum_{t=K}^{T-1} \sqrt{\frac{K}{\tau_t}} + \underbrace{\sqrt{KT} + \sum_{a=1}^K \sum_{n=1}^T \mathbb{E}\left[(U_{a,T,n}^{\text{GM}} - \mu_a - \sqrt{K/T})^+\right]}_{\leq} \\
 &\leq (K-1) + \underbrace{20\sqrt{\frac{\pi}{8}} \sum_{t=K}^{T-1} \sqrt{\frac{K}{\tau_t}}}_{\leq 12.6} + \sqrt{KT} \left(1 + \frac{\pi}{4} + \frac{1}{\sqrt{2}} \int_1^{+\infty} u^{-3/2} \sqrt{\varphi(u)} \, du\right). \quad \square
 \end{aligned}$$

## B. Proofs of the regularity and deviation/concentration results on $\mathcal{K}_{\text{inf}}$

We provide here the proofs of all claims made in Section 4.4 about the  $\mathcal{K}_{\text{inf}}$  function. These proofs are all standard but we occasionally provide simpler or more direct arguments (or slightly refined bounds).

### B.1. Proof of the regularity lemma (Lemma 10)

The proof below is a variation on the proofs that can be found in Honda and Takemura [2015] or earlier references of the same authors.

**Proof:** To prove (19) we lower bound  $\mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon)$ . To that end, given the definition (2), we lower bound  $\text{KL}(\nu, \nu')$  for any fixed probability distribution  $\nu' \in \mathcal{P}[0, 1]$  such that

$$\mathbb{E}(\nu') > \mu - \varepsilon \quad \text{and} \quad \nu' \gg \nu.$$

Since  $\nu'$  has a countable number of atoms, one can pick a real number  $x > \mu$ , arbitrary close to 1, such that  $\delta_x \perp \nu'$  (such that the two probability measures  $\delta_x$  and  $\nu'$  are singular), where  $\delta_x$  is the Dirac distribution at  $x$ . We define

$$\nu'_\alpha = (1 - \alpha)\nu' + \alpha\delta_x \quad \text{where} \quad \alpha = \frac{\varepsilon}{\varepsilon + (x - \mu)} \in (0, 1).$$

The expectation of  $\nu'_\alpha$  satisfies

$$\mathbb{E}(\nu'_\alpha) = (1 - \alpha)\mathbb{E}(\nu') + \alpha x > (1 - \alpha)(\mu - \varepsilon) + \alpha x = \frac{(x - \mu)(\mu - \varepsilon)}{\varepsilon + (x - \mu)} + \frac{\varepsilon x}{\varepsilon + (x - \mu)} = \mu.$$

Since  $\alpha \in (0, 1)$ , we have  $\nu'_\alpha \gg \nu'$ ; therefore,  $\nu'_\alpha \gg \nu' \gg \nu$  and  $\delta_x \perp \nu'$ , which imply the following equalities involving densities (Radon-Nikodym derivatives):

$$\frac{d\nu'}{d\nu'_\alpha} = \frac{1}{1 - \alpha} \quad \text{thus} \quad \frac{d\nu}{d\nu'_\alpha} = \frac{d\nu'}{d\nu'_\alpha} \frac{d\nu}{d\nu'} = \frac{1}{1 - \alpha} \frac{d\nu}{d\nu'}. \quad (48)$$

This allows to compute explicitly the following Kullback-Leibler divergence:

$$\text{KL}(\nu, \nu'_\alpha) = \int_{[0,1]} \ln\left(\frac{d\nu}{d\nu'_\alpha}\right) d\nu = \text{KL}(\nu, \nu') + \ln \frac{1}{1 - \alpha}.$$

Since  $\mathbb{E}(\nu'_\alpha) > \mu$  and by the definition of  $\mathcal{K}_{\text{inf}}$  as an infimum,

$$\mathcal{K}_{\text{inf}}(\nu, \mu) \leq \text{KL}(\nu, \nu'_\alpha) = \text{KL}(\nu, \nu') + \ln \frac{1}{1 - \alpha}.$$

Letting  $x$  go to 1, which implies that  $\alpha$  goes to  $\varepsilon/(1 - \mu + \varepsilon)$ , yields

$$\mathcal{K}_{\text{inf}}(\nu, \mu) \leq \text{KL}(\nu, \nu') + \ln \frac{1 - \mu + \varepsilon}{1 - \mu} = \text{KL}(\nu, \nu') + \ln\left(1 + \frac{\varepsilon}{1 - \mu}\right) \leq \text{KL}(\nu, \nu') + \frac{\varepsilon}{1 - \mu}$$

where we also used  $\ln(1 + u) \leq u$  for all  $u > -1$ . Finally, by taking the infimum in the right-most equation above over all probability distributions  $\nu'$  such that  $\mathbb{E}(\nu') > \mu - \varepsilon$  and  $\nu' \gg \nu$ , we obtain the desired inequality

$$\mathcal{K}_{\text{inf}}(\nu, \mu) \leq \mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon) + \frac{\varepsilon}{1 - \mu}.$$

To prove the second part (20) of Lemma 10, we follow a similar path as above. We lower bound  $\text{KL}(\nu, \nu')$  for any fixed probability distribution  $\nu' \in \mathcal{P}[0, 1]$  such that

$$\mathbb{E}(\nu') > \mu \quad \text{and} \quad \nu' \gg \nu.$$

To that end, we introduce

$$\nu'_\alpha = (1 - \alpha)\nu' + \alpha\nu \quad \text{for} \quad \alpha = \frac{\varepsilon}{\mathbb{E}(\nu') - \mathbb{E}(\nu)} \in (0, 1)$$

where  $\alpha \in (0, 1)$  since  $\mathbb{E}(\nu) \leq \mu - \varepsilon$  by assumption and  $\mathbb{E}(\nu') > \mu$ . These two inequalities also indicate that

$$\mathbb{E}(\nu') - \mathbb{E}(\nu) > \varepsilon \quad \text{thus} \quad \mathbb{E}(\nu'_\alpha) = \mathbb{E}(\nu') - \alpha(\mathbb{E}(\nu') - \mathbb{E}(\nu)) > \mu - \varepsilon \quad (49)$$

so that  $\text{KL}(\nu, \nu'_\alpha) \geq \mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon)$ . Now, thanks to the absolute continuities  $\nu' \gg \nu'_\alpha \gg \nu$ , we have

$$\frac{d\nu}{d\nu'} = \frac{d\nu}{d\nu'_\alpha} \frac{d\nu'_\alpha}{d\nu'} = \frac{d\nu}{d\nu'_\alpha} \left( (1 - \alpha) + \alpha \frac{d\nu}{d\nu'} \right).$$

Therefore, by Fubini's theorem, the Kullback-Leibler divergence between  $\nu$  and  $\nu'$  equals

$$\begin{aligned} \text{KL}(\nu, \nu') &= \int_{[0,1]} \ln\left(\frac{d\nu}{d\nu'}\right) d\nu = \int_{[0,1]} \ln\left(\frac{d\nu}{d\nu'_\alpha}\right) d\nu + \int_{[0,1]} \ln\left((1 - \alpha) + \alpha \frac{d\nu}{d\nu'}\right) d\nu \\ &\geq \int_{[0,1]} \ln\left(\frac{d\nu}{d\nu'_\alpha}\right) d\nu + \alpha \int_{[0,1]} \ln\left(\frac{d\nu}{d\nu'}\right) d\nu \\ &= \text{KL}(\nu, \nu'_\alpha) + \alpha \text{KL}(\nu, \nu') \end{aligned}$$

where we use the concavity of logarithm for the inequality. By Pinsker's inequality together with the data-processing inequality for Kullback-Leibler divergences (see, e.g., Garivier et al., 2018, Lemma 1),

$$\text{KL}(\nu, \nu') \geq \text{KL}\left(\text{Ber}(\mathbb{E}(\nu)), \text{Ber}(\mathbb{E}(\nu'))\right) \geq 2(\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2.$$

Substituting this inequality above, we proved so far

$$\text{KL}(\nu, \nu') \geq \text{KL}(\nu, \nu'_\alpha) + \alpha \text{KL}(\nu, \nu') \geq \text{KL}(\nu, \nu'_\alpha) + 2\alpha(\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2 = \text{KL}(\nu, \nu'_\alpha) + 2\varepsilon(\mathbb{E}(\nu) - \mathbb{E}(\nu'))$$

where we used the definition of  $\alpha$  for the last inequality. By applying the bound (49) and its consequence  $\text{KL}(\nu, \nu'_\alpha) \geq \mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon)$ , we finally get

$$\text{KL}(\nu, \nu') \geq \mathcal{K}_{\text{inf}}(\nu, \mu - \varepsilon) + 2\varepsilon^2.$$

The proof of (20) is concluded by taking the infimum in the left-hand side over the probability distributions  $\nu'$  such that  $\mathbb{E}(\nu') > \mu$  (and  $\nu' \gg \nu$ ).  $\square$

## B.2. A useful tool: a variational formula for $\mathcal{K}_{\text{inf}}$ (statement)

The variational formula below appears in Honda and Takemura [2015] as Theorem 2 (and Lemma 6) and is an essential tool for deriving the deviation and concentration results for the  $\mathcal{K}_{\text{inf}}$ . We state it here (and re-derive it in a direct way in Appendix D) for the sake of completeness.

**Lemma 18** (variational formula for  $\mathcal{K}_{\text{inf}}$ ). *For all  $\nu \in \mathcal{P}[0, 1]$  and all  $0 < \mu < 1$ ,*

$$\mathcal{K}_{\text{inf}}(\nu, \mu) = \max_{0 \leq \lambda \leq 1} \mathbb{E} \left[ \ln \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \right] \quad \text{where } X \sim \nu. \quad (50)$$

Moreover, if we denote by  $\lambda^*$  the value at which the above maximum is reached, then

$$\mathbb{E} \left[ \frac{1}{1 - \lambda^*(X - \mu)/(1 - \mu)} \right] \leq 1. \quad (51)$$

### B.3. Proof of the deviation result (Proposition 13)

The following proof is almost exactly the same as that of Cappé et al. [2013, Lemma 6], except that we correct a small mistake in the constant.

**Proof:** We first upper bound  $\mathcal{K}_{\text{inf}}(\widehat{\nu}_n, \mathbf{E}(\nu))$ : as indicated by the variational formula of Lemma 18, it is a maximum of random variables indexed by  $[0, 1]$ . We provide an upper bound that is a finite maximum. To that end, we fix a real number  $\gamma \in (0, 1)$ , to be determined by the analysis, and let  $S_\gamma$  be the set

$$S_\gamma = \left\{ \frac{1}{2} - \left\lfloor \frac{1}{2\gamma} \right\rfloor \gamma, \dots, \frac{1}{2} - \gamma, \frac{1}{2}, \frac{1}{2} + \gamma, \dots, \frac{1}{2} + \left\lfloor \frac{1}{2\gamma} \right\rfloor \gamma \right\}.$$

The cardinality of this set  $S_\gamma$  is bounded by  $1 + 1/\gamma$ . Lemma 19 below (together with the consequence mentioned after its statement) indicates that for all  $\lambda \in [0, 1]$ , there exists a  $\lambda' \in S_\gamma$  such that for all  $x \in [0, 1]$ ,

$$\ln \left( 1 - \lambda \frac{x - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right) \leq 2\gamma + \ln \left( 1 - \lambda' \frac{x - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right). \quad (52)$$

(The small correction with respect to the original proof is the  $2\gamma$  factor in the inequality above, instead of the claimed  $\gamma$  term therein; this is due to the constraint  $\lambda \leq \lambda' \leq 1/2$  or  $1/2 \leq \lambda' \leq \lambda$  in the statement of Lemma 19.) Now, a combination of the variational formula of Lemma 18 and of the inequality (52) yields a finite maximum as an upper bound on  $\mathcal{K}_{\text{inf}}(\widehat{\nu}_n, \mathbf{E}(\nu))$ :

$$\begin{aligned} \mathcal{K}_{\text{inf}}(\widehat{\nu}_n, \mathbf{E}(\nu)) &= \max_{0 \leq \lambda \leq 1} \frac{1}{n} \sum_{k=1}^n \ln \left( 1 - \lambda \frac{X_k - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right) \\ &\leq 2\gamma + \max_{\lambda' \in S_\gamma} \frac{1}{n} \sum_{k=1}^n \ln \left( 1 - \lambda' \frac{X_k - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right). \end{aligned}$$

In the second part of the proof, we control the deviations of the upper bound obtained. A union bound yields

$$\begin{aligned} \mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\widehat{\nu}_n, \mathbf{E}(\nu)) \geq u \right] &\leq \mathbb{P} \left[ \max_{\lambda' \in S_\gamma} \frac{1}{n} \sum_{k=1}^n \ln \left( 1 - \lambda' \frac{X_k - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right) \geq u - 2\gamma \right] \\ &\leq \sum_{\lambda' \in S_\gamma} \mathbb{P} \left[ \frac{1}{n} \sum_{k=1}^n \ln \left( 1 - \lambda' \frac{X_k - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right) \geq u - 2\gamma \right]. \end{aligned} \quad (53)$$

By the Markov–Chernov inequality, for all  $\lambda' \in [0, 1]$ , we have

$$\begin{aligned} \mathbb{P} \left[ \frac{1}{n} \sum_{k=1}^n \ln \left( 1 - \lambda' \frac{X_k - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right) \geq u - 2\gamma \right] &\leq e^{-n(u-2\gamma)} \mathbb{E} \left[ \prod_{k=1}^n \left( 1 - \lambda' \frac{X_k - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right) \right] \\ &= e^{-n(u-2\gamma)} \prod_{k=1}^n \underbrace{\mathbb{E} \left[ 1 - \lambda' \frac{X_k - \mathbf{E}(\nu)}{1 - \mathbf{E}(\nu)} \right]}_{=1} = e^{-n(u-2\gamma)} \end{aligned}$$

where we used the independence of the  $X_k$ . Substituting in (53) and using the bound  $1 + 1/\gamma$  on the cardinality of  $S_\gamma$ , we get

$$\mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\widehat{\nu}_n, \mathbf{E}(\nu)) \geq u \right] \leq \sum_{\lambda' \in S_\gamma} e^{-n(u-2\gamma)} \leq (1 + 1/\gamma) e^{-n(u-2\gamma)}.$$

Taking  $\gamma = 1/(2n)$  concludes the proof.  $\square$

The proof above relies on the following lemma, which is extracted from Cappé et al. [2013, Lemma 7]. Its elementary proof (not copied here) consists in bounding of derivative of  $\lambda \mapsto \ln(1 - \lambda c)$  and using a convexity argument.

**Lemma 19.** *For all  $\lambda, \lambda' \in [0, 1)$  such that either  $\lambda \leq \lambda' \leq 1/2$  or  $1/2 \leq \lambda' \leq \lambda$ , for all real numbers  $c \leq 1$ ,*

$$\ln(1 - \lambda c) - \ln(1 - \lambda' c) \leq 2|\lambda - \lambda'|.$$

A consequence not drawn by Cappé et al. [2013] is that the lemma above actually also holds for  $\lambda = 1$  and  $\lambda' \in [0, 1)$ . Indeed, by continuity and by letting  $\lambda \rightarrow 1$ , we get from this lemma that for all  $\lambda' \in [1/2, 1)$  and for all real numbers  $c < 1$ ,

$$\ln(1 - c) - \ln(1 - \lambda' c) \leq 2(1 - \lambda').$$

The above inequality is also valid for  $c = 1$  as the left-hand side equals  $-\infty$ .

#### B.4. Proof of the concentration result (Proposition 15)

We recall that Proposition 15—and actually most of its proof below—are similar in spirit to Honda and Takemura [2015, Proposition 11]. However, they are tailored to our needs. The key ingredients in the proof will be the variational formula (50)—again—and Lemma 20 below. This lemma is a concentration result for random variables that are essentially bounded from one side only; it holds for possibly negative  $u$  (there is no lower bound on the  $u$  that can be considered).

**Lemma 20.** *Let  $Z_1, \dots, Z_n$  be i.i.d. random variables such that there exist  $a, b \geq 0$  with*

$$Z_1 \leq a \quad \text{a.s.} \quad \text{and} \quad \mathbb{E}[e^{-Z_1}] \leq b.$$

*Denote  $\gamma = \sqrt{e^a}(16e^{-2b} + a^2)$ . Then  $Z_1$  is integrable and for all real numbers  $u \in (-\infty, \mathbb{E}[Z_1])$ ,*

$$\mathbb{P}\left[\sum_{i=1}^n Z_i \leq nu\right] \leq \begin{cases} \exp(-n\gamma/8) & \text{if } u \leq \mathbb{E}[Z_1] - \gamma/2 \\ \exp(-n(\mathbb{E}[Z_1] - u)^2/(2\gamma)) & \text{if } u > \mathbb{E}[Z_1] - \gamma/2 \end{cases}.$$

##### B.4.1. Proof of Proposition 15 based on Lemma 20

We apply Lemma 18. We denote by  $\lambda^* \in [0, 1]$  a real number achieving the maximum in the variational formula (50) for  $\mathcal{K}_{\text{inf}}(\nu, \mu)$ . We then introduce the random variable

$$Z = \ln\left(1 - \lambda^* \frac{X - \mu}{1 - \mu}\right) \quad \text{where} \quad X \sim \nu$$

and i.i.d. copies  $Z_1, \dots, Z_n$  of  $Z$ . Then,  $\mathcal{K}_{\text{inf}}(\nu, \mu) = \mathbb{E}[Z]$  and by the variational formula (50) again,

$$\mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mu) \geq \frac{1}{n} \sum_{i=1}^n Z_i, \quad \text{therefore,} \quad \mathbb{P}[\mathcal{K}_{\text{inf}}(\hat{\nu}_n, \mu) \leq x] \leq \mathbb{P}\left[\sum_{i=1}^n Z_i \leq nx\right]$$

for all real numbers  $x$ . Now,  $X \geq 0$  and  $\lambda^* \leq 1$ , thus

$$Z \leq \ln\left(1 + \lambda^* \frac{\mu}{1 - \mu}\right) \leq \ln\left(\frac{1}{1 - \mu}\right) \stackrel{\text{def}}{=} a.$$

On the other hand,

$$\mathbb{E}[e^{-Z}] = \mathbb{E}\left[\frac{1}{1 - \lambda^*(X - \mu)/(1 - \mu)}\right] \stackrel{\text{def}}{=} b$$

where  $b \leq 1$  follows from (51). This proves Proposition 15 via Lemma 20, except for the inequality  $e^{-n\gamma/8} \leq e^{-n/4}$  claimed therein. The latter is a consequence of  $\gamma \geq 2$ , as  $\gamma$  is an increasing function of  $\mu > 0$ ,

$$\gamma = \frac{1}{\sqrt{1-\mu}} \left( 16e^{-2} + \ln^2 \left( \frac{1}{1-\mu} \right) \right) > 16e^{-2} > 2$$

**Remark 5.** In the proof of Theorem 3 provided in Section C we will not use Proposition 15 as stated but a stronger result, namely that the bound of Proposition 15 actually holds for the larger quantity

$$\mathbb{P} \left[ \sum_{i=1}^n Z_i \leq nx \right]$$

as is clear from the proof above.

#### B.4.2. Proof of Lemma 20

This lemma is a direct application of the Crámer–Chernov method. We introduce the log-moment generating function  $\Lambda$  of  $Z_1$ :

$$\Lambda : x \mapsto \ln \mathbb{E}[e^{xZ_1}].$$

**Lemma 21.** *The log-moment generating function  $\Lambda$  is well-defined at least on the interval  $[-1, 1]$  and twice differentiable at least on  $(-1, 1)$ , with  $\Lambda'(0) = \mathbb{E}[Z_1]$  and  $\Lambda''(x) \leq \gamma$  for  $x \in [-1/2, 0]$ , where  $\gamma = \sqrt{e^a(16e^{-2}b + a^2)}$  denotes the same constant as in Lemma 20.*

Based on this lemma (proved below), we may resort to a Taylor expansion with a Lagrange remainder and get the bound:

$$\forall x \in [-1/2, 0], \quad \Lambda(x) \leq \Lambda(0) + x \Lambda'(0) + \frac{x^2}{2} \sup_{y \in [-1/2, 0]} \Lambda''(y) \leq x \mathbb{E}[Z_1] + \frac{\gamma}{2} x^2.$$

Therefore, by the Crámer–Chernov method, for all  $x \in [-1/2, 0]$ , the probability of interest is bounded by

$$\begin{aligned} \mathbb{P} \left[ \sum_{i=1}^n Z_i \leq nu \right] &= \mathbb{P} \left[ \prod_{i=1}^n e^{xZ_i} \geq e^{nux} \right] \leq e^{-nux} \left( \mathbb{E}[e^{xZ_1}] \right)^n = \exp \left( -n(ux - \Lambda(x)) \right) \\ &\leq \exp \left( n \left( x^2 \gamma / 2 - x(u - \mathbb{E}[Z_1]) \right) \right). \end{aligned} \quad (54)$$

That is,

$$\mathbb{P} \left[ \sum_{i=1}^n Z_i \leq nu \right] \leq \exp \left( n \min_{x \in [-1/2, 0]} P(x) \right)$$

where we introduced the second-order polynomial function

$$P(x) = x^2 \gamma / 2 - x(u - \mathbb{E}[Z_1]) = \frac{\gamma x}{2} \left( x - 2 \frac{u - \mathbb{E}[Z_1]}{\gamma} \right).$$

The claimed bound is obtained by minimizing  $P$  over  $[-1/2, 0]$  depending on whether  $u > \mathbb{E}[Z_1] - \gamma/2$  or  $u \leq \mathbb{E}[Z_1] - \gamma/2$ , which we do now.

We recall that by assumption,  $u < \mathbb{E}[Z_1]$ . We note that  $P$  is a second-order polynomial function with positive leading coefficient and roots 0 and  $2(u - \mathbb{E}[Z_1])/\gamma < 0$ . Its minimum over the entire real line  $(-\infty, +\infty)$  is thus achieved at the midpoint  $x^* = (u - \mathbb{E}[Z_1])/\gamma < 0$  between these roots. But  $P$

is to be minimized over  $[-1/2, 0]$  only. In the case where  $u > \mathbb{E}[Z_1] - \gamma/2$ , the midpoint  $x^*$  belongs to the interval of interest and

$$\min_{[-1/2, 0]} P = \frac{\gamma x^*}{2} \left( x^* - 2 \frac{u - \mathbb{E}[Z_1]}{\gamma} \right) = -\frac{(u - \mathbb{E}[Z_1])^2}{2\gamma}.$$

Otherwise,  $u - \mathbb{E}[Z_1] \leq -\gamma/2$  and the midpoint  $x^*$  is to the left of  $-1/2$ . Therefore,  $P$  is increasing on  $[-1/2, 0]$ , so that its minimum on this interval is achieved at  $-1/2$ , that is,

$$\min_{[-1/2, 0]} P = P(-1/2) = \frac{\gamma}{8} + \frac{1}{2}(u - \mathbb{E}[Z_1]) \leq \frac{\gamma}{8} - \frac{\gamma}{4} = -\frac{\gamma}{8}.$$

This concludes the proof of Lemma 20. We end this section by proving Lemma 21, which stated some properties of the  $\Lambda$  function.

**Proof: (of Lemma 21)** We will make repeated uses of the fact that  $e^{-Z_1}$  is integrable (by the assumption on  $b$ ), and that so is  $e^{Z_1}$ , as  $e^{Z_1}$  takes bounded values in  $(0, e^a]$ . In particular,  $Z_1$  is integrable, as by Jensen's inequality,

$$\mathbb{E}[|Z_1|] \leq \ln \mathbb{E}[e^{|Z_1|}] \leq \ln(\mathbb{E}[e^{-Z_1}] + \mathbb{E}[e^{Z_1}]) < +\infty.$$

First, that  $\Lambda$  is well-defined over  $[-1, 1]$  follows from the inequality  $e^{xZ_1} \leq e^{Z_1} + e^{-Z_1}$ , which is valid for all  $x \in [-1, 1]$  and whose right-hand side is integrable as already noted above.

Second, that  $\psi : x \mapsto \mathbb{E}[e^{xZ_1}]$  is differentiable at least on  $(-1, 1)$  follows from the fact that  $x \in (-1, 1) \mapsto Z_1 e^{xZ_1}$  is locally dominated by an integrable random variable; indeed, for  $x \in (-1, 1)$ ,

$$|Z_1 e^{xZ_1}| = Z_1 e^{xZ_1} \mathbf{1}_{\{Z_1 \geq 0\}} + Z_1 e^{xZ_1} \mathbf{1}_{\{Z_1 < 0\}} \leq a e^a + \frac{1}{x} \sup_{(-\infty, 0)} f = a e^a + \frac{1}{e x}$$

where  $f(t) = -t e^t$ .

Similarly,  $x \in (-1, 1) \mapsto Z_1^2 e^{xZ_1}$  is also locally dominated by an integrable random variable. Thus,  $\psi$  is twice differentiable at least on  $(-1, 1)$ , with first and second derivatives

$$\psi'(x) = \mathbb{E}[Z_1 e^{xZ_1}] \quad \text{and} \quad \psi''(x) = \mathbb{E}[Z_1^2 e^{xZ_1}].$$

Therefore, so is  $\Lambda = \ln \psi$ , with derivatives

$$\Lambda'(x) = \frac{\psi'(x)}{\psi(x)} = \frac{\mathbb{E}[Z_1 e^{xZ_1}]}{\mathbb{E}[e^{xZ_1}]} \quad \text{and} \quad \Lambda''(x) = \frac{\psi''(x)\psi(x) - (\psi'(x))^2}{\psi(x)^2} \leq \frac{\psi''(x)}{\psi(x)} = \frac{\mathbb{E}[Z_1^2 e^{xZ_1}]}{\mathbb{E}[e^{xZ_1}]}.$$

In particular,  $\Lambda'(0) = \mathbb{E}[Z_1]$ .

Finally, for the bound on  $\Lambda''(x)$ , we note first that  $Z_1 \leq a$  (with  $a \geq 0$ ) and  $x \in [-1/2, 0]$  entail that  $e^{xZ_1} \geq e^{xa} \geq 1/\sqrt{e^a}$ . Second,  $\mathbb{E}[Z_1^2 e^{xZ_1}] \leq 16e^{-2b+a^2}$  follows from replacing  $z$  by  $Z_1$  and taking expectations in the inequality (proved below)

$$\forall x \in [-1/2, 0], \quad z \in (-\infty, a], \quad z^2 e^{xz} \leq 16e^{-2}e^{-z} + a^2. \quad (55)$$

Collecting all elements together, we proved

$$\Lambda''(x) \leq \frac{\mathbb{E}[Z_1^2 e^{xZ_1}]}{\mathbb{E}[e^{xZ_1}]} \leq \sqrt{e^a}(16e^{-2}b + a^2) = \gamma.$$

To see why (55) holds, note that in the case  $z \geq 0$ , since  $x \leq 0$  we have  $z^2 e^{xz} \leq z^2 \leq a^2$ . In the case  $z \leq 0$ , we have (by function study)  $z^2 \leq 16e^{-2-z/2}$ , so that  $z^2 e^{xz} \leq 16e^{-2}e^{(x-1/2)z} \leq 16e^{-2}e^{-z}$  where we used  $x \geq -1/2$  for the final inequality.  $\square$

### C. Proof of Theorem 3 (with the $-\ln \ln T$ term in the regret bound)

We incorporate two refinements to the proof of Theorem 2 in Section 6.2 to obtain Theorem 3 with this improved  $-\ln \ln T$  term. First, the left deviations of the index are controlled with an additional cut on the value of  $U_a(t)$  before using the bound  $U_a(t) \geq U_{a^*}(t)$  that holds when  $A_{t+1} = a$ . This improves the dependency on the parameter  $\delta$  used in the proof; as a consequence,  $\delta = T^{-1/8}$  will be set instead of  $\delta = (\ln T)^{-1/3}$ , which will improve the order of magnitude of second-order terms. Second, to sharpen the bound on the quantity (60), which contains the main logarithmic term, we use a trick introduced in the analysis of the IMED policy by Honda and Takemura [2015, Theorem 5]. Their idea was to deal with the deviations in a more careful way and relate the sum (60) to the behaviour of a biased random walk. Doing so, we obtain a bound of the form  $\kappa W(cT)$ , where  $W$  is Lambert's function, instead of the bound of the form  $\kappa \ln(cT)$  stated in Theorem 2.

We recall that Lambert's function  $W$  is defined, for  $x > 0$ , as the unique solution  $W(x)$  of the equation  $w e^w = x$ , with unknown  $w > 0$ . It is an increasing function satisfying (see, e.g., Hoorfar and Hassani, 2008, Corollary 2.4)

$$\forall x > e, \quad \ln x - \ln \ln x \leq W(x) \leq \ln x - \ln \ln x + \ln(1 + e^{-1}). \quad (56)$$

In particular,  $W(x) = \ln x - \ln \ln x + \mathcal{O}(1)$  as  $x \rightarrow +\infty$ .

What we will exactly prove below is the following. We recall that we assume here  $\mu^* \in (0, 1)$ . Given  $T \geq K/(1 - \mu^*)$ , the KL-UCB-switch algorithm, tuned with the knowledge of  $T$  and the switch function  $f(T, K) = \lfloor (T/K)^{1/5} \rfloor$ , ensures that for all bandit problems  $\underline{\nu}$  over  $[0, 1]$ , for all sub-optimal arms  $a$ , and for all  $\delta > 0$  satisfying

$$\delta < \min \left\{ \mu^*, \frac{\Delta_a}{2}, \frac{1 - \mu^*}{2} \mathcal{K}_{\inf}(\nu_a, \mu^*) \right\}$$

we have

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq 1 \\ &+ \frac{5eK}{(1 - e^{-\Delta_a^2/2})^3} + T e^{-\Delta_a^2 T / (2K)} \\ &+ \frac{K/T}{1 - e^{-\Delta_a^2/8}} \\ &+ \left\lceil \frac{8}{\Delta_a^2} \ln \left( \frac{T}{K} \right) \right\rceil \left( \frac{5eK/T}{(1 - e^{-2\delta^2})^3} + e^{-2\delta^2 T / K} \right) \\ &+ \frac{1}{\mathcal{K}_{\inf}(\nu_a, \mu^*) - \delta / (1 - \mu^*)} \left( W \left( \frac{\ln(1/(1 - \mu^*))}{K} T \right) + \ln(2/(1 - \mu^*)) \right) \\ &\quad + 5 + \frac{1}{1 - e^{-\mathcal{K}_{\inf}(\nu, \mu^*)^2 / (8\gamma_*)}} \\ &+ \frac{1}{1 - e^{-\Delta_a^2/8}}. \end{aligned} \quad (57)$$

We write the bound in this way to match the decomposition of  $\mathbb{E}[N_a(T)]$  appearing in the proof (see page 40). For a choice  $\delta \rightarrow 0$   $T \rightarrow +\infty$ , the previous bound is of the form

$$\mathbb{E}[N_a(T)] \leq \frac{W(c_{\mu^*} T)}{\mathcal{K}_{\inf}(\nu_a, \mu^*) - \delta / (1 - \mu^*)} + \mathcal{O}_T \left( \frac{\ln T}{\delta^6 T} \right) + \mathcal{O}_T((\ln T) e^{-2\delta^2 T / K}) + \mathcal{O}_T(1)$$

where  $c_{\mu^*} = \ln(1/(1 - \mu^*)) / K$ . Based on the first-order approximation  $1/(1 - \varepsilon) = 1 + \varepsilon + \mathcal{O}(\varepsilon)$  as  $\varepsilon \rightarrow 0$  and on the inequalities (56), we get

$$\mathbb{E}[N_a(T)] \leq \frac{\ln T - \ln \ln T}{\mathcal{K}_{\inf}(\nu_a, \mu^*)} (1 + \mathcal{O}_T(\delta)) + \mathcal{O}_T \left( \frac{\ln T}{\delta^6 T} \right) + \mathcal{O}_T((\ln T) e^{-2\delta^2 T / K}) + \mathcal{O}_T(1).$$



The choice  $\delta = T^{-1/8}$  leads to the bound stated in Theorem 3, namely,

$$\mathbb{E}[N_a(T)] \leq \frac{\ln T - \ln \ln T}{\mathcal{K}_{\inf}(\nu_a, \mu^*)} + \mathcal{O}_T(1).$$

We now prove the closed-form bound (57).

**Proof:** As in the proof of Theorem 2, given  $\delta > 0$  sufficiently small, we decompose  $\mathbb{E}[N_a(T)]$ . However, this time we refine the decomposition quite a bit. Instead of simply distinguishing whether  $U_a(t)$  is greater or smaller than  $\mu^* - \delta$ , we add a cutting point at  $(\mu^* + \mu_a)/2$ . In addition, we set a threshold  $n_0 \geq 1$  (to be determined by the analysis) and distinguish whether  $N_a(t) \geq n_0$  or  $N_a(t) \leq n_0 - 1$  when  $U_a(t) < \mu^* - \delta$ , while we keep the integer threshold  $f(T, K)$  in the case  $U_a(t) \geq \mu^* - \delta$ . More precisely,

$$\begin{aligned} \{U_a(t) < \mu^* - \delta\} \cup \{U_a(t) \geq \mu^* - \delta\} &= \{U_a(t) < \mu^* - \delta \text{ and } N_a(t) \geq n_0\} \\ &\cup \{U_a(t) < \mu^* - \delta \text{ and } N_a(t) \leq n_0 - 1\} \\ &\cup \{U_a(t) \geq \mu^* - \delta \text{ and } N_a(t) \leq f(T, K)\} \\ &\cup \{U_a(t) \geq \mu^* - \delta \text{ and } N_a(t) \geq f(T, K) + 1\} \\ &\subseteq \{U_a(t) < (\mu^* + \mu_a)/2 \text{ and } N_a(t) \geq n_0\} \\ &\cup \{(\mu^* + \mu_a)/2 \leq U_a(t) < \mu^* - \delta \text{ and } N_a(t) \geq n_0\} \\ &\cup \{U_a(t) < \mu^* - \delta \text{ and } N_a(t) \leq n_0 - 1\} \\ &\cup \{U_a^{\text{KL}}(t) \geq \mu^* - \delta \text{ and } N_a(t) \leq f(T, K)\} \\ &\cup \{U_a^{\text{M}}(t) \geq \mu^* - \delta \text{ and } N_a(t) \geq f(T, K) + 1\} \end{aligned}$$

where, to get the inclusion, we further cut the first event into two events and we used the definition of the index  $U_a(t)$  to replace it by  $U_a^{\text{KL}}(t)$  or  $U_a^{\text{M}}(t)$  in the last two events.

Hence, by intersecting this partition of the space with the event  $\{A_{t+1} = a\}$  and by slightly simplifying the first and second events of the partition:

$$\begin{aligned} \{A_{t+1} = a\} &\subseteq \{U_a(t) < (\mu^* + \mu_a)/2 \text{ and } A_{t+1} = a\} \\ &\cup \{U_a(t) \geq (\mu^* + \mu_a)/2 \text{ and } A_{t+1} = a \text{ and } N_a(t) \geq n_0\} \\ &\cup \{U_a(t) < \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq n_0 - 1\} \\ &\cup \{U_a^{\text{KL}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(T, K)\} \\ &\cup \{U_a^{\text{M}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \geq f(T, K) + 1\} \end{aligned}$$

Only now do we inject the bound  $U_{a^*}(t) \leq U_a(t)$ , valid when  $A_{t+1} = a$ , as well as a union bound,

to obtain our working decomposition of  $\mathbb{E}[N_a(t)]$ :

$$\mathbb{E}[N_a(T)] \leq 1 + \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}(t) < (\mu^* + \mu_a)/2] \quad (S_1)$$

$$+ \sum_{t=K}^{T-1} \mathbb{P}[U_a(t) \geq (\mu^* + \mu_a)/2 \text{ and } A_{t+1} = a \text{ and } N_a(t) \geq n_0] \quad (S_2)$$

$$+ \sum_{t=K}^{T-1} \mathbb{P}[U_{a^*}(t) < \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq n_0 - 1] \quad (S_3)$$

$$+ \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{KL}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(T, K)] \quad (S_4)$$

$$+ \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{M}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \geq f(T, K) + 1]. \quad (S_5)$$

We call the five sums appearing in the right-hand side  $S_1, S_2, S_3, S_4, S_5$ , respectively and now bound them separately. Most of the efforts will be dedicated to the sum  $S_4$ .

### Bound on $S_5$

As the algorithm considered is the same as in Theorem 2, its analysis is still valid. Fortunately, the  $S_5$  term was already covered in (38): provided that  $\delta < \Delta_a/4$ ,

$$S_5 \leq \frac{1}{1 - e^{-\Delta_a^2/8}}.$$

### Bound on $S_2$

Let

$$n_0 = \left\lceil \frac{8}{\Delta_a^2} \ln\left(\frac{T}{K}\right) \right\rceil. \quad (58)$$

By Pinsker's inequality (8), by definition of the MOSS index, and by our choice of  $n_0$ , we have, when  $N_a(t) \geq n_0$ ,

$$U_a(t) \leq U_a^{\text{M}}(t) = \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \ln_+\left(\frac{T}{KN_a(t)}\right)} \leq \hat{\mu}_a(t) + \underbrace{\sqrt{\frac{1}{2n_0} \ln_+\left(\frac{T}{Kn_0}\right)}}_{\leq \Delta_a/4}. \quad (59)$$

In particular, we get the inclusion

$$\begin{aligned} \{U_a(t) \geq (\mu^* + \mu_a)/2 \text{ and } N_a(t) \geq n_0\} &= \{U_a(t) \geq \mu_a + \Delta_a/2 \text{ and } N_a(t) \geq n_0\} \\ &\subseteq \{\hat{\mu}_a(t) \geq \mu_a + \Delta_a/4 \text{ and } N_a(t) \geq n_0\}. \end{aligned}$$

Thus

$$S_2 \leq \sum_{t=K}^{T-1} \mathbb{P}\left[\hat{\mu}_a(t) \geq \mu_a + \frac{\Delta_a}{4} \text{ and } A_{t+1} = a \text{ and } N_a(t) \geq n_0\right].$$

We now proceed similarly to what we already did on page 21. By a careful application of optional skipping (see Section 4.1), using the fact that all the events  $\{A_{t+1} = a \text{ and } N_a(t) = n\}$  are disjoint

as  $t$  varies, the sum above may be bounded by

$$\sum_{t=K}^{T-1} \mathbb{P} \left[ \hat{\mu}_a(t) \geq \mu_a + \frac{\Delta_a}{4} \text{ and } A_{t+1} = a \text{ and } N_a(t) \geq n_0 \right] \leq \sum_{n \geq n_0} \mathbb{P} \left[ \hat{\mu}_{a,n} \geq \mu_a + \frac{\Delta_a}{4} \right]$$

By a final application of Hoeffding's inequality (Proposition 6, actually not using the maximal form):

$$S_2 \leq \sum_{n=n_0}^T \mathbb{P} \left[ \hat{\mu}_{a,n} \geq \mu_a + \frac{\Delta_a}{4} \right] \leq \sum_{n=n_0}^T e^{-n\Delta_a^2/8} = \frac{e^{-n_0\Delta_a^2/8}}{1 - e^{-\Delta_a^2/8}} \leq \frac{K/T}{1 - e^{-\Delta_a^2/8}}$$

where we substituted the value (58) of  $n_0$ .

### Bounds on $S_1$ and $S_3$

For  $u \in (0, 1)$ , we introduce the event

$$\mathcal{E}_*(u) = \left\{ \exists \tau \in \{K, \dots, T-1\} : U_{a^*}(\tau) < u \right\}$$

so that

$$\{U_{a^*}(t) < (\mu^* + \mu_a)/2\} \subseteq \mathcal{E}_*((\mu^* + \mu_a)/2) \quad \text{and} \quad \{U_{a^*}(t) < \mu^* - \delta\} \subseteq \mathcal{E}_*(\mu^* - \delta).$$

Summing over  $t$ , and using the deterministic control

$$\sum_{t=K}^{T-1} \mathbb{1}_{\{A_{t+1}=a \text{ and } N_a(t) \leq n_0-1\}} \leq n_0$$

for bounding  $S_3$ , we obtain (and this is where it is handy that the  $\mathcal{E}_*$  do not depend on a particular  $t$ )

$$S_1 \leq T \mathbb{P}(\mathcal{E}_*((\mu^* + \mu_a)/2)) \quad \text{and} \quad S_3 \leq n_0 \mathbb{P}(\mathcal{E}_*(\mu^* - \delta))$$

We recall that  $n_0$  was defined in (58). The lemma right below, respectively with  $x = \Delta_a/2$  and  $x = \delta$ , yield the final bounds

$$S_1 \leq \frac{5eK}{(1 - e^{-\Delta_a^2/2})^3} + Te^{-\Delta_a^2 T/(2K)}$$

and

$$S_3 \leq \left\lceil \frac{8}{\Delta_a^2} \ln \left( \frac{T}{K} \right) \right\rceil \left( \frac{5eK/T}{(1 - e^{-2\delta^2})^3} + e^{-2\delta^2 T/K} \right).$$

**Lemma 22.** For all  $x \in (0, \mu^*)$ ,

$$\mathbb{P}(\mathcal{E}_*(\mu^* - x)) = \mathbb{P}[\exists \tau \in \{K, \dots, T-1\} : U_{a^*}(\tau) < \mu^* - x] \leq \frac{eK}{T} \frac{5}{(1 - e^{-2x^2})^3} + e^{-2x^2 T/K}.$$

**Proof:** We first lower bound  $U_{a^*}(\tau)$  depending on whether  $N_{a^*}(\tau) < T/K$  or  $N_{a^*}(\tau) \geq T/K$ . In the first case, we will simply apply Pinsker's inequality (8) to get  $U_{a^*}^{\text{KL}}(\tau) \leq U_{a^*}(\tau)$ . In the second case, since  $T \geq K/(1 - \mu^*) \geq K$ , we have, by definition of  $f(T, K)$ , that  $T/K \geq (T/K)^{1/5} \geq f(T, K)$  and thus, by definition of the  $U_{a^*}(\tau)$  index,  $U_{a^*}(\tau) = U_{a^*}^{\text{M}}(\tau)$ . Now, the  $\ln_+$  in the definition of  $U_{a^*}^{\text{M}}(\tau)$  vanishes when  $N_{a^*}(\tau) \geq T/K$ , so all in all we have  $U_{a^*}(\tau) = \hat{\mu}_{a^*}(\tau)$  when  $N_{a^*}(\tau) \geq T/K$ . Therefore,

by optional skipping (see Section 4.1),

$$\begin{aligned}
 \mathbb{P}\left(\mathcal{E}_*(\mu^* - x)\right) &= \mathbb{P}\left[\exists \tau \in \{K, \dots, T-1\} : U_{a^*}(\tau) < \mu^* - x\right] \\
 &= \mathbb{P}\left[\exists \tau \in \{K, \dots, T-1\} : U_{a^*}(\tau) < \mu^* - x \text{ and } N_{a^*}(\tau) < T/K\right] \\
 &\quad + \mathbb{P}\left[\exists \tau \in \{K, \dots, T-1\} : U_{a^*}(\tau) < \mu^* - x \text{ and } N_{a^*}(\tau) \geq T/K\right] \\
 &\leq \mathbb{P}\left[\exists \tau \in \{K, \dots, T-1\} : U_{a^*}^{\text{KL}}(\tau) < \mu^* - x \text{ and } N_{a^*}(\tau) < T/K\right] \\
 &\quad + \mathbb{P}\left[\exists \tau \in \{K, \dots, T-1\} : \widehat{\mu}_{a^*}(\tau) < \mu^* - x \text{ and } N_{a^*}(\tau) \geq T/K\right] \\
 &\leq \mathbb{P}\left[\exists m \in \{1, \dots, \lfloor T/K \rfloor\} : U_{a^*,m}^{\text{KL}} < \mu^* - x\right] \\
 &\quad + \mathbb{P}\left[\exists m \in \{\lceil T/K \rceil, \dots, T\} : \widehat{\mu}_{a^*,m} < \mu^* - x\right].
 \end{aligned}$$

As in the proof of Corollary 14, by the definition of the  $U_{a^*,m}^{\text{KL}}$  index as some supremum (together with the left-continuity of  $\mathcal{K}_{\text{inf}}$  deriving from Lemma 10), we finally get

$$\begin{aligned}
 \mathbb{P}\left(\mathcal{E}_*(\mu^* - x)\right) &\leq \mathbb{P}\left[\exists m \in \{1, \dots, \lfloor T/K \rfloor\} : \mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,m}, \mu^* - x) > \frac{1}{m} \ln\left(\frac{T}{Km}\right)\right] \\
 &\quad + \mathbb{P}\left[\exists m \in \{\lceil T/K \rceil, \dots, T\} : \widehat{\mu}_{a^*,m} < \mu^* - x\right].
 \end{aligned}$$

The proof is concluded by bounding each probability separately. First, again as in the proof of Corollary 14, we apply Corollary 12 (for the first inequality below) and the deviation inequality of Proposition 13 (for the second inequality below), to see that for all  $x \in (0, \mu^*)$  and  $\varepsilon > 0$ ,

$$\mathbb{P}\left[\mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,m}, \mu^* - x) > \varepsilon\right] \leq \mathbb{P}\left[\mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,m}, \mu^*) > \varepsilon + 2x^2\right] \leq e(2n+1)e^{-n(\varepsilon+2x^2)}.$$

Therefore, by a union bound, the above equation, and the calculations on geometric sums (33) and (34),

$$\begin{aligned}
 &\mathbb{P}\left[\exists m \in \{1, \dots, \lfloor T/K \rfloor\} : \mathcal{K}_{\text{inf}}(\widehat{\nu}_{a^*,m}, \mu^* - x) > \frac{1}{m} \ln\left(\frac{T}{Km}\right)\right] \\
 &\leq \sum_{m=1}^{\lfloor T/K \rfloor} e(2m+1) \frac{Km}{T} e^{-2mx^2} \leq \frac{eK}{T} \sum_{m=1}^{+\infty} m(2m+1) e^{-2mx^2} \leq \frac{eK}{T} \frac{5}{(1 - e^{-2x^2})^3}.
 \end{aligned}$$

Second, by Hoeffding's maximal inequality (Proposition 6),

$$\begin{aligned}
 &\mathbb{P}\left[\exists m \in \{\lceil T/K \rceil, \dots, T\} : \widehat{\mu}_{a^*,m} < \mu^* - x\right] \\
 &= \mathbb{P}\left[\max_{\lceil T/K \rceil \leq m \leq T} \left((1 - \widehat{\mu}_{a^*,m}) - (1 - \mu^*)\right) > x\right] \leq e^{-2\lceil T/K \rceil x^2} \leq e^{-2x^2 T/K}.
 \end{aligned}$$

The proof is concluded by collecting the last two bounds.  $\square$

## Bound on $S_4$

We begin with a now standard use of optional skipping (see Section 4.1), relying on the fact that the events  $\{A_{t+1} = a \text{ and } N_a(t) = n\}$  are disjoint as  $t$  varies:

$$S_4 = \sum_{t=K}^{T-1} \mathbb{P}[U_a^{\text{KL}}(t) \geq \mu^* - \delta \text{ and } A_{t+1} = a \text{ and } N_a(t) \leq f(T, K)] \leq \sum_{n=1}^{f(T, K)} \mathbb{P}[U_{a,n}^{\text{KL}} \geq \mu^* - \delta].$$

We show in this section that

$$\sum_{n=1}^{f(T,K)} \mathbb{P}[U_{a,n}^{\text{KL}} \geq \mu^* - \delta] \leq \frac{1}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \delta/(1 - \mu^*)} \left( W\left(\frac{\ln(1/(1 - \mu^*))}{K} T\right) + \ln(2/(1 - \mu^*)) \right) + 5 + \frac{1}{1 - e^{-\mathcal{K}_{\text{inf}}(\nu, \mu^*)^2/(8\gamma_*)}} \quad (60)$$

where, as in the statement of Proposition 15,

$$\gamma_* = \frac{1}{\sqrt{1 - \mu^*}} \left( 16e^{-2} + \ln^2\left(\frac{1}{1 - \mu^*}\right) \right).$$

To do so, we follow exactly the same method as in the analysis of the IMED policy of Honda and Takemura [2015, Theorem 5]: their idea was to deal with the deviations in a more careful way and relate the sum (60) to the behaviour of a biased random walk.

We start by rewriting the events of interest as

$$\{U_{a,n}^{\text{KL}} \geq \mu^* - \delta\} = \left\{ \mathcal{K}_{\text{inf}}(\hat{\nu}_{a,n}, \mu^* - \delta) \leq \frac{1}{n} \ln\left(\frac{T}{Kn}\right) \right\}$$

where, as in one step of the proof of Lemma 22, we used the definition of  $U_{a,n}^{\text{KL}}$  as well as the left-continuity of  $\mathcal{K}_{\text{inf}}$ . We then follow the same steps as in the proof of Proposition 15 (see Section B.4) and link the deviations in  $\mathcal{K}_{\text{inf}}$  divergence to the ones of a random walk. The variational formula (Lemma 18) for  $\mathcal{K}_{\text{inf}}$  entails the existence of  $\lambda_{a,\delta} \in [0, 1]$  such that

$$\mathcal{K}_{\text{inf}}(\nu_a, \mu^* - \delta) = \mathbb{E} \left[ \ln \left( 1 - \lambda_{a,\delta} \frac{X_a - (\mu^* - \delta)}{1 - (\mu^* - \delta)} \right) \right] \quad \text{where} \quad X_a \sim \nu_a.$$

Note that  $\mathcal{K}_{\text{inf}}(\nu_a, \mu^* - \delta) > 0$  by (7) given that we imposed  $\delta \leq \Delta_a/2$ . We consider i.i.d. copies  $X_{a,1}, \dots, X_{a,n}$  of  $X$  and form the random variables

$$Z_{a,i} = \ln \left( 1 - \lambda_{a,\delta} \frac{X_{a,i} - (\mu^* - \delta)}{1 - (\mu^* - \delta)} \right).$$

By the variational formula (Lemma 18) again, applied this time to  $\mathcal{K}_{\text{inf}}(\hat{\nu}_{a,n}, \mu^* - \delta)$ , we see

$$\mathcal{K}_{\text{inf}}(\hat{\nu}_{a,n}, \mu^* - \delta) \geq \frac{1}{n} \sum_{i=1}^n Z_{a,i}$$

which entails, for each  $n \geq 1$ ,

$$\left\{ \mathcal{K}_{\text{inf}}(\hat{\nu}_{a,n}, \mu^* - \delta) \leq \frac{1}{n} \ln\left(\frac{T}{Kn}\right) \right\} \subseteq \left\{ \sum_{i=1}^n Z_{a,i} \leq \ln\left(\frac{T}{Kn}\right) \right\}. \quad (61)$$

Collecting all previous bounds and inclusions, we proved that the sum of interest (60) is bounded by

$$\begin{aligned} S_4 &\leq \sum_{n=1}^{f(T,K)} \mathbb{P}[U_{a,n}^{\text{KL}} \geq \mu^* - \delta] = \sum_{n=1}^{f(T,K)} \mathbb{P} \left[ \mathcal{K}_{\text{inf}}(\hat{\nu}_{a,n}, \mu^* - \delta) \leq \frac{1}{n} \ln\left(\frac{T}{Kn}\right) \right] \\ &\leq \sum_{n=1}^{f(T,K)} \mathbb{P} \left[ \sum_{i=1}^n Z_{a,i} \leq \ln\left(\frac{T}{Kn}\right) \right] = \mathbb{E} \left[ \sum_{n=1}^{f(T,K)} \mathbb{1}_{\left\{ \sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn)) \right\}} \right] \\ &\leq \mathbb{E} \left[ \sum_{n=1}^T \mathbb{1}_{\left\{ \sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn)) \right\}} \right]. \end{aligned}$$

The last upper bound may seem crude but will be good enough for our purpose.

We may reinterpret

$$\mathbb{E} \left[ \sum_{n=1}^T \mathbb{1}_{\{\sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn))\}} \right]$$

as the expected number of times a random walk with positive drift stays under a decreasing logarithmic barrier. We exploit this interpretation to our advantage by decomposing this sum into the expected hitting time of the barrier and a sum of deviation probabilities for the walk. In what follows,  $\wedge$  denotes the minimum of two numbers. We define the first hitting time  $\tau_a$  of the barrier, if it exists, as

$$\tau_a = \inf \left\{ n \geq 1 : \sum_{i=1}^n Z_{a,i} > \ln \left( \frac{T}{Kn} \right) \right\} \wedge T.$$

The time  $\tau_a$  is bounded by  $T$  and is a stopping time with respect to the filtration generated by the family  $(Z_{a,i})_{1 \leq i \leq n}$ . By distinguishing according to whether or not the condition in the defining infimum of  $\tau_a$  is met for some  $1 \leq n \leq T$ , i.e., whether or not the barrier is hit for  $1 \leq n \leq T$ , we get

$$S_4 \leq \mathbb{E} \left[ \sum_{n=1}^T \mathbb{1}_{\{\sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn))\}} \right] \leq \mathbb{E}[\tau_a] + \mathbb{E} \left[ \sum_{n=\tau_a+1}^T \mathbb{1}_{\{\sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn))\}} \right] \quad (62)$$

where the sum from  $\tau_a + 1$  to  $T$  is void thus null when  $\tau_a = T$  (this is the case, in particular, when the barrier is hit for no  $n \leq T$ ). We now state a lemma, in the spirit of Honda and Takemura [2015, Lemma 18], and will prove it later at the end of this section.

**Lemma 23.** *Let  $(Z_i)_{i \geq 1}$  be a sequence of i.i.d. variables with a positive expectation  $\mathbb{E}[Z_1] > 0$  and such that  $Z_i \leq \alpha$  for some  $\alpha > 0$ . For an integer  $T \geq 1$ , consider the stopping time*

$$\tau \stackrel{\text{def}}{=} \inf \left\{ n \geq 1 : \sum_{i=1}^n Z_i > \ln \left( \frac{T}{Kn} \right) \right\} \wedge T$$

and denote by  $W$  Lambert's function. Then, for all  $T \geq Ke^\alpha$ ,

$$\mathbb{E}[\tau] \leq \frac{W(\alpha T/K) + \alpha + \ln 2}{\mathbb{E}[Z_1]}.$$

The random variables  $Z_{a,i}$  have positive expectation  $\mathcal{K}_{\text{inf}}(\nu_a, \mu^* - \delta) > 0$  and are bounded by  $\alpha = \ln(1/(1 - \mu^*))$ ; indeed, since  $X_{a,i} \geq 0$  and  $\lambda_{a,\delta} \in [0, 1]$ , we have

$$\begin{aligned} Z_{a,i} &= \ln \left( 1 - \lambda_{a,\delta} \frac{X_{a,i} - (\mu^* - \delta)}{1 - (\mu^* - \delta)} \right) \leq \ln \left( 1 + \lambda_{a,\delta} \frac{\mu^* - \delta}{1 - (\mu^* - \delta)} \right) \\ &\leq \ln \left( 1 + \frac{\mu^* - \delta}{1 - (\mu^* - \delta)} \right) = \ln \left( \frac{1}{1 - (\mu^* - \delta)} \right) \leq \ln \left( \frac{1}{1 - \mu^*} \right) \stackrel{\text{def}}{=} \alpha. \end{aligned}$$

In addition, we imposed that  $T > K/(1 - \mu^*) = Ke^\alpha$ . Therefore, Lemma 23 applies and yields the bound

$$\begin{aligned} \mathbb{E}[\tau_a] &\leq \frac{1}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^* - \delta)} \left( W \left( \frac{\ln(1/(1 - \mu^*))}{K} T \right) + \ln(2/(1 - \mu^*)) \right) \\ &\leq \frac{1}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \delta/(1 - \mu^*)} \left( W \left( \frac{\ln(1/(1 - \mu^*))}{K} T \right) + \ln(2/(1 - \mu^*)) \right) \end{aligned}$$

where the second inequality follows by the regularity inequality (19) on  $\mathcal{K}_{\text{inf}}$  (and the denominator therein is still positive thanks to our assumption on  $\delta$ ). All in all, we obtained the first part of the bound (60) and conclude the proof of the latter based on the decomposition (62) by showing that

$$\mathbb{E} \left[ \sum_{n=\tau_a+1}^T \mathbb{1}_{\{\sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn))\}} \right] \leq \beta \stackrel{\text{def}}{=} 5 + \frac{1}{1 - e^{-\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)^2/(8\gamma_*)}}. \quad (63)$$

To that end, note that when  $\tau_a < T$ , we have by definition of  $\tau_a$ ,

$$\ln \left( \frac{T}{K\tau_a} \right) < \sum_{i=1}^{\tau_a} Z_{a,i}$$

The following implication thus holds for any  $n \geq \tau_a$ :

$$\sum_{i=1}^n Z_{a,i} \leq \ln \left( \frac{T}{Kn} \right) \quad \text{implies} \quad \sum_{i=1}^n Z_{a,i} \leq \ln \left( \frac{T}{Kn} \right) \leq \ln \left( \frac{T}{K\tau_a} \right) \leq \sum_{i=1}^{\tau_a} Z_{a,i}. \quad (64)$$

Hence, in this case,

$$\sum_{i=1}^n Z_{a,i} \leq \ln \left( \frac{T}{Kn} \right) \quad \text{implies} \quad \sum_{i=\tau_a+1}^n Z_{a,i} < 0.$$

This, together with a breakdown according to the values of  $\tau_a$  (note that the case  $\tau_a = T$  does not contribute to the expectation) and the independence between  $\{\tau_a = k\}$  and  $Z_{a,k+1}, \dots, Z_{a,T}$ , yields

$$\begin{aligned} & \mathbb{E} \left[ \sum_{n=\tau_a+1}^T \mathbb{1}_{\{\sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn))\}} \right] = \mathbb{E} \left[ \mathbb{1}_{\{\tau_a < T\}} \sum_{n=\tau_a+1}^T \mathbb{1}_{\{\sum_{i=1}^n Z_{a,i} \leq \ln(T/(Kn))\}} \right] \\ & \leq \mathbb{E} \left[ \mathbb{1}_{\{\tau_a < T\}} \sum_{n=\tau_a+1}^T \mathbb{1}_{\{\sum_{i=\tau_a+1}^n Z_{a,i} < 0\}} \right] = \sum_{k=1}^{T-1} \mathbb{E} \left[ \mathbb{1}_{\{\tau_a = k\}} \sum_{n=k+1}^T \mathbb{1}_{\{\sum_{i=k+1}^n Z_{a,i} < 0\}} \right] \\ & = \sum_{k=1}^{T-1} \sum_{n=k+1}^T \mathbb{P}[\tau_a = k] \mathbb{P} \left[ \sum_{i=k+1}^n Z_{a,i} < 0 \right] \\ & = \sum_{k=1}^{T-1} \mathbb{P}[\tau_a = k] \underbrace{\left( \sum_{n=k+1}^T \mathbb{P} \left[ \sum_{i=k+1}^n Z_{a,i} < 0 \right] \right)}_{\text{we show below } \leq \beta, \text{ see (67)}} \leq \beta \end{aligned} \quad (65)$$

where  $\beta$  was defined in (63).

Indeed, we resort to Remark 5 of Section B.4, for the  $n - k$  variables  $Z_{a,k+1}, \dots, Z_{a,n}$  and  $x = 0$ ; we legitimately do so as  $\mu^* - \delta > \mu_a$  by the imposed condition  $\delta < \Delta_a/2$ . Thus, denoting

$$\gamma_{*,\delta} = \frac{1}{\sqrt{1 - (\mu^* - \delta)}} \left( 16e^{-2} + \ln^2 \left( \frac{1}{1 - (\mu^* - \delta)} \right) \right) \leq \gamma_*$$

we have

$$\begin{aligned} \mathbb{P} \left[ \sum_{i=k+1}^n Z_{a,i} \leq 0 \right] & \leq \max \left\{ e^{-(n-k)/4}, \exp \left( -\frac{n-k}{2\gamma_{*,\delta}} \left( \mathcal{K}_{\text{inf}}(\nu_a, \mu^* - \delta) \right)^2 \right) \right\} \\ & \leq e^{-(n-k)/4} + \exp \left( -\frac{n-k}{2\gamma_*} \left( \mathcal{K}_{\text{inf}}(\nu_a, \mu^* - \delta) \right)^2 \right) \\ & \leq e^{-(n-k)/4} + e^{-(n-k)\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)^2/(8\gamma_*)} \end{aligned}$$

where the third inequality follows from (19) and from the imposed condition  $\delta \leq (1 - \mu^*) \mathcal{K}_{\text{inf}}(\nu_a, \mu^*)/2$ :

$$\mathcal{K}_{\text{inf}}(\nu_a, \mu^* - \delta) \geq \mathcal{K}_{\text{inf}}(\nu_a, \mu^*) - \frac{\delta}{1 - \mu^*} \geq \frac{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)}{2}. \quad (66)$$

We finally get, after summation over  $n = k + 1, \dots, T$ ,

$$\sum_{n=k+1}^T \mathbb{P} \left[ \sum_{i=k+1}^n Z_{a,i} \leq 0 \right] \leq \underbrace{\frac{1}{1 - e^{-1/4}}}_{\leq 5} + \frac{1}{1 - e^{-\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)^2/(8\gamma_*)}}, \quad (67)$$

which is the inequality claimed in (65).

It only remains to prove Lemma 23.

**Proof of Lemma 23:** This lemma was almost stated in Honda and Takemura [2015, Lemma 18]: our assumptions and result are slightly different (they are tailored to our needs), which is why we provide below a complete proof, with no significant additional merit compared to the original proof.

We consider the martingale  $(M_n)_{n \geq 0}$  defined by

$$M_n = \sum_{i=1}^n (Z_i - \mathbb{E}[Z_1]).$$

As  $\tau$  is a finite stopping time, Doob's optional stopping theorem indicates that  $\mathbb{E}[M_\tau] = \mathbb{E}[M_0] = 0$ , that is,

$$\mathbb{E}[\tau] \mathbb{E}[Z_1] = \mathbb{E} \left[ \sum_{i=1}^{\tau} Z_i \right].$$

That first step of the proof was exactly similar to the one of Honda and Takemura [2015, Lemma 18]. The idea is now to upper bound the right-hand side of the above equality, which we do by resorting to the very definition of  $\tau$ . An adaptation is needed with respect to the original argument as the value  $\ln(T/(Kn))$  of the barrier varies with  $n$ .

We proceed as follows. Since  $Z_1 \leq \alpha$  and  $T \geq Ke^\alpha$  by assumption, we necessarily have  $\tau \geq 2$ ; using again the boundedness by  $\alpha$ , we have, by definition of  $\tau$ , that

$$\sum_{i=1}^{\tau-1} Z_i \leq \ln \left( \frac{T}{K(\tau-1)} \right)$$

and thus

$$\sum_{i=1}^{\tau-1} Z_i + Z_\tau \leq \ln \left( \frac{T}{K(\tau-1)} \right) + \alpha = \ln \left( \frac{T}{K\tau} \right) + \ln \left( \frac{\tau}{\tau-1} \right) + \alpha \leq \ln \left( \frac{T}{K\tau} \right) + \ln 2 + \alpha.$$

In addition, when  $\tau < T/K$ , and again by definition of  $\tau$ ,

$$\ln \left( \frac{T}{K\tau} \right) < \sum_{i=1}^{\tau} Z_i \leq \tau\alpha \quad \text{thus} \quad 0 < \frac{T}{K\tau} \ln \left( \frac{T}{K\tau} \right) \leq \frac{T\alpha}{K}.$$

Applying the increasing function  $W$  to both sides of the latter inequality, we get, when  $\tau < T/K$ ,

$$\ln \left( \frac{T}{K\tau} \right) \leq W \left( \frac{T\alpha}{K} \right).$$

This inequality also holds when  $\tau \geq T/K$  as the left-hand side then is non-positive, while the right-hand side is positive. Putting all elements together, we successively proved

$$\mathbb{E}[\tau] \mathbb{E}[Z_1] = \mathbb{E} \left[ \sum_{i=1}^{\tau} Z_i \right] \leq W \left( \frac{T\alpha}{K} \right) + \ln 2 + \alpha$$

which concludes the proof.  $\square$



## D. Proof of the variational formula (Lemma 18)

The proof of Honda and Takemura [2015, Theorem 2, Lemma 6] relies on the exhibiting the formula of interest for finitely supported distributions, via KKT conditions, and then taking limits to cover the case of all distributions. We propose a more direct approach that does not rely on discrete approximations of general distributions.

But before we do so, we explain why it is natural to expect to rewrite  $\mathcal{K}_{\text{inf}}$ , which is an infimum, as a maximum. Indeed, given that Kullback-Leibler divergences are given by a supremum,  $\mathcal{K}_{\text{inf}}$  appears as an inf sup, which under some conditions (this is Sion's lemma) is equal to a sup inf.

More precisely, a variational formula for the Kullback-Leibler divergence, see Boucheron et al. [2013, Chapter 4], has it that

$$\text{KL}(\nu, \nu') = \sup \left\{ \mathbb{E}_{\nu}[Y] - \ln \mathbb{E}_{\nu'}[e^Y] : Y \text{ s.t. } \mathbb{E}_{\nu'}[e^Y] < +\infty \right\} \quad (68)$$

where (only in the next few lines) we index the expectation with respect to the assumed distribution of the random variable  $Y$ . In particular, denoting by  $X$  the identity over  $[0, 1]$  and considering, for  $\lambda \in [0, 1]$ , the variables bounded from above

$$Y_{\lambda} = \ln \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \leq \ln \left( 1 + \frac{\lambda \mu}{1 - \mu} \right)$$

we have, for any probability measure  $\nu'$  such that  $\mathbb{E}(\nu') > \mu$ :

$$\ln \mathbb{E}_{\nu'}[e^{Y_{\lambda}}] = \ln \left( \mathbb{E}_{\nu'} \left[ 1 - \lambda \frac{X - \mu}{1 - \mu} \right] \right) = \ln \left( 1 - \lambda \frac{\mathbb{E}(\nu') - \mu}{1 - \mu} \right) \leq 0.$$

Hence, for these distributions  $\nu'$ ,

$$\text{KL}(\nu, \nu') \geq \sup_{\lambda \in [0, 1]} \left\{ \mathbb{E}_{\nu}[Y_{\lambda}] - \ln \mathbb{E}_{\nu'}[e^{Y_{\lambda}}] \right\} \geq \sup_{\lambda \in [0, 1]} \mathbb{E}_{\nu} \left[ \ln \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \right]$$

and by taking the infimum over all distributions  $\nu'$  with  $\mathbb{E}(\nu') > \mu$ :

$$\mathcal{K}_{\text{inf}}(\nu, \mu) \geq \sup_{\lambda \in [0, 1]} \mathbb{E}_{\nu} \left[ \ln \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \right]. \quad (69)$$

**Outline.** We now only need to prove the converse inequality to get the rewriting (50) of Lemma 18, which we will do in Section D.2. Before that, in Section D.1, we prove the second statement of Lemma 18 together with several useful facts for the proof provided in Section D.2, including the fact that the supremum in the right-hand side of (69) is achieved. We conclude in Section D.3 with an alternative (sketch of) proof of the inequality (69), not relying on the variational formula (68) for the Kullback-Leibler divergences.

### D.1. A function study

Let  $X$  denote a random variable with distribution  $\nu \in \mathcal{P}[0, 1]$ . We recall that  $\mu \in (0, 1)$ . The following function is well defined:

$$H : \lambda \in [0, 1] \mapsto \mathbb{E} \left[ \ln \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \right] \in \mathbb{R} \cup \{-\infty\}.$$

Indeed, since  $X \in [0, 1]$ , the random variable  $\ln(1 - \lambda(X - \mu)/(1 - \mu))$  is bounded from above by  $\ln(1 + \lambda\mu/(1 - \mu))$ . Hence,  $H$  is well defined. For  $\lambda \in [0, 1]$ , the considered random variable is bounded

from below by  $\ln(1 - \lambda)$ , hence  $H$  takes finite values. For  $\lambda = 1$ , we possibly have that  $H(1)$  equals  $-\infty$  (this is the case in particular when  $\nu\{1\} > 0$ ).

We begin by a study of the function  $H$ .

**Lemma 24.** *The function  $H$  is continuous and strictly concave on  $[0, 1]$ , differentiable at least on  $[0, 1)$ , and its derivative  $H'(1)$  can be defined at 1, with  $H'(1) \in \mathbb{R} \cup \{-\infty\}$ . We have the closed-form expression: for all  $\lambda \in [0, 1]$ ,*

$$H'(\lambda) = -\mathbb{E} \left[ \left( \frac{X - \mu}{1 - \mu} \right) \frac{1}{1 - \lambda \frac{X - \mu}{1 - \mu}} \right] = \frac{1}{\lambda} \left( 1 - \mathbb{E} \left[ \frac{1}{1 - \lambda \frac{X - \mu}{1 - \mu}} \right] \right). \quad (70)$$

It reaches a unique maximum over  $[0, 1]$ , denoted by  $\lambda^*$ ,

$$\arg \max_{0 \leq \lambda \leq 1} H(\lambda) = \{\lambda^*\}$$

at which  $H'(\lambda^*) = 0$  if  $\lambda^* \in [0, 1)$  and  $H'(\lambda^*) \geq 0$  if  $\lambda^* = 1$ .

Moreover, under the additional condition  $\mathbb{E}(\nu) < \mu$ ,

$$\mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] = 1 \quad \text{if } \lambda^* \in [0, 1) \quad \text{and} \quad \mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] = \mathbb{E} \left[ \frac{1 - \mu}{1 - X} \right] \leq 1 \quad \text{if } \lambda^* = 1$$

where we have in particular  $\nu\{1\} = 0$  in the latter case  $\lambda^* = 1$ .

Note that  $\mathcal{K}_{\text{inf}}(\nu, \mu) = 0$  when  $\mu \leq \mathbb{E}(\nu)$ . In this case, necessarily  $\lambda^* = 0$  (there is a unique maximum) and we still have

$$\mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] = 1.$$

This concludes the proof of the statement (51) of Lemma 18.

**Proof:** For the continuity of  $H$ , we note that the discussion before the statement of the lemma entails that the random variables  $\ln(1 - \lambda(X - \mu)/(1 - \mu))$  are uniformly bounded on ranges of the form  $[0, \lambda_0]$  for  $\lambda_0 < 1$ . By a standard continuity theorem under the integral sign, this proves that  $H$  is continuous on  $[0, 1)$ . For the continuity at 1, we separate the  $H(\lambda)$  and  $H(1)$  into two pieces, for which monotone convergences take place:

$$\begin{aligned} \lim_{\lambda \rightarrow 1} \mathbb{E} \left[ \ln \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \mathbf{1}_{\{X \in [0, \mu]\}} \right] &= \mathbb{E} \left[ \ln \left( \frac{1 - X}{1 - \mu} \right) \mathbf{1}_{\{X \in [0, \mu]\}} \right] \\ \lim_{\lambda \rightarrow 1} \mathbb{E} \left[ \ln \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \mathbf{1}_{\{X \in (\mu, 1]\}} \right] &= \mathbb{E} \left[ \ln \left( \frac{1 - X}{1 - \mu} \right) \mathbf{1}_{\{X \in (\mu, 1]\}} \right] \end{aligned}$$

where the first expectation is finite (but the second may equal  $-\infty$ ).

The strict concavity of  $H$  on  $[0, 1]$  follows from the one of  $\ln$  on  $(0, 1]$  and from the continuity of  $H$  on  $[0, 1]$ .

For  $\lambda \in [0, 1)$ , we get, by legitimately differentiating under the expectation,

$$H'(\lambda) = -\mathbb{E} \left[ \left( \frac{X - \mu}{1 - \mu} \right) \frac{1}{1 - \lambda \frac{X - \mu}{1 - \mu}} \right] = \frac{1}{\lambda} \left( 1 - \mathbb{E} \left[ \frac{1}{1 - \lambda \frac{X - \mu}{1 - \mu}} \right] \right).$$

Indeed as long as  $\lambda < 1$ , the random variables in the expectations are uniformly bounded on ranges of the form  $[0, \lambda_0]$  for  $\lambda_0 < 1$ , so that we may invoke a standard differentiation theorem under the

integral sign. A similar argument of double monotone convergences as above shows that  $H'(\lambda)$  has a limit value as  $\lambda \rightarrow 1$ , with

$$\lim_{\lambda \rightarrow 1} H'(\lambda) = -\mathbb{E} \left[ \frac{X - \mu}{1 - X} \right].$$

By a standard limit theorem on derivatives, when the above value is finite,  $H$  is differentiable at 1 and  $H'(1)$  equals the limit above; otherwise,  $H$  is not differentiable at 1 but we still denote  $H'(1) = -\infty$ .

Since  $H$  is strictly concave on  $[0, 1]$  and continuous, it reaches its maximum exactly once on  $[0, 1]$ . Now, under the condition  $\mu < \mathbb{E}(\nu) < 1$ , we have

$$H'(0) = -\frac{\mathbb{E}(\nu) - \mu}{1 - \mu} > 0.$$

As  $H$  is concave,  $H'$  is decreasing: either  $H'(1) \geq 0$  and  $H$  reaches its maximum at  $\lambda^* = 1$ , or  $H'(1) < 0$  and  $H$  reaches its maximum on the open interval  $(0, 1)$ . It may be proved (by a standard continuity theorem under the integral sign) that  $H'$  is continuous on  $[0, 1)$ , that is, that  $H$  is continuously differentiable on  $[0, 1)$ . In the case  $H'(1) < 0$ , the derivative at the maximum therefore satisfies  $H'(\lambda^*) = 0$ . Substituting the expression for  $H'(\lambda^*)$  concludes the proof.  $\square$

## D.2. Proof of $\leq$ in the equality (50)

We keep the notation introduced in the previous section. To prove this inequality, by the rewriting of  $\mathcal{K}_{\text{inf}}(\nu, \mu)$  stated in Corollary 11, it is enough to show that there exists a probability measure  $\nu'$  on  $[0, 1]$  such that  $\mathbb{E}(\nu') \geq \mu$  and  $\nu \ll \nu'$  and

$$\text{KL}(\nu, \nu') \leq \mathbb{E} \left[ \ln \left( 1 - \lambda^* \frac{X - \mu}{1 - \mu} \right) \right] \quad (71)$$

Given the definition of the KL divergence, it suffices to find a probability measure  $\nu'$  on  $[0, 1]$  such that  $\mathbb{E}(\nu') \geq \mu$  and  $\nu \ll \nu'$  and

$$\frac{d\nu}{d\nu'}(x) = 1 - \lambda^* \frac{x - \mu}{1 - \mu} \quad \nu\text{-a.s.} \quad (72)$$

It can be shown (proof omitted as this statement is only given to explain the intuition behind the proof) that

$$\frac{d\nu}{d\nu'} > 0 \quad \nu\text{-a.s.} \quad \text{with} \quad \frac{d\nu'_{\text{ac}}}{d\nu} = \left( \frac{d\nu}{d\nu'} \right)^{-1} \quad \nu\text{-a.s.} \quad (73)$$

where  $\nu'_{\text{ac}}$  denotes the absolute part of  $\nu'$  with respect to  $\nu$ . This is why we introduce the measure  $\nu'$  on  $[0, 1]$  defined by

$$d\nu'(x) = \underbrace{\frac{1}{1 - \lambda^* \frac{x - \mu}{1 - \mu}}}_{\geq 0} d\nu(x) + \left( 1 - \mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] \right) d\delta_1(x)$$

where  $\delta_1$  denotes the Dirac point-mass distribution at 1 and where  $X$  denotes a random variable with distribution  $\nu$ . The measure  $\nu'$  is a probability measure as by Lemma 24,

$$\mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] \leq 1.$$

Now, we show first that  $\nu \ll \nu'$  with the density (72). We do so by distinguishing two cases. If  $\lambda^* \in [0, 1)$ , then by the last statement of Lemma 24, the probability measure  $\nu'$  is actually defined by

$$d\nu'(x) = \underbrace{\frac{1}{1 - \lambda^* \frac{x-\mu}{1-\mu}}}_{>0} d\nu(x)$$

and the strict positivity underlined in the equality above ensures the desired result by a standard theorem on Radon-Nikodym derivatives. In that case,  $\nu$  and  $\nu'$  are actually equivalent measures:  $\nu \ll \nu'$  and  $\nu' \ll \nu$ . If  $\lambda^* = 1$ , then again by Lemma 24, we know that  $\nu$  does not put any probability mass at 1. The strict positivity of  $f(x) = 1 - (x - \mu)/(1 - \mu)$  on  $[0, 1)$  and the fact that  $\nu\{1\} = 0$  ensure the first equality below: for all Borel sets  $A$  of  $[0, 1]$ ,

$$\nu(A) = \int \mathbf{1}_A f \frac{1}{f} d\nu = \int \mathbf{1}_A f \left( \frac{1}{f} d\nu + r d\delta_1 \right) = \int \mathbf{1}_A f d\nu'$$

while the second equality follows from  $f(1) = 0$  and the third equality is by definition of  $\nu'$ . Put differently,  $\nu \ll \nu'$  with the density  $f$  claimed in (72). In that case,  $\nu \ll \nu'$  but  $\nu'$  is not necessarily absolutely continuous with respect to  $\nu$ .

We conclude this proof by showing that  $\mathbb{E}(\nu') \geq \mu$ . We recall that Lemma 24 indicates that

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{X - \mu}{1 - \mu} \right) \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] &= -H'(\lambda^*) \\ \mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] &= 1 - \lambda^* H'(\lambda^*) \end{aligned}$$

where  $X$  denotes a random variable with distribution  $\nu$  and where both expectations are well defined (possibly with values  $+\infty$  when  $\lambda^* = 1$ ). Therefore,

$$\begin{aligned} \mathbb{E}(\nu') &= \mathbb{E} \left[ \overbrace{\frac{X}{1 - \lambda^* \frac{X - \mu}{1 - \mu}}}^{\text{"}\nu \text{ part of } \nu'"} + \overbrace{\left( 1 - \mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] \right)}^{\text{"}\delta_1 \text{ part of } \nu'"} \right] \\ &= (1 - \mu) \mathbb{E} \left[ \left( \frac{X - \mu}{1 - \mu} \right) \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] + \mu \mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] + \left( 1 - \mathbb{E} \left[ \frac{1}{1 - \lambda^* \frac{X - \mu}{1 - \mu}} \right] \right) \\ &= -(1 - \mu) H'(\lambda^*) + \mu(1 - \lambda^* H'(\lambda^*)) + \lambda^* H'(\lambda^*) = \mu - ((1 - \mu)(1 - \lambda^*) H'(\lambda^*)) \end{aligned}$$

where the first equality is justified in the case  $\lambda^* = 1$  by the same arguments of monotone convergence as in the proof of Lemma 24. All in all, we have  $\mathbb{E}(\nu') \geq \mu$  as desired if and only if  $(1 - \lambda^*) H'(\lambda^*) \leq 0$ . This is the case as we actually have  $(1 - \lambda^*) H'(\lambda^*) = 0$  in all cases, i.e., whether  $\lambda^* = 1$  or  $\lambda^* \in [0, 1)$ .

### D.3. Alternative proof of $\geq$ in the equality (50)

We use the notation of Sections D.1 and D.2 and prove the desired inequality (69), that is, the  $\geq$  part of the equality (50), without resorting to the variational formula (68) for the Kullback-Leibler divergences. Actually, we only provide a sketch of proof and omit proofs of some facts about Radon-Nikodym derivatives.

Let  $\nu'' \in \mathcal{P}[0, 1]$  be such that  $\mathbb{E}(\nu'') > \mu$  and  $\nu \ll \nu''$ ; with no loss of generality, we assume that  $\text{KL}(\nu, \nu'') < +\infty$ . By definition of  $\nu'$ , the divergence  $\text{KL}(\nu, \nu')$  equals the maximum of the continuous function  $H$  over  $[0, 1]$  and therefore also satisfies  $\text{KL}(\nu, \nu') < +\infty$ . We denote by  $\mathbb{L}_1(\nu)$  the set of  $\nu$ -integrable random variables. That these divergences are finite means that

$$\left| \ln \frac{d\nu}{d\nu'} \right| \in \mathbb{L}_1(\nu) \quad \text{and} \quad \left| \ln \frac{d\nu}{d\nu''} \right| \in \mathbb{L}_1(\nu).$$

Hence,

$$\text{KL}(\nu, \nu'') - \text{KL}(\nu, \nu') = - \int \left( \ln \frac{d\nu}{d\nu'} - \ln \frac{d\nu}{d\nu''} \right) d\nu.$$

Now, by (72),

$$\ln \frac{d\nu}{d\nu'}(x) = \ln \left( 1 - \lambda^* \frac{x - \mu}{1 - \mu} \right) \quad \nu\text{-a.s.}$$

and by (73),

$$- \ln \frac{d\nu}{d\nu''} = \ln \frac{d\nu''_{\text{ac}}}{d\nu}(x) \quad \nu\text{-a.s.}$$

so that

$$\begin{aligned} \text{KL}(\nu, \nu'') - \text{KL}(\nu, \nu') &= - \int \ln \left( \left( 1 - \lambda^* \frac{x - \mu}{1 - \mu} \right) \frac{d\nu''_{\text{ac}}}{d\nu}(x) \right) d\nu(x) \\ &\geq - \ln \left( \int \underbrace{\left( 1 - \lambda^* \frac{x - \mu}{1 - \mu} \right)}_{\geq 0} \underbrace{\frac{d\nu''_{\text{ac}}}{d\nu}(x) d\nu(x)}_{d\nu''_{\text{ac}}(x)} \right) \\ &\geq - \ln \left( \int \underbrace{\left( 1 - \lambda^* \frac{x - \mu}{1 - \mu} \right) d\nu''(x)}_{\leq 1 \text{ as } E(\nu'') > \mu} \right) \geq 0 \end{aligned}$$

where Jensen's inequality provided the first inequality, while the second one followed by increasing the integral in the logarithm. Taking the infimum over distributions  $\nu'' \in \mathcal{P}[0, 1]$  with  $E(\nu'') > \mu$  and  $\nu \ll \nu''$  and  $\text{KL}(\nu, \nu'') < +\infty$ , we proved

$$\mathcal{K}_{\text{inf}}(\nu, \mu) - \text{KL}(\nu, \nu') \geq 0$$

which was the desired result.