

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

Antoine Barrier^{1,2} · Aurélien Garivier^{1,3} · Gilles Stoltz²

1. UMPA, ÉNS de Lyon

2. LMO, Université Paris-Saclay

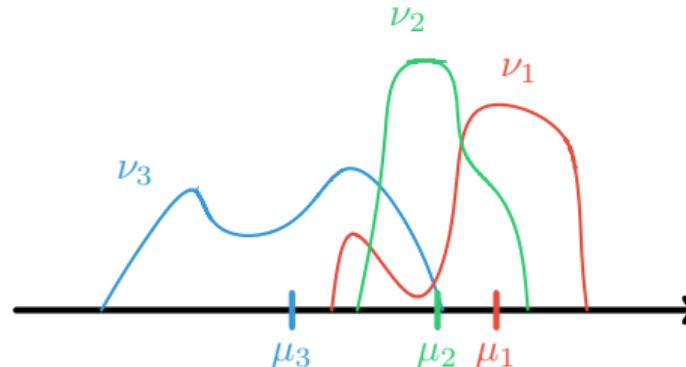
3. LIP, ÉNS de Lyon



34th International Conference on Algorithmic Learning Theory
February 21, 2023 · Singapore

Bandit algorithms

- Model \mathcal{D} : set of distributions
(e.g. Gaussian or Bernoulli distributions, exponential model, $\mathcal{P}[0, 1], \dots$)
- K arms: K unknown distributions $\underline{\nu} = (\nu_1, \dots, \nu_K)$ of \mathcal{D} , $E(\nu_a) = \mu_a$



- At round t :
 - choose an arm A_t depending on previous observations
 - observe a new independent reward $Y_t \sim \nu_{A_t}$

$$a^*(\underline{\nu}) = \operatorname{argmax}_{1 \leq a \leq K} \mu_a$$

$$\nu^* = \nu_{a^*}(\underline{\nu})$$

$$\mu^* = \mu_{a^*}(\underline{\nu})$$

Regret minimization: optimal non-parametric approach

$$\text{Minimize } R_T(\underline{\nu}) = \sum_{t=1}^T \mu^* - \mathbb{E}_{\underline{\nu}}[\mu_{A_t}] = \sum_{\substack{1 \leq a \leq K \\ a \neq a^*(\underline{\nu})}} \mathbb{E}_{\underline{\nu}}[N_a(T)] \Delta_a$$

with $N_a(T)$: # of pulls of arm a at time T ; $\Delta_a = \mu^* - \mu_a$: gap of arm a

→ Non-parametric inequalities based on \mathcal{K}_{\inf} : for $\nu \in \mathcal{D}$ and $x \in \mathbb{R}$

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

Lower bound [Lai and Robbins, 1985, Burnetas and Katehakis, 1996]

$$\forall a \neq a^*(\underline{\nu}), \quad \liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{\ln T} \geq \frac{1}{\mathcal{K}_{\inf}^>(\nu_a, \mu^*)}$$

Upper bounds Gap-based (e.g. UCB) and then $\mathcal{K}_{\inf}^>$ -based (e.g. KL-UCB-switch) guarantees

Best-Arm Identification (BAI) with a fixed confidence

Return an estimate \hat{a}_τ of $a^*(\underline{\nu})$ after some stopping time τ

Minimize $\mathbb{E}_{\underline{\nu}}[\tau]$ such that $\mathbb{P}_{\underline{\nu}}(\hat{a}_\tau \neq a^*(\underline{\nu})) \leq \delta$

Lower bound [Garivier and Kaufmann, 2016, Jourdan et al., 2022]

$$\mathbb{E}_{\underline{\nu}}[\tau] \geq T(\underline{\nu}) \ln(1/\delta)$$

$$T(\underline{\nu})^{-1} = \sup_{w \in \Sigma_K} \min_{a \neq a^*(\underline{\nu})} \inf_{x \in [\mu_a, \mu^*]} w_{a^*(\underline{\nu})} \mathcal{K}_{\inf}^<(\nu^*, x) + w_a \mathcal{K}_{\inf}^>(\nu_a, x)$$

Upper bound Asymptotically ($\delta \rightarrow 0$)

- optimal strategies for exponential models
- almost optimal strategies for non-parametric models

→ Non-parametric theory based on \mathcal{K}_{\inf}

Best-Arm Identification (BAI) with a fixed budget

Gap-based approaches [Audibert et al., 2010]

Minimize $\mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu}))$ given budget T

Order arms by their means: $\mu_{(1)} > \mu_{(2)} > \mu_{(3)} > \dots > \mu_{(K)}$

Lower bound For the model $\mathcal{D} = \mathcal{B}_{[p, 1-p]} = \{\text{Ber}(x) : x \in [p, 1-p]\}$:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq -\frac{5}{p(1-p)} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}$$

Upper bound For the model $\mathcal{D} = \mathcal{P}[0, 1]$ (or subGaussian): the successive-rejects strategy satisfies:

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

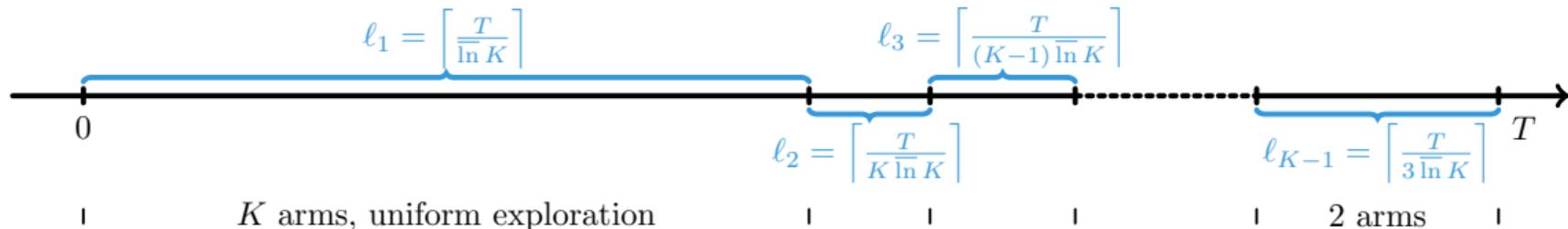
See also [Karnin et al., 2013, Kaufmann et al., 2016, Carpentier and Locatelli, 2016]

- Gap-based theory for specific models
- What about non-parametric bounds?

The successive-rejects strategy

- Exploration is split in $K - 1$ phases of various lengths $\ell_1, \dots, \ell_{K-1}$.
- At each phase, all remaining arms are pulled uniformly and the arm with worst empirical mean is removed.

Phase lengths of [Audibert et al., 2010]:



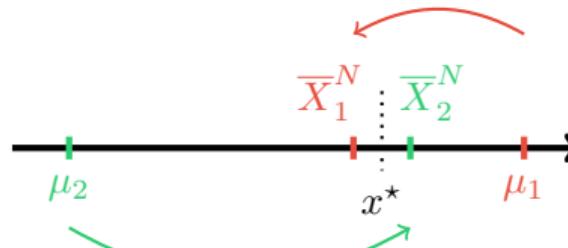
$$\limsup_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \leq -\frac{1}{\ln K} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}$$

→ Proof based on Hoeffding's inequality
to bound errors when comparing empirical means

Non-parametric upper bound

$$\mathcal{D} = \mathcal{P}[0, 1]$$

$$\Delta = \mu_1 - \mu_2 > 0$$



Hoeffding's inequality: $\frac{1}{N} \log \mathbb{P}\left(\overline{X}_2^N \geq \overline{X}_1^N\right) \leq -\Delta^2$

Cramér-Chernoff bound separately for \overline{X}_1^N and \overline{X}_2^N with the **best cut-off** x^* :

$$\limsup_{N \rightarrow +\infty} \frac{1}{N} \log \mathbb{P}\left(\overline{X}_2^N \geq \overline{X}_1^N\right) \leq -\left(\mathcal{L}_{\inf}^>(x^*, \nu_2) + \mathcal{L}_{\inf}^<(x^*, \nu_1)\right) \stackrel{\text{def}}{=} -\mathcal{L}(\nu_2, \nu_1)$$

where, for $\nu \in \mathcal{D}$ and $x \in \mathbb{R}$:

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

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

Non-parametric upper bound

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

where arms are ordered such that

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

Pinsker's inequality: $\mathcal{L}(\nu_2, \nu_1) \geq (\mu_1 - \mu_2)^2$

→ Generalization and improvement of the bound of [Audibert et al., 2010]

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

Non-parametric lower bounds

Recipe For any bandit problem $\underline{\lambda}$ such that $a^*(\underline{\lambda}) \neq a^*(\underline{\nu})$:

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

Difficulty Control of $\frac{\mathbb{E}_{\underline{\lambda}}[N_a(T)]}{T}$ when $\lambda_a \neq \nu_a$.

Gap-based bound for a Bernoulli model [Audibert et al., 2010]

→ Careful construction of $\underline{\lambda}$ with consistency assumption

Adaptation to non-parametric models?

→ Natural hypothesis on the strategies

Weak monotonicity

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

Strong monotonicity

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

for all $a \in [K]$

→ satisfied by the successive-rejects strategy

Non-parametric lower bounds

For a general model \mathcal{D} , under some assumptions including **weak** monotonicity

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

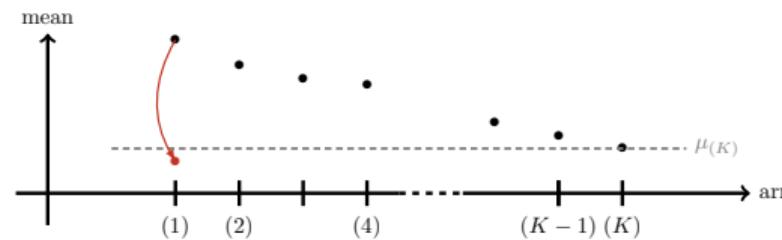
Sketch of proof (case $k = K$)

- ① Weak monotonicity on $\underline{\lambda}$ only differing from $\underline{\nu}$ at arm (1) with $E(\underline{\lambda}_{(1)}) < \mu_{(K)}$

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

- ② Take infimum on $E(\underline{\lambda}_{(1)}) < \mu_{(K)}$

$$\inf_{\lambda_{(1)} \in \mathcal{D} \text{ s.t. } E(\lambda_{(1)}) < \mu_{(K)}} \text{KL}(\lambda_{(1)}, \nu^*) = \mathcal{L}_{\text{inf}}^<(\mu_{(K)}, \nu^*)$$



Non-parametric lower bounds

→ Generalization and improvement of the bound of [Audibert et al., 2010]

For the model $\mathcal{D} = \mathcal{B}_{[p, 1-p]} = \{\text{Ber}(x) : x \in [p, 1-p]\}$:

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(\hat{a}_T \neq a^*(\underline{\nu})) \geq -\frac{5}{p(1-p)} \min_{2 \leq k \leq K} \frac{\Delta_{(k)}^2}{k}$$

Improvement towards cutoff x^*

For a general model \mathcal{D} and a consistent and **strongly monotonous** strategy

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

