

# Multi-armed bandit problems: a statistical view, focused on lower bounds

Gilles Stoltz

Laboratoire de mathématiques d'Orsay



Joint works with **Hédi Hadiji**, now at CentraleSupélec,  
and with ENS Lyon colleagues: **Aurélien Garviev**, **Pierre Ménard**, **Antoine Barrier**

# $K$ -armed stochastic bandits

Framework, possible objectives, index strategies

$K$  probability distributions  $\nu_1, \dots, \nu_K$   
with expectations  $\mu_1, \dots, \mu_K$

$$\longrightarrow \mu^* = \max_{a \in [K]} \mu_a$$

At each round  $t = 1, 2, \dots,$

1. Statistician picks **arm**  $A_t \in [K]$
2. She gets a reward  $Y_t$  drawn according to  $\nu_{A_t}$
3. This is the **only feedback** she receives

→ **Exploration–exploitation dilemma**  
estimate the  $\nu_a$  **vs.** get high rewards  $Y_t$

**Link with UQ?** Emmanuel Vazquez told me:

Conceptually                      arms  $\leftrightarrow$  parameters of numerical experiments

Technically                        leverage bandit techniques to study EI strategy

Setting: at round each round  $t \geq 1$ , pick arm  $A_t \in [K]$ , get and observe  $Y_t \sim \nu_{A_t}$

Objective #1: Maximize cumulative rewards  $\leftrightarrow$   
Minimize **pseudo-regret**

$$\begin{aligned} R_T &= \sum_{t=1}^T (\mu^* - \mathbb{E}[Y_t]) = \sum_{t=1}^T (\mu^* - \mathbb{E}[\mu_{A_t}]) \\ &= \sum_{a \in [K]} \left( (\mu^* - \mu_a) \mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}} \right] \right) = \sum_{a \in [K]} (\mu^* - \mu_a) \mathbb{E}[N_a(T)] \end{aligned}$$

$\leftrightarrow$  Control the  $\mathbb{E}[N_a(T)]$

Objective #2: **Identify best arm**  $\leftrightarrow$  Minimize  $\mathbb{P} \left( I_T \notin \arg \max_{a \in [K]} \mu_a \right)$

Model:  $\nu_1, \dots, \nu_K$  are distributions over  $[0, 1]$

A classical strategy: **UCB** [upper confidence bound]

Auer, Cesa-Bianchi and Fisher [2002]

For  $t \geq K$ , pick  $A_{t+1} \in \arg \max_{a \in [K]} \left\{ \hat{\mu}_a(t) + \sqrt{\frac{2 \ln t}{N_a(t)}} \right\}$

**Exploitation**: cf. empirical mean  $\hat{\mu}_a(t)$

**Exploration**: cf.  $\sqrt{2 \ln t / N_a(t)}$  favors arms  $a$  not pulled often

**Suboptimal** regret bounds of two types

– **Distribution-dependent** bound:  $R_T \lesssim \sum_{a: \mu_a < \mu^*} \frac{8 \ln T}{\mu^* - \mu_a}$

– **Distribution-free** bound:  $\sup_{\nu_1, \dots, \nu_K} R_T \lesssim \sqrt{8KT \ln T}$

Model:  $\nu_1, \dots, \nu_K$  are distributions over  $[0, 1]$

Another **index-based** strategy:

**MOSS** [minimax optimal strategy in the stochastic setting]

Audibert and Bubeck [2009]

For  $t \geq K$ , pick  $A_{t+1} \in \arg \max_{a \in [K]} \left\{ \hat{\mu}_a(t) + \sqrt{\frac{1}{N_a(t)} \ln_+ \frac{T}{KN_a(t)}} \right\}$

$\ln_+ = \max\{\ln, 0\}$ ; there exist anytime versions

Distribution-free regret bounds  $\sup_{\nu_1, \dots, \nu_K} R_T$  of optimal order  $\sqrt{KT}$

– **Upper** bound:  $49\sqrt{KT}$  for MOSS

– **Lower** bound:  $(1/20)\sqrt{KT}$

Auer, Cesa-Bianchi, Freund and Schapire [2002]

Model:  $\nu_1, \dots, \nu_K$  are distributions over  $[0, 1]$

## KL-UCB strategy

Honda and Takemura [2015]; Cappé, Garivier, Maillard, Munos, Stoltz [2013];  
Garivier, Hadiji, Ménard, Stoltz [2022]

Key quantity  $\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) = \inf \{ \text{KL}(\nu_a, \nu'_a) : E(\nu'_a) > \mu^* \}$

Indices  $U_a(t) = \sup \left\{ \mu \in [0, 1] : \mathcal{K}_{\text{inf}}(\hat{\nu}_a(t), \mu) \leq \frac{\varphi(t, N_a(t))}{N_a(t)} \right\}$

Typically,  $\varphi(t, N_a(t))$  of order  $\ln t$ ; for  $t \geq K$ , pick  $A_{t+1} \in \arg \max_{a \in [K]} U_a(t)$

Optimal distribution-dependent regret bounds:

$$\sum_{a: \mu_a < \mu^*} \frac{\mu^* - \mu_a}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} \ln T - \Theta(\ln \ln T)$$

For lower bounds: Lai and Robbins [1985]; Burnetas and Katehakis [1996]; Garivier, Ménard and Stoltz [2019]

Model:  $\nu_1, \dots, \nu_K$  are distributions over  $[0, 1]$

## KL-UCB-Switch strategy

Garivier, Hadiji, Ménard, Stoltz [2022]

Index strategy of the form: for each arm  $a \in [K]$ , use

- KL-UCB index if  $N_a(t) \leq (t/K)^5$
- MOSS index if  $N_a(t) \geq (t/K)^5$

Optimal bounds of the two types:

- **Distribution-dependent** bound, with the  $-\Theta(\ln \ln T)$  term
- **Distribution-free** bound: 
$$\sup_{\nu_1, \dots, \nu_K} R_T \lesssim K + 23\sqrt{KT}$$



Model:  $\nu_1, \dots, \nu_K$  are distributions over  $[0, 1]$

**Summary** / Reviewed index strategies **all** of the form:

For  $t \geq K$ , pick  $A_{t+1} \in \arg \max_{a \in [K]} \left\{ \hat{\mu}_a(t) + \text{expl}(t, N_a(t)) \right\}$

Possibly with fancy, or null, **exploration** bonuses  $\text{expl}(t, N_a(t))$

**Exploitation**: cf. empirical mean  $\hat{\mu}_a(t)$

Various bounds achieved, depending on how  $\text{expl}(t, N_a(t))$  is set

– Optimal **distribution-dependent** bounds:

$$\sum_{a: \mu_a < \mu^*} \frac{\mu^* - \mu_a}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} \ln T - \Theta(\ln \ln T)$$

– Optimal **distribution-free** bounds:  $\sup_{\nu_1, \dots, \nu_K} R_T = \Theta(\sqrt{KT})$

**Proofs** for upper bounds: control  $\mathbb{E}[N_a(T)]$

# Proofs of the regret lower bounds on $[0, 1]$

(At least, high-level ideas...)

## Proof ideas for the lower bounds

Strategy  $\psi$ : maps  $H_t = (Y_1, \dots, Y_t) \mapsto A_{t+1} = \psi_t(H_t)$

Change of measure: compare distributions of  $H_T$   
under  $\underline{\nu} = (\nu_1, \dots, \nu_K)$  vs.  $\underline{\nu}' = (\nu'_1, \dots, \nu'_K)$

**Fundamental inequality:** performs an **implicit** change of measure

Reference: Lai and Robbins [1985], Auer et al. [2002], Garivier et al. [2019]

For all  $Z$  taking values in  $[0, 1]$  and  $\sigma(H_T)$ -measurable

$$\text{(chain rule)} \quad \sum_{a \in [K]} \mathbb{E}_{\underline{\nu}}[N_a(T)] \text{KL}(\nu_a, \nu'_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{H_T}, \mathbb{P}_{\underline{\nu}'}^{H_T})$$

$$\text{(data-proc. ineq.)} \quad \geq \text{kl}(\mathbb{E}_{\underline{\nu}}[Z], \mathbb{E}_{\underline{\nu}'}[Z])$$

where  $\text{kl}(p, q) = \text{KL}(\text{Ber}(p), \text{Ber}(q))$

**Later use:**  $\underline{\nu}'$  only differs from  $\underline{\nu}$  at some  $a$ , with  $Z = N_a(T)/T$

Distribution-free lower bound, for distributions over  $[0, 1]$ 

Problem  $\underline{\nu}_0 = (\text{Ber}(1/2))_{a \in [K]}$  vs.  $\underline{\nu}_k = (\text{Ber}(1/2 + \varepsilon \mathbb{I}_{\{a=k\}}))_{a \in [K]}$

$$R_T \stackrel{\text{def}}{=} \sum_{a \neq k} \varepsilon \mathbb{E}_{\underline{\nu}_k} [N_a(T)] = T\varepsilon \left( 1 - \mathbb{E}_{\underline{\nu}_k} [N_k(T)/T] \right)$$

Thus,  $\sup_{\underline{\nu}} R_T \geq \sup_{\varepsilon \in (0,1)} \max_{k \in [K]} T\varepsilon \left( 1 - \mathbb{E}_{\underline{\nu}_k} [N_k(T)/T] \right)$

Fundamental inequality,

with  $Z = N_k(T)/T$

+ Pinsker's inequality

and  $k \in [K]$  such that  $\mathbb{E}_{\underline{\nu}_0} [N_k(T)/T] \leq 1/K$

$$\underbrace{\mathbb{E}_{\underline{\nu}_0} [N_k(T)]}_{\leq T/K} \underbrace{\text{KL}(\text{Ber}(1/2), \text{Ber}(1/2 + \varepsilon))}_{= -\ln(1-4\varepsilon^2)/2 \leq 2.5\varepsilon^2}$$

$$\geq \text{kl}(\mathbb{E}_{\underline{\nu}_0} [Z], \mathbb{E}_{\underline{\nu}_k} [Z]) \geq 2 \left( \mathbb{E}_{\underline{\nu}_k} [N_0(T)/T] - \mathbb{E}_{\underline{\nu}_k} [N_k(T)/T] \right)^2$$

Thus,  $\sup_{\underline{\nu}} R_T \geq \sup_{\varepsilon \in (0,1/4)} T\varepsilon \left( 1 - 1/K - \varepsilon \sqrt{1.25 T/K} \right) \geq \Theta(\sqrt{KT})$

Distribution-dependent bound: 
$$R_T = \sum_{a \in [K]} (\mu^* - \mu_a) \mathbb{E}_{\underline{\nu}}[N_a(T)]$$

We lower bound each  $\mathbb{E}_{\underline{\nu}}[N_a(T)]$  for a fixed  $a$  with  $\mu_a < \mu^*$ ; let  $\nu'_a$  with  $\mu_a > \mu^*$

Problems  $\underline{\nu} = (\nu_a)_{a \in [K]}$  vs.  $\underline{\nu}' = (\nu_1, \dots, \nu_{a-1}, \nu'_a, \nu_{a+1}, \dots, \nu_K)$

Fundamental inequality

on “good” strategies

& lower bound on kl

$\forall \alpha \in (0, 1], \mathbb{E}[N_k(T)] = o(T^\alpha)$  for subopt.  $k$

$\text{kl}(p, q) \geq (1-p) \ln(1/(1-q)) - \ln 2$

$$\begin{aligned} \mathbb{E}_{\underline{\nu}}[N_a(T)] \text{KL}(\nu_a, \nu'_a) &\geq \text{kl} \left( \overbrace{\mathbb{E}_{\underline{\nu}}[N_a(T)/T], \mathbb{E}_{\underline{\nu}'}[N_a(T)/T]}^{=o(1)} \right) \\ &\gtrsim \ln \left( 1 / (1 - \mathbb{E}_{\underline{\nu}'}[N_a(T)/T]) \right) \end{aligned}$$

Since  $\mathbb{E}_{\underline{\nu}'}[N_a(T)/T] = 1 - \sum_{k \neq a} \mathbb{E}_{\underline{\nu}'}[N_k(T)/T] \gtrsim 1 - T^{\alpha-1}$ , we get:

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{KL}(\nu_a, \nu'_a) \gtrsim \ln T^{1-\alpha}$$

**Distribution-dependent bound:** 
$$R_T = \sum_{a \in [K]} (\mu^* - \mu_a) \mathbb{E}_{\underline{\nu}}[N_a(T)]$$

We lower bound each  $\mathbb{E}_{\underline{\nu}}[N_a(T)]$  for a fixed  $a$  with  $\mu_a < \mu^*$ ; let  $\nu'_a$  with  $\mu_a > \mu^*$

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{KL}(\nu_a, \nu'_a) \gtrsim \ln T^{1-\alpha}, \quad \text{that is,} \quad \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{KL}(\nu_a, \nu'_a)}{\ln T} \gtrsim 1 - \alpha \rightarrow 1$$

Therefore, “good” strategies can ensure, at best:

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{\ln T} \geq \sup_{\nu'_a: \mu'_a > \mu^*} \frac{1}{\text{KL}(\nu_a, \nu'_a)} \stackrel{\text{def}}{=} \frac{1}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)}$$

By summing over suboptimal arms:

$$\liminf_{T \rightarrow \infty} \frac{R_T}{\ln T} \geq \sum_{a \in [K]} \frac{\mu^* - \mu_a}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)}$$

Note: **general proof**, valid for any model  $\mathcal{D}$

# Adaptation to the range

## Bounded but unknown range

Reference for this part of the talk: Hadiji and Stoltz [2020]

That is, model: 
$$\mathcal{D} = \bigcup_{m, M: m < M} \mathcal{D}_{m, M}$$

where  $\mathcal{D}_{m, M}$  set of distributions  $\nu$  over a given interval  $[m, M]$   
Before, we were only dealing with  $\mathcal{D}_{0, 1}$

What changes?

Same distribution-free lower bound:

$\Theta((M - m)\sqrt{KT})$  by rescaling

No worsening due to ignorance of the range

Different distribution-dependent lower bound:

$R_T / \ln T \rightarrow +\infty$  as  $\mathcal{K}_{\text{inf}}(\nu_a, \mu^*, \mathcal{D}) = 0$

But any rate  $\gg \ln T$  may be achieved



## Focus on the UCB strategy

With a known range  $[m, M]$ , reads (knowledge of the **range is key!**)

$$A_{t+1} \in \arg \max_{a \in [K]} \left\{ \hat{\mu}_a(t) + (M - m) \sqrt{\frac{2 \ln t}{N_a(t)}} \right\}$$

Extension to an unknown range:

$$A_{t+1} \in \arg \max_{a \in [K]} \left\{ \hat{\mu}_a(t) + \sqrt{\frac{\varphi(t)}{N_a(t)}} \right\}$$

where  $\ln t \ll \varphi(t) \ll t$ ; eventually,  $\sqrt{\varphi(t)} \geq (M - m) \sqrt{2 \ln t}$

Guarantee: for all bandit problems  $\nu_1, \dots, \nu_K$  in  $\mathcal{D}$ ,

$$\limsup \frac{R_T}{\varphi(T)} < +\infty$$

$\Phi_{\text{dep}} = \varphi$  is the corresponding **distribution-dependent rate for adaptation** to the range

## Distribution-free rate for adaptation to the range

$\Phi_{\text{free}} : \mathbb{N} \rightarrow (0, +\infty)$  such that

$$\forall m < M,$$

$$\forall \nu_1, \dots, \nu_K \text{ in } \mathcal{D}_{m,M},$$

$$\forall T \geq 1,$$

$$R_T \leq (M - m)\Phi_{\text{free}}(T)$$

By the lower bound proved for  $[m, M] = [0, 1]$ :

$$\Phi_{\text{free}}(T) \geq \Theta(\sqrt{KT})$$

AdaHedge on estimated payoffs + mixing achieves

$$\Phi_{\text{free}}(T) \approx 7(M - m)\sqrt{TK \ln K}$$

Reference for AdaHedge: Cesa-Bianchi, Mansour, Stoltz [2005, 2007] and De Rooij, van Erven, Grünwald, Koolen [2014]

Note:  $\sqrt{\ln K}$  shaved off (with different strategy) when  $M$  is known

## AdaHedge on estimated payoffs + mixing

Randomized strategy:  $A_t \sim \mathbf{p}_t$  for  $t \geq K + 1$

Unbiased estimated payoffs:  $\hat{X}_{t,a} = \frac{Y_t - C}{p_{t,a}} + C$

where  $C$  is the average of the payoffs in the first  $K$  rounds

AdaHedge:  $q_{t+1,a} \propto \exp\left(-\eta_t \sum_{s=K+1}^t \hat{X}_{s,a}\right)$

Mixing:  $\mathbf{p}_t = (1 - \gamma_t)\mathbf{q}_t + \gamma_t \mathbf{1}/K$

Strategy actually built for adversarial payoffs  
(= arbitrary sequences)

## What about simultaneous bounds?

Reminder for known range  $[0, 1]$ :  $\ln T$  and  $\sqrt{T}$  rates for regret upper bounds

Theorem: If  $\Phi_{\text{free}}(T) \ll T$  then  $\Phi_{\text{dep}}(T) \times \Phi_{\text{free}}(T) \geq \Theta(T)$

Example:  $\Phi_{\text{free}}(T) = \Theta(\sqrt{T})$  now forces  $\Phi_{\text{dep}}(T) \geq \Theta(\sqrt{T})$

→ We finally exhibit some heavy **price for adaptation!**

Proof: by the fundamental inequality  
+ lack of upper end on payoffs in  $\mathcal{D}$

AdaHedge on estimated payoffs + mixing simultaneously achieves

$$\Phi_{\text{free}}(T) = \Theta(\sqrt{T}) \quad \text{and} \quad \Phi_{\text{dep}}(T) = \Theta(\sqrt{T})$$

Analysis heavily based on Seldin and Lugosi [2017]

Actually, all pairs  $\Phi_{\text{free}}(T) = \Theta(T^\alpha)$  and  $\Phi_{\text{dep}}(T) = \Theta(T^{1-\alpha})$   
with  $\alpha \in [1/2, 1)$  may be achieved, by setting the mixing factor properly

FYI—Proof of “If  $\Phi_{\text{free}}(T) \ll T$  then  $\Phi_{\text{dep}}(T) \times \Phi_{\text{free}}(T) \geq \Theta(T)$ ”Based on fundamental inequality + lack of upper end on payoffs in  $\mathcal{D}$ We lower bound each  $\mathbb{E}_{\underline{\nu}}[N_a(T)]$  for a fixed  $a$  with  $\mu_a < \mu^*$ 

Problems  $\underline{\nu}, \underline{\nu}'$  only differing at  $\nu'_a = (1 - \varepsilon)\nu_a + \varepsilon \delta_{\mu_a + c/\varepsilon}$   
 such that  $\nu_a \perp \delta_{\mu_a + c/\varepsilon}$  and  $\mu'_a > \mu^*$

$f = \frac{d\nu_a}{d\nu'_a} = \frac{1}{1 - \varepsilon}$  so that  $\text{KL}(\nu_a, \nu'_a) = \mathbb{E}_{\nu_a}[\ln f] \approx \varepsilon$

Fundamental inequality and  $\text{kl}(p, q) \gtrsim (1 - p) \ln(1/(1 - q))$ 

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \overbrace{\text{KL}(\nu_a, \nu'_a)}^{\approx \varepsilon} \geq \text{kl} \left( \overbrace{\mathbb{E}_{\underline{\nu}}[N_a(T)/T]}^{=o(1)}, \mathbb{E}_{\underline{\nu}'}[N_a(T)/T] \right)$$

$$\gtrsim \ln \left( 1 / (1 - \mathbb{E}_{\underline{\nu}'}[N_a(T)/T]) \right)$$

Indeed:  $(\mu^* - \mu_a) \mathbb{E}_{\underline{\nu}}[N_a(T)] \leq R_T(\underline{\nu}) \leq (M - m) \Phi_{\text{free}}(T) \ll T$ Similarly:  $\ln \left( 1 / (1 - \mathbb{E}_{\underline{\nu}'}[N_a(T)/T]) \right) \gtrsim \ln(c' \Phi_{\text{free}}(T) / (T\varepsilon))$ As:  $(\mu'_a - \mu^*) (T - \mathbb{E}_{\underline{\nu}'}[N_a(T)]) \leq R_T(\underline{\nu}') \leq (M + c/\varepsilon - m) \Phi_{\text{free}}(T)$ Picking  $\varepsilon \sim \Phi_{\text{free}}(T)/T$ :  $(\Phi_{\text{free}}(T)/T) \mathbb{E}_{\underline{\nu}}[N_a(T)] \gtrsim \text{cst}$

## Best-arm identification

With fixed budget and for possibly non-parametric models

## Objective #2: BAI with fixed budget $T$

Reference for this final part of the talk: Barrier, Garivier, and Stoltz [2022]

Bandit problem  $\underline{\nu} = (\nu_1, \dots, \nu_K)$  with unique optimal arm  $a^*(\underline{\nu})$

where optimality is in expectation:  $\mu_{a^*} = \max_{a \in [K]} \mu_a$

$T$  rounds, where arms  $A_t$  are pulled, rewards  $Y_t$  are obtained

Then: issue a **recommendation**  $I_T \in [K]$

Goal: upper and lower bound  $\mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu}))$

Note—BAI with fixed confidence  $\delta$  well understood

Track and Stop strategies, by Aurélien Garivier, Emilie Kaufmann, and co-authors

Typical strategy: **Successive rejects**

By Audibert and Bubeck [2010], with analysis based on **Hoeffding's** inequality

$K - 1$  regimes, and in each regime  $r = 1, \dots, K - 1$ :

- Denote by  $S_r$  the set of arms not dropped so far;  $S_1 = [K]$
- Play each arm in  $S_r$  an equal number of times
- **Drop arm** with smallest average payoff since the beginning (not smallest average payoff in regime  $r$ )

By carefully setting regimes (based on  $T$  and  $K$ ): when  $\mathcal{D} \subseteq \mathcal{P}_{0,1}$ ,

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

where  $\mu_{a^*} = \mu_{(1)} > \mu_{(2)} \geq \dots \geq \mu_{(K)}$

→ Called a **gap-based bound**



## Lower bounds?

Gap-based approach by Audibert and Bubeck [2010]

Studied  $\mathcal{D}_p = \{\text{Ber}(x) : x \in [p, 1 - p]\}$  for  $p > 0$

Methodology actually extends to models  $\mathcal{D}$  such that

$$\forall \nu, \nu' \text{ in } \mathcal{D}, \quad \text{KL}(\nu, \nu') \leq C_{\mathcal{D}} (\mathbb{E}(\nu) - \mathbb{E}(\nu'))^2$$

For instance,  $C_{\mathcal{D}_p} = 1/(2p(1-p))$

Careful and explicit analysis leading to: for all strategies,

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

Difference:  $-5 C_{\mathcal{D}}$  vs.  $-1/\bar{\ln}K$  in front of the min

## New non-parametric approach

Key quantities: note the reverse order in the KL compared to  $\mathcal{K}_{\text{inf}}$

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

and

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

Analysis of **Successive rejects** based on Cramér-Chernoff bounds:

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

where  $\mathcal{L}(\nu_{(k)}, \nu_{a^*}) = \inf_{x \in [\mu_{(k)}, \mu_{a^*}]} \left\{ \mathcal{L}_{\text{inf}}^{\geq}(x, \nu_{(k)}) + \mathcal{L}_{\text{inf}}^{\leq}(x, \nu_{a^*}) \right\}$

By Pinsker's inequality: yields the gap-based upper bound for  $\mathcal{D} = \mathcal{P}_{0,1}$

Special case  $x = \mu_{(k)}$  for the lower bounds

## New non-parametric approach: simple lower bound

Alternative problem  $\underline{\nu}'$  differing from  $\underline{\nu}$  only at  $a^*(\underline{\nu})$ ,  
with distribution  $\zeta$  s.t.  $\mathbb{E}(\zeta) < \mu(K)$

$$q_T \stackrel{\text{def}}{=} \mathbb{P}_{\underline{\nu}'}(I_T \neq a^*(\underline{\nu})) \geq \mathbb{P}_{\underline{\nu}'}(I_T = a^*(\underline{\nu}')) \xrightarrow{T \rightarrow +\infty} 1,$$

$$\text{while } p_T \stackrel{\text{def}}{=} \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \xrightarrow{T \rightarrow +\infty} 0$$

Fundamental inequality:

$$-\frac{1}{T} \ln p_T \sim \frac{\text{KL}(\text{Ber}(p_T), \text{Ber}(q_T))}{T} \leq \overbrace{\frac{\mathbb{E}_{\underline{\nu}'}[N_{a^*(\underline{\nu})}]}{T}}^{\leq 1/K} \text{KL}(\zeta, \nu_{a^*})$$

Hence,

$$\liminf_{T \rightarrow +\infty} \frac{1}{T} \ln \mathbb{P}_{\underline{\nu}}(I_T \neq a^*(\underline{\nu})) \geq -\frac{1}{K} \overbrace{\inf_{\mathbb{E}(\zeta) < \mu(K)} \text{KL}(\zeta, \nu_{a^*})}^{\mathcal{L}_{\text{inf}}^<(\mu(K), \nu_{a^*})}$$

Pruning argument to get the min over  $k \in \{2, \dots, K\}$