

# On the notion of Complexity in the Stochastic Multi-armed Bandit Problems

Pierre Ménard

Supervised by Aurélien Garivier and Gilles Stoltz

July 3, 2018

Chapters addressed in this presentation:

- Explore First, Exploit Next: The True Shape of Regret in Bandit Problems
- kl-UCB Algorithms for Exponential Families
- KL-UCB Algorithms for Bounded Rewards
- Thresholding Bandit for Dose-ranging: The Impact of Monotonicity
- Fano's inequality for random variables

# K-armed bandit problem: parametric setting

Bernoulli rewards:

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$



...



...



Game: for each round  $1 \leq t \leq T$ :

1. Player pulls arm  $A_t \in \{1, \dots, K\}$ .
2. He gets a reward  $Y_t \sim \mathcal{B}(\mu_{A_t})$ .

## Regret

Player wants to maximize

$$\mathbb{E} \left[ \sum_{t=1}^T Y_t \right],$$

equivalently, minimize his regret

$$R_T = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T Y_t \right],$$

where  $\mu^* = \max_{a=1,\dots,K} \mu_a$ .

## Regret

Player wants to maximize

$$\mathbb{E} \left[ \sum_{t=1}^T Y_t \right],$$

equivalently, minimize his regret

$$R_T = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T Y_t \right],$$

where  $\mu^* = \max_{a=1, \dots, K} \mu_a$ .

Chain rule

$$R_T = \sum_{a=1}^K (\mu^* - \mu_a) \mathbb{E}[N_a(T)]$$

where  $N_a(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}}$ .

## Regret

Player wants to maximize

$$\mathbb{E} \left[ \sum_{t=1}^T Y_t \right],$$

equivalently, minimize his regret

$$R_T = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T Y_t \right],$$

where  $\mu^* = \max_{a=1, \dots, K} \mu_a$ .

Chain rule

$$R_T = \sum_{a=1}^K (\mu^* - \mu_a) \mathbb{E}[N_a(T)] (\sim T \text{ worst case}),$$

where  $N_a(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}}$ .

## Ideas of strategy

- First idea: pull an arm uniformly at random at each round.  
⇒ Exploration      ⇒  $R_T \sim T$

## Ideas of strategy

- First idea: pull an arm uniformly at random at each round.  
 $\Rightarrow$  Exploration  $\Rightarrow R_T \sim T$
- Second idea: pull the current best empirical arm,

$$A_{t+1} = \operatorname{argmax}_{a \in \{1, \dots, K\}} \hat{\mu}_{a, N_a(t)} \quad \hat{\mu}_{a, N_a(t)} = \sum_{s=1}^t Y_s \mathbb{I}_{A_t=a} / N_a(t)$$

$$\Rightarrow \text{Exploitation} \quad \Rightarrow R_T \sim T$$



## Ideas of strategy

- First idea: pull an arm uniformly at random at each round.

⇒ Exploration      ⇒  $R_T \sim T$

- Second idea: pull the current best empirical arm,

$$A_{t+1} = \operatorname{argmax}_{a \in \{1, \dots, K\}} \hat{\mu}_{a, N_a(t)} \quad \hat{\mu}_{a, N_a(t)} = \sum_{s=1}^t Y_s \mathbb{I}_{A_t=a} / N_a(t)$$

⇒ Exploitation      ⇒  $R_T \sim T$

⇒ Exploration-Exploitation tradeoff

⇒  $R_T \sim \log(T)$

# UCB algorithm

---

## Algorithm 1: UCB

---

**Initialization:** Play each arm once.

**For**  $t = K$  to  $T - 1$ , **do**

1. Compute for each arm  $a$  the upper confidence bound

$$U_a^{\text{UCB}}(t) = \underbrace{\hat{\mu}_{a, N_a(t)}}_{\text{Exploitation}} + \underbrace{\sqrt{\frac{\log(T)}{2N_a(t)}}}_{\text{Exploration}}$$

2. Play  $A_{t+1} \in \operatorname{argmax}_{a \in \{1, \dots, K\}} U_a^{\text{UCB}}(t)$ .
-

## Upper Confident Bound

$X_1, \dots, X_n$  i.i.d.  $\sim \mathcal{B}(\mu)$  with  $\hat{\mu}_n = \sum_{k=1}^n X_k/n$

Hoeffding inequality for  $x < \mu$

$$\mathbb{P}(\hat{\mu}_n < x) \leq e^{-2n(x-\mu)^2}.$$

With probability at least  $1 - \delta$

$$\mu \leq \hat{\mu}_n + \sqrt{\frac{\log(1/\delta)}{2n}}.$$

## Upper Confident Bound

$X_1, \dots, X_n$  i.i.d.  $\sim \mathcal{B}(\mu)$  with  $\hat{\mu}_n = \sum_{k=1}^n X_k/n$

Hoeffding inequality for  $x < \mu$

$$\mathbb{P}(\hat{\mu}_n < x) \leq e^{-2n(x-\mu)^2}.$$

With probability at least  $1 - \delta$

$$\mu \leq \hat{\mu}_n + \sqrt{\frac{\log(1/\delta)}{2n}}.$$

UCB index  $\delta = 1/T$

$$U_a^{\text{UCB}}(t) = \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{\log(T)}{2N_a(t)}}$$

## Upper Confident Bound

$X_1, \dots, X_n$  i.i.d.  $\sim \mathcal{B}(\mu)$  with  $\hat{\mu}_n = \sum_{k=1}^n X_k/n$

Hoeffding inequality for  $x < \mu$

$$\mathbb{P}(\hat{\mu}_n < x) \leq e^{-2n(x-\mu)^2}.$$

With probability at least  $1 - \delta$

$$\mu \leq \hat{\mu}_n + \sqrt{\frac{\log(1/\delta)}{2n}}.$$

UCB index  $\delta = 1/T$

$$U_a^{\text{UCB}}(t) = \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{\log(T)}{2N_a(t)}}$$

# Regret bound

## Theorem

For the UCB algorithm, for all  $a$  such that  $\mu^* - \mu_a > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{1}{2(\mu^* - \mu_a)^2} \log(T) + o(\log(T)),$$

# Regret bound

## Theorem

For the UCB algorithm, for all  $a$  such that  $\mu^* - \mu_a > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{1}{2(\mu^* - \mu_a)^2} \log(T) + o(\log(T)),$$

therefore (Chain rule)

$$R_T \leq \sum_{a: \mu^* > \mu_a} \frac{1}{2(\mu^* - \mu_a)^2} \log(T) + o(\log(T)).$$

## Regret bound

### Theorem

For the UCB algorithm, for all  $a$  such that  $\mu^* - \mu_a > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{1}{2(\mu^* - \mu_a)^2} \log(T) + o(\log(T)),$$

therefore (Chain rule)

$$R_T \leq \sum_{a: \mu^* > \mu_a} \frac{1}{2(\mu^* - \mu_a)^2} \log(T) + o(\log(T)).$$

Is that the best we can do?  $\Rightarrow$  Lower bound



# Kullback-Leibler divergence

For two probability distributions  $P$  and  $Q$

$$\text{KL}(P, Q) = \begin{cases} \int \log\left(\frac{dP}{dQ}\right) dQ & \text{if } P \ll Q \\ +\infty & \text{else.} \end{cases}$$

Example with Bernoulli

$$\text{kl}(p, q) := \text{KL}(\mathcal{B}(p), \mathcal{B}(q)) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$$

## An asymptotic lower bound

Strategy which always pulls the same arm  $\Rightarrow$  assumptions on the strategy.

## An asymptotic lower bound

Strategy which always pulls the same arm  $\Rightarrow$  assumptions on the strategy.

### Definition

A strategy is consistent if for all bandit problems  $\nu$ , for all suboptimal arms  $a$ , i.e., for all arms  $a$  such that  $\mu^* - \mu_a > 0$ , it satisfies  $\mathbb{E}[N_a(T)] = o(T^\alpha)$  for all  $0 < \alpha \leq 1$ .

## An asymptotic lower bound

Strategy which always pulls the same arm  $\Rightarrow$  assumptions on the strategy.

### Definition

A strategy is consistent if for all bandit problems  $\nu$ , for all suboptimal arms  $a$ , i.e., for all arms  $a$  such that  $\mu^* - \mu_a > 0$ , it satisfies  $\mathbb{E}[N_a(T)] = o(T^\alpha)$  for all  $0 < \alpha \leq 1$ .

### Theorem (Asymptotic lower bound from Lai & Robbins)

*For all consistent strategies, for all suboptimal arms  $a$ ,*

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\log T} \geq \frac{1}{\text{kl}(\mu_a, \mu^*)}.$$

## Sketch of proof 1/2

a suboptimal arm ( $\mu^* - \mu_a > 0$ ).

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\underline{\nu}' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

Information at time  $t$ :  $Y^{1:t} = (Y_1, \dots, Y_t)$ .

## Sketch of proof 1/2

a suboptimal arm ( $\mu^* - \mu_a > 0$ ).

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\underline{\nu}' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

Information at time  $t$ :  $Y^{1:t} = (Y_1, \dots, Y_t)$ .

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{kl}(\mu_a, \mu'_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{Y_{1:T}}, \mathbb{P}_{\underline{\nu}'}^{Y_{1:T}})$$

Chain rule

## Sketch of proof 1/2

a suboptimal arm ( $\mu^* - \mu_a > 0$ ).

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\underline{\nu}' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

Information at time  $t$ :  $Y^{1:t} = (Y_1, \dots, Y_t)$ .

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{kl}(\mu_a, \mu'_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{Y_{1:T}}, \mathbb{P}_{\underline{\nu}'}^{Y_{1:T}})$$

$$\text{contraction of entropy} \quad \geq \text{KL}(\mathbb{P}_{\underline{\nu}}^{N_a(T)/T}, \mathbb{P}_{\underline{\nu}'}^{N_a(T)/T})$$

## Sketch of proof 1/2

a suboptimal arm ( $\mu^* - \mu_a > 0$ ).

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\underline{\nu}' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

Information at time  $t$ :  $Y^{1:t} = (Y_1, \dots, Y_t)$ .

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{ kl}(\mu_a, \mu'_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{Y_{1:T}}, \mathbb{P}_{\underline{\nu}'}^{Y_{1:T}})$$

contraction of entropy

$$\geq \text{KL}(\mathbb{P}_{\underline{\nu}}^{N_a(T)/T}, \mathbb{P}_{\underline{\nu}'}^{N_a(T)/T})$$

projection

$$\geq \text{kl}\left(\mathbb{E}_{\underline{\nu}}[N_a(T)]/T, \mathbb{E}_{\underline{\nu}'}[N_a(T)]/T\right)$$



## Sketch of proof 1/2

a suboptimal arm ( $\mu^* - \mu_a > 0$ ).

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\underline{\nu}' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

Information at time  $t$ :  $Y^{1:t} = (Y_1, \dots, Y_t)$ .

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{kl}(\mu_a, \mu'_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{Y_{1:T}}, \mathbb{P}_{\underline{\nu}'}^{Y_{1:T}})$$

contraction of entropy

$$\geq \text{KL}(\mathbb{P}_{\underline{\nu}}^{N_a(T)/T}, \mathbb{P}_{\underline{\nu}'}^{N_a(T)/T})$$

projection

$$\geq \text{kl}\left(\mathbb{E}_{\underline{\nu}}[N_a(T)]/T, \mathbb{E}_{\underline{\nu}'}[N_a(T)]/T\right)$$

$$\text{kl}(p, q) \geq p \log(1/q) - \log(2) \geq \left(1 - \mathbb{E}_{\underline{\nu}}[N_a(T)]/T\right) \log \frac{T}{T - \mathbb{E}_{\underline{\nu}'}[N_a(T)]} - \log(2)$$

## Sketch of proof 1/2

a suboptimal arm ( $\mu^* - \mu_a > 0$ ).

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\underline{\nu}' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

Information at time  $t$ :  $Y^{1:t} = (Y_1, \dots, Y_t)$ .

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{kl}(\mu_a, \mu'_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{Y_{1:T}}, \mathbb{P}_{\underline{\nu}'}^{Y_{1:T}})$$

contraction of entropy

$$\geq \text{KL}(\mathbb{P}_{\underline{\nu}}^{N_a(T)/T}, \mathbb{P}_{\underline{\nu}'}^{N_a(T)/T})$$

projection

$$\geq \text{kl}\left(\mathbb{E}_{\underline{\nu}}[N_a(T)]/T, \mathbb{E}_{\underline{\nu}'}[N_a(T)]/T\right)$$

$$\text{Consistent} \geq \left(1 - \underbrace{\mathbb{E}_{\underline{\nu}}[N_a(T)]/T}_{o(1)}\right) \log \frac{T}{\underbrace{T - \mathbb{E}_{\underline{\nu}'}[N_a(T)]}_{O(T^\alpha)}} - \log(2)$$

## Sketch of proof 1/2

a suboptimal arm ( $\mu^* - \mu_a > 0$ ).

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\underline{\nu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\underline{\nu}' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

Information at time  $t$ :  $Y^{1:t} = (Y_1, \dots, Y_t)$ .

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{kl}(\mu_a, \mu'_a) = \text{KL}(\mathbb{P}_{\underline{\nu}}^{Y_{1:T}}, \mathbb{P}_{\underline{\nu}'}^{Y_{1:T}})$$

contraction of entropy

$$\geq \text{KL}(\mathbb{P}_{\underline{\nu}}^{N_a(T)/T}, \mathbb{P}_{\underline{\nu}'}^{N_a(T)/T})$$

projection

$$\geq \text{kl}\left(\mathbb{E}_{\underline{\nu}}[N_a(T)]/T, \mathbb{E}_{\underline{\nu}'}[N_a(T)]/T\right)$$

$$\gtrsim (1 - \alpha) \log(T) - \log(2)$$

For all  $\alpha \in (0, 1]$ :

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\nu} [N_a(T)]}{\log T} \geq \frac{1 - \alpha}{\text{kl}(\mu_a, \mu'_a)}.$$

For all  $\alpha \in (0, 1]$ :

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\nu} [N_a(T)]}{\log T} \geq \frac{1 - \alpha}{\text{kl}(\mu_a, \mu'_a)}.$$

Very generic tools, can be used to prove :

- minimax/problem dependent lower bound on the Bayesian posterior concentration rates
- minimax lower bound for sparse adversarial bandits
- lower bounds in best arm identification

## Sub-optimality of UCB

UCB

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\log T} \leq \frac{1}{2(\mu_a - \mu^*)^2},$$

Lower bound

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\log T} \geq \frac{1}{\text{kl}(\mu_a, \mu^*)}.$$

Pinsker inequality

$$\text{kl}(\mu_a, \mu^*) \geq 2(\mu_a - \mu^*)^2$$

## Chernoff Bound

$X_1, \dots, X_n$  i.i.d.  $\sim \mathcal{B}(\mu)$  with  $\hat{\mu}_n = \sum_{k=1}^n X_k/n$

Chernoff inequality for  $x < \mu$

$$\mathbb{P}(\hat{\mu}_n < x) \leq e^{-nkl(x, \mu)} \underset{\text{ Pinsker }}{\leq} e^{-2n(x-\mu)^2}$$

## Chernoff Bound

$X_1, \dots, X_n$  i.i.d.  $\sim \mathcal{B}(\mu)$  with  $\hat{\mu}_n = \sum_{k=1}^n X_k/n$

Chernoff inequality for  $x < \mu$

$$\mathbb{P}(\hat{\mu}_n < x) \leq e^{-n \text{kl}(x, \mu)} \stackrel{\text{ Pinsker }}{\leq} e^{-2n(x-\mu)^2}$$

Inverting for  $u = \text{kl}(x, \mu)$

$$\mathbb{P}(\hat{\mu}_n < \mu \text{ and } \text{kl}(\hat{\mu}_n, \mu) > u) \leq e^{-nu}$$

New upper confidence bound, with probability at least  $1 - \delta$

$$\hat{\mu}_n \geq \mu \text{ or } \text{kl}(\hat{\mu}_n, \mu) \leq \frac{\log(1/\delta)}{n}$$

$$\mu \leq \sup \left\{ \mu' \geq \hat{\mu}_n : \text{kl}(\hat{\mu}_n, \mu') \leq \frac{\log(1/\delta)}{n} \right\}$$



## Get the right constant: kl-UCB algorithm

### kl-UCB algorithm:

Play each arm once, then for  $K \leq t \leq T - 1$  play according to the index:

$$U_a^{\text{kl}}(t) = \sup \left\{ \mu \geq \hat{\mu}_a(t) : N_a(t) \text{kl}(\hat{\mu}_a(t), \mu) \leq \log(T) \right\}.$$

## Get the right constant: kl-UCB algorithm

### kl-UCB algorithm:

Play each arm once, then for  $K \leq t \leq T - 1$  play according to the index:

$$U_a^{\text{kl}}(t) = \sup \left\{ \mu \geq \hat{\mu}_a(t) : N_a(t) \text{kl}(\hat{\mu}_a(t), \mu) \leq \log(T) \right\}.$$

### Regret upper bound:

For the kl-UCB algorithm, for all  $a$  such that  $\mu^* - \mu_a > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{1}{\text{kl}(\mu_a, \mu^*)} \log(T) + o(\log(T)),$$

Lower bound

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\log T} \geq \frac{1}{\text{kl}(\mu_a, \mu^*)}.$$

# Minimax optimality

## Theorem (Minimax lower bound)

$$\sup_{\underline{\nu}} R_{T, \underline{\nu}} \geq C\sqrt{KT}.$$

An algorithm is **minimax** optimal if

$$R_T \leq C'\sqrt{KT}.$$

UCB algorithm is not minimax optimal

$$R_T \leq C''\sqrt{KT \log(T)}.$$

# MOSS-algorithm

## MOSS index:

$$U_a^M(t) = \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{\log_+ \left( T / (KN_a(t)) \right)}{2N_a(t)}}.$$

MOSS algorithm is minimax optimal, i.e.

$$R_T \leq 17\sqrt{KT} + K,$$

but is not asymptotically optimal (as UCB)...

## A minimax and asymptotically optimal algorithm

**kl-UCB<sup>++</sup> index:**

$$U_a^{\text{kl}^{++}}(t) = \sup \left\{ \mu \geq \hat{\mu}_a(t) : N_a(t) \text{kl}(\hat{\mu}_a(t), \mu) \leq \log_+ \left( T / (KN_a(t)) \right) \right\}.$$

## A minimax and asymptotically optimal algorithm

**kl-UCB<sup>++</sup> index:**

$$U_a^{\text{kl}^{++}}(t) = \sup \left\{ \mu \geq \hat{\mu}_a(t) : N_a(t) \text{kl}(\hat{\mu}_a(t), \mu) \leq \log_+ \left( T / (KN_a(t)) \right) \right\}.$$

kl-UCB<sup>++</sup> algorithm is **minimax** and **asymptotically** optimal:

**Minimax optimality:**

$$R_T \leq 17\sqrt{KT} + K.$$

**Asymptotic optimality:** For any sub-optimal arm  $a$ ,

$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{\text{kl}(\mu_a, \mu^*)} + o(\log(T)).$$

## Second order optimality

With a stronger assumption:

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

refinement of the Lai and Robbins lower bound

$$\mathbb{E}[N_a(T)] \geq \frac{\log T}{\text{kl}(\mu_a, \mu^*)} - O(\log \log(T)).$$

## Second order optimality

With a stronger assumption:

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

refinement of the Lai and Robbins lower bound

$$\mathbb{E}[N_a(T)] \geq \frac{\log T}{\text{kl}(\mu_a, \mu^*)} - O(\log \log(T)).$$

For the kl-UCB<sup>++</sup> algorithm, for any sub-optimal arm  $a$ ,

$$\mathbb{E}[N_a(T)] \leq \frac{\log(T) - \log \log(T)}{\text{kl}(\mu_a, \mu^*)} + O(1).$$



## An intuition for the exploration function

Bandit problems  $\underline{\nu}$  and  $\underline{\nu}'$  with  $\underline{\nu} = \underline{\nu}'$  except  $\mu'_a > \mu^*$ . Lower bound:

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{kl}(\mu_a, \mu'_a) \geq \left(1 - \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{T}\right) \log \frac{T}{\sum_{b \neq a} \mathbb{E}_{\underline{\nu}'}[N_b(T)]} - \ln 2.$$

## An intuition for the exploration function

Bandit problems  $\underline{\nu}$  and  $\underline{\nu}'$  with  $\underline{\nu} = \underline{\nu}'$  except  $\mu'_a > \mu^*$ . Lower bound:

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \text{kl}(\mu_a, \mu'_a) \geq \left(1 - \frac{\mathbb{E}_{\underline{\nu}}[N_a(T)]}{T}\right) \log \frac{T}{\sum_{b \neq a} \mathbb{E}_{\underline{\nu}'}[N_b(T)]} - \ln 2.$$

Approximations:

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \approx N_a(T) \left( \approx \frac{\log(T)}{\text{kl}(\mu_a, \mu^*)} \right), \quad \mu_a \approx \hat{\mu}_a(T), \quad \sum_{b \neq a} \mathbb{E}_{\underline{\nu}'}[N_b(T)] \approx KN_a(T).$$

## An intuition for the exploration function

Bandit problems  $\underline{\nu}$  and  $\underline{\nu}'$  with  $\underline{\nu} = \underline{\nu}'$  except  $\mu'_a > \mu^*$ . Lower bound:

$$N_a(T) \text{kl}(\hat{\mu}_a(T), \mu'_a) \gtrsim \log \frac{T}{\sum_{b \neq a} \mathbb{E}_{\underline{\nu}'}[N_b(T)]}.$$

Approximations:

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \approx N_a(T) \left( \approx \frac{\log(T)}{\text{kl}(\mu_a, \mu^*)} \right), \quad \mu_a \approx \hat{\mu}_a(T), \quad \sum_{b \neq a} \mathbb{E}_{\underline{\nu}'}[N_b(T)] \approx KN_a(T).$$

## An intuition for the exploration function

Bandit problems  $\underline{\nu}$  and  $\underline{\nu}'$  with  $\underline{\nu} = \underline{\nu}'$  except  $\mu'_a > \mu^*$ . Lower bound:

$$N_a(T) \text{kl}(\hat{\mu}_a(T), \mu'_a) \gtrsim \log \frac{T}{\sum_{b \neq a} \mathbb{E}_{\underline{\nu}'}[N_b(T)]}.$$

Approximations:

$$\mathbb{E}_{\underline{\nu}}[N_a(T)] \approx N_a(T) \left( \approx \frac{\log(T)}{\text{kl}(\mu_a, \mu^*)} \right), \quad \mu_a \approx \hat{\mu}_a(T), \quad \sum_{b \neq a} \mathbb{E}_{\underline{\nu}'}[N_b(T)] \approx KN_a(T).$$

Index:

$$\begin{aligned} U_a^{\text{kl}++}(T) &= \sup \left\{ \mu \geq \hat{\mu}_a(T) : N_a(T) \text{kl}(\hat{\mu}_a(T), \mu) \leq \log \frac{T}{KN_a(T)} \right\} \\ &= \inf \left\{ \mu \geq \hat{\mu}_a(T) : N_a(T) \text{kl}(\hat{\mu}_a(T), \mu) \geq \log \frac{T}{KN_a(T)} \right\}, \end{aligned}$$

# K-armed bandit problem: non-parametric setting

Bounded rewards:  $\nu_a \in \mathcal{P}[0, 1]$



Game: for each round  $1 \leq t \leq T$ :

1. Player pulls arm  $A_t \in \{1, \dots, K\}$ .
2. He gets a reward  $Y_t \sim \nu_{A_t}$ .

$$\mu_a = E(\nu_a)$$

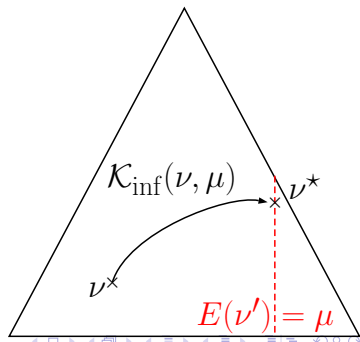
## Lower bound

### Theorem (Asymptotic lower)

For all consistent strategies, for all arms  $a$  such that  $\mu^* - E(\nu_a) > 0$ ,

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

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



## Lower bound

### Theorem (Asymptotic lower)

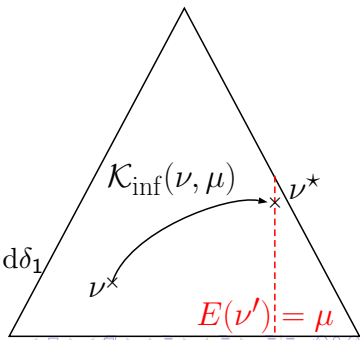
For all consistent strategies, for all arms  $a$  such that  $\mu^* - E(\nu_a) > 0$ ,

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

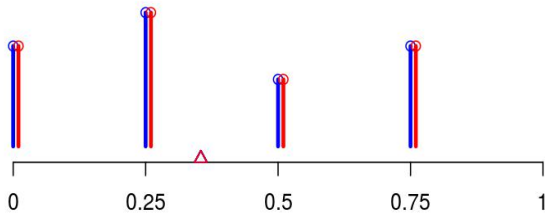
$$\begin{aligned} \mathcal{K}_{\text{inf}}(\nu, \mu) &:= \inf \{ \text{KL}(\nu, \nu') : E(\nu') > \mu \} \\ &= \text{KL}(\nu, \nu^*) \end{aligned}$$

$\nu^*$  of the form, for a certain  $\lambda^* \in [0, 1]$

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

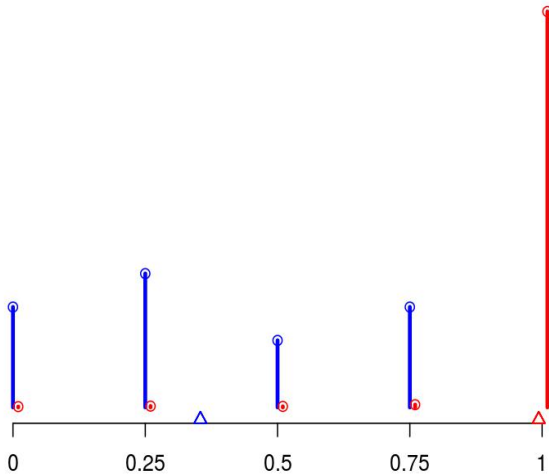


$$E(\nu) \leq E(\nu^*) = \mu$$





$$E(\nu) \leq E(\nu^*) = \mu$$



## Sub-optimality of kl-UCB

kl-UCB

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\log T} \leq \frac{1}{\text{kl}(\mu_a, \mu^*)},$$

Lower bound

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

Pseudo-Pinsker inequality

$$\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) \geq \text{kl}(E(\nu_a), \mu^*)$$

## Sub-optimality of kl-UCB

kl-UCB

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\log T} \leq \frac{1}{\text{kl}(\mu_a, \mu^*)},$$

Lower bound

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

Pseudo-Pinsker inequality

$$\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) \geq \text{kl}(E(\nu_a), \mu^*)$$

Reduction to kl for Bernoulli:

$$\mathcal{K}_{\text{inf}}(\mathcal{B}(\mu_a), \mu^*) = \text{kl}(\mu_a, \mu^*)$$

## Sub-optimality of kl-UCB

kl-UCB

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\log T} \leq \frac{1}{\text{kl}(\mu_a, \mu^*)},$$

Lower bound

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

Pseudo-Pinsker inequality

$$\mathcal{K}_{\text{inf}}(\nu_a, \mu^*) \geq \text{kl}(E(\nu_a), \mu^*)$$

$$\underbrace{\inf \left\{ \text{KL}(\nu, \nu') : E(\nu') > \mu \right\}}_{=\mathcal{K}_{\text{inf}}(\nu, \mu)} \geq \underbrace{\inf \left\{ \text{KL}(\nu'', \nu') : E(\nu') > \mu, E(\nu'') = E(\nu) \right\}}_{=\text{kl}(E(\nu), \mu)}$$

## Index ?

Move from empirical mean  $\hat{\mu}_n$  to empirical distribution  $\hat{\nu}_n = 1/n \sum_{k=1}^n \delta_{X_k}$

New index

$$U_a^{\text{kl}}(t) = \sup \left\{ \mu' \geq \hat{\mu}_a(t) : \mu' \in [0, 1], \text{kl}(\hat{\mu}_a(t), \mu') \leq \frac{\log(T)}{N_a(t)} \right\}$$

$$U_a^{\text{KL}}(t) = \sup \left\{ E\nu' \geq E(\hat{\nu}_a(t)) : \nu' \in \mathcal{P}[0, 1], \text{KL}(\hat{\nu}_a(t), \nu') \leq \frac{\log(T)}{N_a(t)} \right\}$$

## Index ?

Move from empirical mean  $\hat{\mu}_n$  to empirical distribution  $\hat{\nu}_n = 1/n \sum_{k=1}^n \delta_{X_k}$

New index

$$U_a^{\text{kl}}(t) = \sup \left\{ \mu' \geq \hat{\mu}_a(t) : \mu' \in [0, 1], \text{kl}(\hat{\mu}_a(t), \mu') \leq \frac{\log(T)}{N_a(t)} \right\}$$

$$U_a^{\text{KL}}(t) = \sup \left\{ E\nu' \geq E(\hat{\nu}_a(t)) : \nu' \in \mathcal{P}[0, 1], \text{KL}(\hat{\nu}_a(t), \nu') \leq \frac{\log(T)}{N_a(t)} \right\}$$

$$= \sup \left\{ \mu' : \mu' \in [0, 1], \mu' \geq \hat{\mu}_a(t), \mathcal{K}_{\text{inf}}(\hat{\nu}_a(t), \mu') \leq \frac{\log(T)}{N_a(t)} \right\}.$$

# An asymptotically optimal algorithm

## KL-UCB index:

$$U_a^{\text{KL}}(t) = \sup \left\{ \mu' \geq \hat{\mu}_a(t) : N_a(t) \mathcal{K}_{\text{inf}}(\hat{\nu}_a(t), \mu') \leq \log(T) \right\}.$$

## An asymptotically optimal algorithm

### KL-UCB index:

$$U_a^{\text{KL}}(t) = \sup \left\{ \mu' \geq \hat{\mu}_a(t) : N_a(t) \mathcal{K}_{\text{inf}}(\hat{\nu}_a(t), \mu') \leq \log(T) \right\}.$$

**Upper bound on the regret:** for all  $a$  such that  $\mu^* - E(\mu_a) > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{1}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} \log(T) + o(\log(T)),$$

Lower bound

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



## Non-parametric upper confidence bound

$X_1, \dots, X_n$  i.i.d.  $\sim \nu$  with  $\hat{\nu}_n = \sum_{k=1}^n \delta_{X_k} / n$ .

### Deviations of kl

$$\mathbb{P}\left(\hat{\mu}_n < E(\nu) \text{ and } \text{kl}(\hat{\mu}_n, E(\nu)) > u\right) \leq e^{-nu}$$

### Deviations of $\mathcal{K}_{\text{inf}}$

$$\mathbb{P}\left(\mathcal{K}_{\text{inf}}(\hat{\nu}_n, E(\nu)) > u\right) \leq e(n+3)e^{-nu}.$$

# Non-parametric upper confidence bound

$X_1, \dots, X_n$  i.i.d.  $\sim \nu$  with  $\hat{\nu}_n = \sum_{k=1}^n \delta_{X_k} / n$ .

## Deviations of kl

$$\mathbb{P}\left(\hat{\mu}_n < E(\nu) \text{ and } \text{kl}(\hat{\mu}_n, E(\nu)) > u\right) \leq e^{-nu}$$

## Deviations of $\mathcal{K}_{\text{inf}}$

$$\mathbb{P}\left(\mathcal{K}_{\text{inf}}(\hat{\nu}_n, E(\nu)) > u\right) \leq e^{(n+3)u} e^{-nu}.$$

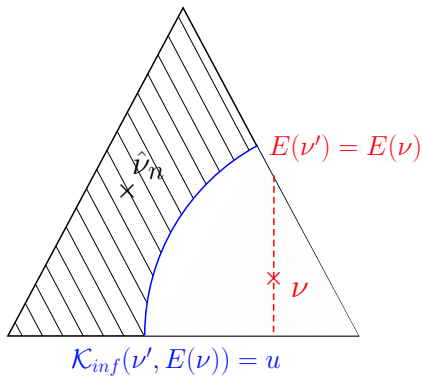
Open question: remove the factor  $(n+3)$  ?

## Variational formula

$$\mathcal{K}_{\text{inf}}(\nu, \mu) = \max_{0 \leq \lambda \leq 1} \mathbb{E}_{\nu} \left[ \underbrace{\log \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right)}_{:= f_{\lambda}(X)} \right]$$

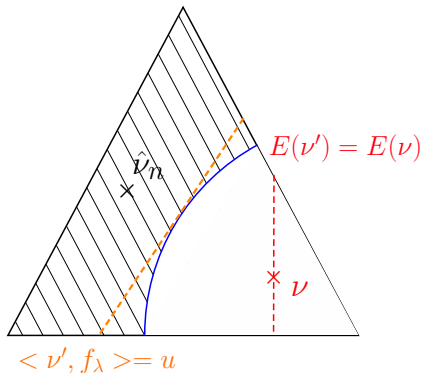
## Variational formula

$$\mathcal{K}_{\text{inf}}(\hat{\nu}_n, E(\nu)) = \max_{0 \leq \lambda \leq 1} \mathbb{E}_{\hat{\nu}_n} \left[ \underbrace{\log \left( 1 - \lambda \frac{X - E(\nu)}{1 - E(\nu)} \right)}_{:= f_\lambda(X)} \right]$$



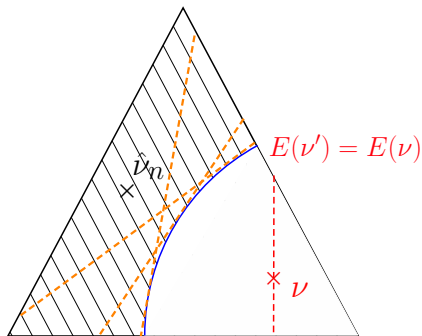
## Variational formula

$$\mathcal{K}_{\text{inf}}(\hat{\nu}_n, E(\nu)) = \max_{0 \leq \lambda \leq 1} \langle \hat{\nu}_n, f_\lambda \rangle$$



## Variational formula

$$\mathcal{K}_{\text{inf}}(\hat{\nu}_n, E(\nu)) = \max_{0 \leq \lambda \leq 1} \langle \hat{\nu}_n, f_\lambda \rangle$$



## Variational formula: Worst family

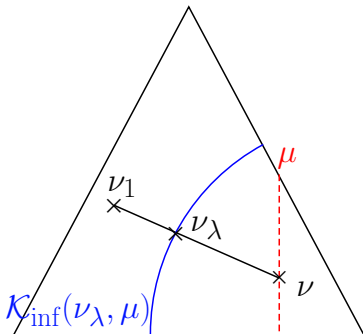
$$\mathcal{K}_{\text{inf}}(\nu, \mu) = \max_{0 \leq \lambda \leq 1} \mathbb{E}_{\nu} \left[ \log \left( 1 - \lambda \frac{X - \mu}{1 - \mu} \right) \right].$$

If  $E(\nu) = \mu$ . Convex family of probability distributions:  $\frac{d\nu_{\lambda}}{d\nu} = \left( 1 - \lambda \frac{x - \mu}{1 - \mu} \right)$

$$\nu_{\lambda} = \lambda\nu_1 + (1 - \lambda)\nu$$

Worst family for  $\nu$ :

$$\mathcal{K}_{\text{inf}}(\nu_{\lambda}, \mu) = \text{KL}(\nu_{\lambda}, \nu)$$



## Variational formula: Worst family

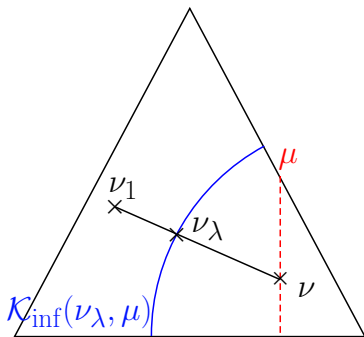
$$\mathcal{K}_{\text{inf}}(\nu, \mu) = - \min_{0 \leq \lambda \leq 1} \text{KL}(\nu, \nu_\lambda) = 0.$$

If  $E(\nu) = \mu$ . Convex family of probability distributions:  $\frac{d\nu_\lambda}{d\nu} = \left(1 - \lambda \frac{x-\mu}{1-\mu}\right)$

$$\nu_\lambda = \lambda\nu_1 + (1-\lambda)\nu$$

Worst family for  $\nu$ :

$$\mathcal{K}_{\text{inf}}(\nu_\lambda, \mu) = \text{KL}(\nu_\lambda, \nu)$$





## Thompson Sampling for arms in $\mathcal{P}[0, 1]$ ?

$\pi_a^0 = \mathcal{D}(\alpha, H)$  Dirichlet process with  $H \in \mathcal{P}[0, 1]$ .

---

### Algorithm 2: Thompson sampling

---

**Parameter:** prior  $\Pi^0 = (\pi_1^0, \dots, \pi_K^0)$

**For**  $t = 0$  to  $T - 1$ , **do**

1. **For**  $a = 1$  to  $K$ , **do**

    Sample  $\nu_a(t) \sim \pi_a^t$ .

2. Play  $A_t \in \operatorname{argmax}_{a \in \{1, \dots, K\}} E(\nu_a(t))$ , update the posterior  $\Pi^{t+1}$ .

---

Is Thompson Sampling asymptotically optimal ?

## A minimax and asymptotically optimal algorithm

**KL-UCB<sup>++</sup> index:**

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

## A minimax and asymptotically optimal algorithm

**KL-UCB<sup>++</sup> index:**

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

**KL-UCB-switch index:**

$$U_a^{\text{KL-S}}(t) := \begin{cases} U_a^{\text{KL}^{++}}(t) & \text{si } N_a(t) \leq \lfloor (T/K)^{1/5} \rfloor \\ U_a^{\text{M}}(t) & \text{si } N_a(t) > \lfloor (T/K)^{1/5} \rfloor \end{cases},$$

## A minimax and asymptotically optimal algorithm

**KL-UCB<sup>++</sup> index:**

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

**KL-UCB-switch index:**

$$U_a^{\text{KL-S}}(t) := \begin{cases} U_a^{\text{KL}^{++}}(t) & \text{si } N_a(t) \leq \lfloor (T/K)^{1/5} \rfloor \\ U_a^{\text{M}}(t) & \text{si } N_a(t) > \lfloor (T/K)^{1/5} \rfloor \end{cases},$$

KL-UCB-switch algorithm is **minimax** and **asymptotically** optimal:

**Minimax optimality:**

$$R_T \leq 25\sqrt{KT} + K.$$

**Asymptotic optimality:** For any sub-optimal arm  $a$

$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} + o(\log(T))$$

## A minimax and asymptotically optimal algorithm

**KL-UCB<sup>++</sup> index:**

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

**KL-UCB-switch index:**

$$U_a^{\text{KL-S}}(t) := \begin{cases} U_a^{\text{KL}^{++}}(t) & \text{si } N_a(t) \leq \lfloor (T/K)^{1/5} \rfloor \\ U_a^{\text{M}}(t) & \text{si } N_a(t) > \lfloor (T/K)^{1/5} \rfloor \end{cases},$$

Open question, does the same theorem hold for KL-UCB<sup>++</sup> ?

**Minimax optimality:**

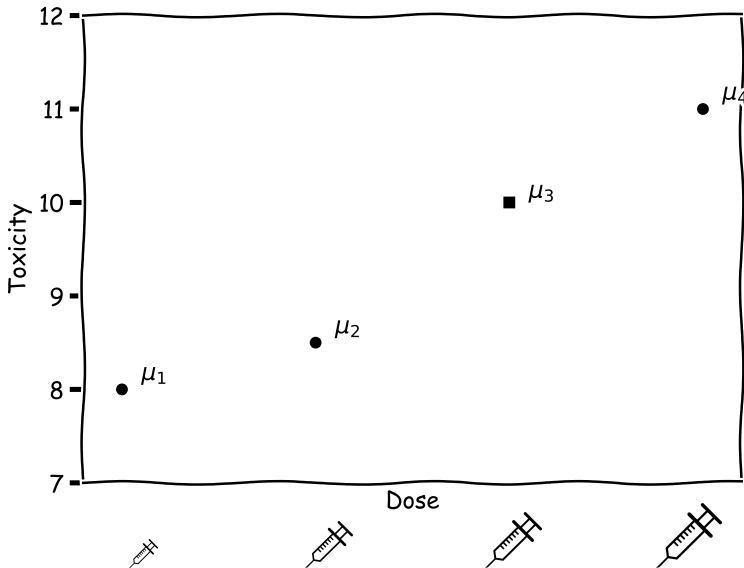
$$R_T \leq 25\sqrt{KT} + K.$$

**Asymptotic optimality:** For any sub-optimal arm  $a$

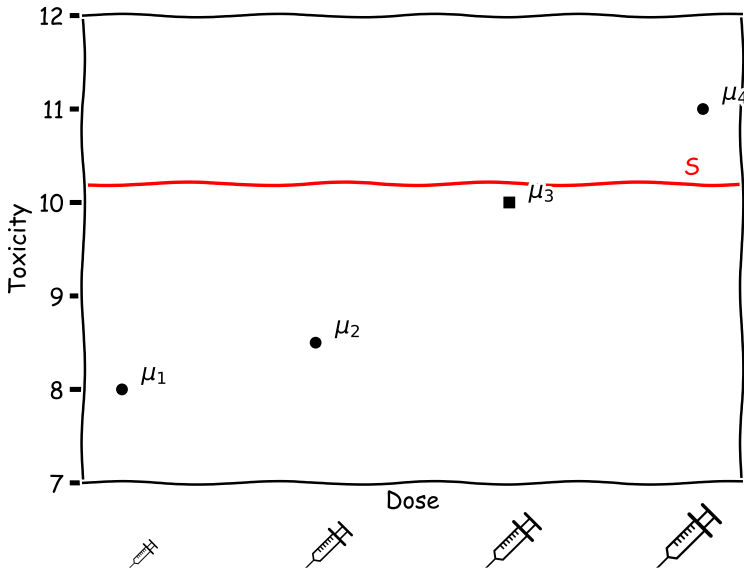
$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)} + o(\log(T))$$

Merci pour votre attention !

# Dose-ranging

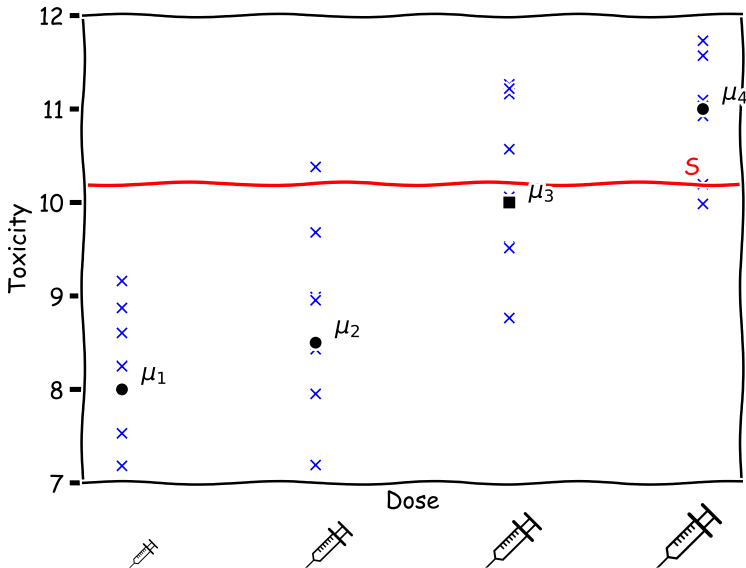


# Dose-ranging





## Dose-ranging



## Best arm identification

$$\boldsymbol{\mu} = (\mathcal{N}(\mu_1, 1) \quad \cdots \quad \mathcal{N}(\mu_a, 1) \quad \cdots \quad \mathcal{N}(\mu_K, 1))$$



...



...



## Best arm identification

$$\boldsymbol{\mu} = \begin{pmatrix} \mathcal{N}(\mu_1, 1) & \cdots & \mathcal{N}(\mu_a, 1) & \cdots & \mathcal{N}(\mu_K, 1) \\ \mu_1 & \cdots & \mu_a & \cdots & \mu_K \end{pmatrix}$$

Optimal arm (dose):  $a_{\boldsymbol{\mu}}^* \in \arg \min_{1 \leq a \leq K} |\mu_a - S|$

## Best arm identification

$$\begin{aligned} \boldsymbol{\mu} &= (\mathcal{N}(\mu_1, 1) \quad \cdots \quad \mathcal{N}(\mu_a, 1) \quad \cdots \quad \mathcal{N}(\mu_K, 1)) \\ &\sim \quad [\mu_1 \quad \cdots \quad \mu_a \quad \cdots \quad \mu_K] \end{aligned}$$

Optimal arm (dose):  $a_{\boldsymbol{\mu}}^* \in \arg \min_{1 \leq a \leq K} |\mu_a - S|$

**Game:** while  $t < \tau$ :

1. Player pulls arm (dose)  $A_t \in \{1, \dots, K\}$ .
2. He gets an observation (toxicity)  $Y_t \sim \mathcal{N}(\mu_{A_t}, 1)$ .

Predict best arm  $\hat{a}_\tau$ .

## Best arm identification

$$\begin{aligned} \mu &= (\mathcal{N}(\mu_1, 1) \quad \cdots \quad \mathcal{N}(\mu_a, 1) \quad \cdots \quad \mathcal{N}(\mu_K, 1)) \\ &\sim \begin{bmatrix} \mu_1 & \cdots & \mu_a & \cdots & \mu_K \end{bmatrix} \end{aligned}$$

Optimal arm (dose):  $a_\mu^* \in \arg \min_{1 \leq a \leq K} |\mu_a - S|$

**Game:** while  $t < \tau$ :

1. Player pulls arm (dose)  $A_t \in \{1, \dots, K\}$ .
2. He gets an observation (toxicity)  $Y_t \sim \mathcal{N}(\mu_{A_t}, 1)$ .

Predict best arm  $\hat{a}_\tau$ .

**$\delta$ -correct algorithm.** ( $\mathcal{F}_t = \sigma(Y_1, \dots, Y_t)$  information available at step  $t$ )

- a **sampling rule**  $(A_t)_{t \geq 1}$ , where  $A_t$  is  $\mathcal{F}_{t-1}$ -measurable;
- a **stopping rule**  $\tau$ , stopping time with respect to the filtration  $(\mathcal{F}_t)_{t \geq 1}$ ;
- a  $\mathcal{F}_\tau$ -measurable **decision rule**  $\hat{a}_\tau$ ;

An algorithm is  **$\delta$ -correct** if  $\mathbb{P}_\mu(\hat{a}_{\tau_\delta} \neq a_\mu^*) \leq \delta$  and  $\mathbb{P}_\mu(\tau_\delta < +\infty) = 1$ .

**Goal:** find a  $\delta$ -correct algorithm that minimize  $\mathbb{E}_{\mu}[\tau_{\delta}]$ .

→ lower bound on  $\mathbb{E}_{\mu}[\tau_{\delta}]$ ?

## Lower bound

$$\mathcal{M} = \{\boldsymbol{\mu} \in \mathbb{R}^K : \mathbf{a}_{\boldsymbol{\mu}}^* \text{ is unique}\} \quad \mathcal{I} = \{\boldsymbol{\mu} \in \mathcal{M} : \mu_1 < \dots < \mu_K\}$$

Alternative set for  $\mathcal{S} \in \{\mathcal{M}, \mathcal{I}\}$ :  $\text{Alt}(\boldsymbol{\mu}, \mathcal{S}) := \{\boldsymbol{\lambda} \in \mathcal{S} : \mathbf{a}_{\boldsymbol{\lambda}}^* \neq \mathbf{a}_{\boldsymbol{\mu}}^*\}$ .

## Lower bound

$$\mathcal{M} = \{\boldsymbol{\mu} \in \mathbb{R}^K : a_{\boldsymbol{\mu}}^* \text{ is unique}\} \quad \mathcal{I} = \{\boldsymbol{\mu} \in \mathcal{M} : \mu_1 < \dots < \mu_K\}$$

Alternative set for  $\mathcal{S} \in \{\mathcal{M}, \mathcal{I}\}$ :  $\text{Alt}(\boldsymbol{\mu}, \mathcal{S}) := \{\boldsymbol{\lambda} \in \mathcal{S} : a_{\boldsymbol{\lambda}}^* \neq a_{\boldsymbol{\mu}}^*\}$ .

### Theorem

Let  $\mathcal{S} \in \{\mathcal{M}, \mathcal{I}\}$ . For all  $\delta$ -correct algorithm,

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau_{\delta}]}{\log(1/\delta)} \geq T_{\mathcal{S}}^*(\boldsymbol{\mu}),$$

where the characteristic time is

$$T_{\mathcal{S}}^*(\boldsymbol{\mu})^{-1} = \sup_{\boldsymbol{\omega} \in \Sigma_K} \inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu}, \mathcal{S})} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2},$$

where  $\Sigma_K$  is the simplex of dimension  $K - 1$ .



## Key quantities

Characteristic time:

$$T_S^*(\boldsymbol{\mu})^{-1} = \sup_{\omega \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\boldsymbol{\mu}, \mathcal{S})} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2}.$$

Optimal weights:

$$\omega^*(\boldsymbol{\mu}) := \operatorname{argmax}_{\omega \in \Sigma_K} \inf_{\lambda \in (\boldsymbol{\mu}, \mathcal{S})} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2}.$$

where  $\Sigma_K$  simplex of dimension  $K - 1$ .

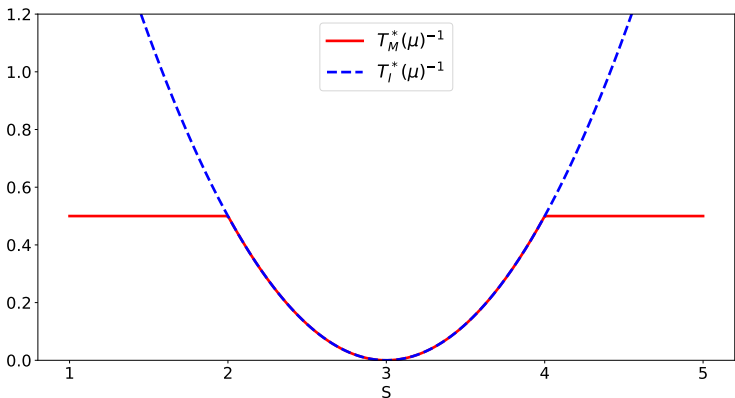
## Two arms

If  $K = 2$ ,

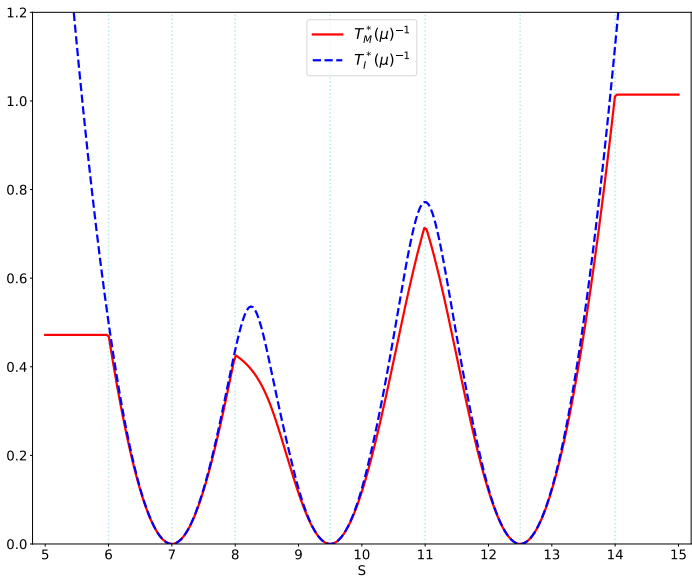
$$T_{\mathcal{I}}^*(\boldsymbol{\mu})^{-1} = (2S - \mu_1 - \mu_2)^2 / 8$$

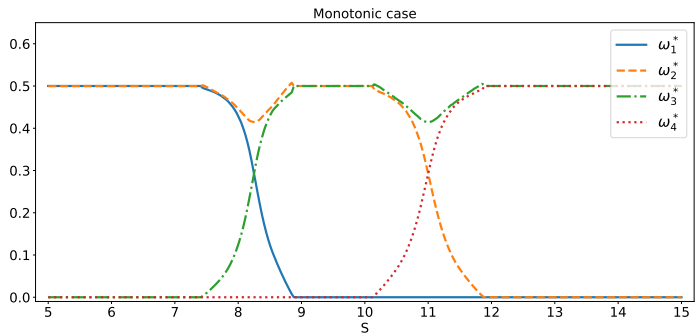
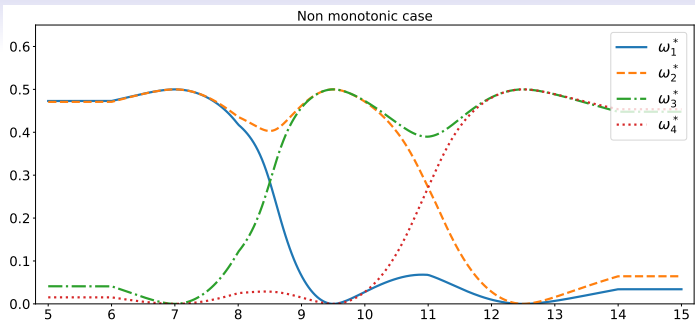
$$T_{\mathcal{M}}^*(\boldsymbol{\mu})^{-1} = \min\left((2S - \mu_1 - \mu_2)^2, (\mu_1 - \mu_2)^2\right) / 8.$$

$\boldsymbol{\mu} = [2, 4]$



K arms:  $\mu = [6, 8, 11, 14]$





# Asymptotically optimal algorithm

---

## Algorithm 3: Direct-tracking

---

### Sampling rule

$$A_{t+1} \in \begin{cases} \operatorname{argmin}_{a \in U_t} N_a(t) \text{ if one } N_a(t) \text{ "too small"} & (\text{forced exploration}) \\ \operatorname{argmax}_{1 \leq a \leq K} \omega_a^*(\hat{\mu}(t)) - N_a(t)/t & (\text{direct tracking}) \end{cases}$$

### Stopping rule

$$\tau_\delta = \inf \left\{ t \in \mathbb{N}^* : \hat{\mu}(t) \in \mathcal{M} \text{ and } \inf_{\lambda \in \operatorname{Alt}(\hat{\mu}(t), \mathcal{S})} \sum_{a=1}^K N_a(t) \frac{(\hat{\mu}_a(t) - \lambda_a)^2}{2} > \beta(t, \delta) \right\}$$

### Decision rule

$$\hat{a}_\tau \in \operatorname{argmin}_{1 \leq a \leq K} |\hat{\mu}_a(\tau) - \mathcal{S}|.$$

## Theorem (Asymptotic optimality)

For  $S \in \{\mathcal{I}, \mathcal{M}\}$ , for  $\beta(t, \delta) = \log(tC/\delta) + (3K + 2) \log \log(tC/\delta)$   
Direct-tracking is  $\delta$ -**correct on  $S$**  and **asymptotically optimal**, i.e.

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\log(1/\delta)} \leq T_S^*(\mu).$$

where

$$C := e^{K+1} \left(\frac{2}{K}\right)^K (2(3K + 2))^{3K} \frac{4}{\log(3)}.$$

## Theorem (Asymptotic optimality)

For  $\mathcal{S} \in \{\mathcal{I}, \mathcal{M}\}$ , for  $\beta(t, \delta) = \log(tC/\delta) + (3K + 2) \log \log(tC/\delta)$   
 Direct-tracking is  $\delta$ -correct on  $\mathcal{S}$  and **asymptotically optimal**, i.e.

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\log(1/\delta)} \leq T_{\mathcal{S}}^*(\mu).$$

One key tool :

## Theorem (Concentration inequality for KL divergences)

For all  $\delta \geq (K + 1)$  and  $t \in \mathbb{N}^*$  we have

$$\mathbb{P}\left(\sum_{a=1}^K N_a(t) \frac{(\hat{\mu}_a(t) - \mu_a)^2}{2} \geq \delta\right) \leq e^{K+1} \left(\frac{2\delta(\delta \log(t) + 1)}{K}\right)^K e^{-\delta}.$$

## Practical implementation

How to compute the optimal weights ?

$$\omega^*(\boldsymbol{\mu}) = \operatorname{argmax}_{\omega \in \Sigma_K} \inf_{\lambda \in \mathcal{Alt}(\boldsymbol{\mu}, \mathcal{S})} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2}.$$

**Non-monotonous case:**  $\mathcal{S} = \mathcal{M}$

→ boils down to solve one scalar equation → **Fast** !



**Monotonous case:**  $\mathcal{S} = \mathcal{I}$

$$\mathcal{I}_b := \{\lambda \in \mathcal{I}, a_\lambda^* = b\}$$

$$F : \omega \mapsto \inf_{\lambda \in \text{Alt}(\mu, \mathcal{I})} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2} = \min_{b \neq a_\mu^*} \inf_{\lambda \in \mathcal{I}_b} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2}, \quad (1)$$

$F$  is **concave**  $\rightarrow$  **sub-gradient ascent** on the simplex.

**Monotonous case:**  $\mathcal{S} = \mathcal{I}$

$$\mathcal{I}_b := \{\lambda \in \mathcal{I}, a_\lambda^* = b\}$$

$$F : \omega \mapsto \inf_{\lambda \in \text{Alt}(\mu, \mathcal{I})} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2} = \min_{b \neq a_\mu^*} \inf_{\lambda \in \mathcal{I}_b} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2}, \quad (1)$$

$F$  is **concave**  $\rightarrow$  **sub-gradient ascent** on the simplex.

Optimal alternative in  $\bar{\mathcal{I}}_b$ :

$$\lambda^b := \arg \min_{\lambda \in \bar{\mathcal{I}}_b} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2},$$

Sub-gradient of  $F$  at  $\omega$

$$\partial F(\omega) = \text{Conv}_{b \in B_{\text{Opt}}} \left[ \frac{(\mu_a - \lambda_a^b)^2}{2} \right]_{a \in \{1, \dots, K\}},$$

where  $\text{Conv}$  denotes the convex hull and  $B_{\text{Opt}}$  the set of points that attain the minimum in (1).

## How to compute $\lambda^b$ ?

$$\lambda^b := \arg \min_{\lambda \in \bar{I}_b} \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda)^2}{2},$$

~ to compute the **Unimodal regression**  $\hat{\lambda}^b$  of  $\mu'$  with mode at  $b$ :

$$\hat{\lambda}^b = \arg \min_{\substack{\lambda'_1 \leq \dots \leq \lambda'_b \\ \lambda'_K \leq \dots \leq \lambda'_b}} \sum_{a=1}^K \omega_a \frac{(\mu'_a - \lambda'_a)^2}{2}.$$

Unimodal regression can be efficiently computed via isotonic regressions (Pool Adjacent Violators Algorithm) in  $O(K)$ .

→ sub-gradient in  $O(K^2)$ .

## Open questions

- Extend to any one-exponential family ? For example

$$\boldsymbol{\mu} = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K)),$$

- Compute the sub-gradient in  $O(K)$ .
- Find algorithm with better theoretical/practical properties ("finite  $\delta$ " bound).