



**HAL**  
open science

# On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits

Antoine Barrier, Aurélien Garivier, Gilles Stoltz

► **To cite this version:**

Antoine Barrier, Aurélien Garivier, Gilles Stoltz. On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits. 2022. hal-03792668v1

**HAL Id: hal-03792668**

**<https://hal.science/hal-03792668v1>**

Preprint submitted on 30 Sep 2022 (v1), last revised 31 Jan 2023 (v2)

**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.

# On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits

**Antoine Barrier**

*Univ. Lyon, ENS de Lyon, UMPA UMR 5669, Lyon, France*

*Université Paris-Saclay, CNRS, Laboratoire de mathématiques d’Orsay, 91405, Orsay, France*

ANTOINE.BARRIER@ENS-LYON.FR

**Aurélien Garivier**

*Univ. Lyon, ENS de Lyon, UMPA UMR 5669, LIP UMR 5668, Lyon, France*

AURELIEN.GARIVIER@ENS-LYON.FR

**Gilles Stoltz**

*Université Paris-Saclay, CNRS, Laboratoire de mathématiques d’Orsay, 91405, Orsay, France*

GILLES.STOLTZ@UNIVERSITE-PARIS-SACLAY.FR

## Abstract

We lay the foundations of a non-parametric theory of best-arm identification in multi-armed bandits with a fixed budget  $T$ . We consider general, possibly non-parametric, models  $\mathcal{D}$  for distributions over the arms; an overarching example is the model  $\mathcal{D} = \mathcal{P}[0, 1]$  of all probability distributions over  $[0, 1]$ . We propose upper bounds on the average log-probability of misidentifying the optimal arm based on information-theoretic quantities that we name  $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$  and  $\mathcal{L}_{\text{inf}}^{>}(\cdot, \nu)$  and that correspond to infima over Kullback-Leibler divergences between some distributions in  $\mathcal{D}$  and a given distribution  $\nu$ . This is made possible by a refined analysis of the successive-rejects strategy of [Audibert et al. \(2010\)](#). We finally provide lower bounds on the same average log-probability, also in terms of the same new information-theoretic quantities; these lower bounds are larger when the (natural) assumptions on the considered strategies are stronger. All these new upper and lower bounds generalize existing bounds based, e.g., on gaps between distributions.

**Keywords:** Multi-armed bandits, best-arm identification, non-parametric models, Kullback-Leibler divergences, information-theoretic bounds

## 1. Introduction and brief literature review

We consider a class  $\mathcal{D}$  of distributions over  $\mathbb{R}$  with finite first moments, which we refer to as the model  $\mathcal{D}$ . A  $K$ -armed bandit problem in  $\mathcal{D}$  is a  $K$ -tuple  $\underline{\nu} = (\nu_1, \dots, \nu_K)$  of distributions in  $\mathcal{D}$ . We denote by  $(\mu_1, \dots, \mu_K)$  the  $K$ -tuple of their expectations. An agent sequentially interacts with  $\underline{\nu}$ : at each step  $t \geq 1$ , she selects an arm  $A_t$  and receives a reward  $Y_t$  drawn from the distribution  $\nu_{A_t}$ . This is the only feedback that she obtains.

While regret minimization has been vastly studied (see [Lattimore and Szepesvári, 2020](#)), another relevant objective is *best-arm identification*, that is, identifying the distribution with highest expectation. In the fixed-confidence setting, this identification is performed under the constraint that a given confidence level  $1 - \delta$  is respected, while minimizing the expected number of pulls of the arms. This setting is fairly well understood (see [Lattimore and Szepesvári, 2020](#), Chapter 33 for a review). A turning point in this literature was achieved by [Garivier and Kaufmann \(2016\)](#), who provided matching upper and lower bounds on the expected number of pulls of the arms in the case of canonical one-parameter exponential families. Since then, improvements have been made in several directions, including for example non-asymptotic bounds ([Degenne et al., 2019](#)) and the

problem of  $\varepsilon$ -best-arm identification (Garivier and Kaufmann, 2021); however, no generalization to non-parametric models has been considered yet, to the best of our knowledge.

**Best-arm identification with a fixed budget.** The *fixed-budget setting* is much less understood. Therein, the total number  $T$  of pulls of the arms is fixed. After these  $T$  pulls, a strategy must issue a recommendation  $I_T$ . Assuming that  $\underline{\nu}$  contains a unique optimal distribution  $\nu^*$  of index  $a^*(\underline{\nu})$ , one aims at minimizing  $\mathbb{P}(I_T \neq a^*(\underline{\nu}))$ . We are interested in (upper and lower) bounds that hold for all problems  $\underline{\nu}$  in  $\mathcal{D}$ , possibly under the restriction that they only contain a unique optimal arm. It may be straightforwardly seen that the probability of error can decay exponentially fast—for instance, by uniformly exploring the arms (pulling each of them about  $T/K$  times) and recommending the one with the largest empirical average. This is why the literature focuses on upper and lower bound functions  $\ell \leq U < 0$  of the typical form: *for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ , with a unique optimal arm,*

$$\ell(\underline{\nu}) \leq \liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq \limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq U(\underline{\nu}) < 0,$$

$$\text{or, put differently,} \quad \exp\left(\ell(\underline{\nu}) T(1+o(1))\right) \leq \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq \exp\left(U(\underline{\nu}) T(1+o(1))\right).$$

This problem is generally considered more difficult than the fixed-confidence setting, and even for parametric models like canonical one-parameter exponential models, no strategy with matching upper and lower bounds (i.e., no optimal strategy) is known so far. Indeed, even the possibility of obtaining such bounds is disputed.

**Earlier approaches.** So far, four main approaches were considered for the problem of best-arm identification with a fixed budget. *First*, the early approach by Audibert et al. (2010) relies on gaps: we define the gap  $\Delta_a$  of arm  $a$  as the difference  $\mu^* - \mu_a$  between the largest expectation  $\mu^*$  in  $\underline{\nu}$  and the expectation of the distribution  $\nu_a$ . They introduce a successive-rejects strategy and provide gap-based upper bounds for sub-Gaussian models, based on Hoeffding’s inequality. They however propose a lower bound only in the case of a Bernoulli model, not for larger, non-parametric, models. This lower bound was further discussed by Carpentier and Locatelli (2016). *A second series of approaches* (see, e.g., Kaufman et al., 2016) focused on Gaussian bandits with fixed variances, but results do not seem to be easily generalized to other models as they often rely on specific facts (related, among others, to the symmetry of the Kullback-Leibler divergence). *A third approach*, led by Russo (2016, 2020), considered canonical one-parameter exponential families, but a for a different target probability. Namely, a Bayesian setting is considered and the quality of a strategy is measured as the posterior probability of identifying the best arm. An optimal non-gap-based complexity is exhibited, together with optimal strategies matching this complexity. However, Komiyama (2022) argue that such an approach is specific to the Bayesian case and is not suited to the frequentist case that we consider. *A fourth approach* is to focus on the case of  $K = 2$  arms, see, e.g., Kaufman et al. (2016). The non-parametric bounds obtained therein do not enjoy any obvious generalization to the case of  $K \geq 3$  arms. By considering very specific models, Kato et al. (2022) constructed a strategy that is optimal (only) in the regime where the gap between the 2 arms is small—yet, this gap-based approach does not, by nature, go in the direction of non-parametric bounds.

We will provide more details concerning some of these approaches while presenting and discussing our main results, in Section 2.2; see also Appendix E.

**Content and outline of this article.** We focus our attention on general upper and lower bounds, holding for all possible models  $\mathcal{D}$ , including non-parametric models, and valid for any number  $K$  of arms. Put differently, we target a high degree of generality. While admittedly not exhibiting matching upper and lower bounds, we show that the same information-theoretic quantities  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^>$  are at stake in these upper and lower bounds. They are defined, in Section 2, as infima of Kullback-Leibler divergences and provide a quantification of the difficulty of the identification in terms of the geometry of information of the problem. We also provide in this section an overview of our results, which we carefully compare to existing bounds (restated therein, occasionally with some improvements). We study in Section 3 the classical successive-rejects strategy and provide an improved analysis, not relying on gaps through Hoeffding’s lemma. Section 4 exhibits several possible lower bounds, which are inversely larger to the strength of the assumptions made on the strategies. These lower bounds generalize known lower bounds in the literature, like the lower bound for Bernoulli models by [Audibert et al. \(2010\)](#), but hold for arbitrary models. They share some similar flavor with the lower bounds by [Lai and Robbins \(1985\)](#) and [Burnetas and Katehakis \(1996\)](#) for the cumulative regret.

## 2. Overview of the results and more extended literature review

Before being able to actually provide a formal summary of our results, we introduce new quantifications of the difficulty of a bandit problem in terms of geometry of the information. These quantifications should be reminiscent of the  $\mathcal{K}_{\text{inf}}$  quantity that appears in the optimal bounds on the cumulative regret; it was introduced by [Honda and Takemura \(2015\)](#), see also [Garivier et al. \(2022\)](#). Denoting by  $\nu$  the distribution of interest, this  $\mathcal{K}_{\text{inf}}$  is defined as an infimum over divergences of the form  $\text{KL}(\nu, \zeta)$ . For our objectives, the arguments in the KL are in reverse order, and we are rather interested in infima over divergences of the form  $\text{KL}(\zeta, \nu)$ .

Except for very specific models (e.g., the model  $\mathcal{D}_{\sigma^2}$  of Gaussian distributions with a fixed variance  $\sigma^2 > 0$ ), the Kullback-Leibler divergence is not symmetric, i.e.,  $\text{KL}(\zeta, \nu)$  and  $\text{KL}(\nu, \zeta)$  differ in general. Specific best-arm-identification results were obtained by [Kaufman et al. \(2016\)](#) for the model  $\mathcal{D}_{\sigma^2}$ , based on the Bretagnolle-Huber inequality; they indicate that the sum of the inverse squared gaps would be driving both the lower bound and upper bound functions  $\ell$  and  $U$ . However, a close look at the proof reveals that they heavily rely (among others) on the symmetry of KL for this model. There is therefore absolutely no hope to provide any generalization beyond the Gaussian case. Appendix E.2 provides further details and discussions on this matter.

### 2.1. The key new quantities: $\mathcal{L}_{\text{inf}}^<$ and $\mathcal{L}_{\text{inf}}^{\leq}$ , as well as $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$

In this article, we only consider models  $\mathcal{D}$  whose distributions all admit an expectation. We denote by  $\mathbb{E}(\zeta)$  the expectation of a distribution  $\zeta \in \mathcal{D}$ . For a distribution  $\nu \in \mathcal{D}$  and a real number  $x \in \mathbb{R}$ , we then introduce

$$\begin{aligned} \mathcal{L}_{\text{inf}}^<(x, \nu) &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) < x \} \\ \text{and } \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) \leq x \}, \end{aligned}$$

where KL denotes the Kullback-Leibler divergence and with the usual convention that the infimum of an empty set equals  $+\infty$ . Symmetrically, by considering rather distributions  $\zeta$  with expectations larger than  $x$ , we define

$$\begin{aligned} \mathcal{L}_{\text{inf}}^>(x, \nu) &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) > x \} \\ \text{and } \mathcal{L}_{\text{inf}}^{\geq}(x, \nu) &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } \mathbb{E}(\zeta) \geq x \}. \end{aligned}$$

We state some general properties on these quantities in Appendix A—among others, that  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^{\leq}$ , as well as  $\mathcal{L}_{\text{inf}}^>$  and  $\mathcal{L}_{\text{inf}}^{\geq}$ , are almost identical for the model  $\mathcal{P}[0, 1]$ . The same holds for canonical one-parameter exponential models, as discussed in Appendix C.3.

Lower bounds will be typically expressed with  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^>$  quantities, while upper bounds will rely on  $\mathcal{L}_{\text{inf}}^{\leq}$  and  $\mathcal{L}_{\text{inf}}^{\geq}$  quantities.

## 2.2. Overview of the results

The paper provides new and more general (possibly non-parametric) bounds based on information-theoretic quantities, namely the quantities  $\mathcal{L}_{\text{inf}}^<$ ,  $\mathcal{L}_{\text{inf}}^{\leq}$ ,  $\mathcal{L}_{\text{inf}}^>$  and  $\mathcal{L}_{\text{inf}}^{\geq}$ . In particular, we consider a version of Chernoff information defined, for  $\nu, \nu'$  in  $\mathcal{D}$  with  $\mathbb{E}(\nu') < \mathbb{E}(\nu)$ , as

$$\mathcal{L}(\nu', \nu) = \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ \mathcal{L}_{\text{inf}}^{\geq}(x, \nu') + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \right\}.$$

Given a bandit problem  $\underline{\nu}$  with a unique optimal distribution denoted by  $\nu^*$ , we may rank the arms  $a$  in non-decreasing order of  $\mathcal{L}(\nu_a, \nu^*)$ , i.e., consider the permutation  $\sigma$  such that

$$0 = \mathcal{L}(\nu_{\sigma_1}, \nu^*) < \mathcal{L}(\nu_{\sigma_2}, \nu^*) \leq \dots \leq \mathcal{L}(\nu_{\sigma_{K-1}}, \nu^*) \leq \mathcal{L}(\nu_{\sigma_K}, \nu^*).$$

*Our first main result* (Corollary 3 together with Lemma 4) considers models  $\mathcal{D}$  like  $\mathcal{D} = \mathcal{P}[0, 1]$ , the set of all probability distributions over  $[0, 1]$ , or  $\mathcal{D} = \mathcal{D}_{\text{exp}}$ , any canonical one-parameter exponential family. It studies the successive-rejects strategy, introduced by Audibert et al. (2010), for which arms are rejected one by one at the end of phases of uniform exploration, and states that this strategy is such that for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}$  with a unique optimal arm,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\ln K} \min_{2 \leq k \leq K} \frac{\mathcal{L}(\nu_{\sigma_k}, \nu^*)}{k}, \quad (1)$$

where  $\overline{\ln K} \approx \ln K$ . The key for this result (Lemma 1, of independent interest) is a grid-based application of the Cramér-Chernoff bound to control  $\mathbb{P}(\overline{X}_N \leq \overline{Y}_N)$ , where  $\overline{X}_N$  and  $\overline{Y}_N$  are averages of two independent  $N$ -samples. This approach can be used to analyze similar algorithms, like sequential halving (Karnin et al., 2013).

The corresponding lower bounds are stated rather in terms of  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^>$  quantities, but Appendix A explains why, except in a single pathological case,  $\mathcal{L}(\nu', \nu)$  could be alternatively defined with  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^>$  instead of  $\mathcal{L}_{\text{inf}}^{\leq}$  and  $\mathcal{L}_{\text{inf}}^{\geq}$ . We actually state several lower bounds in Section 4, that are larger as the assumptions on the strategies considered are more restrictive; as usual, there is a trade-off between the strength of a lower bound and its generality. However, all assumptions considered remain rather mild and are satisfied by successive-rejects-type strategies: e.g., one may restrict the attention to strategies such that for all bandit problems, the arm associated with the smallest expectation is pulled less than a fraction  $1/K$  of the time. *Our second main result* (Theorem 12) gives the following alma matter lower bound: for all bandit problems  $\underline{\nu}$  with no two same expectations,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq -\min_{2 \leq k \leq K} \inf_{x \in [\mu_{(k)}, \mu_{(k-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{k-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{k} \right\}, \quad (2)$$

where we considered the order statistics in reverse order,  $\mu_{(1)} > \mu_{(2)} > \mu_{(3)} > \dots > \mu_{(K)}$ , and where  $\nu_{(a)}$  denotes the distribution with expectation  $\mu_{(a)}$ .

This lower bound does not match the exhibited upper bound, mainly because the infima in (2) are only taken on restricted ranges  $[\mu_{(k)}, \mu_{(k-1)})$  and not on the entire intervals  $[\mu_{(k)}, \mu^*]$  as in (1). Still, we argue that quantities defined as infima over  $x$  of  $\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)}) + \mathcal{L}_{\text{inf}}^<(x, \nu^*)$  should measure how difficult a best-arm-identification problem is under a fixed budget. *This is the main insight of this article.*

We now compare our general bounds to existing bounds, for sub-Gaussian models and for exponential models. To do so, we will sometimes consider the following consequence of the lower bound (2), obtained by picking  $x = \mu_{(k)}$ :

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \frac{\mathcal{L}_{\text{inf}}^<(\mu_{(k)}, \nu^*)}{k}. \quad (3)$$

**Comparison to the gap-based approaches.** Audibert et al. (2010) propose an analysis of the successive-rejects strategy based on Hoeffding's inequality, stating that for all bandit problems in  $\mathcal{P}[0, 1]$  with a unique optimal arm,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \leq - \frac{1}{\ln K} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}, \quad (4)$$

where we recall the definition of the gaps  $\Delta_{(k)} = \mu^* - \mu_{(k)}$ . This bound is a consequence of (Corollary 3, a slightly more general form of the) bound (1), given Pinsker's inequality (18):

$$\mathcal{L}(\nu_{(k)}, \nu^*) \geq \inf_{x \in [\mu_{(k)}, \mu^*]} \left\{ 2(x - \mu_{(k)})^2 + 2(x - \mu^*)^2 \right\} = (\mu^* - \mu_{(k)})^2 = \Delta_{(k)}^2. \quad (5)$$

We remark that the bound (4) and the lower bound on  $\mathcal{L}(\nu_{(k)}, \nu^*)$  may actually be extended to the model of  $\sigma^2$ -sub-Gaussian distributions, up to dividing the bound (4) by a factor  $4\sigma^2$ . We do not discuss the UCB-E algorithm of Audibert et al. (2010), as its performance and analysis crucially depend on a tuning parameter set with some knowledge of the gaps.

Audibert et al. (2010) also propose a carefully constructed lower bound for the model  $\mathcal{D}_p = \{\text{Ber}(x) : x \in [p, 1-p]\}$  of Bernoulli distributions  $\text{Ber}(x)$  with parameters  $x$  in  $[p, 1-p]$  for some  $p \in (0, 1/2)$ . A key inequality in their proof follows from the Kullback-Leibler  $-\chi^2$ -divergence bound:

$$\forall x, y \in [p, 1-p], \quad \text{KL}(\text{Ber}(x), \text{Ber}(y)) \leq \frac{(x-y)^2}{2p(1-p)}.$$

Their construction may actually be generalized to models  $\mathcal{D}$  with  $C_{\mathcal{D}} > 0$  such that for all  $\nu, \nu'$  in  $\mathcal{D}$ , one has  $\text{KL}(\nu, \nu') \leq C_{\mathcal{D}}(\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2$ . This is a property that clearly holds for some exponential problems: on top of the restricted Bernoulli model discussed above, for which  $C_{\mathcal{D}_p} = 1/(2p(1-p))$ , we may cite the model  $\mathcal{D}_{\sigma^2}$  of Gaussian distributions with variance  $\sigma^2$ , for which  $C_{\mathcal{D}_{\sigma^2}} = 1/(2\sigma^2)$ . For models enjoying the existence of such a constant  $C_{\mathcal{D}}$ , (a straightforward modification of) the analysis by Audibert et al. (2010) entails that for any  $\underline{\nu}$  in  $\mathcal{D}$ ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq -5C_{\mathcal{D}} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}. \quad (6)$$

As by the very assumption on the model,  $\mathcal{L}_{\text{inf}}^<(\mu_{(k)}, \nu^*) \leq C_{\mathcal{D}} \Delta_{(k)}^2$ , the lower bound (3) implies the stated lower bound (6), with an improved constant factor.

The lower bound (6) and the upper bound (4) differ in particular by a factor of about  $\approx \ln K$ . [Carpentier and Locatelli \(2016\)](#) discuss this gap in the case of a restricted Bernoulli model and improve the lower bound (6) by a factor of  $\ln K$ , but not simultaneously for all bandit problems  $\underline{\nu}$  (as we aim for); they obtain the improvement just for one bandit problem  $\underline{\nu}$ . Their lower bound result (formally stated and discussed in [Appendix E.1](#)) is therefore of a totally different nature. More results on how and when given lower bounds with a given complexity measure may, or may not, be improved were stated by [Komiyama et al. \(2022\)](#).

**Discussion of the non-parametric bound for  $K = 2$  arms of [Kaufman et al. \(2016\)](#).** It turns out that the existing literature offered so far a non-parametric bound, in the case of  $K = 2$  arms. Namely, in a general, possibly non-parametric model  $\mathcal{D}$ , [Kaufman et al. \(2016, Theorem 12\)](#) stated a lower bound for all 2-armed bandit problems  $\underline{\nu} = (\nu_1, \nu_2)$ :

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \inf_{\substack{\lambda \text{ in } \mathcal{D}: \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < \mathbb{E}(\lambda_{w_*(\underline{\nu})})}} \max \left\{ \text{KL}(\lambda_{w_*(\underline{\nu})}, \nu_{w_*(\underline{\nu})}), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\}, \quad (7)$$

where  $w_*(\underline{\nu})$  denotes the suboptimal arm in  $\underline{\nu}$  and where the infimum is over all alternative bandit problems  $(\lambda_1, \lambda_2)$  in  $\mathcal{D}$  with inverse order on the expectations compared to  $\underline{\nu}$ . We note (see the proof of [Theorem 13](#)) that we may actually rewrite this lower bound in a more readable way, in terms of  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^>$  quantities, illustrating once again that these quantities are key in measuring the complexity of best-arm-identification under a fixed budget:

$$\begin{aligned} \inf_{\substack{\lambda \text{ in } \mathcal{D}: \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < \mathbb{E}(\lambda_{w_*(\underline{\nu})})}} \max \left\{ \text{KL}(\lambda_{w_*(\underline{\nu})}, \nu_{w_*(\underline{\nu})}), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\} \\ = \inf_{x \in [\mu_{w_*(\underline{\nu})}, \mu^*]} \left\{ \max \left\{ \mathcal{L}_{\text{inf}}^>(x, \nu_{w_*(\underline{\nu})}), \mathcal{L}_{\text{inf}}^<(x, \nu^*) \right\} \right\}. \end{aligned} \quad (8)$$

The proof technique of [Kaufman et al. \(2016\)](#) may be applied in a pairwise fashion to generalize the above lower bound for 2 arms into a lower bound for  $K \geq 2$  arms, stated in [Theorem 13](#): for all  $\underline{\nu}$  in  $\mathcal{D}$  with a unique optimal arm,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{k \neq a^*(\underline{\nu})} \inf_{x \in [\mu_k, \mu^*]} \left\{ \max \left\{ \mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*) \right\} \right\}.$$

We however do not claim that this is a deep and interesting bound, as it only involves pairwise comparisons with the best arm. In particular, we lack divisions by the ranks of the arms, as in (2).

**Bounds for exponential models.** We denote by  $\mathcal{D}_{\text{exp}}$  the model corresponding to a canonical one-parameter exponential family with expectations defined on an open interval  $\mathcal{M}$  (see [Appendix C.3](#) for a reminder on this matter). For such a model, we denote by  $d$  the mean-parameterized Kullback-Leibler divergence. By continuity of  $d$ , we have that for all  $\nu$  in  $\mathcal{D}_{\text{exp}}$  and for all  $x \in \mathcal{M}$ ,

$$\forall x \leq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^<(x, \nu) = \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = d(x, \mathbb{E}(\nu)), \quad (9)$$

$$\text{and } \forall x \geq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^>(x, \nu) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu) = d(x, \mathbb{E}(\nu)). \quad (10)$$

All bounds above then admit simple reformulations in terms of  $d$ . The version  $\mathcal{L}$  of Chernoff information we introduced may also be mean-parameterized: for  $\mu' < \mu$ ,

$$L(\mu', \mu) = \min_{x \in [\mu', \mu]} \left\{ d(x, \mu') + d(x, \mu) \right\}.$$



**ALGORITHM: SUCCESSIVE-REJECTS STRATEGY**

**Parameters:**  $K$  arms, budget  $T$ , lengths  $\ell_1, \dots, \ell_{K-1} \geq 1$  with  $\ell_1 + \dots + \ell_{K-1} = T$

**Initialization:**  $S_0 = \{1, \dots, K\}$

**For each regime:**  $r \in \{1, \dots, K-1\}$ :

1. For each arm  $a \in S_{r-1}$ 
  - (a) Pull it  $\lfloor \ell_r / (K-1+r) \rfloor$  times
  - (b) Compute the empirical average  $\bar{X}_a^r$  of the payoffs obtained in this regime and in the previous regimes
2. Drop the arm  $a_r$  with smallest average (ties broken arbitrarily):

$$S_r = S_{r-1} \setminus \{a_r\}, \quad \text{where} \quad a_r \in \operatorname{argmin}_{a \in S_{r-1}} \bar{X}_a^r$$

**Output:** Recommend arm  $I_T$ , where  $S_{K-1} = \{I_T\}$

The original definition of the Chernoff information  $D(\mu', \mu)$  is the value  $d(y, \mu)$  for  $y \in [\mu', \mu]$  such that  $d(y, \mu') = d(y, \mu)$ . This is the quantity at stake in (8): given that  $d(\cdot, \mu')$  and  $d(\cdot, \mu)$  are respectively increasing and decreasing on  $[\mu', \mu]$ ,

$$\min_{x \in [\mu', \mu]} \max\{d(x, \mu'), d(x, \mu)\} = D(\mu', \mu).$$

Therefore,  $D(\mu', \mu) \leq L(\mu', \mu) \leq 2D(\mu', \mu)$ . This justifies that we called  $L$  (and therefore  $\mathcal{L}$ ) a version of Chernoff information.

### 3. Upper bound: successive-rejects strategy, with an improved analysis

We consider the successive-rejects strategy introduced by Audibert et al. (2010), for  $K$  arms and a budget  $T$ . Regime lengths are set beforehand; they are denoted by  $\ell_1, \dots, \ell_{K-1} \geq 1$  and satisfy  $\ell_1 + \dots + \ell_{K-1} = T$ . The strategy maintains a list of candidate arms, starting with all arms, i.e.,  $S_0 = \{1, \dots, K\}$ . At the end of each regime  $r \in \{1, \dots, K-1\}$ , it drops an arm to get  $S_r$ , while during the regime  $r$ , it operates with the  $K-r+1$  arms in  $S_{r-1}$ .

More precisely, during regime  $r \in \{1, \dots, K-1\}$ , the strategy draws  $\lfloor \ell_r / (K-r+1) \rfloor$  times each arm in  $S_{r-1}$  (and does not use the few remaining time steps, if there are some). At the end of each regime  $r$ , the strategy computes the empirical averages  $\bar{X}_a^r$  of the payoffs obtained by each arm  $a \in S_{r-1}$  since the beginning; i.e.,  $\bar{X}_a^r$  is an average over

$$N_r = \lfloor \ell_1 / K \rfloor + \dots + \lfloor \ell_r / (K-r+1) \rfloor$$

i.i.d. realizations of  $\nu_a$ . It then drops the arm  $a_r$  with smallest empirical average (ties broken arbitrarily). This description is summarized in the algorithm box.



### 3.1. General analysis

The key quantities for the general analysis will be the logarithmic moment-generating function  $\phi_\nu$  of a distribution  $\nu \in \mathcal{D}$ , and its Fenchel-Legendre transform  $\phi_\nu^*$ :

$$\forall \lambda \in \mathbb{R}, \quad \phi_\nu(\lambda) = \ln \int_{\mathbb{R}} e^{\lambda x} d\nu(x) \quad \text{and} \quad \forall x \in \mathbb{R}, \quad \phi_\nu^*(x) = \sup_{\lambda \in \mathbb{R}} \{\lambda x - \phi_\nu(\lambda)\}. \quad (11)$$

Based on them, we can now define, for all  $\nu, \nu' \in \mathcal{D}$  with  $E(\nu') < E(\nu)$ ,

$$\Phi(\nu', \nu) \stackrel{\text{def}}{=} \inf_{x \in [E(\nu'), E(\nu)]} \{\phi_{\nu'}^*(x) + \phi_\nu^*(x)\}.$$

The following simple lemma shows that  $\Phi$  plays a significant role for bounding the probability that two sample averages are in inverse order compared to the expectations of the underlying distributions. It supersedes the use of Hoeffding's inequality in [Audibert et al. \(2010\)](#).

**Lemma 1** *Fix  $\nu$  and  $\nu'$  in  $\mathcal{D}$ , with respective expectations  $\mu = E(\nu) > \mu' = E(\nu')$ . For all  $N \geq 1$ , let  $\bar{X}_N$  and  $\bar{Y}_N$  be the averages of  $N$ -samples with respective distributions  $\nu$  and  $\nu'$ . Then,*

$$\limsup_{N \rightarrow +\infty} \frac{1}{N} \ln \mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq - \inf_{x \in [\mu', \mu]} \{\phi_{\nu'}^*(x) + \phi_\nu^*(x)\} \stackrel{\text{def}}{=} -\Phi(\nu', \nu).$$

**Proof sketch** The fact  $\bar{X}_N \leq \bar{Y}_N$  entails the existence of  $x$  such that  $\bar{X}_N \leq x \leq \bar{Y}_N$ . By independence, together with two applications of the Cramér-Chernoff bound (recalled in [Appendix B.1](#)),

$$\mathbb{P}(\bar{X}_N \leq x \leq \bar{Y}_N) = \mathbb{P}(\bar{X}_N \leq x) \mathbb{P}(x \leq \bar{Y}_N) \leq \exp(-N \phi_\nu^*(x)) \exp(-N \phi_{\nu'}^*(x)).$$

The technical issue is then to deal with some union over  $x$  of the events  $\{\bar{X}_N \leq x \leq \bar{Y}_N\}$ . We do so with a sequence of finite grids, with vanishing steps, and use lower-semi-continuity arguments to obtain an infimum over an interval based on a sequence of finite minima. A complete proof is to be found in [Appendix B.2](#). ■

The main performance upper bound is stated below in terms of  $\Phi$ , that is, in terms of Fenchel-Legendre transforms of logarithmic moment-generating functions. [Section 3.2](#) will later explain why and when the latter may be replaced by  $\mathcal{L}_{\text{inf}}^{\leq}$  and  $\mathcal{L}_{\text{inf}}^{\geq}$  quantities, leading to a rewriting  $\Phi = \mathcal{L}$  and to the bound claimed in [\(1\)](#).

**Theorem 2** *Fix  $K \geq 2$  and a model  $\mathcal{D}$ . Consider a sequence of successive-rejects strategies, indexed by  $T$ , such that  $N_r/T \rightarrow \gamma_r > 0$  as  $T \rightarrow +\infty$  for all  $r \in \{1, \dots, K-1\}$ . Let  $\underline{\nu}$  be a bandit problem in  $\mathcal{D}$  with a unique optimal arm and, for each  $r \in \{1, \dots, K-1\}$ , let  $\mathcal{A}_r$  be a subset of arms of cardinality  $r$  that does not contain  $a^*(\underline{\nu})$ . Then*

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq - \min_{1 \leq r \leq K-1} \left\{ \gamma_r \min_{k \in \mathcal{A}_r} \Phi(\nu_k, \nu^*) \right\}.$$

**Proof sketch** A complete proof may be found in [Appendix B.3](#); it mimics the analysis by [Audibert et al. \(2010\)](#), the main modification being the substitution of Hoeffding's inequality by the bound of [Lemma 1](#). We have  $I_T \neq a^*(\underline{\nu})$  if and only if  $a^*(\underline{\nu})$  is rejected in some regime, i.e.,

$$\{I_T \neq a^*(\underline{\nu})\} = \bigcup_{r=1}^{K-1} \{a_r = a^*(\underline{\nu})\} \subseteq \bigcup_{r=1}^{K-1} \left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{X}_{a^*(\underline{\nu})}^r \leq \bar{X}_k^r \right\}.$$

By optional skipping (see [Doob, 1953](#), Chapter III, Theorem 5.2, p. 145) and by the fact that by the pigeonhole principle, the (random) set  $S_{r-1}$  necessarily contains one element of the deterministic set  $\mathcal{A}_r$ ,

$$\mathbb{P}\left(a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{X}_{a^*(\underline{\nu})}^r \leq \bar{X}_k^r\right) \leq \sum_{k \in \mathcal{A}_r} \mathbb{P}\left(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r\right),$$

where, for all  $a$ , the  $\bar{Y}_a^r$  are the averages of independent  $N_r$ -samples distributed according to  $\nu_a$ . The proof is concluded by Lemma 1 and the fact that a sum of exponentially fast decaying quantities is driven by its largest term.  $\blacksquare$

We conclude this subsection by stating the bound of Theorem 2 for the regime lengths suggested by [Audibert et al. \(2010\)](#), namely,  $\ell_1 = T/\bar{\ln} K$  and for  $r \in \{2, \dots, K-1\}$ ,

$$\ell_r = \frac{T}{(K-r+2)\bar{\ln} K}, \quad \text{where} \quad \bar{\ln} K = \frac{1}{2} + \sum_{k=2}^K \frac{1}{k}. \quad (12)$$

We also consider lower bounds  $f(\nu_k, \nu^*)$  on the  $\Phi(\nu_k, \nu^*)$ . We may of course use  $f = \Phi$  but sometimes, it is handy to rely on more readable lower bounds. For instance, in the case of the  $\mathcal{P}[0, 1]$  model, Hoeffding's inequality entails that

$$\phi_\nu^*(x) \geq 2(x - \mathbb{E}(\nu))^2, \quad \text{so that} \quad \Phi(\nu_k, \nu^*) \geq \Delta_k^2 \stackrel{\text{def}}{=} f(\nu_k, \nu^*); \quad (13)$$

see more details in Appendix B.4. Such bounds hold more generally in models consisting of sub-Gaussian distributions.

We now order the arms into  $\sigma_1, \dots, \sigma_K$  based on  $f$ , namely, we let  $\sigma_1 = a^*(\underline{\nu})$  and

$$0 = f(\nu_{\sigma_1}, \nu^*) < f(\nu_{\sigma_2}, \nu^*) \leq \dots \leq f(\nu_{\sigma_{K-1}}, \nu^*) \leq f(\nu_{\sigma_K}, \nu^*), \quad (14)$$

and we take  $\mathcal{A}_r = \{\sigma_{K-r+1}, \dots, \sigma_K\}$ . We obtain immediately the following corollary, for which a detailed proof may be found, for the sake of completeness, in Appendix B.4.

**Corollary 3** *Fix  $K \geq 2$ , a model  $\mathcal{D}$ , and consider a lower bound  $f$  on  $\Phi$ . The sequence of successive-rejects strategies based on the regime lengths (12) ensures, that for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}$  with a unique optimal arm,*

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\bar{\ln} K} \min_{2 \leq k \leq K} \frac{f(\nu_{\sigma_k}, \nu^*)}{k},$$

where arms were reordered as in (14).

### 3.2. On links between $\Phi$ and the quantities $\mathcal{L}_{\text{inf}}^<$ , $\mathcal{L}_{\text{inf}}^{\leq}$ , $\mathcal{L}_{\text{inf}}^>$ and $\mathcal{L}_{\text{inf}}^{\geq}$

The Fenchel-Legendre transform  $\phi_\nu^*$  of the logarithmic moment-generating function of  $\nu$  admits a classical (see, e.g., [Boucheron et al., 2013](#), Exercice 4.13) dual formulation in terms of infima of Kullback-Leibler divergences. The following lemma, proved in Appendix C.2, reveals that these infima correspond to  $\mathcal{L}_{\text{inf}}^{\leq}$  and  $\mathcal{L}_{\text{inf}}^{\geq}$  for the model  $\mathcal{P}[0, 1]$  of distributions supported on  $[0, 1]$ .

**Lemma 4** Consider the model  $\mathcal{D} = \mathcal{P}[0, 1]$ . For all  $\nu \in \mathcal{P}[0, 1]$ ,

$$\forall x \leq E(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\inf}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq E(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\inf}^{\geq}(x, \nu).$$

Based on this lemma, we have the following rewriting, which is useful to reinterpret the quantities appearing in Theorem 2 and Corollary 3:  $\Phi(\nu', \nu) = \mathcal{L}(\nu', \nu)$  for the model  $\mathcal{P}[0, 1]$ , i.e.,

$$\inf_{x \in [E(\nu'), E(\nu)]} \{\phi_{\nu'}^*(x) + \phi_\nu^*(x)\} = \inf_{x \in [E(\nu'), E(\nu)]} \{\mathcal{L}_{\inf}^{\geq}(x, \nu') + \mathcal{L}_{\inf}^{\leq}(x, \nu)\}. \quad (15)$$

For canonical one-parameter exponential models  $\mathcal{D}_{\text{exp}}$ , a slightly weaker version of Lemma 4, only holding for  $x$  corresponding to expectations in  $\mathcal{D}_{\text{exp}}$  and provided in Appendix C.3, similarly shows (15), i.e.,  $\Phi = \mathcal{L}$ . Conditions on general models for  $\Phi = \mathcal{L}$  to hold are discussed in Appendix C.4.

#### 4. Lower bounds

In most of this section, we restrict our attention to generic  $K$ -armed bandit problems  $\underline{\nu}$ , that are such that  $\mu_j \neq \mu_k$  for  $j \neq k$ . In particular, the best arm  $a^*(\underline{\nu})$  is unique.

**Definition of a strategy, and of a (doubly-indexed) sequence of strategies.** A strategy  $(\psi, \varphi)$  depends on the budget  $T$  and the number  $K$  of arms; it consists of a sampling scheme  $\psi = (\psi_t)_{1 \leq t \leq T}$  and a recommendation function  $\varphi$ . At each round  $t \in \{1, \dots, T\}$ , the strategy picks an arm  $A_t$ , possibly at random using an auxiliary randomization  $U_{t-1}$ . Given this choice  $A_t$ , the strategy observes a payoff  $Y_t$  drawn at random according to  $\nu_{A_t}$ , independently from the past. For  $t \geq 2$ , the choice  $A_t$  is therefore a measurable function  $A_t = \psi_t(H_t)$  of the history  $H_t = (U_0, Y_1, \dots, Y_{t-1}, U_{t-1})$ , while  $A_1 = \psi_1(H_0)$ , where  $H_0 = U_0$ . At round  $T$ , the strategy recommends the arm  $I_T = \varphi(H_T)$ .

For our lower bounds, we will consider sequences of strategies, either only indexed by  $T \geq 1$  given a value of  $K \geq 2$ , or doubly indexed by  $T$  and  $K$ . These sequences will also be assumed to be “reasonable” in the sense below.

**Consistent (or exponentially consistent) sequences of strategies.** The probability  $\mathbb{P}(I_T \neq a^*(\underline{\nu}))$  of misidentifying the unique optimal arm may vanish asymptotically (and even vanish exponentially fast) for all bandit problems—in not too large a model  $\mathcal{D}$ , as illustrated in Section 3. We will therefore only be interested in such sequences of strategies, called (exponentially) consistent. In the sequel and for extra clarity, we index the probabilities by the ambient bandit problem  $\underline{\nu}$  considered.

**Definition 5** Fix  $K \geq 2$ . A sequence of strategies indexed by  $T \geq 1$  is consistent, respectively, exponentially consistent, on a model  $\mathcal{D}$  if for all generic problems  $\underline{\nu}$  in  $\mathcal{D}$ ,

$$\mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \xrightarrow{T \rightarrow +\infty} 0, \quad \text{respectively,} \quad \limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) < 0.$$

By extension, a doubly-indexed sequence of strategies is (exponentially) consistent if for all  $K \geq 2$ , the associated sequences of strategies are so.

**The fundamental inequality.** We denote by 
$$N_a(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}}$$

the number of times arm  $a$  was pulled in the  $T$  exploration rounds of a given strategy with budget  $T \geq 1$ . The fundamental inequality by Garivier et al. (2019), together with the very definition of consistency, yields in a straightforward manner our building block for lower bounds. Details of the derivation are provided in Appendix D.1, for the sake of completeness.

**Lemma 6** Fix  $K \geq 2$  and a model  $\mathcal{D}$ . Consider a consistent sequence of strategies on  $\mathcal{D}$ , and two generic bandit problems  $\underline{\nu}$  and  $\underline{\lambda}$  in  $\mathcal{D}$  such that  $a^*(\underline{\lambda}) \neq a^*(\underline{\nu})$ . Then

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \limsup_{T \rightarrow +\infty} \sum_{a=1}^K \frac{\mathbb{E}_{\underline{\lambda}}[N_a(T)]}{T} \text{KL}(\lambda_a, \nu_a).$$

#### 4.1. A lower bound revisiting and extending the one by Audibert et al. (2010)

The focus of this subsection is to establish the lower bound (3), from which we derived the gap-based lower bound (4) by Audibert et al. (2010). The lower bound (3) is smaller than the lower bound to be exhibited in the next subsection, but it comes with less restrictive assumptions on the behaviors of the sequences of strategies considered.

Firstly, we only consider sequences of strategies—actually, sequences of sampling schemes—that do not pull too often the worst arm, and which we will refer to as being balanced against the worst arm. Successive-rejects-type strategies sample the worst arm less than other arms in expectations, and hence, are indeed balanced against the worst arm. To define this constraint formally, we denote by  $w_*(\underline{\nu})$  the index of the unique worst arm of a generic bandit problem  $\underline{\nu}$ .

**Definition 7** A doubly-indexed sequence of strategies is balanced against the worst arm on a model  $\mathcal{D}$  if for all  $K \geq 2$ , for all generic  $K$ -armed bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ ,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \mathbb{E}_{\underline{\nu}}[N_{w_*(\underline{\nu})}(T)] \leq \frac{1}{K}.$$

A second constraint is related to bandit subproblems. We say that  $\underline{\nu}'$  is a subproblem of a  $K$ -armed bandit problem  $\underline{\nu}$  if  $\underline{\nu}' = (\nu_a)_{a \in \mathcal{A}}$  for a subset  $\mathcal{A} \subseteq \{1, \dots, K\}$  of cardinality greater than or equal to 2; we denote by  $\underline{\nu}' \subseteq \underline{\nu}$  this fact. We say in addition that  $\underline{\nu}'$  and  $\underline{\nu}$  feature the same optimal arm if  $\nu'_{a^*(\underline{\nu}')} = \nu_{a^*(\underline{\nu})}$ . It should be easier to identify the best arm in  $\underline{\nu}'$  than in  $\underline{\nu}$ , in the sense below, and this defines the fact that a strategy cleverly exploits pruning of suboptimal arms. Again, successive-rejects-type strategies naturally satisfy this constraint.

**Definition 8** A doubly-indexed sequence of strategies cleverly exploits pruning of suboptimal arms on a model  $\mathcal{D}$  if for all generic bandit problems  $\underline{\nu}$  in  $\mathcal{D}$  with  $K \geq 2$  arms, for all subproblems  $\underline{\nu}' \subseteq \underline{\nu}$  featuring the same optimal arm,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq \liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}'}(I_T \neq a^*(\underline{\nu}')).$$

We use again the order statistics  $\mu_{w_*(\underline{\nu})} = \mu_{(K)} < \mu_{(K-1)} < \dots < \mu_{(1)} = \mu_{a^*(\underline{\nu})}$ .

**Theorem 9** Fix a model  $\mathcal{D}$ . Consider a doubly-indexed sequence of strategies that is consistent, balanced against the worst arm on  $\mathcal{D}$ , and that cleverly exploits the pruning of suboptimal arms on  $\mathcal{D}$ . For all generic bandit problems  $\underline{\nu}$  in  $\mathcal{D}$  with  $K \geq 2$  arms,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \frac{\mathcal{L}_{\text{inf}}^<(\mu_{(k)}, \nu^*)}{k}.$$

**Proof sketch** The bound is proved for  $k = K$  by considering alternative bandit problems  $\underline{\lambda}$  differing from  $\underline{\nu}$  only at arm  $a^*(\underline{\nu})$ , where  $\nu^*$  is replaced by distributions  $\zeta \in \mathcal{D}$  with  $\text{E}(\zeta) < \mu_{(K)}$ . For  $\underline{\lambda}$ , the arm  $a^*(\underline{\nu})$  is the worst arm, and is therefore pulled less than a fraction  $1/K$  of the time,

asymptotically and on average, as the strategy is balanced against the worst arm. An application of Lemma 6 concludes the case  $k = K$ . The extension to  $k \leq K - 1$  is obtained by clever exploitation of the pruning of suboptimal arms. A complete proof may be found in Appendix D.2. ■

#### 4.2. A larger lower bound, for a more restrictive class of strategies

In this section, we derive a slightly stronger version of the lower bound (2). This lower bound is larger than the bound exhibited in the previous subsection but relies on stronger assumptions on the strategies considered. Namely, we introduce an assumption of monotonicity, which extends Definition 7 to provide frequency constraints on each arm  $a \in \{1, \dots, K\}$ .

**Definition 10** Fix  $K \geq 2$ . A sequence of strategies is monotonous on a model  $\mathcal{D}$  if for all generic problems  $\underline{\nu}$  in  $\mathcal{D}$ , for all arms  $a \in \{1, \dots, K\}$ ,

$$\limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\nu}}[N_{(a)}(T)]}{T} \leq \frac{1}{a},$$

where arms are ordered such that  $\mu_{(1)} > \mu_{(2)} > \dots > \mu_{(K)}$ .

This condition is satisfied as soon as a given arm is not pulled more often, asymptotically and on average, than better-performing arms (note that Definition 10 is slightly weaker than this). Successive-rejects-type strategies naturally satisfy this requirement.

We also rely on the following assumption on the model  $\mathcal{D}$ , which essentially indicates that there is “no gap” in  $\mathcal{D}$ . Once again, the model  $\mathcal{P}[0, 1]$  and canonical one-parameter exponential models  $\mathcal{D}_{\text{exp}}$  all satisfy this mild requirement (see Appendix D.3 for the immediate details).

**Definition 11** A model  $\mathcal{D}$  is normal if for all  $\nu \in \mathcal{D}$ , for all  $x \geq E(\nu)$ ,

$$\begin{aligned} \forall \varepsilon > 0, \quad \mathcal{L}_{\text{inf}}^>(x, \nu) &\stackrel{\text{def}}{=} \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } E(\zeta) > x \} \\ &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > E(\zeta) > x \}. \end{aligned}$$

**Theorem 12** Fix  $K \geq 2$  and a normal model  $\mathcal{D}$ . Consider a sequence of strategies which is consistent and monotonous on  $\mathcal{D}$ . For all generic bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \min_{2 \leq j \leq k} \inf_{x \in [\mu_{(j)}, \mu_{(j-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{j-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{j} \right\}.$$

**Proof sketch** A complete proof may be found in Appendix D.4. For triplets  $(k, j, x)$  satisfying the stated requirements, we consider an alternative problem  $\underline{\lambda}$  differing from the original bandit problem  $\underline{\nu}$  at the best arm (1) and at the  $k$ -th best arm ( $k$ ), for which we pick distributions such that  $E(\lambda_{(1)}) < x < E(\lambda_{(k)}) < \mu_{(j-1)}$ . Then arm (1) is at best the  $j$ -th best arm of  $\underline{\lambda}$ , while arm ( $k$ ) is exactly the  $j-1$ -th best arm of  $\underline{\lambda}$ . By monotonicity and Lemma 6, we obtain

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \left( \frac{\text{KL}(\lambda_{(k)}, \nu_{(k)})}{j-1} + \frac{\text{KL}(\lambda_{(1)}, \nu^*)}{j} \right). \quad (16)$$

We get  $-\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})/(j-1) - \mathcal{L}_{\text{inf}}^<(x, \nu^*)/j$  as a lower bound by taking (separate) suprema of the lower bound (16) over  $E(\lambda_{(1)}) < x$  and  $x < E(\lambda_{(k)}) < \mu_{(j-1)}$ , where the  $< \mu_{(j-1)}$  constraint disappears thanks to normality of the model. ■

### 4.3. A general lower bound, valid for any strategy

The previous subsections illustrated what may be achieved under restrictions—though natural restrictions—on the classes of strategies considered. For the sake of completeness, we also provide a lower bound relying on no other restriction than consistency; it extends the lower bound (7) exhibited by Kaufman et al. (2016) for  $K = 2$  arms, and is formulated in terms of  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^>$ . A proof of the following theorem may be found in Appendix D.5.

**Theorem 13** *Fix  $K \geq 2$  and a model  $\mathcal{D}$ . Consider a consistent sequence of strategies on  $\mathcal{D}$ . For all generic bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ ,*

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{k \neq a^*(\underline{\nu})} \inf_{x \in [\mu_k, \mu^*]} \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\}.$$

## Acknowledgments

We thank Hédi Hadiji for pointers relative to the equality between  $\phi^*$  and  $d$  in the case of exponential models  $\mathcal{D}_{\text{exp}}$ .

## References

- J.-Y. Audibert, S. Bubeck, and R. Munos. Best arm identification in multi-armed bandits. In *Proceedings of the 23th Conference on Learning Theory (COLT 2010)*, 2010.
- S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- J. Bretagnolle and C. Huber. Estimation des densités: risque minimax. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 47, 1979.
- A.N. Burnetas and M.N. Katehakis. Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*, 17(2):122–142, 1996.
- A. Carpentier and A. Locatelli. Tight (lower) bounds for the fixed budget best arm identification bandit problem. In *Proceedings of the 29th Conference on Learning Theory (COLT 2016)*, volume 49, pages 590–604. PMLR, 2016.
- Y. Chow and H. Teicher. *Probability Theory*. Springer, 1988.
- R. Degenne, W. Koolen, and P. Ménard. Non-asymptotic pure exploration by solving games. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- J.L. Doob. *Stochastic Processes*. Wiley Publications in Statistics. John Wiley & Sons, 1953.
- A. Garivier and E. Kaufmann. Optimal best arm identification with fixed confidence. In *Proceedings of the 29th Conference on Learning Theory (COLT 2016)*, volume 49, pages 998–1027. PMLR, 2016.
- A. Garivier and E. Kaufmann. Nonasymptotic sequential tests for overlapping hypotheses applied to near-optimal arm identification in bandit models. *Sequential Analysis*, 40(1):61–96, 2021.
- A. Garivier, P. Ménard, and G. Stoltz. Explore first, exploite next: the true shape of regret in bandit problems. *Mathematics of Operations Research*, 44(2):377–399, 2019.
- A. Garivier, H. Hadiji, P. Ménard, and G. Stoltz. KL-UCB-switch: optimal regret bounds for stochastic bandits from both a distribution-dependent and a distribution-free viewpoints. *Journal of Machine Learning Research*, 23(179):1–66, 2022.
- 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.
- Z. Karnin, T. Koren, and O. Somekh. Almost optimal exploration in multi-armed bandits. In *Proceedings of the 30th International Conference on Machine Learning (ICML 2013)*, volume 28, pages 1238–1246. PMLR, 2013.



- M. Kato, K. Ariu, M. Imaizumi, M. Nomura, and C. Qin. Optimal best arm identification in two-armed bandits with a fixed budget under a small gap, 2022. Preprint, arXiv:2201.04469.
- E. Kaufman, O. Cappé, and A. Garivier. On the complexity of best-arm identification in multi-armed bandit models. *Journal of Machine Learning Research*, 17(1):1–42, 2016.
- J. Komiyama. Suboptimal performance of the Bayes optimal algorithm in frequentist best arm identification, 2022. Preprint, arXiv:2202.05193.
- J. Komiyama, T. Tsuchiya, and J. Honda. Globally optimal algorithms for fixed-budget best arm identification, 2022. Preprint, arXiv:2206.04646.
- T.L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- E.L. Lehmann and G. Casella. *Theory of Point Estimation*. Springer Texts in Statistics. Springer, 2nd edition, 1998.
- D. Russo. Simple Bayesian algorithms for best arm identification. In *Proceedings of the 29th Conference on Learning Theory (COLT 2016)*, volume 49, pages 1417–1418. PMLR, 2016.
- D. Russo. Simple Bayesian algorithms for best arm identification. *Operations Research*, 68(6): 1625–1647, 2020.

## Supplementary material for “On Best-Arm Identification with a Fixed Budget in Non-Parametric Multi-Armed Bandits ”

The appendices of this article contain the following elements.

- Appendix [A](#) states and proves some basic properties on quantities  $\mathcal{L}_{\text{inf}}^<$ ,  $\mathcal{L}_{\text{inf}}^{\leq}$ ,  $\mathcal{L}_{\text{inf}}^>$ , and  $\mathcal{L}_{\text{inf}}^{\geq}$  that were introduced in Section [2.1](#).
- Appendix [B](#) provides the proofs for the first part of the analysis of the successive-rejects strategy, namely, the general analysis in terms of  $\Phi$ , to be found in Section [3.1](#).
- Appendix [C](#) provides the proofs for the second part of the analysis of the successive-rejects strategy, namely, the rewriting of  $\Phi$  as  $\mathcal{L}$  that was the key contribution of Section [3.2](#).
- Appendix [D](#) is related to the lower bounds of Section [4](#), and provides detailed proofs thereof.
- Appendix [E](#) contains additional elements on the literature review of Sections [1](#) and [2](#); it states and discusses some important existing lower bounds.

## Appendix A. Properties of the $\mathcal{L}_{\text{inf}}^<$ , $\mathcal{L}_{\text{inf}}^{\leq}$ , $\mathcal{L}_{\text{inf}}^>$ , and $\mathcal{L}_{\text{inf}}^{\geq}$ quantities

We separate the list of properties in two categories: general properties, that hold for all models  $\mathcal{D}$ , in Appendix A.1; specific properties for the model  $\mathcal{D} = \mathcal{P}[0, 1]$ , in Appendix A.2. It is also worth noting that the  $\mathcal{L}_{\text{inf}}^<$ ,  $\mathcal{L}_{\text{inf}}^{\leq}$ ,  $\mathcal{L}_{\text{inf}}^>$ , and  $\mathcal{L}_{\text{inf}}^{\geq}$  quantities admit a simple rewriting in the case of canonical one-parameter exponential models  $\mathcal{D}_{\text{exp}}$ , as the mean-parameterized Kullback-Leibler divergence  $d$ , see Appendix C.3. Properties in this case thus follow from classical properties of  $d$ .

### A.1. General properties

We state some properties for  $\mathcal{L}_{\text{inf}}^<$ , that all also hold for  $\mathcal{L}_{\text{inf}}^{\leq}$ ; the corresponding properties for  $\mathcal{L}_{\text{inf}}^>$  and  $\mathcal{L}_{\text{inf}}^{\geq}$  are deduced by symmetry.

The function  $\mathcal{L}_{\text{inf}}^<(\cdot, \nu)$  is non-increasing and satisfies  $\mathcal{L}_{\text{inf}}^<(x, \nu) = 0$  for all  $x > \mathbb{E}(\nu)$ , as can be seen by taking  $\zeta = \nu$ . Also, whenever  $\mathcal{D}$  is convex, the function  $\mathcal{L}_{\text{inf}}^<$  is jointly convex over  $\mathbb{R} \times \mathcal{D}$ , as indicated in the lemma below. In particular,  $x \mapsto \mathcal{L}_{\text{inf}}^<(x, \nu)$  is continuous on the interior of its domain (the set where it takes finite values).

**Lemma 14** *When  $\mathcal{D}$  is a convex model, all four functions  $\mathcal{L}_{\text{inf}}^<$ ,  $\mathcal{L}_{\text{inf}}^{\leq}$ ,  $\mathcal{L}_{\text{inf}}^>$ , and  $\mathcal{L}_{\text{inf}}^{\geq}$  are jointly convex over  $\mathbb{R} \times \mathcal{D}$ .*

**Proof** We provide the proof for  $\mathcal{L}_{\text{inf}}^<$ , and it may be adapted in a straightforward manner for the other functions.

We set two distributions  $\nu$  and  $\nu'$  of  $\mathcal{D}$ , two expectation levels  $\mu$  and  $\mu'$  in  $\mathbb{R}$ , and a weight  $\lambda \in (0, 1)$ . We want to prove that

$$\mathcal{L}_{\text{inf}}^<(\lambda\mu + (1 - \lambda)\mu', \lambda\nu + (1 - \lambda)\nu') \leq \lambda\mathcal{L}_{\text{inf}}^<(\mu, \nu) + (1 - \lambda)\mathcal{L}_{\text{inf}}^<(\mu', \nu'). \quad (17)$$

The desired inequality holds whenever  $\mathcal{L}_{\text{inf}}^<(\mu, \nu) = +\infty$  or  $\mathcal{L}_{\text{inf}}^<(\mu', \nu') = +\infty$ . Otherwise, assuming that both  $\mathcal{L}_{\text{inf}}^<(\mu, \nu)$  and  $\mathcal{L}_{\text{inf}}^<(\mu', \nu')$  are finite, we set  $\delta > 0$  (which we will ultimately let converge to 0) and pick  $\zeta$  and  $\zeta'$  in  $\mathcal{D}$  such that  $\mathbb{E}(\zeta) < \mu$  and  $\mathbb{E}(\zeta') < \mu'$ , as well as

$$\text{KL}(\zeta, \nu) \leq \mathcal{L}_{\text{inf}}^<(\mu, \nu) + \delta \quad \text{and} \quad \text{KL}(\zeta', \nu') \leq \mathcal{L}_{\text{inf}}^<(\mu', \nu') + \delta.$$

Then, by joint convexity of the Kullback-Leibler divergence:

$$\begin{aligned} \lambda\mathcal{L}_{\text{inf}}^<(\mu, \nu) + (1 - \lambda)\mathcal{L}_{\text{inf}}^<(\mu', \nu') + \delta &\geq \lambda\text{KL}(\zeta, \nu) + (1 - \lambda)\text{KL}(\zeta', \nu') \\ &\geq \text{KL}(\lambda\zeta + (1 - \lambda)\zeta', \lambda\nu + (1 - \lambda)\nu') \\ &\geq \mathcal{L}_{\text{inf}}^<(\lambda\mu + (1 - \lambda)\mu', \lambda\nu + (1 - \lambda)\nu'), \end{aligned}$$

where for the last inequality, we used the definition of  $\mathcal{L}_{\text{inf}}^<$  as an infimum and the fact that by convexity, the distribution  $\lambda\zeta + (1 - \lambda)\zeta'$  belongs to  $\mathcal{D}$ , with expectation larger than  $\lambda\mu + (1 - \lambda)\mu'$ . The desired convexity inequality (17) follows by letting  $\delta \rightarrow 0$ .  $\blacksquare$

## A.2. Specific properties for $\mathcal{D} = \mathcal{P}[0, 1]$

We now consider only the model  $\mathcal{P}[0, 1]$  of all distributions over  $[0, 1]$ .

Since we are considering distributions over  $[0, 1]$ , the data-processing inequality for Kullback-Leibler divergences ensures (see, e.g., [Garivier et al., 2019](#), Lemma 1) that for all  $\zeta \in \mathcal{P}[0, 1]$ ,

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

where  $\text{Ber}(p)$  denotes the Bernoulli distribution with parameter  $p$  and where we applied Pinsker's inequality for Bernoulli distributions. Therefore, taking the infimum over distributions  $\zeta \in \mathcal{P}[0, 1]$  with  $\mathbb{E}(\zeta) < x$ ,

$$\forall x \leq \mathbb{E}(\nu), \quad \mathcal{L}_{\text{inf}}^<(x, \nu) \geq 2(\mathbb{E}(\nu) - x)^2. \quad (18)$$

We denote by  $m(\nu) = \min(\text{Supp}(\nu)) \geq 0$  the minimum of the closed support  $\text{Supp}(\nu)$  of  $\nu$ ; that is,  $m(\nu)$  is the largest value such that  $\text{Supp}(\nu) \subseteq [m(\nu), 1]$ . We will refer to  $m(\nu)$  as the lower end of the support of  $\nu$ . Though we will not need it immediately, we also define the upper end of the support of  $\nu$  as  $M(\nu) = \max(\text{Supp}(\nu)) \leq 1$ ; by symmetry, it will be considered when studying  $\mathcal{L}_{\text{inf}}^>$  and  $\mathcal{L}_{\text{inf}}^{\geq}$  instead of  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^{\leq}$ .

The lemma below states that the functions  $\mathcal{L}_{\text{inf}}^<(\cdot, \nu)$  and  $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$  coincide, except maybe at  $m(\nu)$ . One may wonder what happens at  $x = m(\nu)$ . We denote by  $\nu\{m(\nu)\}$  the probability mass assigned by  $\nu$  to the point  $m(\nu)$ . It follows from the second part the lemma below that  $\mathcal{L}_{\text{inf}}^<(m(\nu), \nu) = \mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu)$  if and only if  $\{m(\nu)\}$  is not an atom of  $\nu$ .

**Lemma 15** *For the model  $\mathcal{D} = \mathcal{P}[0, 1]$ , we have, on the one hand,*

$$\forall \mu \neq m(\nu), \quad \mathcal{L}_{\text{inf}}^<(\mu, \nu) = \mathcal{L}_{\text{inf}}^{\leq}(\mu, \nu), \quad (19)$$

*and on the other hand, at  $\mu = m(\nu)$ ,*

$$\ln \frac{1}{\nu\{m(\nu)\}} = \mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu) \leq \mathcal{L}_{\text{inf}}^<(m(\nu), \nu) = +\infty. \quad (20)$$

**Proof** To prove (19), we first identify the interior of the domain of  $\mathcal{L}_{\text{inf}}^<$ .

Distributions  $\zeta$  such that  $\mathbb{E}(\zeta) < m(\nu)$  cannot be absolutely continuous with respect to  $\nu$ ; otherwise, they would also give a null probability to values strictly smaller than  $m(\nu)$ , which contradicts the assumption  $\mathbb{E}(\zeta) < m(\nu)$ . Hence  $\text{KL}(\zeta, \nu) = +\infty$  for these distributions. It follows that  $\mathcal{L}_{\text{inf}}^<(\mu, \nu) = \mathcal{L}_{\text{inf}}^{\leq}(\mu, \nu) = +\infty$  for  $\mu < m(\nu)$ ; we note in passing that we also have  $\mathcal{L}_{\text{inf}}^<(m(\nu), \nu) = +\infty$ .

For  $\mu > m(\nu)$ , we take  $\varepsilon > 0$  with  $m(\nu) + \varepsilon < \mu$  and have, by definition of the support of a measure, that  $[m(\nu), m(\nu) + \varepsilon]$  has a positive  $\nu$ -measure denoted by  $\kappa$ . The distribution  $\zeta$  given by  $\nu$  conditioned to the interval  $[m(\nu), m(\nu) + \varepsilon]$  is absolutely continuous with respect to  $\nu$ , with density  $d\zeta/d\nu = 1/\kappa$  on  $[m(\nu), m(\nu) + \varepsilon]$ , and 0 elsewhere; therefore,  $\text{KL}(\zeta, \nu) = \ln(1/\kappa) < +\infty$  and  $\mathcal{L}_{\text{inf}}^<(\mu, \nu) < +\infty$ .

The interior of the domain of  $\mu \mapsto \mathcal{L}_{\text{inf}}^<(\mu, \nu)$  is therefore  $(m(\nu), +\infty)$ , and we recall that  $\mathcal{L}_{\text{inf}}^<(\cdot, \nu)$  is continuous on this interval. We fix some  $\mu > m(\nu)$ . For all  $\varepsilon > 0$ , by the very definitions of all quantities as infima of nested sets, we have

$$\mathcal{L}_{\text{inf}}^<(\mu, \nu - \varepsilon) \leq \mathcal{L}_{\text{inf}}^{\leq}(\mu, \nu) \leq \mathcal{L}_{\text{inf}}^<(\mu, \nu).$$

Letting  $\varepsilon \rightarrow 0$ , we get, by a sandwich argument, that  $\mathcal{L}_{\inf}^{\leq}(\mu, \nu) = \mathcal{L}_{\inf}^{<}(\mu, \nu)$ . This concludes the proof of (19).

We turn our attention to (20). We already showed above that  $\mathcal{L}_{\inf}^{<}(m(\nu), \nu) = +\infty$ . Now, to compute  $\mathcal{L}_{\inf}^{\leq}(\mu, \nu)$ , we wonder which are the distributions  $\zeta$  that are absolutely continuous with respect to  $\nu$ , and thus, give a null probability to values strictly smaller than  $m(\nu)$ , and are also such that  $E(\zeta) \leq m(\nu)$ : at most one such distribution exists, the Dirac mass at  $m(\nu)$ , denoted by  $\delta_{m(\nu)}$ . We then distinguish the cases  $\nu\{m(\nu)\} > 0$  and  $\nu\{m(\nu)\} = 0$  to establish, respectively, the equalities

$$\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = \text{KL}(\delta_{m(\nu)}, \nu) = \ln \frac{1}{\nu\{m(\nu)\}} \quad \text{and} \quad \mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = +\infty = \ln \frac{1}{\nu\{m(\nu)\}}.$$

In both cases, the first equality in (20) is proved, which concludes the proof.  $\blacksquare$

We also have the following result, which is the most important and useful one, as it discussed the quantity that appears in the upper bounds on the average log-probability of misidentification of the optimal arm; see Corollary 3 together with Lemma 4.

**Lemma 16** *Let  $\nu, \nu' \in \mathcal{P}[0, 1]$  with  $\mu = E(\nu) > E(\nu') = \mu'$ . Then*

$$\inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu') = \inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{<}(x, \nu) + \mathcal{L}_{\inf}^{>}(x, \nu')$$

*if and only if either  $m(\nu) \neq M(\nu')$  or  $\nu\{m(\nu)\} \times \nu'\{M(\nu')\} = 0$ .*

**Remark 17** *In other words, the only case for which the two infima differ is when  $m(\nu) = M(\nu')$ , i.e., the upper end of the support of  $\nu'$  equals the lower end of the support of  $\nu$ , and both  $\nu$  and  $\nu'$  admit this common value as an atom.*

**Proof** The first lines of the proof of Lemma 15 show that  $\mathcal{L}_{\inf}^{\leq}(x, \nu) = \mathcal{L}_{\inf}^{<}(x, \nu) = +\infty$  for  $x < m(\nu)$ . We can symmetrically show that  $\mathcal{L}_{\inf}^{\geq}(x, \nu') = \mathcal{L}_{\inf}^{>}(x, \nu') = +\infty$  for  $x > M(\nu')$ . Therefore,  $\mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu')$  and  $\mathcal{L}_{\inf}^{<}(x, \nu) + \mathcal{L}_{\inf}^{>}(x, \nu')$  are infinite whenever  $x$  lies outside of  $[m(\nu), M(\nu')]$ . This implies that

$$\begin{aligned} \inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu') &= \inf_{x \in [\mu', \mu] \cap [m(\nu), M(\nu')]} \mathcal{L}_{\inf}^{\leq}(x, \nu) + \mathcal{L}_{\inf}^{\geq}(x, \nu') \\ \text{and} \quad \inf_{x \in [\mu', \mu]} \mathcal{L}_{\inf}^{<}(x, \nu) + \mathcal{L}_{\inf}^{>}(x, \nu') &= \inf_{x \in [\mu', \mu] \cap [m(\nu), M(\nu')]} \mathcal{L}_{\inf}^{<}(x, \nu) + \mathcal{L}_{\inf}^{>}(x, \nu'). \end{aligned}$$

We now split the analysis according to how large the interval  $\mathcal{I}$  is, where

$$\mathcal{I} = [\mu', \mu] \cap [m(\nu), M(\nu')] = \left[ \max\{\mu', m(\nu)\}, \min\{\mu, M(\nu')\} \right].$$

*Case 1:  $\mathcal{I}$  is empty.* In that case, the two infima are over an empty set and both equal  $+\infty$ .

*Case 2:  $\mathcal{I}$  has a non-empty interior.* When  $a \neq b$ , the infimum of a convex function over a closed interval  $[a, b]$  equals the infimum over  $(a, b)$ , whether the function takes finite or infinite values at  $a$  and  $b$ . Now, the interior of  $\mathcal{I} = [a, b]$  equals

$$(a, b) = \left( \max\{\mu', m(\nu)\}, \min\{\mu, M(\nu')\} \right) = (\mu', \mu) \cap (m(\nu), M(\nu'))$$

and does not contain neither  $m(\nu)$  nor  $M(\nu')$ . By Lemma 15, the functions  $\mathcal{L}_{\inf}^{\leq}(\cdot, \nu)$  and  $\mathcal{L}_{\inf}^{\leq}(\cdot, \nu)$  coincide on  $\mathbb{R} \setminus \{m(\nu)\}$ . It may be similarly shown that  $\mathcal{L}_{\inf}^{\geq}(\cdot, \nu')$  and  $\mathcal{L}_{\inf}^{\geq}(\cdot, \nu')$  coincide on  $\mathbb{R} \setminus \{M(\nu')\}$ . In particular, the functions  $\mathcal{L}_{\inf}^{\leq}(\cdot, \nu) + \mathcal{L}_{\inf}^{\geq}(\cdot, \nu')$  and  $\mathcal{L}_{\inf}^{\leq}(\cdot, \nu) + \mathcal{L}_{\inf}^{\geq}(\cdot, \nu')$  coincide on the interior of  $\mathcal{I}$ . Their infima over the interior of  $\mathcal{I}$ , which, by convexity, are equal to the infima over  $\mathcal{I}$ , are therefore equal.

*Case 3:  $\mathcal{I}$  is a singleton.* This case arises if and only if  $m(\nu) = M(\nu')$ , as by definition,  $m(\nu) \leq \mu$  and  $M(\nu') \geq \mu'$ . We then have  $\mathcal{I} = \{m(\nu)\} = \{M(\nu')\}$ , and both infima are equal to the values of the sums at  $m(\nu) = M(\nu')$ . By Lemma 15 and by symmetric results for  $\mathcal{L}_{\inf}^{\geq}$  and  $\mathcal{L}_{\inf}^{\leq}$ , on the one hand,

$$\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) = \mathcal{L}_{\inf}^{\geq}(M(\nu'), \nu') = +\infty,$$

and on the other hand,

$$\mathcal{L}_{\inf}^{\leq}(m(\nu), \nu) + \mathcal{L}_{\inf}^{\geq}(M(\nu'), \nu') = \ln \frac{1}{\nu\{m(\nu)\}} + \ln \frac{1}{\nu'\{M(\nu')\}}.$$

We get the desired equality if and only if either  $\nu\{m(\nu)\} = 0$  or  $\nu'\{M(\nu')\} = 0$ . ■

## Appendix B. General analysis of successive-rejects in terms of $\Phi$

This appendix is devoted to the technical elements omitted in the general analysis of the successive-rejects strategy presented in Section 3.1.

### B.1. The Cramér-Chernoff bound

In this section, we recall, for the sake of completeness, the highly classical Cramér-Chernoff bound. With the notation introduced in Section 3, it states that, for an  $N$ -sample  $X_1, \dots, X_N$ , distributed according to  $\nu$  and of average denoted by  $\bar{X}_N$ ,

$$\forall x \leq \mathbb{E}(\nu), \quad \mathbb{P}(\bar{X}_N \leq x) \leq \exp(-N \phi_\nu^*(x)), \quad (21)$$

$$\text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \mathbb{P}(\bar{X}_N \geq x) \leq \exp(-N \phi_\nu^*(x)). \quad (22)$$

**Proof** For all  $\lambda < 0$ , by Markov's inequality first and then by independence,

$$\begin{aligned} \mathbb{P}(\bar{X}_N \leq x) &= \mathbb{P}(e^{\lambda \bar{X}_N} \geq e^{\lambda x}) \leq e^{-\lambda x} \mathbb{E}[e^{\lambda \bar{X}_N}] = e^{-\lambda x} \left( \mathbb{E}[e^{\lambda X_1/N}] \right)^N \\ &= \exp(-\lambda x + N \phi_\nu(\lambda/N)) = \exp(-N(\lambda' x - \phi_\nu(\lambda'))), \end{aligned}$$

where  $\lambda' = \lambda/N$ . The bound also holds for  $\lambda = \lambda' = 0$  given that  $\phi_\nu(0) = 0$ . Optimizing over  $\lambda \leq 0$  (or, equivalently, over  $\lambda' \leq 0$ ), we proved so far

$$\mathbb{P}(\bar{X}_N \leq x) \leq \exp\left(-N \sup_{\lambda \leq 0} \{\lambda x - \phi_\nu(\lambda)\}\right).$$

Now, by Jensen's inequality,

$$\forall \lambda \in \mathbb{R}, \quad \phi_\nu(\lambda) = \ln \mathbb{E}[e^{\lambda X}] \geq \lambda \mathbb{E}[X] = \lambda \mathbb{E}(\nu); \quad (23)$$

therefore, for  $x \leq \mathbb{E}(\nu)$ ,

$$\forall \lambda \geq 0, \quad \lambda x - \phi_\nu(\lambda) \leq \lambda(x - \mathbb{E}(\nu)) \leq 0.$$

In particular,

$$0 = -\phi_\nu(0) \leq \sup_{\lambda \leq 0} \{\lambda x - \phi_\nu(\lambda)\} = \sup_{\lambda \in \mathbb{R}} \{\lambda x - \phi_\nu(\lambda)\} \stackrel{\text{def}}{=} \phi_\nu^*(x). \quad (24)$$

This concludes the proof of (21). The bound (22) follows by symmetry.  $\blacksquare$

We also note, in passing, that Jensen's inequality entails, for  $x = \mathbb{E}(\nu)$ , that

$$\forall \lambda \in \mathbb{R}, \quad \lambda \mathbb{E}(\nu) - \phi_\nu(\lambda) \leq \lambda(\mathbb{E}(\nu) - \mathbb{E}(\nu)) = 0,$$

thus showing that  $\phi_\nu^*(\mathbb{E}(\nu)) = 0$ . The property (24) and its counterpart for  $x \geq \mathbb{E}(\nu)$  and  $\lambda \geq 0$  actually show that  $\phi_\nu^*$  is non-increasing on  $(-\infty, \mathbb{E}(\nu)]$  and non-decreasing on  $[\mathbb{E}(\nu), +\infty)$ .



## B.2. Proof of Lemma 1

We first restate the lemma, for the convenience of the reader.

**Lemma 1** Fix  $\nu$  and  $\nu'$  in  $\mathcal{D}$ , with respective expectations  $\mu = \mathbb{E}(\nu) > \mu' = \mathbb{E}(\nu')$ . For all  $N \geq 1$ , let  $\bar{X}_N$  and  $\bar{Y}_N$  be the averages of  $N$ -samples with respective distributions  $\nu$  and  $\nu'$ . Then,

$$\limsup_{N \rightarrow +\infty} \frac{1}{N} \ln \mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq - \inf_{x \in [\mu', \mu]} \{ \phi_{\nu'}^*(x) + \phi_{\nu}^*(x) \} \stackrel{\text{def}}{=} -\Phi(\nu', \nu).$$

**Proof** The proof consists in two parts. We first show that for any finite grid  $\mathcal{G} = \{g_2, \dots, g_{G-1}\}$  in  $(\mu', \mu)$ , to which we add the points  $g_1 = \mu'$  and  $g_G = \mu$ , we have

$$\limsup_{N \rightarrow +\infty} \frac{1}{N} \ln \mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq - \min \left\{ \phi_{\nu}^*(\mu'), \min_{2 \leq j \leq G-1} \{ \phi_{\nu'}^*(g_{j-1}) + \phi_{\nu}^*(g_j) \}, \phi_{\nu'}^*(\mu) \right\}. \quad (25)$$

Indeed, by identifying, when  $\bar{X}_N$  and  $\bar{Y}_N$  belong to  $[\mu', \mu]$ , in which interval  $[g_{j-1}, g_j]$  lies  $\bar{X}_N$ , we note that

$$\{ \bar{X}_N \leq \bar{Y}_N \} \subseteq \{ \bar{X}_N \leq \mu' \} \cup \{ \bar{Y}_N \geq \mu \} \cup \bigcup_{j=2}^{G-1} \{ \bar{Y}_N \geq g_{j-1} \text{ and } \bar{X}_N \leq g_j \}.$$

First, by independence and by the Cramér-Chernoff inequalities (21) and (22),

$$\mathbb{P}(\bar{Y}_N \geq g_{j-1} \text{ and } \bar{X}_N \leq g_j) = \mathbb{P}(\bar{Y}_N \geq g_{j-1}) \mathbb{P}(\bar{X}_N \leq g_j) \leq \exp\left(-N(\phi_{\nu'}^*(g_{j-1}) + \phi_{\nu}^*(g_j))\right).$$

Second, again by the Cramér-Chernoff inequalities,

$$\mathbb{P}(\bar{X}_N \leq \mu') \leq \exp(-N \phi_{\nu}^*(\mu')) \quad \text{and} \quad \mathbb{P}(\bar{Y}_N \geq \mu) \leq \exp(-N \phi_{\nu'}^*(\mu)).$$

By a union bound,

$$\mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq \exp(-N \phi_{\nu}^*(\mu')) + \exp(-N \phi_{\nu'}^*(\mu)) + \sum_{j=2}^{G-1} \exp\left(-N(\phi_{\nu'}^*(g_{j-1}) + \phi_{\nu}^*(g_j))\right).$$

The stated bound (25) follows by identifying the (finitely many) terms with the smallest rate in the exponent.

In the second part of the proof, we note that the bound (25) holds for any finite grid in  $(\mu', \mu)$ , and we consider a sequence

$$\mathcal{G}^{(n)} = \{g_2^{(n)}, \dots, g_{G_n-1}^{(n)}\}$$

of such finite grids. In particular,

$$\limsup_{N \rightarrow +\infty} \frac{1}{N} \ln \mathbb{P}(\bar{X}_N \leq \bar{Y}_N) \leq - \min \left\{ \phi_{\nu}^*(\mu'), \max_{n \geq 1} S_n, \phi_{\nu'}^*(\mu) \right\},$$

$$\text{where} \quad S_n \stackrel{\text{def}}{=} \min_{2 \leq j \leq G_n-1} \left\{ \phi_{\nu'}^*(g_{j-1}^{(n)}) + \phi_{\nu}^*(g_j^{(n)}) \right\}.$$

To obtain the claimed bound, given that (see the end of Appendix B.1)

$$\phi_{\nu}^*(\mu) = 0 = \phi_{\nu'}^*(\mu'),$$

it suffices to show that

$$\max_{n \geq 1} S_n \geq \inf_{x \in [\mu', \mu]} \{ \phi_{\nu'}^*(x) + \phi_{\nu}^*(x) \}.$$

To that end, we assume that the steps  $\varepsilon_n$  of the grids  $\mathcal{G}^{(n)}$ , which are defined as

$$\varepsilon_n \stackrel{\text{def}}{=} \max_{2 \leq j \leq G_n} |g_j^{(n)} - g_{j-1}^{(n)}|,$$

vanish asymptotically, i.e.,  $\varepsilon_n \rightarrow 0$ . For each grid  $\mathcal{G}^{(n)}$ , we denote by  $x_n^* \in (\mu', \mu)$  the argument of the minimum in the definition of  $S_n$ . As a consequence, for each  $n \geq 1$ ,

$$S_n = \phi_{\nu'}^*(x_n^* - \varepsilon_n^*) + \phi_{\nu}^*(x_n^*),$$

for some  $0 < \varepsilon_n^* \leq \varepsilon_n$ . The quantity  $x_n^* - \varepsilon_n^*$  denotes the point in the grid that is right before  $x_n^*$ , and it belongs to  $[\mu', \mu]$ . We note that we also have  $\varepsilon_n^* \rightarrow 0$ . In the compact interval  $[\mu', \mu]$ , the Bolzano-Weierstrass property ensures the existence of a converging subsequence: there exists  $x_{\infty}^* \in [\mu', \mu]$  and a sequence  $(n_k)_{k \geq 1}$  of integers such that

$$x_{n_k}^* \xrightarrow[k \rightarrow +\infty]{} x_{\infty}^*, \quad \text{which also entails} \quad x_{n_k}^* - \varepsilon_{n_k}^* \xrightarrow[k \rightarrow +\infty]{} x_{\infty}^*.$$

Now, the functions  $\phi_{\nu}^*$ , respectively,  $\phi_{\nu'}^*$ , are lower semi-continuous, as the suprema over  $\lambda \in \mathbb{R}$  of the continuous functions  $x \mapsto \lambda x - \varphi_{\nu}(\lambda)$ , respectively,  $x \mapsto \lambda x - \varphi_{\nu'}(\lambda)$ . Therefore, by these lower semi-continuities,

$$\begin{aligned} \max_{n \geq 1} S_n &\geq \liminf_{k \rightarrow +\infty} \phi_{\nu'}^*(x_{n_k}^* - \varepsilon_{n_k}^*) + \phi_{\nu}^*(x_{n_k}^*) \geq \phi_{\nu'}^*(x_{\infty}^*) + \phi_{\nu}^*(x_{\infty}^*) \\ &\geq \inf_{x \in [\mu', \mu]} \{ \phi_{\nu'}^*(x) + \phi_{\nu}^*(x) \}. \end{aligned}$$

This concludes the proof. ■

### B.3. Proof of Theorem 2

The proof mimics the analysis by [Audibert et al. \(2010\)](#), the main modification being the substitution of Hoeffding's inequality by the bound of Lemma 1.

**Proof** We recall that for  $r \in \{1, \dots, K-1\}$ , we denoted by  $N_r = \lfloor \ell_1/K \rfloor + \dots + \lfloor \ell_r/(K-r+1) \rfloor$  the total number of times an arm still considered in regime  $r$ , i.e., belonging to  $S_{r-1}$ , was pulled in regimes 1 to  $r$ . For each arm  $a$ , we denote by  $\bar{Y}_a^r$  the average of a  $N_r$ -sample distributed according to  $\nu_a$ . By optional skipping (see [Doob, 1953](#), Chapter III, Theorem 5.2, p. 145, or [Chow and Teicher, 1988](#), Section 5.3 for a more recent reference), we may assume, with no loss of generality, that for each  $r \in \{1, \dots, K-1\}$ ,

$$\text{on the event } \{a \in S_{r-1}\}, \quad \bar{X}_a^r = \bar{Y}_a^r. \quad (26)$$

We fix a bandit problem  $\underline{\nu}$  with a unique optimal arm  $a^*(\underline{\nu})$ . The successive-rejects strategy fails if (and only) if it rejects  $a^*(\underline{\nu})$  in ones of the regimes. This corresponds to the event

$$\{I_T \neq a^*(\underline{\nu})\} = \bigcup_{r=1}^{K-1} \{a_r = a^*(\underline{\nu})\} \subseteq \bigcup_{r=1}^{K-1} \left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{X}_{a^*(\underline{\nu})}^r \leq \bar{X}_k^r \right\}.$$

(We have an inclusion because ties are broken arbitrarily.) By optional skipping (26),

$$\begin{aligned} \bigcup_{r=1}^{K-1} \left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{X}_{a^*(\underline{\nu})}^r \leq \bar{X}_k^r \right\} \\ = \bigcup_{r=1}^{K-1} \left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r \right\}. \end{aligned}$$

Recall that the set  $S_{r-1}$  is a random set; dealing with it therefore requires some care. On the event of interest,  $S_{r-1}$  contains  $K - r + 1$  elements, among which  $a^*(\underline{\nu})$ . The set  $\mathcal{A}_r$  is of cardinality  $r$  and does not contain  $a^*(\underline{\nu})$ . By the pigeonhole principle,  $S_{r-1}$  thus necessarily contains one arm in  $\mathcal{A}_r$ . As a consequence, for each regime  $r \in \{1, \dots, K - 1\}$ ,

$$\left\{ a^*(\underline{\nu}) \in S_{r-1} \text{ and } \forall k \in S_{r-1}, \bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r \right\} \subseteq \bigcup_{k \in \mathcal{A}_r} \left\{ \bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r \right\}.$$

Summarizing the inclusions above, taking unions bounds, and upper bounding the obtained sum in a crude way, we proved so far

$$\mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq \sum_{r=1}^{K-1} \sum_{k \in \mathcal{A}_r} \mathbb{P}\left(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r\right) \leq K^2 \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \mathbb{P}\left(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r\right),$$

or equivalently,

$$\begin{aligned} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) &\leq \frac{2}{T} \ln K + \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \frac{1}{T} \ln \mathbb{P}\left(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r\right) \\ &= \frac{2}{T} \ln K + \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \frac{N_r}{T} \frac{1}{N_r} \ln \mathbb{P}\left(\bar{Y}_{a^*(\underline{\nu})}^r \leq \bar{Y}_k^r\right). \end{aligned}$$

As  $N_r/T \rightarrow \gamma_r > 0$  as  $T \rightarrow +\infty$ , we may apply Lemma 1, together with an exchange between the lim sup and the maximum over a finite number of quantities. We obtain

$$\begin{aligned} \limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) &\leq \max_{1 \leq r \leq K-1} \max_{k \in \mathcal{A}_r} \left\{ \gamma_r \left( -\Phi(\nu_k, \nu^*) \right) \right\} \\ &= - \min_{1 \leq r \leq K-1} \left\{ \gamma_r \min_{k \in \mathcal{A}_r} \Phi(\nu_k, \nu^*) \right\}. \end{aligned}$$

This concludes the proof. ■

#### B.4. Proof of Corollary 3 and of the bound (13) on $\Phi$

In this final subsection, we provide two series of proofs: first, a proof of Corollary 3; and then a proof of the bound  $\Phi(\nu_k, \nu^*) \geq \Delta_k^2$  stated as (13).

**Proof of Corollary 3.** To apply Theorem 2, we need only to show that the regime lengths of (12) are such that  $N_r/T$  converges to a positive value, and to identify this limit value  $\gamma_r$ . As  $N_1 =$

$\lfloor \ell_1/K \rfloor$ , where  $\ell_1 = T/\overline{\ln K}$ , we immediately have  $N_1/T \rightarrow \gamma_1 = 1/(K \overline{\ln K}) > 0$ . For  $r \in \{2, \dots, K-1\}$ ,

$$\begin{aligned} \frac{N_r}{T} &= \sum_{p=1}^r \frac{1}{T} \left\lfloor \frac{\ell_p}{K} \right\rfloor = \frac{1}{T} \left( \left\lfloor \frac{T}{K \overline{\ln K}} \right\rfloor + \sum_{p=2}^r \left\lfloor \frac{T}{(K-p+1)(K-p+2) \overline{\ln K}} \right\rfloor \right) \\ &\xrightarrow{T \rightarrow +\infty} \gamma_r \stackrel{\text{def}}{=} \frac{1}{\overline{\ln K}} \left( \frac{1}{K} + \sum_{p=2}^r \frac{1}{K-p+1} - \frac{1}{K-p+2} \right) = \frac{1}{(K-r+1) \overline{\ln K}}. \end{aligned}$$

The bound of Theorem 2 reads:

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\overline{\ln K}} \min_{1 \leq r \leq K-1} \left\{ \frac{1}{K-r+1} \min_{k \in \mathcal{A}_r} \Phi(\nu_k, \nu^*) \right\}.$$

It implies, in terms of lower bounds  $f(\nu_k, \nu^*) \leq \Phi(\nu_k, \nu^*)$ ,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}(I_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\overline{\ln K}} \min_{1 \leq r \leq K-1} \left\{ \frac{1}{K-r+1} \min_{k \in \mathcal{A}_r} f(\nu_k, \nu^*) \right\}. \quad (27)$$

The permutation  $\sigma$  in (14) and the sets  $\mathcal{A}_r = \{\sigma_{K-r+1}, \dots, \sigma_K\}$  were exactly picked, for each  $r \in \{1, \dots, K-1\}$ , to minimize

$$\min_{k \in \mathcal{B}_r} f(\nu_k, \nu^*)$$

over sets  $\mathcal{B}_r$  abiding by the indicated constraints: being of cardinal  $r$  and not containing the optimal arm  $a^*(\underline{\nu}) = \sigma_1$ . We get

$$\min_{k \in \mathcal{A}_r} f(\nu_k, \nu^*) = \min_{K-r+1 \leq k \leq K} f(\nu_{\sigma_k}, \nu^*) = f(\nu_{\sigma_{K-r+1}}, \nu^*),$$

which, together with (27), yields the stated bound, up to replacing  $K-r+1$  with  $r \in \{1, \dots, K-1\}$  by  $k \in \{2, \dots, K\}$ :

$$-\frac{1}{\overline{\ln K}} \min_{1 \leq r \leq K-1} \left\{ \frac{1}{K-r+1} f(\nu_{\sigma_{K-r+1}}, \nu^*) \right\} = -\frac{1}{\overline{\ln K}} \min_{2 \leq k \leq K} \left\{ \frac{1}{k} f(\nu_{\sigma_k}, \nu^*) \right\}. \quad \blacksquare$$

We now move to the proof of the bound (13) on  $\Phi$ , when the model is  $\mathcal{D} = \mathcal{P}[0, 1]$ . We restate it here for the convenience of the reader:

$$\phi_\nu^*(x) \geq 2(x - \mathbb{E}(\nu))^2, \quad \text{so that} \quad \Phi(\nu_k, \nu^*) \geq \Delta_k^2 \stackrel{\text{def}}{=} f(\nu_k, \nu^*).$$

For the ease of exposition, the path followed in Section 2 to show that  $\Phi(\nu_k, \nu^*) \geq \Delta_k^2$  was to first note that  $\Phi = \mathcal{L}$  when  $\mathcal{D} = \mathcal{P}[0, 1]$  (see Lemma 4) and then use Pinsker's inequality (5). We provide here a slightly more direct but equivalent approach, based on Hoeffding's inequality.

**Proof of the bound (13) on  $\Phi$ .** When  $\nu \in \mathcal{P}[0, 1]$ , Hoeffding's inequality exactly states that

$$\begin{aligned} \forall \lambda \in \mathbb{R}, \quad \phi_\nu(\lambda) &\leq \lambda \mathbb{E}(\nu) + \frac{\lambda^2}{8}, \\ \text{so that} \quad \forall x \in \mathbb{R}, \quad \phi_\nu^*(x) &\geq \sup_{\lambda \in \mathbb{R}} \left\{ \lambda(x - \mathbb{E}(\nu)) - \frac{\lambda^2}{8} \right\} = 2(x - \mathbb{E}(\nu))^2. \end{aligned}$$

This corresponds to the first part of (13).

For its second part, we consider a pair  $\nu, \nu'$  of distributions in  $\mathcal{P}[0, 1]$ , we set any  $x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]$ , and we apply twice the bound of the first part to get

$$\phi_{\nu'}^*(x) + \phi_{\nu}^*(x) \geq 2(x - \mathbb{E}(\nu'))^2 + 2(x - \mathbb{E}(\nu))^2.$$

From the definition of  $\Phi$ , it follows that

$$\Phi(\nu', \nu) \geq \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ 2(x - \mathbb{E}(\nu'))^2 + 2(x - \mathbb{E}(\nu))^2 \right\} = (\mathbb{E}(\nu') - \mathbb{E}(\nu))^2.$$

This corresponds to the second part of (13). ■

### Appendix C. Rewriting of $\phi_\nu^*$ as some $\mathcal{L}_{\text{inf}}^{\leq}$ or $\mathcal{L}_{\text{inf}}^{\geq}$ , i.e., of $\Phi$ as $\mathcal{L}$

We use the notation of Sections 2.1 and 3 and discuss conditions on models guaranteeing that  $\Phi = \mathcal{L}$ , i.e., that (15) holds. We do so for  $\mathcal{D} = \mathcal{P}[0, 1]$  in Section C.2 and for canonical one-parameter exponential families in Section C.3; based on these two examples, we provide a set of conditions for general models, in Section C.4. A building block of these results is that for all models  $\mathcal{D}$ , the functions  $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$  and  $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$  dominate  $\phi_\nu^*$  defined in (11); we prove this in Section C.1.

This rewriting of  $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$  or  $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$  as  $\phi_\nu^*$  claimed, e.g., by Lemma 4, can be seen as a counterpart to a similar rewriting of the  $\mathcal{K}_{\text{inf}}$  as the supremum of a function of  $\lambda \in [0, 1]$ . More precisely, for  $\nu \in \mathcal{P}[0, 1]$  and  $x \in [0, 1]$ ,

$$\mathcal{K}_{\text{inf}}(\nu, x) = \inf \{ \text{KL}(\nu, \zeta) : \zeta \in \mathcal{P}[0, 1] \text{ s.t. } \mathbb{E}(\zeta) < x \},$$

and Honda and Takemura (2015, Theorem 2)—see also Garivier et al., 2022, Lemma 18—show that

$$\mathcal{K}_{\text{inf}}(\nu, x) = \sup_{0 \leq \lambda \leq 1} \mathbb{E} \left[ \ln \left( 1 - \lambda \frac{Y - x}{1 - x} \right) \right],$$

where  $Y$  is a random variable distributed according to  $\nu$ . In both cases, for  $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$  or  $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$ , and for  $\mathcal{K}_{\text{inf}}$ , being able to rewrite the infimum of a set of Kullback-Leibler divergences as a supremum is not unexpected, as a Kullback-Leibler divergence can be formulated as a supremum, see (28), and exchanges between the infimum and the supremum may be justified (e.g., through Sion's lemma, when applicable).

#### C.1. $\mathcal{L}_{\text{inf}}^{\leq}(\cdot, \nu)$ and $\mathcal{L}_{\text{inf}}^{\geq}(\cdot, \nu)$ dominate $\phi_\nu^*$

This domination is a consequence of a variational formula for the Kullback-Leibler divergence (28).

**Lemma 18** *For all models  $\mathcal{D}$ , for all  $\nu \in \mathcal{D}$ ,*

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_\nu^*(x) \leq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \phi_\nu^*(x) \leq \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

In the next sections we will obtain the converse inequalities for the specific models mentioned above. The proofs are immediate adaptations of a rather standard result, stated, among others, but in a slightly different form (and for the model  $\mathcal{D}$  of all real-valued distributions with a first moment), by Boucheron et al. (2013, Exercise 4.13).

**Proof** We rely on a key variational formula for the Kullback-Leibler divergence, see Boucheron et al. (2013, Chapter 4): for all distributions  $\nu, \nu'$  over  $\mathbb{R}$ ,

$$\begin{aligned} \text{KL}(\nu', \nu) &= \sup \left\{ \mathbb{E}_{\nu'}[Y] - \ln \mathbb{E}_\nu[e^Y] : \text{r.v. } Y \in \mathbb{L}^1(\nu') \text{ s.t. } \mathbb{E}_\nu[e^Y] < +\infty \right\}, \\ &= \sup \left\{ \mathbb{E}_{\nu'}[Y] - \ln \mathbb{E}_\nu[e^Y] : \text{r.v. } Y \in \mathbb{L}^1(\nu') \right\}, \end{aligned} \quad (28)$$

where  $\mathbb{E}_\nu$  and  $\mathbb{E}_{\nu'}$  indicate that expectations are relative to  $\nu$  and  $\nu'$ , respectively. In particular, when  $\nu$  and  $\nu'$  lie in  $\mathcal{D}$ , those two distributions admit a finite first moment, hence all random variables

$Y = \lambda \text{id}_{\mathbb{R}}$  are  $\nu'$ -integrable, where  $\text{id}_{\mathbb{R}}$  denotes the identity function on  $\mathbb{R}$  and where  $\lambda \in \mathbb{R}$ . They satisfy  $\mathbb{E}_{\nu'}[Y] = \lambda \mathbb{E}(\nu')$ . A consequence of (28) is therefore

$$\text{KL}(\nu', \nu) \geq \sup_{\lambda \in \mathbb{R}} \left\{ \lambda \mathbb{E}(\nu') - \ln \mathbb{E}_{\nu} [e^{\lambda \text{id}_{\mathbb{R}}}] \right\} = \phi_{\nu}^*(\mathbb{E}(\nu')). \quad (29)$$

Using the variations of  $\phi_{\nu}^*$  indicated at the end of Appendix B.1, we see that

$$\phi_{\nu}^*(\mathbb{E}(\nu')) \geq \phi_{\nu}^*(x) \quad \text{when } \mathbb{E}(\nu') \leq x \leq \mathbb{E}(\nu) \quad \text{or} \quad \mathbb{E}(\nu') \geq x \geq \mathbb{E}(\nu).$$

Therefore, by taking the infimum in (29) over  $\nu' \in \mathcal{D}$  either with  $\mathbb{E}(\nu') \leq x$  or  $\mathbb{E}(\nu') \geq x$ , we proved the claimed inequalities.  $\blacksquare$

### C.2. The case of $\mathcal{P}[0, 1]$

In this section, we focus on the model  $\mathcal{P}[0, 1]$  and prove that the inequalities of Lemma 18 are in fact equalities, as claimed by Lemma 4.

**Lemma 4** *Consider the model  $\mathcal{D} = \mathcal{P}[0, 1]$ . For all  $\nu \in \mathcal{P}[0, 1]$ ,*

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_{\nu}^*(x) = \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \phi_{\nu}^*(x) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

The lemma holds for all  $x \in \mathbb{R}$ , that is, even outside of the  $[0, 1]$  interval, though the proof reveals that when  $x$  is smaller than the lower end  $m(\nu)$  of the support of  $\nu$ , we actually have  $\phi_{\nu}^*(x) = \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = +\infty$ . The counterpart statement  $\phi_{\nu}^*(x) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu) = +\infty$  holds for  $x$  larger than the upper end  $M(\nu)$  of the support of  $\nu$ . The pieces of notation  $m(\nu)$  and  $M(\nu)$  were formally defined in Appendix A.2.

**Proof** Note first that by Lemma 18, it suffices to prove that

$$\forall x \leq \mathbb{E}(\nu), \quad \phi_{\nu}^*(x) \geq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq \mathbb{E}(\nu), \quad \phi_{\nu}^*(x) \geq \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

We only deal with the first inequality, namely  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \phi_{\nu}^*(x)$  for  $x \leq \mathbb{E}(\nu)$ , and obtain the other one by symmetry.

In the case  $x = \mathbb{E}(\nu)$ , we have  $\phi_{\nu}^*(\mathbb{E}(\nu)) = 0$ , as stated at the end of Section B.1, and  $\mathcal{L}_{\text{inf}}^{\leq}(\mathbb{E}(\nu), \nu) = 0$  as can be seen by taking  $\zeta = \nu$  in the defining infimum. We therefore only consider  $x < \mathbb{E}(\nu)$  in the sequel. We will rely on the standard fact that, by Hölder's inequality, the logarithmic moment-generating function

$$\phi_{\nu} : \lambda \in \mathbb{R} \mapsto \ln \mathbb{E}_{\nu} [e^{\lambda \text{id}_{[0,1]}}],$$

is convex, where  $\text{id}_{[0,1]}$  denotes the identity function on  $[0, 1]$ . Also, by two applications of a standard theorem of differentiation under the integral, given that  $\nu$  is supported by  $[0, 1]$ , we have that  $\phi_{\nu}$  is continuously differentiable over  $\mathbb{R}$ , with derivative

$$\phi'_{\nu} : \lambda \in \mathbb{R} \mapsto \frac{\mathbb{E}_{\nu} [\text{id}_{[0,1]} e^{\lambda \text{id}_{[0,1]}}]}{\mathbb{E}_{\nu} [e^{\lambda \text{id}_{[0,1]}}]}.$$



By convexity of  $\phi_\nu$ , this derivative is non-decreasing. Therefore, the limit of  $\phi'_\nu$  at  $-\infty$  exists; we denote it by  $\ell$  and have that a priori  $\ell \in \{-\infty\} \cup \mathbb{R}$ . We now prove that actually,

$$\lim_{\lambda \rightarrow -\infty} \phi'_\nu(\lambda) \stackrel{\text{def}}{=} \ell = m(\nu). \quad (30)$$

On the one hand, by definition of  $m(\nu)$ , we have  $\text{id}_{[0,1]} \geq m(\nu)$   $\nu$ -a.s., which entails  $\phi'_\nu(\lambda) \geq m(\nu)$  for all  $\lambda \in \mathbb{R}$ , and hence,  $\ell \geq m(\nu)$ . On the other hand, as  $\phi'_\nu$  is non-decreasing, it is always larger than its limit  $\ell$  at  $-\infty$ :

$$\forall \lambda \in \mathbb{R}, \quad \phi'_\nu(\lambda) \geq \ell, \quad \text{thus,} \quad \mathbb{E}_\nu \left[ (\text{id}_{[0,1]} - \ell) e^{\lambda \text{id}_{[0,1]}} \right] \geq 0, \quad (31)$$

$$\text{or} \quad \mathbb{E}_\nu \left[ (\text{id}_{[0,1]} - \ell) e^{\lambda(\text{id}_{[0,1]} - \ell)} \right] \geq 0. \quad (32)$$

The last inequality and limit arguments as  $\lambda \rightarrow -\infty$  impose that  $\text{id}_{[0,1]} - \ell \geq 0$   $\nu$ -a.s., which in turn entails that  $\ell \leq m(\nu)$ . This concludes the proof of (30).

The various properties exhibited above for  $\phi_\nu$ , including the fact that the derivative  $\phi'_\nu$  takes values in  $[m(\nu), +\infty)$ , entail that the function

$$\Lambda : \lambda \in \mathbb{R} \mapsto \lambda x - \phi_\nu(\lambda)$$

is concave, continuously differentiable, with a non-increasing derivative  $\Lambda'$  taking values in the interval  $(-\infty, x - m(\nu)]$  and with limit  $x - m(\nu)$  at  $-\infty$ .

We split the analysis of the case  $x < \mathbb{E}(\nu)$  into three sub-cases, depending on the respective positions of  $x$  and  $m(\nu)$ , and recall that we want to show that  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \phi_\nu^*(x)$ .

*Case 1:  $x > m(\nu)$ .* By Jensen's inequality (23) and given that we consider  $x < \mathbb{E}(\nu)$ , the limit of  $\Lambda$  at  $+\infty$  equals  $-\infty$ . The limit of  $\Lambda$  at  $-\infty$  also equals  $-\infty$ , as the derivative  $\Lambda'$  has limit  $x - m(\nu) > 0$  at  $-\infty$ . By concavity of  $\Lambda$  and the fact that  $\Lambda'$  is continuous, this implies the existence of some  $\lambda^* \in \mathbb{R}$  such that

$$\Lambda'(\lambda^*) = x - \phi'_\nu(\lambda^*) = 0 \quad \text{and} \quad \phi_\nu^*(x) = \sup_{\lambda \in \mathbb{R}} \{\Lambda(\lambda)\} = \Lambda(\lambda^*).$$

Denoting by  $\zeta_{\lambda^*}$  the distribution absolutely continuous with respect to  $\nu$  with density

$$\frac{d\zeta_{\lambda^*}}{d\nu} = \frac{e^{\lambda^* \text{id}_{[0,1]}}}{\mathbb{E}_\nu [e^{\lambda^* \text{id}_{[0,1]}}]} = e^{\lambda^* \text{id}_{[0,1]} - \phi_\nu(\lambda^*)},$$

we have  $\mathbb{E}_{\zeta_{\lambda^*}}[\text{id}_{[0,1]}] = \mathbb{E}(\zeta_{\lambda^*}) = \phi'_\nu(\lambda^*) = x$ . Therefore, by definition of  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$  and of the Kullback-Leibler divergence,

$$\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \text{KL}(\zeta_{\lambda^*}, \nu) = \mathbb{E}_{\zeta_{\lambda^*}} \left[ \ln \frac{d\zeta_{\lambda^*}}{d\nu} \right] = \lambda^* \mathbb{E}_{\zeta_{\lambda^*}}[\text{id}_{[0,1]}] - \phi_\nu(\lambda^*) = \Lambda(\lambda^*) = \phi_\nu^*(x).$$

*Case 2:  $x = m(\nu)$ .* In that case,  $\Lambda' \rightarrow 0$  at  $-\infty$  and  $\Lambda'$  is non-increasing, thus  $\Lambda' \leq 0$  on  $\mathbb{R}$  and  $\Lambda$  is non-increasing on  $\mathbb{R}$ . Thus,

$$\phi_\nu^*(m(\nu)) = \sup_{\lambda \in \mathbb{R}} \{\Lambda(\lambda)\} = \lim_{\lambda \rightarrow -\infty} \Lambda(\lambda) = \lim_{\lambda \rightarrow -\infty} -\ln \mathbb{E}_\nu \left[ e^{\lambda(\text{id}_{[0,1]} - m(\nu))} \right].$$

By monotone convergence based on  $\text{id}_{[0,1]} - m(\nu) \geq 0$   $\nu$ -a.s.,

$$\lim_{\lambda \rightarrow -\infty} -\ln \mathbb{E}_\nu \left[ e^{\lambda(\text{id}_{[0,1]} - m(\nu))} \right] = -\ln \nu\{m(\nu)\},$$

whether  $\nu\{m(\nu)\}$  is positive or null. Moreover, the very end of the proof of Lemma 15 shows that

$$\mathcal{L}_{\text{inf}}^{\leq}(m(\nu), \nu) = -\ln \nu\{m(\nu)\}.$$

We therefore have  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = \phi_\nu^*(x)$  in this case.

*Case 3:  $x < m(\nu)$ .* In that case, as  $\Lambda' \rightarrow x - m(\nu) < 0$  at  $-\infty$ , we get that  $\Lambda \rightarrow +\infty$  at  $-\infty$  and  $\phi_\nu^*(x) = +\infty$ . No distribution  $\zeta \in \mathcal{P}[0, 1]$  with  $\mathbb{E}(\zeta) \leq x$ , if some exists, can be absolutely continuous with respect to  $\nu$ , as  $x < m(\nu)$  imposes that  $\zeta$  puts some probability mass to the left of the support of  $\nu$ . Therefore,  $\text{KL}(\zeta, \nu) = +\infty$ . All in all,  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$  appears as the infimum of either an empty set or of  $+\infty$  values, so that  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = +\infty$ . In this case as well,  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) = \phi_\nu^*(x)$ , both being equal to  $+\infty$ .  $\blacksquare$

### C.3. The case of canonical one-parameter exponential models $\mathcal{D}_{\text{exp}}$

In this section, we consider a model  $\mathcal{D}_{\text{exp}}$  of distributions given by a canonical one-parameter exponential family and will prove that the target equality (15) is satisfied by  $\mathcal{D}_{\text{exp}}$ . Before we do so, let us recall briefly some definitions and properties of exponential families; more detail (including proofs of the stated properties) may be found in Lehmann and Casella (1998).

**Canonical one-parameter exponential families.** We fix a reference real measure  $\rho$  and consider the set, called the natural parameter space,

$$\Theta = \left\{ \theta \in \mathbb{R} : \int_{\mathbb{R}} \exp(\theta y) d\rho(y) < +\infty \right\}.$$

We assume that  $\Theta$  is an open interval (the model is said to be regular) and we consider the model  $\mathcal{D}_{\text{exp}} = \{\nu_\theta : \theta \in \Theta\}$ , where, for  $\theta \in \Theta$ , the distribution  $\nu_\theta$  is absolutely continuous with respect to  $\rho$ , with density

$$\frac{d\nu_\theta}{d\rho} = \exp(\theta \text{id}_{\mathbb{R}} - b(\theta)), \quad (33)$$

for a twice differentiable function  $b : \Theta \rightarrow \mathbb{R}$ . Due to density constraints, we get, for all  $\theta \in \Theta$ ,

$$\underbrace{\int_{\mathbb{R}} \exp(\theta y - b(\theta)) d\rho(y)}_{=d\nu_\theta(y)} = \mathbb{E}_{\nu_\theta}[1] = 1 \quad \text{or, equivalently,} \quad b(\theta) = \ln \int_{\mathbb{R}} e^{\theta y} d\rho(y). \quad (34)$$

It can be seen that  $b$  is strictly convex and that  $\mathbb{E}(\nu_\theta) = b'(\theta)$  for all  $\theta \in \Theta$ . Thus,  $b'$  is a one-to-one mapping between  $\Theta$  and the set  $\mathcal{M}$  of the expectations of the distributions of  $\mathcal{D}_{\text{exp}}$ ; the set  $\mathcal{M}$  is an open interval of  $\mathbb{R}$ , by continuity of  $b'$ , whose bounds are denoted by  $\mu_-$  and  $\mu^+$ . In particular, a distribution of  $\mathcal{D}_{\text{exp}}$  is characterised by its mean. We can also parameterize the Kullback-Leibler divergence function by the expectations: we set, for all  $\theta_1, \theta_2 \in \Theta$ ,

$$d(\mathbb{E}(\nu_{\theta_1}), \mathbb{E}(\nu_{\theta_2})) \stackrel{\text{def}}{=} \text{KL}(\nu_{\theta_1}, \nu_{\theta_2}).$$

This defines a divergence  $d$  which is strictly convex and differentiable on the open set  $\mathcal{M} \times \mathcal{M}$ . In particular,  $d$  is continuous, is such that  $d(\mu, \mu') = 0$  if and only if  $\mu = \mu'$ , and, for all  $\mu \in \mathcal{M}$ , both  $d(\mu, \cdot)$  and  $d(\cdot, \mu)$  are decreasing on  $(\mu_-, \mu]$ , and increasing on  $[\mu, \mu^+)$ .

In the following, we extend  $d$  to  $\mathbb{R} \times \mathbb{R}$  by  $+\infty$  values outside of  $\mathcal{M} \times \mathcal{M}$ .

**Links between  $\mathcal{L}_{\text{inf}}^<$ ,  $\mathcal{L}_{\text{inf}}^{\leq}$ ,  $\mathcal{L}_{\text{inf}}^>$ ,  $\mathcal{L}_{\text{inf}}^{\geq}$  and  $d$ .** A direct application of the continuity and monotonicity properties of  $d$  is that all functions  $\mathcal{L}_{\text{inf}}^<$ ,  $\mathcal{L}_{\text{inf}}^{\leq}$ ,  $\mathcal{L}_{\text{inf}}^>$ ,  $\mathcal{L}_{\text{inf}}^{\geq}$  easily rewrite as  $d$ , as stated in equalities (9) and (10). For instance, for  $\nu \in \mathcal{D}_{\text{exp}}$  and  $x \leq \mathbb{E}(\nu)$ , we have, when  $x \in \mathcal{M}$ :

$$\mathcal{L}_{\text{inf}}^<(x, \nu) = \inf_{\mu < x} \{d(\mu, \nu)\} = \lim_{\substack{\mu \rightarrow x \\ \mu < x}} d(\mu, \nu) = d(x, \nu);$$

and when  $x \notin \mathcal{M}$ , by convention,  $\mathcal{L}_{\text{inf}}^<(x, \nu) = +\infty = d(x, \nu)$ .

We are now able to prove that all canonical one-parameter exponential models satisfy Equation (15), which rewrites as follows, thanks to (9) and (10): for all  $\nu, \nu' \in \mathcal{D}_{\text{exp}}$  with  $\mathbb{E}(\nu') < \mathbb{E}(\nu)$ ,

$$\inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ \phi_{\nu'}^*(x) + \phi_{\nu}^*(x) \right\} = \inf_{x \in [\mathbb{E}(\nu'), \mathbb{E}(\nu)]} \left\{ d(x, \mathbb{E}(\nu')) + d(x, \mathbb{E}(\nu)) \right\}. \quad (35)$$

To obtain this result, we will not prove an exact counterpart of Lemma 4: we will rather only focus on points  $x$  belonging to  $\mathcal{M}$ , not to all  $x \in \mathbb{R}$ ; see (36) below. In particular, we avoid the cases  $x \in \{\mu_-, \mu^+\}$ , the boundary points of  $\mathcal{M}$ , at which the equality of  $d(\cdot, \mathbb{E}(\nu))$  and  $\phi_{\nu}^*$  does not seem to hold in general.

**Proof of (35).** It suffices to show that, for all  $\nu \in \mathcal{D}_{\text{exp}}$  and  $x \in \mathcal{M}$ ,

$$\phi_{\nu}^*(x) \stackrel{\text{def}}{=} \sup_{\lambda \in \mathbb{R}} \{ \lambda x - \phi_{\nu}(\lambda) \} = d(x, \mathbb{E}(\nu)). \quad (36)$$

Now, by Lemma 18, we only need to show that  $\phi_{\nu}^*(x) \geq d(x, \mathbb{E}(\nu))$ , and to prove this inequality, we will justify the existence of  $\lambda^* \in \mathbb{R}$  such that

$$d(x, \mathbb{E}(\nu)) = \lambda^* x - \phi_{\nu}(\lambda^*). \quad (37)$$

Let  $\theta_1 \in \Theta$  be such that  $\nu = \nu_{\theta_1}$  and  $\theta_2 = (b')^{-1}(x) \in \Theta$  be such that  $\mathbb{E}(\nu_{\theta_2}) = x$ . We will prove (37) with  $\lambda^* \stackrel{\text{def}}{=} \theta_2 - \theta_1$ . Using the model density (33), we note that  $\nu_{\theta_2}$  is absolutely continuous with respect to  $\nu_{\theta_1}$  and compute, by definition of the Kullback-Leibler divergence,

$$\begin{aligned} d(x, \mathbb{E}(\nu)) &= \text{KL}(\nu_{\theta_2}, \nu_{\theta_1}) = \mathbb{E}_{\nu_{\theta_2}} \left[ \ln \frac{d\nu_{\theta_2}}{d\nu_{\theta_1}} \right] = \mathbb{E}_{\nu_{\theta_2}} \left[ (\theta_2 - \theta_1) \text{id}_{\mathbb{R}} - (b(\theta_2) - b(\theta_1)) \right] \\ &= (\theta_2 - \theta_1) \mathbb{E}(\nu_{\theta_2}) - (b(\theta_2) - b(\theta_1)) = \lambda^* x - (b(\theta_2) - b(\theta_1)). \end{aligned} \quad (38)$$

To obtain (37), it only remains to show that

$$b(\theta_2) - b(\theta_1) = \phi_{\nu}(\lambda^*). \quad (39)$$

Using (34) at  $\theta_2$  and again the density (33) at  $\theta_1$ , we obtain

$$\begin{aligned} b(\theta_2) &= \ln \int_{\mathbb{R}} e^{\theta_2 y} d\rho(y) = \ln \int_{\mathbb{R}} e^{(\theta_2 - \theta_1)y} \underbrace{e^{\theta_1 y - b(\theta_1)}}_{=d\nu_{\theta_1}(y) = d\nu(y)} d\rho(y) + b(\theta_1) \\ &= \ln \int_{\mathbb{R}} e^{\lambda^* y} d\nu(y) + b(\theta_1) = \phi_{\nu}(\lambda^*) + b(\theta_1), \end{aligned} \quad (40)$$

which gives (39) and concludes the proof of (37).

**Remark 19** A more direct approach bypassing Lemma 18 can be followed with  $\mathcal{D}_{\text{exp}}$  models, along the following lines. The result (40) can be generalized into

$$\forall \theta \in \Theta, \quad \phi_\nu(\theta - \theta_1) = b(\theta) - b(\theta_1).$$

As  $b$  is differentiable on  $\Theta$ , the function  $\phi_\nu$  is also differentiable; at  $\lambda^* = \theta_2 - \theta_1$ , we have

$$\phi'_\nu(\lambda^*) = \phi'_\nu(\theta_2 - \theta_1) = b'(\theta_2) = x.$$

Thus, the derivative of the strictly concave function

$$\Lambda : \lambda \in \mathbb{R} \mapsto \lambda x - \phi_\nu(\lambda)$$

vanishes at  $\lambda^*$ , which is therefore the argument of its maximum:  $\phi_\nu^*(x) = \Lambda(\lambda^*)$ . The latter equality, together with the closed-form calculation (38), shows (36).

#### C.4. Condition for general models

In this section, we extend the equality (15) to more general models; we did so by analyzing which conditions were actually required in the proof of Lemma 4. We recall that the pieces of notation  $m(\nu)$  and  $M(\nu)$  for the lower and upper ends of the support of a distribution  $\nu$  were introduced in Appendix A.2 and lie in  $\mathbb{R} \cup \{-\infty\}$  and  $\mathbb{R} \cup \{+\infty\}$  respectively in all generality. We also remind the reader that all considered models contain distributions with finite first moments.

**Lemma 20** Assume that a model  $\mathcal{D}$  is such that for all distributions  $\nu \in \mathcal{D}$  and all  $\lambda \in \mathbb{R}$ , on the one hand, the quantity  $\phi_\nu(\lambda)$  is well-defined and finite, and, on the other hand, the distribution  $\nu_\lambda$  with density

$$\frac{d\nu_\lambda}{d\nu} = \frac{e^{\lambda \text{id}_{\mathbb{R}}}}{\mathbb{E}_\nu[e^{\lambda \text{id}_{\mathbb{R}}}]}$$
 with respect to  $\nu$

belongs to  $\mathcal{D}$ . Assume also that  $\delta_x$ , the Dirac mass at  $x$ , belongs to  $\mathcal{D}$  whenever there exists  $\nu \in \mathcal{D}$  with  $x \in \{m(\nu), M(\nu)\}$  and  $\nu\{x\} > 0$ .

Then, for all  $\nu \in \mathcal{D}$ ,

$$\forall x \leq E(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq E(\nu), \quad \phi_\nu^*(x) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

**Proof** By Lemma 18, we only need to prove that

$$\forall x \leq E(\nu), \quad \phi_\nu^*(x) \geq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \quad \text{and} \quad \forall x \geq E(\nu), \quad \phi_\nu^*(x) \geq \mathcal{L}_{\text{inf}}^{\geq}(x, \nu).$$

We will follow the proof of Lemma 4 for the model  $\mathcal{P}[0, 1]$  (see Appendix C.2), and show that the same proof applies to  $\mathcal{D}$ , up to a few changes in the justifications. In particular, we only treat the first case, showing that

$$\forall x \leq E(\nu), \quad \phi_\nu^*(x) \geq \mathcal{L}_{\text{inf}}^{\leq}(x, \nu), \tag{41}$$

and obtain the other by symmetry. For  $x = E(\nu)$ , we show, as in the  $\mathcal{P}[0, 1]$  case, that

$$\phi_\nu^*(E(\nu)) = 0 = \mathcal{L}_{\text{inf}}^{\leq}(E(\nu), \nu).$$

Before proving (41) for  $x < E(\nu)$ , we show that  $\phi_\nu$  satisfies the same required properties than in the  $\mathcal{P}[0, 1]$  case. As  $\phi_\nu$  takes finite values on  $\mathbb{R}$ , the moment-generating function of  $\nu$  is a power series with infinite radius of convergence. In particular,  $\phi_\nu$  is differentiable over  $\mathbb{R}$  and, as Hölder's inequality still entails that  $\phi_\nu$  is convex, the derivative  $\phi'_\nu$  is non-decreasing. We now prove that the limit of the latter at  $-\infty$  (which we recall exists and belongs to  $\{-\infty\} \cup \mathbb{R}$  by monotonicity) equals  $m(\nu)$ :

$$\lim_{\lambda \rightarrow -\infty} \phi'_\nu(\lambda) \stackrel{\text{def}}{=} \ell = m(\nu). \quad (42)$$

If  $m(\nu) \in \mathbb{R}$ , the same proof as in the  $\mathcal{P}[0, 1]$  case applies. Otherwise,  $m(\nu) = -\infty$  and we prove that  $\ell$  cannot be finite. Indeed, if  $\ell \in \mathbb{R}$ , then using that  $\phi'_\nu$  is always larger than its limit  $\ell$  by non-decreasing, and following the same analysis as in (31)–(32), we obtain  $\text{id}_{\mathbb{R}} - \ell \geq 0$   $\nu$ -a.s., which in turn entails that  $\ell \leq m(\nu)$  and contradicts the fact that  $m(\nu) = -\infty$ . This concludes the proof of (42).

Finally, we prove that

$$\phi'_\nu : \lambda \in \mathbb{R} \mapsto \frac{\mathbb{E}_\nu[\text{id}_{\mathbb{R}} e^{\lambda \text{id}_{\mathbb{R}}}]}{\mathbb{E}_\nu[e^{\lambda \text{id}_{\mathbb{R}}}]},$$

by applying a standard theorem of differentiation under the integral. To obtain the domination on a subset  $(\lambda_-, \lambda_+)$  of  $\mathbb{R}$ , we write that

$$\forall \lambda \in (\lambda_-, \lambda_+), \quad |\text{id}_{\mathbb{R}} e^{\lambda \text{id}_{\mathbb{R}}}| \leq |\text{id}_{\mathbb{R}}| \left( e^{\lambda_- \text{id}_{\mathbb{R}}} + e^{\lambda_+ \text{id}_{\mathbb{R}}} \right) \stackrel{\text{def}}{=} h$$

and get that  $h$  is  $\nu$ -integrable whenever distributions  $\nu_{\lambda_-}$  and  $\nu_{\lambda_+}$  admit a finite first moment, which holds by assumption as those distributions belong to  $\mathcal{D}$ .

We now get all necessary properties to follow the proof of Lemma 4. Similarly to this proof, we split the analysis of the case  $x < E(\nu)$  into three sub-cases, depending on the respective positions of  $x$  and  $m(\nu)$ .

*Case 1:  $x > m(\nu)$ .* The properties of  $\phi_\nu$  ensure, as in the corresponding  $\mathcal{P}[0, 1]$  case, the existence of  $\lambda^*$  such that

$$\phi_\nu^*(x) = \lambda^* x - \phi_\nu(\lambda^*).$$

Besides, we get, again as in the proof of Lemma 4, that  $E(\nu_{\lambda^*}) = \phi'_\nu(\lambda^*) = x$  and compute that

$$\text{KL}(\nu_{\lambda^*}, \nu) = \lambda^* x - \phi_\nu(\lambda^*).$$

This gives (41) noting that  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu) \leq \text{KL}(\delta_{m(\nu)}, \nu)$ , as by assumptions  $\nu_{\lambda^*} \in \mathcal{D}$ .

*Case 2:  $x = m(\nu)$ .* As in the proof of this case in Lemma 4, we prove that

$$\phi_\nu^*(m(\nu)) = -\ln \nu\{m(\nu)\},$$

and get

$$\text{KL}(\delta_{m(\nu)}, \nu) = -\ln \nu\{m(\nu)\}.$$

This also concludes this case, as  $\delta_{m(\nu)} \in \mathcal{D}$  by assumption.

*Case 3:  $x < m(\nu)$ .* Similarly to the corresponding  $\mathcal{P}[0, 1]$  case, we show that both  $\mathcal{L}_{\text{inf}}^{\leq}(x, \nu)$  and  $\phi_\nu^*(x)$  are infinite.  $\blacksquare$

## Appendix D. Proofs for lower bounds (Section 4)

This section provides the detailed proofs that were omitted when stating our various lower bounds in Section 4.

### D.1. Proof of Lemma 6

We restate it for the convenience of the reader.

**Lemma 6** Fix  $K \geq 2$  and a model  $\mathcal{D}$ . Consider a consistent sequence of strategies on  $\mathcal{D}$ , and two generic bandit problems  $\underline{\nu}$  and  $\underline{\lambda}$  in  $\mathcal{D}$  such that  $a^*(\underline{\lambda}) \neq a^*(\underline{\nu})$ . Then

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \limsup_{T \rightarrow +\infty} \sum_{a=1}^K \frac{\mathbb{E}_{\underline{\lambda}}[N_a(T)]}{T} \text{KL}(\lambda_a, \nu_a).$$

**Proof** The considered sequence of strategies being consistent on  $\mathcal{D}$ , and as  $a^*(\underline{\lambda}) \neq a^*(\underline{\nu})$ ,

$$\begin{aligned} q_T &\stackrel{\text{def}}{=} \mathbb{P}_{\underline{\lambda}}(I_T \neq a^*(\underline{\nu})) \geq \mathbb{P}_{\underline{\lambda}}(I_T = a^*(\underline{\lambda})) \xrightarrow{T \rightarrow +\infty} 1, \\ \text{while } p_T &\stackrel{\text{def}}{=} \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \xrightarrow{T \rightarrow +\infty} 0. \end{aligned}$$

Note that we introduced above short-hand notation  $p_T$  and  $q_T$ .

The fundamental inequality for lower bounds in bandit problems (which is a consequence of the chain rule and of the data-processing inequality for Kullback-Leibler divergences, see [Garivier et al., 2019](#)), applied for  $Z = \mathbb{I}_{\{I_T \neq a^*(\underline{\nu})\}}$ , exactly states here that

$$\sum_{a=1}^K \mathbb{E}_{\underline{\lambda}}[N_a(T)] \text{KL}(\lambda_a, \nu_a) \geq \text{KL}(\text{Ber}(q_T), \text{Ber}(p_T)), \quad (43)$$

where we recall that  $\text{Ber}(p)$  refers to the Bernoulli distribution with parameter  $p$ . Given the asymptotics of  $p_T$  and  $q_T$ ,

$$\text{KL}(\text{Ber}(q_T), \text{Ber}(p_T)) = q_T \ln \frac{q_T}{p_T} + (1 - q_T) \ln \frac{1 - q_T}{1 - p_T} \sim -\ln p_T \quad \text{as } T \rightarrow +\infty.$$

Put differently,

$$\frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \sim -\frac{\text{KL}(\text{Ber}(q_T), \text{Ber}(p_T))}{T}.$$

Combining this limit behavior with the previous inequality leads to the stated result, namely:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \limsup_{T \rightarrow +\infty} \sum_{a=1}^K \frac{\mathbb{E}_{\underline{\lambda}}[N_a(T)]}{T} \text{KL}(\lambda_a, \nu_a). \quad \blacksquare$$

## D.2. Proof of Theorem 9

We restate it for the convenience of the reader.

**Theorem 9** *Fix a model  $\mathcal{D}$ . Consider a doubly-indexed sequence of strategies that is consistent, balanced against the worst arm on  $\mathcal{D}$ , and that cleverly exploits the pruning of suboptimal arms on  $\mathcal{D}$ . For all generic bandit problems  $\underline{\nu}$  in  $\mathcal{D}$  with  $K \geq 2$  arms,*

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \frac{\mathcal{L}_{\text{inf}}^<(\mu_{(k)}, \nu^*)}{k}.$$

**Proof** The proof consists of two steps. The first step is to prove that for a generic bandit problem  $\underline{\nu}$  in  $\mathcal{D}$  with  $K \geq 2$  arms, we have,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \frac{\mathcal{L}_{\text{inf}}^<(\mu_{(K)}, \nu^*)}{K}. \quad (44)$$

In the second step, we use this lower bound and the very definition of the clever exploitation of the pruning of suboptimal arms to get the claimed bound.

**Step 1: lower bound (44).** We follow a well-established methodology and consider an alternative bandit problem only differing from  $\underline{\nu}$  at one arm, namely, at the best arm. To do so, we set some distribution  $\zeta \in \mathcal{D}$  with  $\mathbb{E}(\zeta) < \mu_{(K)}$ , if some exists, and define the bandit problem  $\underline{\lambda} = (\lambda_1, \dots, \lambda_K)$  as

$$\lambda_a = \begin{cases} \zeta & \text{if } a = a^*(\underline{\nu}), \\ \nu_a & \text{if } a \neq a^*(\underline{\nu}). \end{cases}$$

Observe that  $\underline{\lambda}$  is also a generic bandit problem in  $\mathcal{D}$ , that  $a^*(\underline{\nu})$  is the worst arm in  $\underline{\lambda}$  (and also that the second best arm of  $\underline{\nu}$  is the optimal arm in  $\underline{\lambda}$ , but we will not use this specific fact). Therefore, Lemma 6 yields, as  $\underline{\lambda}$  and  $\underline{\nu}$  only differ at arm  $a^*(\underline{\nu})$ ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_{a^*(\underline{\nu})}(T)]}{T} \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu^*),$$

where we recall that  $\nu^* = \nu_{a^*(\underline{\nu})}$ . Given that  $a^*(\underline{\nu})$  is the worst arm of  $\underline{\lambda}$ , and since by assumption, the sequence of strategies is balanced against the worst arm,

$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \mathbb{E}_{\underline{\lambda}}[N_{a^*(\underline{\nu})}(T)] \leq \frac{1}{K},$$

proving that

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \frac{\text{KL}(\zeta, \nu^*)}{K}.$$

The claimed inequality (44) follows from taking the supremum in the right-hand side over distributions  $\zeta \in \mathcal{D}$  with  $\mathbb{E}(\zeta) < \mu_{(K)}$ .



**Step 2: clever exploitation of pruning.** For each  $k \in \{2, \dots, K-1\}$ , define  $\underline{\nu}'_{1:k}$  as the subproblem of  $\underline{\nu}$  obtained by keeping the  $k$  best arms and dropping the  $K-k$  worse arms. Use the definition of clever exploitation of pruning of suboptimal arms and apply (44) to  $\underline{\nu}'_{1:k}$  to get

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq \liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}'_{1:k}}(I_T \neq a^*(\underline{\nu}'_{1:k})) \geq -\frac{\mathcal{L}_{\text{inf}}^<(\mu^{(k)}, \nu^*)}{k}.$$

Taking the maximum of all lower bounds exhibited as  $k$  varies between 2 and  $K$ , we proved the claimed result.  $\blacksquare$

### D.3. Proof of the normality of the models $\mathcal{P}[0, 1]$ and $\mathcal{D}_{\text{exp}}$

In this section, we show that  $\mathcal{P}[0, 1]$  and canonical one-parameter exponential models are normal. We focus first on  $\mathcal{P}[0, 1]$ .

**Proposition 21**  *$\mathcal{P}[0, 1]$  is a normal model.*

**Proof** We fix  $\nu \in \mathcal{P}[0, 1]$ , a real  $x \geq \mathbb{E}(\nu)$  and consider a positive  $\varepsilon$ . We recall that the pieces of notation  $m(\nu)$  and  $M(\nu)$  for the lower and upper ends of the support of a distribution  $\nu$  were introduced in Appendix A.2. If  $x \geq M(\nu)$ , there exists no distribution  $\zeta$  absolutely continuous with respect to  $\nu$  and such that  $\mathbb{E}(\zeta) > x$ , hence the considered infima are infinite:

$$\mathcal{L}_{\text{inf}}^>(x, \nu) = +\infty = \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x\}.$$

Assume now that  $\mathbb{E}(\nu) \leq x < M(\nu)$  and note that this case only occurs when  $\mathbb{E}(\nu) < M(\nu)$ , i.e. when  $\nu$  is not a Dirac distribution. It is clear that

$$\mathcal{L}_{\text{inf}}^>(x, \nu) \leq \inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x\}.$$

To prove the other inequality  $\geq$ , we will first remind why

$$\mathcal{L}_{\text{inf}}^>(x, \nu) = \text{KL}(\zeta_{\lambda^*}, \nu), \quad (45)$$

where  $\zeta_{\lambda^*}$  is a distribution of  $\mathcal{D}$  of mean  $x$  introduced in the proof of Lemma 4. As the expectation of  $\zeta_{\lambda^*}$  does not belong to  $(x, x + \varepsilon)$ , we will then need to prove that

$$\inf\{\text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x\} \leq \text{KL}(\zeta_{\lambda^*}, \nu). \quad (46)$$

In the remaining of the proof, we rely on notation and results proved in the first case (by symmetry) of the proof of Lemma 4, and for  $\lambda \in \mathbb{R}$ , we define  $\zeta_{\lambda}$  the distribution absolutely continuous with respect to  $\nu$  with density

$$\frac{d\zeta_{\lambda}}{d\nu} = \frac{e^{\lambda \text{id}_{[0,1]}}}{\mathbb{E}_{\nu}[e^{\lambda \text{id}_{[0,1]}}]} = e^{\lambda \text{id}_{[0,1]} - \phi_{\nu}(\lambda)},$$

On the one hand, by Lemma 15 (or more precisely its formulation for  $\mathcal{L}_{\text{inf}}^>$  and  $\mathcal{L}_{\text{inf}}^{\geq}$  instead of  $\mathcal{L}_{\text{inf}}^<$  and  $\mathcal{L}_{\text{inf}}^{\leq}$ ), we know that

$$\mathcal{L}_{\text{inf}}^>(x, \nu) = \mathcal{L}_{\text{inf}}^{\geq}(x, \nu),$$

while, on the other hand, the proof of the first case of Lemma 4 ensures the existence of  $\lambda^* \in \mathbb{R}$  such that  $\mathbb{E}(\zeta_{\lambda^*}) = \phi'_\nu(\lambda^*) = x$  and

$$\mathcal{L}_{\inf}^{\geq}(x, \nu) = \text{KL}(\zeta_{\lambda^*}, \nu).$$

We thus obtained (45) and move to Equation (46).

We noticed that  $\nu$  is not a Dirac distribution, as  $\mathbb{E}(\nu) \leq x < M(\nu)$ . This entails, by Hölder's inequality, that  $\phi_\nu$  is strictly convex, hence  $\phi'_\nu$  is an increasing function. Considering distributions  $\zeta_\lambda$  for which  $\lambda > \lambda^*$ , this implies that

$$\mathbb{E}_{\zeta_\lambda}[\text{id}_{[0,1]}] = \mathbb{E}(\zeta_\lambda) = \phi'_\nu(\lambda) > x = \mathbb{E}(\zeta_{\lambda^*}).$$

As  $\phi'_\nu$  is a continuous function (see, again, Lemma 4), we also get

$$\lim_{\substack{\lambda \rightarrow \lambda^* \\ \lambda > \lambda^*}} \mathbb{E}(\zeta_\lambda) = \mathbb{E}(\zeta_{\lambda^*}),$$

so that there exists  $\lambda_\varepsilon > \lambda^*$  such that

$$\forall \lambda \in (\lambda^*, \lambda_\varepsilon), \quad \mathbb{E}(\zeta_\lambda) \in (x, x + \varepsilon).$$

As a consequence,

$$\inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x \} \leq \inf_{\lambda \in (\lambda_\varepsilon, \lambda^*)} \{ \text{KL}(\zeta_\lambda, \nu) \}, \quad (47)$$

and we compute that, for  $\lambda \in (\lambda^*, \lambda_\varepsilon)$

$$\begin{aligned} \text{KL}(\zeta_\lambda, \nu) &= \mathbb{E}_{\zeta_\lambda} \left[ \ln \frac{d\zeta_\lambda}{d\nu} \right] = \lambda \mathbb{E}_{\zeta_\lambda}[\text{id}_{[0,1]}] - \phi_\nu(\lambda) = \lambda \phi'_\nu(\lambda) - \phi_\nu(\lambda) \\ &\xrightarrow{\lambda \rightarrow \lambda^*} \lambda^* \phi'_\nu(\lambda^*) - \phi_\nu(\lambda^*) = \text{KL}(\zeta_{\lambda^*}, \nu), \end{aligned}$$

by continuity of  $\phi_\nu$  and  $\phi'_\nu$ . Combining this limit behaviour with inequality (47) leads to (46).  $\blacksquare$

We now consider canonical one-parameter exponential models, for which normality is easily obtained by the rewriting of  $\mathcal{L}_{\inf}^{\geq}$  as  $d$ .

**Proposition 22** *All canonical one-parameter exponential models  $\mathcal{D}_{\text{exp}}$  are normal.*

**Proof** We fix  $\nu \in \mathcal{D}_{\text{exp}}$ , a real  $x \geq \mathbb{E}(\nu)$  and consider a positive  $\varepsilon$ . As in the  $\mathcal{P}[0, 1]$  case, the required equality holds (both terms are infinite) when  $x \geq M(\nu)$ .

Assume now that  $x < M(\nu)$ . By equality (10), we recall that

$$\mathcal{L}_{\inf}^{\geq}(x, \nu) = d(x, \mathbb{E}(\nu)).$$

Using the basic properties of  $d$  (see Appendix C.3), namely that  $d(\cdot, \mathbb{E}(\nu))$  is continuous at  $x$  and is non-decreasing on  $(x, +\infty) \subset [\mathbb{E}(\nu), +\infty)$ , this gives

$$\begin{aligned} \mathcal{L}_{\inf}^{\geq}(x, \nu) &= \lim_{y \rightarrow x^+} d(y, \mathbb{E}(\nu)) \\ &= \inf_{y \in (x, x+\varepsilon)} \{ d(y, \mathbb{E}(\nu)) \} \\ &= \inf \{ \text{KL}(\zeta, \nu) : \zeta \in \mathcal{D} \text{ s.t. } x + \varepsilon > \mathbb{E}(\zeta) > x \}. \end{aligned}$$

$\blacksquare$

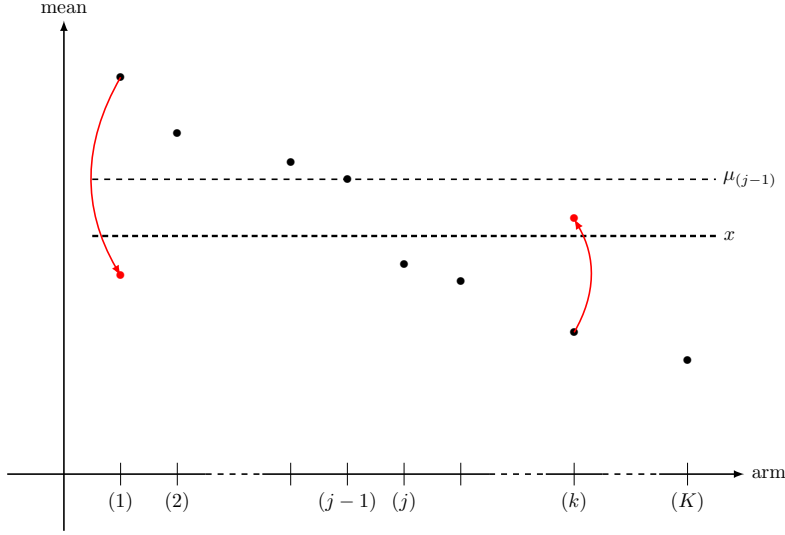


Figure 1: Original bandit problem  $\underline{\nu}$  (in dark) and modifications made to arms (1) and (k) to obtain an alternative bandit problem  $\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})$  (in red): in  $\underline{\lambda}$ , arm (k) is the  $j - 1$ -th best arm, while arm (1) =  $a^*(\underline{\nu})$  is at best the  $j$ -th best arm.

#### D.4. Proof of Theorem 12

We restate it for the convenience of the reader.

**Theorem 12** Fix  $K \geq 2$  and a normal model  $\mathcal{D}$ . Consider a sequence of strategies which is consistent and monotonous on  $\mathcal{D}$ . For all generic bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{2 \leq k \leq K} \min_{2 \leq j \leq k} \inf_{x \in [\mu_{(j)}, \mu_{(j-1)})} \left\{ \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{j-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu_{(1)})}{j} \right\}.$$

**Proof** We fix a generic bandit  $\underline{\nu}$  in  $\mathcal{D}$  and consider the following sets of alternative bandit problems, indexed by triplets  $(k, j, x)$  satisfying  $2 \leq k \leq K$  and  $2 \leq j \leq k$ , as well as  $x \in [\mu_{(j)}, \mu_{(j-1)})$ :

$$\text{Alt}_{k,j,x}(\underline{\nu}) = \left\{ \underline{\lambda} \text{ in } \mathcal{D} : \mathbb{E}(\lambda_{(1)}) < x < \mathbb{E}(\lambda_{(k)}) < \mu_{(j-1)} \text{ and } \lambda_a = \nu_a \text{ for } a \notin \{(1), (k)\} \right\};$$

in particular, an alternative problem  $\underline{\lambda}$  in  $\text{Alt}_{k,j,x}(\underline{\nu})$  only differ from the original bandit problem  $\underline{\nu}$  at the best arm (1) and at the  $k$ -th best arm (k). Given  $x \in [\mu_{(j)}, \mu_{(j-1)})$  and  $\mathbb{E}(\lambda_{(1)}) < x$ , arm (1) is at best the  $j$ -th best arm of  $\underline{\lambda}$ , but it can be possibly worse. Similarly, the same condition on  $x$  and the fact that  $x < \mathbb{E}(\lambda_{(k)})$  implies that arm (k) is exactly the  $j - 1$ -th best arm of  $\underline{\lambda}$ . Both facts are illustrated on Figure 1. Thus, by monotonicity of the strategy,

$$\limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_{(k)}(T)]}{T} \leq \frac{1}{j-1} \quad \text{and} \quad \limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_{(1)}(T)]}{T} \leq \frac{1}{j}.$$

Given that the optimal arm in  $\underline{\lambda}$  is different from the optimal arm (1) of  $\underline{\nu}$ , Lemma 6 may be applied; together with the two upper bounds above, it yields

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \left( \frac{\text{KL}(\lambda_{(k)}, \nu_{(k)})}{j-1} + \frac{\text{KL}(\lambda_{(1)}, \nu_{(1)})}{j} \right).$$

We can now take the infimum over all bandit problems  $\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})$  and obtain the following lower bound, where we define a quantity  $\mathcal{I}_{k,j,x}(\underline{\nu})$ :

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \inf_{\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})} \left\{ \frac{\text{KL}(\lambda_{(k)}, \nu_{(k)})}{j-1} + \frac{\text{KL}(\lambda_{(1)}, \nu^*)}{j} \right\} \stackrel{\text{def}}{=} -\mathcal{I}_{k,j,x}(\underline{\nu}).$$

We prove below that

$$\mathcal{I}_{k,j,x}(\underline{\nu}) = \frac{\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})}{j-1} + \frac{\mathcal{L}_{\text{inf}}^<(x, \nu^*)}{j}, \quad (48)$$

from which the lower bound claimed in Theorem 12 will follow, by taking the supremum of  $-\mathcal{I}_{k,j,x}(\underline{\nu})$  first over  $x \in [\mu_{(j)}, \mu_{(j-1)})$ , then the maximum over  $2 \leq j \leq k$ , and finally, the maximum over  $2 \leq k \leq K$ .

We now prove (48). The infimum over  $\underline{\lambda} \in \text{Alt}_{k,j,x}(\underline{\nu})$  may be split into two separate infima, respectively over  $\lambda_{(k)}$  and  $\lambda_{(1)}$ ; given that each term of the sum of KL only depends either on  $\lambda_{(k)}$ , or on  $\lambda_{(1)}$ , but not on both, we may write

$$\begin{aligned} \mathcal{I}_{k,j,x}(\underline{\nu}) &= \inf_{\substack{\lambda_{(1)}, \lambda_{(k)} \in \mathcal{D}: \\ \mathbb{E}(\lambda_{(1)}) < x \\ x < \mathbb{E}(\lambda_{(k)}) < \mu_{(j-1)}}} \left\{ \frac{\text{KL}(\lambda_{(k)}, \nu_{(k)})}{j-1} + \frac{\text{KL}(\lambda_{(1)}, \nu^*)}{j} \right\} \\ &= \frac{1}{j-1} \underbrace{\inf_{\substack{\lambda_{(k)} \in \mathcal{D}: \\ x < \mathbb{E}(\lambda_{(k)}) < \mu_{(j-1)}}} \text{KL}(\lambda_{(k)}, \nu_{(k)})}_{=\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})} + \frac{1}{j} \underbrace{\inf_{\substack{\lambda_{(1)} \in \mathcal{D}: \\ \mathbb{E}(\lambda_{(1)}) < x}} \text{KL}(\lambda_{(1)}, \nu^*)}_{=\mathcal{L}_{\text{inf}}^<(x, \nu^*)}, \end{aligned}$$

where we obtain  $\mathcal{L}_{\text{inf}}^<(x, \nu^*)$  by definition while we rely on the normality of the model to obtain  $\mathcal{L}_{\text{inf}}^>(x, \nu_{(k)})$ : we do so with  $\varepsilon = \mu_{(j-1)} - x$ , which is indeed positive as we considered  $x < \mu_{(j-1)}$ . ■

### D.5. Proof of Theorem 13

We restate it for the convenience of the reader.

**Theorem 13** Fix  $K \geq 2$  and a model  $\mathcal{D}$ . Consider a consistent sequence of strategies on  $\mathcal{D}$ . For all generic bandit problems  $\underline{\nu}$  in  $\mathcal{D}$ ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \min_{k \neq a^*(\underline{\nu})} \inf_{x \in [\mu_k, \mu^*]} \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\}.$$

**Proof** Let  $\underline{\nu}$  be a generic bandit problem. We fix  $k \neq a^*(\underline{\nu})$  and  $x \in [\mu_k, \mu^*]$ , and prove that

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \max\{\mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*)\},$$

from which the stated lower bound follows, by taking suprema. To do so, we consider the set of alternative bandit problems

$$\text{Alt}_{k,x}(\underline{\nu}) = \left\{ \underline{\lambda} \text{ in } \mathcal{D} : \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < x < \mathbb{E}(\lambda_k) \text{ and } \lambda_a = \nu_a \text{ for } a \notin \{a^*(\underline{\nu}), k\} \right\};$$

it is composed of bandit problems, only differing from  $\underline{\nu}$  at arms  $a^*(\underline{\nu})$  and  $k$ , and for which arm  $k$  is better than arm  $a^*(\underline{\nu})$ , with associated expectations separated by  $x$ . In particular, the optimal arm in  $\underline{\lambda}$  is different from the optimal arm  $a^*(\underline{\nu})$  of  $\underline{\nu}$ . Lemma 6 may therefore be applied; it states that

$$\begin{aligned} & \liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \\ & \geq - \limsup_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\lambda}}[N_k(T)]}{T} \text{KL}(\lambda_k, \nu_k) + \frac{\mathbb{E}_{\underline{\lambda}}[N_{a^*(\underline{\nu})}(T)]}{T} \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \\ & \geq - \max \left\{ \text{KL}(\lambda_k, \nu_k), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\}, \end{aligned}$$

where we used, for the second inequality, the crude upper bound  $N_k(T) + N_{a^*(\underline{\nu})}(T) \leq T$ . Taking the supremum of the obtained lower bound over all  $\underline{\lambda} \in \text{Alt}_{k,x}(\underline{\nu})$  leads to the following inequality, where we define the short-hand notation  $\mathcal{I}_{k,x}(\underline{\nu})$ :

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq - \inf_{\underline{\lambda} \in \text{Alt}_{k,x}(\underline{\nu})} \max \left\{ \text{KL}(\lambda_k, \nu_k), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\} \stackrel{\text{def}}{=} -\mathcal{I}_{k,x}(\underline{\nu}).$$

The proof is concluded below by showing that  $\mathcal{I}_{k,x}(\underline{\nu}) = \max \left\{ \mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*) \right\}$ .

As in the proof of Theorem 12 (see Appendix D.4), we use a separation of the infima, in the abstract form, for two functions  $f$  and  $g$ ,

$$\inf_{u,v} \max \{ f(u), g(v) \} = \max \left\{ \inf_u f(u), \inf_v g(v) \right\}.$$

Here, by definition of  $\text{Alt}_{k,x}(\underline{\nu})$ ,

$$\begin{aligned} \mathcal{I}_{k,x}(\underline{\nu}) &= \inf_{\substack{\lambda_{a^*(\underline{\nu})}, \lambda_k \in \mathcal{D} \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < x \\ \mathbb{E}(\lambda_k) > x}} \max \left\{ \text{KL}(\lambda_k, \nu_k), \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\} \\ &= \max \left\{ \inf_{\substack{\lambda_k \in \mathcal{D} \\ \mathbb{E}(\lambda_k) > x}} \text{KL}(\lambda_k, \nu_k), \inf_{\substack{\lambda_{a^*(\underline{\nu})} \in \mathcal{D} \\ \mathbb{E}(\lambda_{a^*(\underline{\nu})}) < x}} \text{KL}(\lambda_{a^*(\underline{\nu})}, \nu_{a^*(\underline{\nu})}) \right\} \\ &= \max \left\{ \mathcal{L}_{\text{inf}}^>(x, \nu_k), \mathcal{L}_{\text{inf}}^<(x, \nu^*) \right\}, \end{aligned}$$

which concludes the proof. ■

## Appendix E. Additional comments for the literature review

This appendix is devoted to additional discussions concerning the fixed-budget literature. More precisely, we discuss in detail two gap-based lower bounds that we believe are somewhat detached from the spirit of the article, namely, the minimax lower bound of [Carpentier and Locatelli \(2016\)](#) (Appendix E.1) and the Bretagnolle-Huber technique (Appendix E.2).

### E.1. The minimax lower bound of [Carpentier and Locatelli \(2016\)](#)

[Carpentier and Locatelli \(2016\)](#) proved the following non-asymptotic minimax lower bound: for the model  $\mathcal{D}_{1/4}$  of Bernoulli distributions with parameters in  $[1/4, 3/4]$ , for a given consistent strategy, and each large enough budget  $T$ ,

$$\begin{aligned} \exists \underline{\nu} \text{ in } \mathcal{D}_{1/4}, \quad \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) &\geq -\frac{400}{\ln K} \left( \sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2} \right)^{-1} + o(1) \\ &\geq -\frac{400}{\ln K} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k} + o(1), \end{aligned} \quad (49)$$

where the  $o(1)$  is with respect to  $T \rightarrow +\infty$ .

**Different nature: a non-uniform lower bound.** This result is different in nature from the lower bounds considered in this article, as we now discuss. First and foremost, it improves the lower bound (6) of [Audibert et al. \(2010\)](#) for only one (unspecified) bandit problem  $\underline{\nu}$  (belonging to a known collection of  $K$  bandit problems). This is in strong contrast with the instance-dependent lower bounds (bounds holding simultaneously for all bandit problems of the model) presented in this article.

Second, the result is stated for the restricted Bernoulli model  $\mathcal{D}_{1/4}$  and does not seem to be easily generalized beyond this model or beyond similar models (e.g., the model  $\mathcal{D}_{\sigma^2}$  of Gaussian distributions with a fixed variance  $\sigma^2 > 0$ ). This is because of the truly gap-based arguments used (see the second part of the proof below).

**Asymptotic statement and proof.** The bound (49) may actually be stated for a sequence of strategies (as we do for all other bounds in this article) thanks to a straightforward adaptation of its proof (relying on the pigeonhole principle). More precisely, the counterpart of (49) would be that there exists an increasing sequence of budgets  $(T_n)_{n \in \mathbb{N}}$  such that

$$\exists \underline{\nu} \text{ in } \mathcal{D}_{1/4}, \quad \liminf_{n \rightarrow +\infty} \frac{1}{T_n} \ln \mathbb{P}_{\underline{\nu}}(I_{T_n} \neq a^*(\underline{\nu})) \geq -\frac{400}{\ln K} \left( \sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2} \right)^{-1}. \quad (50)$$

We are actually able to slightly improve the numerical factor 400 into 100/3 by using a sharper change-of-measure argument (namely, Lemma 6) than the original argument by [Carpentier and Locatelli \(2016\)](#).

**Proposition 23** Fix  $K \geq 3$ , consider the model  $\mathcal{D}_{1/4}$  of Bernoulli distributions with parameters in  $[1/4, 3/4]$ , and a consistent sequence of strategies on  $\mathcal{D}_{1/4}$ . Then, there exists an increasing

sequence of budgets  $(T_n)_{n \in \mathbb{N}}$  such that

$$\exists \underline{\nu} \text{ in } \mathcal{D}_{1/4}, \quad \liminf_{n \rightarrow +\infty} \frac{1}{T_n} \ln \mathbb{P}_{\underline{\nu}}(I_{T_n} \neq a^*(\underline{\nu})) \geq -\frac{100}{3 \ln K} \left( \sum_{a \neq a^*(\underline{\nu})} \frac{1}{\Delta_a^2} \right)^{-1}.$$

**Proof** We begin the proof by introducing a collection of bandit problems  $(\underline{\nu}^{(k)})_{1 \leq k \leq K}$  in  $\mathcal{D}_{1/4}$  and associated notation. Let  $\underline{\nu}^{(1)}$  be a Bernoulli bandit problem such that the mean of arm 1 is  $p_1 = 1/2$  and the mean of arm  $k \in \{2, \dots, K\}$  is  $p_k \in [1/4, 1/2)$ . Let  $\Delta_k = 1/2 - p_k$  denotes the gap of arm  $k$  in  $\underline{\nu}^{(1)}$ . For  $k \in \{2, \dots, K\}$ , we set  $\underline{\nu}^{(k)}$  the bandit problem obtained by changing the arm  $k$  of  $\underline{\nu}^{(1)}$  to a Bernoulli distribution with mean  $1 - p_k$ , hence  $a^*(\underline{\nu}^{(k)}) = k$ , and we define

$$H(\underline{\nu}^{(k)}) = \sum_{a \neq k} \frac{1}{(\Delta_a^{(k)})^2},$$

where  $\Delta_a^{(k)}$  is the gap of arm  $a$  in  $\underline{\nu}^{(k)}$ . As for  $a \neq k$ , we notice that

$$\Delta_a^{(k)} = (1 - p_k) - p_a = (1/2 - p_k) + (1/2 - p_a) = \Delta_k + \Delta_a,$$

we have

$$H(\underline{\nu}^{(k)}) = \sum_{a \neq k} \frac{1}{(\Delta_k + \Delta_a)^2}. \quad (51)$$

Finally we define

$$H^* = \sum_{k=2}^K \frac{1}{\Delta_k^2 H(\underline{\nu}^{(k)})}. \quad (52)$$

The proof consists of two steps. Firstly, we show that there exists  $k \in \{2, \dots, K\}$  and an increasing sequence of budgets  $(T_n)_{n \in \mathbb{N}}$  such that

$$\liminf_{n \rightarrow +\infty} \frac{1}{T_n} \ln \mathbb{P}_{\underline{\nu}^{(k)}}(I_{T_n} \neq k) \geq -\frac{10}{H^*} \frac{1}{H(\underline{\nu}^{(k)})}. \quad (53)$$

Then, in a second step, we set particular values of the  $(p_k)_{2 \leq k \leq K}$  and show that

$$H^* \geq \frac{3}{10} \ln K. \quad (54)$$

**Step 1: lower bound (53) for general values of  $(p_k)_k$ .** We observe that, for each  $T \in \mathbb{N}^*$ ,

$$\sum_{k=2}^K \frac{\mathbb{E}_{\underline{\nu}^{(1)}}[N_k(T)]}{T} \leq 1 = \frac{1}{H^*} \sum_{k=2}^K \frac{1}{\Delta_k^2 H(\underline{\nu}^{(k)})} = \sum_{k=2}^K \frac{1}{H^* \Delta_k^2 H(\underline{\nu}^{(k)})},$$

hence there exists an arm  $k_T \in \{2, \dots, K\}$  such that

$$\frac{\mathbb{E}_{\underline{\nu}^{(1)}}[N_{k_T}(T)]}{T} \leq \frac{1}{H^* \Delta_{k_T}^2 H(\underline{\nu}^{(k_T)})}. \quad (55)$$

As there is a finite number of arms, by the pigeonhole principle, there exists an increasing sequence  $(T_n)_{n \in \mathbb{N}}$  of budgets such that all arms  $(k_{T_n})_{n \in \mathbb{N}}$  are the same. Letting  $k$  denotes this arm, (55) gives

$$\limsup_{n \rightarrow +\infty} \frac{\mathbb{E}_{\nu^{(1)}}[N_k(T_n)]}{T_n} \leq \frac{1}{H^* \Delta_k^2 H(\underline{\nu}^{(k)})},$$

and applying a slightly modified version of Lemma 6 (where we only consider the sequence of budgets  $(T_n)_{n \in \mathbb{N}}$  with bandit problems  $\underline{\nu}^{(k)}$  and  $\underline{\nu}^{(1)}$  respectively, we obtain

$$\begin{aligned} \liminf_{n \rightarrow +\infty} \frac{1}{T_n} \ln \mathbb{P}_{\underline{\nu}^{(k)}}(I_{T_n} \neq k) &\geq - \limsup_{n \rightarrow +\infty} \frac{\mathbb{E}_{\nu^{(1)}}[N_k(T_n)]}{T_n} \times \text{KL}(\text{Ber}(1 - p_k), \text{Ber}(p_k)) \\ &\geq - \frac{\text{KL}(\text{Ber}(1 - p_k), \text{Ber}(p_k))}{H^* \Delta_k^2 H(\underline{\nu}^{(k)})} \\ &\geq - \frac{10}{H^*} \frac{1}{H(\underline{\nu}^{(k)})}, \end{aligned}$$

where, in the last inequality, we used that for all  $x \in [1/4, 1/2)$ ,

$$\text{KL}(\text{Ber}(1 - x), \text{Ber}(x)) \leq 10 \left( \frac{1}{2} - x \right)^2,$$

which can be checked analytically. This concludes the proof of (53).

**Step 2: control (54) of  $H^*$  for specific values of  $(p_k)_{2 \leq k \leq K}$ .** We proceed as [Carpentier and Locatelli \(2016\)](#). We set, for  $k \in \{2, \dots, K\}$ ,

$$p_k = \frac{1}{2} - \frac{1}{4} \frac{k}{K} \quad \text{or, equivalently,} \quad \Delta_k = \frac{1}{4} \frac{k}{K},$$

and show first that  $\Delta_k^2 H(\nu^{(k)}) \leq 2k$ . Indeed, by (51):

$$\Delta_k^2 H(\nu^{(k)}) = \Delta_k^2 \sum_{a \neq k} \frac{1}{(\Delta_k + \Delta_a)^2} = \sum_{a < k} \frac{\Delta_k^2}{(\Delta_k + \Delta_a)^2} + \sum_{a > k} \frac{\Delta_k^2}{(\Delta_k + \Delta_a)^2},$$

and, lower bounding  $\Delta_k + \Delta_a$  by  $\Delta_k$  in the first sum, and by  $\Delta_a$  in the second, we obtain

$$\Delta_k^2 H(\nu^{(k)}) \leq k - 1 + \sum_{a > k} \left( \frac{\Delta_k}{\Delta_a} \right)^2 = k - 1 + \sum_{a > k} \frac{k^2}{a^2} = k - 1 + k^2 \left( \frac{1}{k} - \frac{1}{K} \right) \leq 2k.$$

We finally get (54) by plugging this inequality into the definition (52) of  $H^*$ :

$$H^* = \sum_{k=2}^K \frac{1}{\Delta_k^2 H(\underline{\nu}^{(k)})} \geq \sum_{k=2}^K \frac{1}{2k} = \frac{1}{2} (\ln(K+1) - \ln 2) \geq \frac{3}{10} \ln(K),$$

the latter inequality being easily verified for  $K \geq 3$ . ■



## E.2. The Bretagnolle-Huber technique

An alternative (non-asymptotic) method to obtain lower bounds consists in using, together with the data-processing inequality and the chain rule for the Kullback-Leibler divergence, the Bretagnolle-Huber inequality (Bretagnolle and Huber, 1979), which states that, for  $p, q \in [0, 1]$ ,

$$p + 1 - q \geq \frac{1}{2} \exp\left(-\text{KL}(\text{Ber}(p), \text{Ber}(q))\right). \quad (56)$$

**Lower bound of Kaufman et al. (2016).** The method was used by Kaufman et al. (2016) on the model  $\mathcal{D}_{\sigma^2}$  of Gaussian distributions with a fixed variance  $\sigma^2 > 0$ . We state here an asymptotic version of their result, that will be generalized to all models, possibly non-parametric, in Proposition 24. The statement relies on a measure of complexity

$$C(\underline{\nu}) = \sum_{a \neq a^*(\underline{\nu})} \frac{2\sigma^2}{\Delta_a^2}.$$

It reads: for all strategies and for all bandit problems  $\underline{\nu}$  in  $\mathcal{D}_{\sigma^2}$  with a unique optimal arm, there exists a set of alternative bandit instances  $(\underline{\nu}^{(k)})_{k \neq a^*(\underline{\nu})}$  in  $\mathcal{D}_{\sigma^2}$  such that the best arm of  $\underline{\nu}^{(k)}$  is arm  $k$  and, denoting  $\underline{\nu}^{a^*(\underline{\nu})} = \underline{\nu}$ ,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \left( \max_{1 \leq k \leq K} \left\{ \mathbb{P}_{\underline{\nu}^{(k)}}(I_T \neq k) \right\} \right) \geq -4C(\underline{\nu})^{-1} \geq -\frac{2}{\sigma^2} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}, \quad (57)$$

$$\text{and } \forall k \neq a^*(\underline{\nu}), \quad C(\underline{\nu}^{(k)}) \leq C(\underline{\nu}).$$

(The adaptation to a bound for a sequence of strategies comes at a small cost: the original bound proposed by Kaufman et al. (2016) for a given budget  $T$  only involved two bandit problems,  $\underline{\nu}$  and a single alternative problem  $\underline{\lambda}$ .)

**Generalization to possibly non-parametric models.** The Bretagnolle-Huber methodology readily extends to general, possibly non-parametric, models. It leads to the following bound. However, the issue is the lack of interpretability of that bound, as we discuss after the statement of the proposition.

**Proposition 24** *Fix  $K \geq 2$ , a model  $\mathcal{D}$ , and a sequence of strategies. Let  $\underline{\nu}$  be a bandit problem in  $\mathcal{D}$  with a unique optimal arm. Consider, for each  $k \neq a^*(\underline{\nu})$ , a distribution  $\zeta_k \in \mathcal{D}$  such that  $E(\zeta_k) > \mu^*$ . Denoting by  $\underline{\nu}^{(k)}$ , for  $k \neq a^*(\underline{\nu})$ , the bandit problem obtained from  $\underline{\nu}$  by changing arm  $k$  to distribution  $\zeta_k$ , and by  $\underline{\nu}^{(a^*(\underline{\nu}))}$  the original bandit problem  $\underline{\nu}$ , we have*

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \left( \max_{1 \leq k \leq K} \left\{ \mathbb{P}_{\underline{\nu}^{(k)}}(I_T \neq k) \right\} \right) \geq - \left( \sum_{a \neq a^*(\underline{\nu})} \frac{1}{\text{KL}(\nu_a, \zeta_a)} \right)^{-1}.$$

Before proving this result, we provide a few comments concerning the (lack of) interpretability of this result and explain how to derive the bound (57) for Gaussian models.

*Lack of interpretability of the bound for general models.* To derive an interesting and interpretable bound from this result, one needs to choose carefully the distributions  $\zeta_k$ : there is a tradeoff between obtaining the best possible bound by choosing  $\zeta_k$  close to  $\nu_k$  in terms of Kullback-Leibler

divergences, and controlling the maximum of the misidentification probabilities. In particular, when  $E(\zeta_k)$  is close to  $\mu^*$ , the probability of misidentification under  $\underline{\nu}^{(k)}$  will be larger, so that

$$\max_{1 \leq k \leq K} \left\{ \mathbb{P}_{\underline{\nu}^{(k)}}(I_T \neq k) \right\} \gg \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})).$$

That is, the bound will become uninformative on the targeted quantity  $\mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu}))$ . This trade-off seems to be unsolvable in general, unless there exist some specific properties for the Kullback-Leibler divergence of the model, as we illustrate now.

*The case of Gaussian models  $\mathcal{D}_{\sigma^2}$ .* For the model  $\mathcal{D}_{\sigma^2}$ , we obtain (57) and the desired inequalities on complexities by considering  $\zeta_k = \mathcal{N}(\mu^* + \Delta_k, \sigma^2)$ , the Gaussian distribution of mean  $\mu^* + \Delta_k$  and variance  $\sigma^2$ . Indeed, we recall that

$$\forall \nu, \zeta \in \mathcal{D}_{\sigma^2}, \quad \text{KL}(\zeta, \nu) = \frac{(E(\nu) - E(\zeta))^2}{2\sigma^2} = \text{KL}(\nu, \zeta).$$

Note a key property (actually stronger than symmetry), which will turn out to be useful in the calculations below: the Kullback-Leibler divergence only depends on the expectation gaps between the two distributions. Now, with the distributions  $\zeta_k$  defined above, on the one hand, the bound of Proposition 24 rewrites

$$\sum_{a \neq a^*(\underline{\nu})} \frac{1}{\text{KL}(\nu_a, \zeta_a)} = \sum_{a \neq a^*(\underline{\nu})} \frac{2\sigma^2}{\left( \underbrace{E(\nu_a)}_{\mu^* - \Delta_a} - \underbrace{E(\zeta_a)}_{\mu^* + \Delta_a} \right)^2} = \frac{C(\underline{\nu})}{4},$$

and on the other hand, for  $k \neq a^*(\underline{\nu})$ , as the best arm of  $\underline{\nu}^{(k)}$  is  $k$ , with associated expectation  $\mu^* + \Delta_k$ ,

$$\begin{aligned} C(\underline{\nu}^{(k)}) &= \sum_{a \neq k} \frac{2\sigma^2}{(\mu^* + \Delta_k - \mu_a)^2} = \frac{2\sigma^2}{\Delta_k^2} + \sum_{a \notin \{k, a^*(\underline{\nu})\}} \frac{2\sigma^2}{(\mu^* + \Delta_k - \mu_a)^2} \\ &\leq \frac{2\sigma^2}{\Delta_k^2} + \sum_{a \notin \{k, a^*(\underline{\nu})\}} \frac{2\sigma^2}{(\mu^* - \mu_a)^2} = \sum_{a \neq a^*(\underline{\nu})} \frac{2\sigma^2}{\Delta_a^2} = C(\underline{\nu}). \end{aligned}$$

As underlined above, the calculations led are highly specific to the Gaussian model and exploit the gap-based rewriting of the Kullback-Leibler divergence. They would only extend to models for which gap-based rewritings of (or upper and lower bounds on) the Kullback-Leibler divergence would be available.

*Reverse order of the arguments in the KL.* We observe that the bound of Proposition 24 involves Kullback-Leibler divergences with arguments in reverse order compared to the lower bounds presented in Section 4. Indeed, taking the infimum of the lower bound over distributions  $\zeta_k$  such that  $E(\zeta_k) > \mu^*$  would lead to a complexity in terms of the  $\mathcal{K}_{\text{inf}}^>(\nu_k, \mu^*)$ , where

$$\mathcal{K}_{\text{inf}}^>(\nu, x) \stackrel{\text{def}}{=} \inf \{ \text{KL}(\nu, \zeta) : \zeta \in \mathcal{D} \text{ s.t. } E(\zeta) > x \},$$

rather than in terms of the  $\mathcal{L}_{\text{inf}}^>(\mu^*, \nu_k)$ . Given all bounds presented in this article, it does not seem that this would be the correct notion of complexity for the fixed-budget best-arm identification.

**Proof of Proposition (24).** We will prove that for all  $(u_b)_{b \neq a^*(\underline{\nu})} \in \Sigma_{K-1}$ , where  $\Sigma_{K-1}$  is the simplex of  $\mathbb{R}^{K-1}$ , we get

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \left( \max_{1 \leq k \leq K} \left\{ \mathbb{P}_{\underline{\nu}^{(k)}}(I_T \neq k) \right\} \right) \geq - \max_{b \neq a^*(\underline{\nu})} \{u_b \text{KL}(\nu_b, \zeta_b)\}. \quad (58)$$

The result will follow by choosing  $(u_b)_{b \neq a^*(\underline{\nu})}$  so as to maximize this lower bound, that is taking

$$u_b = \left( \sum_{a \neq a^*(\underline{\nu})} \frac{1}{\text{KL}(\nu_a, \zeta_a)} \right)^{-1} \times \frac{1}{\text{KL}(\nu_b, \zeta_b)}.$$

We now fix  $(u_b)_{b \neq a^*(\underline{\nu})} \in \Sigma_{K-1}$  and prove (58). Consider a given budget  $T$  and let  $b \neq a^*(\underline{\nu})$ . We get, as  $a^*(\underline{\nu}) \neq b$  first, and by the Bretagnolle-Huber inequality (56) then,

$$\begin{aligned} \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) + \mathbb{P}_{\underline{\nu}^{(b)}}(I_T \neq b) &\geq \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) + \mathbb{P}_{\underline{\nu}^{(b)}}(I_T = a^*(\underline{\nu})) \\ &\geq \frac{1}{2} \exp\left(-\text{KL}(\text{Ber}(p_T), \text{Ber}(q_T))\right), \end{aligned}$$

where  $p_T \stackrel{\text{def}}{=} \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu}))$  and  $q_T \stackrel{\text{def}}{=} \mathbb{P}_{\underline{\nu}^{(b)}}(I_T \neq a^*(\underline{\nu}))$ . Applying the combination (43) of the data-compressing inequality and the chain rule for Kullback-Leibler divergences leads to, as  $\underline{\nu}$  and  $\underline{\nu}^{(b)}$  only differ at arm  $b$ ,

$$\mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) + \mathbb{P}_{\underline{\nu}^{(b)}}(I_T \neq b) \geq \frac{1}{2} \exp\left(-\mathbb{E}_{\underline{\nu}}[N_b(T)] \text{KL}(\nu_b, \zeta_b)\right).$$

As  $\max(u, v) \geq (u + v)/2$ , we obtained so far

$$\begin{aligned} \max_{1 \leq k \leq K} \left\{ \mathbb{P}_{\underline{\nu}^{(k)}}(I_T \neq k) \right\} &\geq \max\left\{ \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})), \mathbb{P}_{\underline{\nu}^{(b)}}(I_T \neq b) \right\} \\ &\geq \frac{1}{4} \exp\left(-\mathbb{E}_{\underline{\nu}}[N_b(T)] \text{KL}(\nu_b, \zeta_b)\right). \quad (59) \end{aligned}$$

This bound holds for any  $b \neq a^*(\underline{\nu})$ . In particular, as

$$\sum_{b \neq a^*(\underline{\nu})} \frac{\mathbb{E}_{\underline{\nu}}[N_b(T)]}{T} \leq 1 = \sum_{b \neq a^*(\underline{\nu})} u_b,$$

we know that there exists  $b^* \neq a^*(\underline{\nu})$  such that  $\mathbb{E}_{\underline{\nu}}[N_{b^*}(T)] \leq T u_{b^*}$ , and, applying (59) with  $b^*$ ,

$$\max_{1 \leq k \leq K} \left\{ \mathbb{P}_{\underline{\nu}^{(k)}}(I_T \neq k) \right\} \geq \frac{1}{4} \exp\left(-T u_{b^*} \text{KL}(\nu_{b^*}, \zeta_{b^*})\right) \geq \frac{1}{4} \exp\left(-T \max_{b \neq a^*(\underline{\nu})} \{u_b \text{KL}(\nu_b, \zeta_b)\}\right),$$

or, to put it differently,

$$\frac{1}{T} \ln \left( \max_{1 \leq k \leq K} \left\{ \mathbb{P}_{\underline{\nu}^{(k)}}(I_T \neq k) \right\} \right) \geq \frac{1}{T} \ln \frac{1}{4} - \max_{b \neq a^*(\underline{\nu})} \{u_b \text{KL}(\nu_b, \zeta_b)\}.$$

We obtain (58) by taking the lim inf in  $T$  on that inequality. ■