

Efficient Stopping Rules for Bandit Pure Exploration

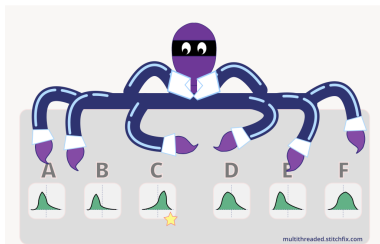
Emilie Kaufmann



Algorithmic Statistics Workshop
Oxford, November 2025

The Multi Armed Bandit (MAB) model

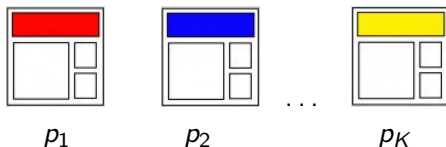
- K unknown distributions ν_1, \dots, ν_K called *arms*
- a time t , select an arm A_t and collect an observation $X_t \sim \nu_{A_t}$



Sequential strategy / algorithm : A_{t+1} can depend on:

- previous observation $A_1, X_1, \dots, A_t, X_t$
- some external randomization $U_t \sim \mathcal{U}([0, 1])$
- some knowledge about the possible distributions: $\nu_a \in \mathcal{D}$

Example: A/B/n Testing



p_a : probability that a visitor seeing version a buys a product

For the t -th visitor:

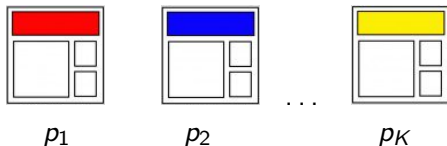
- choose a version A_t to display
- observe $X_t = 1$ if a product is bought, 0 otherwise

Objective 1: Maximizing Rewards

- observation = reward
- maximize $\mathbb{E}[\sum_{t=1}^T X_t]$ for some (possibly unknown) T

Regret minimization in bandits: UCB, Thompson Sampling...

Example: A/B/n Testing



p_a : probability that a visitor seeing version a buys a product

For the t -th visitor:

- choose a version A_t to display
- observe $X_t = 1$ if a product is bought, 0 otherwise

Objective 2: Pure Exploration

- **identify** quickly some interesting arms
- e.g. $a_\star = \arg \max_a p_a$ (best arm identification)

This talk: a generic recipe for **pure exploration**

Possible bandit models: $\boldsymbol{\nu} = (\nu_1, \dots, \nu_K) \in \mathcal{B}$

(e.g. independent sub-Gaussian arms, or Bernoulli arms)

Possible vectors of arms means $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K) \in \mathcal{M}$

Identification task

Given a **correct answer** function

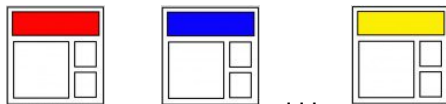
$$\begin{aligned} i_\star : \mathcal{M} &\longrightarrow \mathcal{I} \\ \boldsymbol{\mu} &\mapsto i_\star(\boldsymbol{\mu}) \end{aligned}$$

find a correct answer with high probability.

Examples of correct answers

- Best Arm Identification

[Even-Dar et al., 2006]

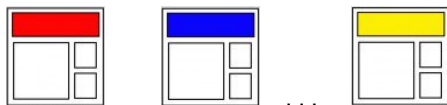


$$i_*(\mu) = \arg \max_{a \in [K]} \mu_a$$

Examples of correct answers

- Best Arm Identification

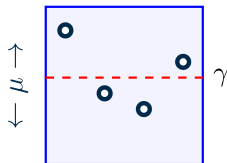
[Even-Dar et al., 2006]



$$i_*(\mu) = \arg \max_{a \in [K]} \mu_a$$

- Threshold-based questions: which means are below γ ?

[Locatelli et al., 2016]

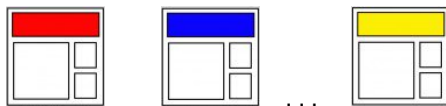


$$i_*(\mu) = (\mathbb{1}(\mu_1 > \gamma), \dots, \mathbb{1}(\mu_K > \gamma)) \in \{0, 1\}^K$$

Examples of correct answers

- Best Arm Identification

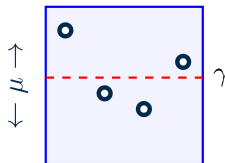
[Even-Dar et al., 2006]



$$i_*(\mu) = \arg \max_{a \in [K]} \mu_a$$

- Threshold-based questions: is there a mean below γ ?

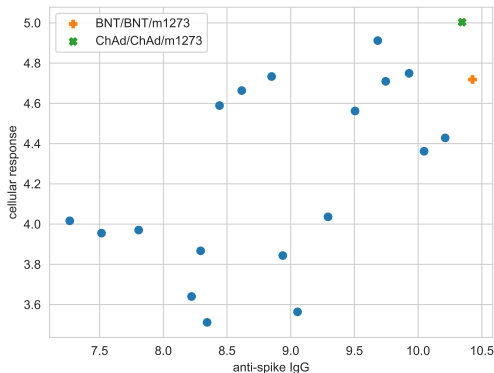
[Kaufmann et al., 2018]



$$i_*(\mu) = \mathbb{1}(\min_i \mu_i < \gamma) \in \{0, 1\}$$

Examples of correct answers

- Pareto Set Identification [Auer et al., 2016]



- arms are multi-variate distributions
- $i_{\star}(\mu)$ is the Pareto Set of the means

Pure Exploration with Fixed Confidence

An algorithm is made of:

- a **sampling rule** $A_t \in [K]$: what is the next arm to explore?
- get a new observation $X_t \sim \nu_{A_t}$
- a **recommendation rule** \hat{i}_t : a guess for the correct answer
- a **stopping rule** τ : when to stop the data collection?

Definition

An algorithm is **δ -correct** if, for all $\mu \in \mathcal{M}$, $\mathbb{P}_\mu(\hat{i}_\tau \neq i_\star(\mu)) \leq \delta$.

Goal: a δ -correct algorithm with small **sample complexity** $\mathbb{E}_\mu[\tau]$

- 1 (Optimal) Pure Exploration: A General Recipe
- 2 Best Arm Identification
- 3 Pareto Set Identification

- 1 (Optimal) Pure Exploration: A General Recepte
- 2 Best Arm Identification
- 3 Pareto Set Identification

A lower bound on the sample complexity

Setting: independent arms, parametrized by their means

$$d(\mu, \mu') := \text{KL}(\nu_\mu, \nu_{\mu'})$$

Theorem

[Garivier and Kaufmann, 2016]

For any δ -correct algorithm,

$$\mathbb{E}_\mu[\tau] \geq T^*(\mu) \ln \left(\frac{1}{3\delta} \right),$$

where

$$T^*(\mu)^{-1} = \sup_{\mathbf{w} \in \Delta_K} \inf_{\boldsymbol{\lambda} \in \text{Alt}(i_*(\mu))} \left(\sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right).$$

with

$$\begin{aligned} \Delta_K &= \left\{ \mathbf{w} \in [0, 1]^K : \sum_{i=1}^K w_i = 1 \right\} \\ \text{Alt}(i) &= \left\{ \boldsymbol{\lambda} \in \mathcal{M} : i_*(\boldsymbol{\lambda}) \neq i \right\} \end{aligned}$$

Optimal proportions

$$T^*(\mu)^{-1} = \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(i_*(\mu))} \left(\sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right).$$

- $N_a(t) = \sum_{s=1}^t \mathbb{1}(A_s = a)$: number of selections of arm a

The proof of the lower bound further suggests that the vector

$$\left(\frac{\mathbb{E}_{\mu}[N_1(\tau)]}{\mathbb{E}_{\mu}[\tau]}, \dots, \frac{\mathbb{E}_{\mu}[N_K(\tau)]}{\mathbb{E}_{\mu}[\tau]} \right)$$

should belong to

$$w^*(\mu) = \operatorname{argmax}_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(\mu)} \left(\sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right)$$

→ algorithmic strategy: let's make this happen!

The GLR stopping rule

Given a candidate best arm i , the (log) **Generalized Likelihood Ratio statistic** associated to

$$\mathcal{H}_0 = (\mu \in \text{Alt}(i)) \text{ against } \mathcal{H}_1 : (\mu \notin \text{Alt}(i))$$

is

$$\begin{aligned} Z_i(t) &= \log \frac{\sup_{\lambda \in \mathcal{M}} \ell(X_1, \dots, X_t; \lambda)}{\sup_{\lambda \in \text{Alt}(i)} \ell(X_1, \dots, X_t; \lambda)} \\ &= \inf_{\lambda \in \text{Alt}(i)} \log \frac{\ell(X_1, \dots, X_t; \hat{\mu}(t))}{\ell(X_1, \dots, X_t; \lambda)} \\ &= \inf_{\lambda \in \text{Alt}(i)} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a) \end{aligned}$$

for **exponential families** (Bernoulli, Gaussian with known variance, etc.)

Idea: stop the first time that one of the $Z_i(t)$ is large enough

A stopping rule aligned with the lower bound

GLR stopping rule

Given a threshold function $\beta(t, \delta)$:

$$\tau_\delta = \inf \left\{ t \in \mathbb{N} : \inf_{\lambda \in \text{Alt}(\hat{i}_t^*)} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a) \geq \beta(t, \delta) \right\}$$

with the recommendation rule $\hat{i}_t^* = i_*(\hat{\mu}(t))$

→ reminiscent of

$$T^*(\mu)^{-1} = \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(i_*(\mu))} \left(\sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right).$$

- if $\frac{N_a(t)}{t} \simeq w_a^*(\mu)$ and $\beta(t, \delta) \simeq \log(1/\delta)$, we get

$$\tau_\delta \simeq T_*(\mu) \log(1/\delta)$$

Converging to the optimal proportions

- Introducing $U_t = \{a : N_a(t) < \sqrt{t}\}$,

$$A_{t+1} \in \begin{cases} \operatorname{argmin}_{a \in U_t} N_a(t) & \text{if } U_t \neq \emptyset \quad (\text{forced exploration}) \\ \operatorname{argmax}_{1 \leq a \leq K} \left[w_a^*(\hat{\mu}(t)) - \frac{N_a(t)}{t} \right] & (\text{tracking}) \end{cases}$$

Lemma

Assume that

- for all $\mu \in \mathcal{M}$, $|w^*(\mu)| = 1$ (unique optimal allocation)
- $\mu \mapsto w^*(\mu)$ is continuous in all $\mu \in \mathcal{M}$

Under the **Tracking sampling rule**,

$$\mathbb{P}_\mu \left(\lim_{t \rightarrow \infty} \frac{N_a(t)}{t} = w_a^*(\mu) \right) = 1.$$

An asymptotically optimal algorithm

Theorem [Garivier and Kaufmann, 2016, Kaufmann and Koolen, 2021]

When the arm distributions belong to a one-dimensional exponential family, the Track-and-Stop strategy, that uses

- the Tracking sampling rule
- the GLR stopping rule with

$$\beta(t, \delta) \simeq \ln(1/\delta) + \ln \ln(1/\delta) + K \ln(\ln(t))$$

- and the recommendation rule $\hat{i}_t = i_*(\hat{\mu}(t))$

is δ -correct for every $\delta \in]0, 1[$ and satisfies

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\ln(1/\delta)} = T^*(\mu).$$

Calibration of the Stopping Rule

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}}(\tau < \infty, \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu})) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu}), \inf_{\lambda \in \text{Alt}(\hat{i}_{\tau}^*)} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \end{aligned}$$

Needed:

- a time uniform deviation inequality
- where the deviations are measured with KL-divergence
- and aggregated over arms

Solution: (a product of) e-processes

Calibration of the Stopping Rule

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}}(\tau < \infty, \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu})) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu}), \inf_{\lambda \in \text{Alt}(\hat{i}_{\tau}^*)} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \end{aligned}$$

Needed:

- a **time uniform** deviation inequality
- where the deviations are measured with KL-divergence
- and aggregated over arms

Solution: (a product of) e-processes

Calibration of the Stopping Rule

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}}(\tau < \infty, \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu})) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu}), \inf_{\lambda \in \text{Alt}(\hat{i}_{\tau}^*)} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \end{aligned}$$

Needed:

- a time uniform deviation inequality
- where the deviations are measured with KL-divergence
- and aggregated over arms

Solution: (a product of) e-processes

Calibration of the Stopping Rule

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}}(\tau < \infty, \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu})) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \hat{i}_{\tau}^* \neq i_{\star}(\boldsymbol{\mu}), \inf_{\lambda \in \text{Alt}(\hat{i}_{\tau}^*)} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P}_{\boldsymbol{\mu}} \left(\exists t \in \mathbb{N} : \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \end{aligned}$$

Needed:

- a time uniform deviation inequality
- where the deviations are measured with KL-divergence
- and aggregated over arms

Solution: (a product of) e-processes

Calibration of the Stopping Rule

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log(N_a(t)))$$

Step 1: $e^{\lambda X_a(t)}$ is (almost) an e-process

$$\forall \lambda \in \Lambda : M_a^\lambda(t) \geq e^{\lambda X_a(t) - g(\lambda)}$$

where $M_a^\lambda(t)$ is a **test martingale** and g a correction function.

Calibration of the Stopping Rule

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log(N_a(t)))$$

Step 1: $e^{\lambda X_a(t)}$ is (almost) an e-process

$$\forall \lambda \in \Lambda : M_a^\lambda(t) \geq e^{\lambda X_a(t) - g(\lambda)}$$

where $M_a^\lambda(t)$ is a **test martingale** and g a correction function.

① Link $X_a(t)$ to the martingale $W_a^\eta(t) = e^{\eta S_a(t) - \phi_{\mu_a}(\eta) N_a(t)}$

$$S_a(t) = \sum_{s=1}^t X_s \mathbf{1}(A_s = a) \quad \phi_{\mu_a} = \log \mathbb{E}_{X \sim \nu_{\mu_a}} [e^{\lambda X}]$$

[Robbins, 1970]

Calibration of the Stopping Rule

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log(N_a(t)))$$

Step 1: $e^{\lambda X_a(t)}$ is (almost) an e-process

$$\forall \lambda \in \Lambda : M_a^\lambda(t) \geq e^{\lambda X_a(t) - g(\lambda)}$$

where $M_a^\lambda(t)$ is a **test martingale** and g a correction function.

① Link $X_a(t)$ to the martingale $W_a^\eta(t) = e^{\eta S_a(t) - \phi_{\mu_a}(\eta) N_a(t)}$

$$S_a(t) = \sum_{s=1}^t X_s \mathbb{1}(A_s = a) \quad \phi_{\mu_a} = \log \mathbb{E}_{X \sim \nu_{\mu_a}} [e^{\lambda X}]$$

[Robbins, 1970]

For exponential families:

If $N_a(t) \in [(1 + \xi)^{i-1}, (1 + \xi)^i]$, there exists some $\eta \in \{\eta_i^\pm(x)\}$:

$$\{N_a(t)d(\hat{\mu}_a(t), \mu_a) \geq x\} \subseteq \left\{W_t^\eta(t) \geq e^{\frac{x}{1+\xi}}\right\}$$

Calibration of the Stopping Rule

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log(N_a(t)))$$

Step 1: $e^{\lambda X_a(t)}$ is (almost) an e-process

$$\forall \lambda \in \Lambda : M_a^\lambda(t) \geq e^{\lambda X_a(t) - g(\lambda)}$$

where $M_a^\lambda(t)$ is a **test martingale** and g a correction function.

① Link $X_a(t)$ to the martingale $W_a^\eta(t) = e^{\eta S_a(t) - \phi_{\mu_a}(\eta) N_a(t)}$

$$S_a(t) = \sum_{s=1}^t X_s \mathbb{1}(A_s = a) \quad \phi_{\mu_a} = \log \mathbb{E}_{X \sim \nu_{\mu_a}} [e^{\lambda X}]$$

[Robbins, 1970]

For exponential families:

The martingale $Z_a^{(x)}(t) = \sum_i \frac{c}{i^2} \left[W_a^{\eta_i^- (x)}(t) + W_a^{\eta_i^+ (x)}(t) \right]$ satisfies

$$\{X_a(t) - f(\xi) \geq x\} \subseteq \left\{ Z_a^{(x)}(t) \geq e^{\frac{x}{1+\xi}} \right\}$$

Calibration of the Stopping Rule

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log(N_a(t)))$$

Step 1: $e^{\lambda X_a(t)}$ is (almost) an e-process

$$\forall \lambda \in \Lambda : M_a^\lambda(t) \geq e^{\lambda X_a(t) - g(\lambda)}$$

where $M_a^\lambda(t)$ is a **test martingale** and g a correction function.

① Link $X_a(t)$ to the martingale $W_a^\eta(t) = e^{\eta S_a(t) - \phi_{\mu_a}(\eta) N_a(t)}$

$$S_a(t) = \sum_{s=1}^t X_s \mathbb{1}(A_s = a) \quad \phi_{\mu_a} = \log \mathbb{E}_{X \sim \nu_{\mu_a}} [e^{\lambda X}]$$

[Robbins, 1970]

For exponential families:

The martingale $Z_a^{(x)}(t) = \sum_i \frac{c}{i^2} \left[W_a^{\eta_i^- (x)}(t) + W_a^{\eta_i^+ (x)}(t) \right]$ satisfies

$$\left\{ e^{\lambda X_a(t) - \lambda f(\xi)} \geq z \right\} \subseteq \left\{ \tilde{Z}_a^{(\lambda, z)}(t) \geq 1 \right\}$$

Calibration of the Stopping Rule

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3\log(1 + \log(N_a(t)))$$

Step 1: $e^{\lambda X_a(t)}$ is (almost) an e-process

$$\forall \lambda \in \Lambda : M_a^\lambda(t) \geq e^{\lambda X_a(t) - g(\lambda)}$$

where $M_a^\lambda(t)$ is a **test martingale** and g a correction function.

① Link $X_a(t)$ to the martingale $W_a^\eta(t) = e^{\eta S_a(t) - \phi_{\mu_a}(\eta)N_a(t)}$

$$S_a(t) = \sum_{s=1}^t X_s \mathbb{1}(A_s = a) \quad \phi_{\mu_a} = \log \mathbb{E}_{X \sim \nu_{\mu_a}} [e^{\lambda X}]$$

[Robbins, 1970]

For exponential families:

The final test martingale is

$$M_a^\lambda(t) \propto 1 + \int_1^\infty \tilde{Z}_a^{(\lambda, z)}(t) dz$$

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log(N_a(t)))$$

Step 2: Product martingales

$\forall \lambda \in \Lambda : M^\lambda(t) = \prod_{a \in [K]} M_a^\lambda(t)$ is still a test martingale

→ Chernoff method + Ville's inequality

$$X_a(t) = N_a(t)d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log(N_a(t)))$$

Step 2: Product martingales

$$\forall \lambda \in \Lambda : M^\lambda(t) = \prod_{a \in [K]} M_a^\lambda(t) \text{ is still a test martingale}$$

→ Chernoff method + Ville's inequality

$$\begin{aligned} \mathbb{P} \left(\exists t \in \mathbb{N} : \sum_{a=1}^K X_a(t) > x \right) &\leq \mathbb{P} \left(\exists t \in \mathbb{N} : e^{\sum_{a=1}^K (\lambda X_a(t) - g(\lambda))} > e^{\lambda x - K g(\lambda)} \right) \\ &\leq \mathbb{P} \left(\exists t \in \mathbb{N} : M(t) > e^{\lambda x - K g(\lambda)} \right) \\ &\leq e^{-\lambda x + K g(\lambda)} \end{aligned}$$

Then optimize over λ :

$$\mathbb{P} \left(\exists t \in \mathbb{N} : \sum_{a=1}^K X_a(t) > K \min_{\lambda \in \Lambda} \frac{g(\lambda) + \log(1/\delta)/K}{\lambda} \right) \leq \delta.$$

Correctness [Kaufmann and Koolen, 2021]

When the arm distributions belong to a one-dimensional exponential family, there exists a threshold such that

$$\beta(t, \delta) \simeq \log(1/\delta) + \log \log(1/\delta) + K \log \log(t)$$

for which, $\mathbb{P}_{\mu}(\tau < \infty, \hat{i}_{\tau} \neq i_{\star}(\mu)) \leq \delta$.

(the factor K may be reduced for some particular identification tasks)

Wait! Can we actually implement it?

Track-and-Stop requires the computation in every round t of the “minimal distance”

$$\inf_{\lambda \in \text{Alt}(\hat{\mu}(t))} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a)$$

for checking the stopping rule, and

$$\arg \max_{w \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\hat{\mu}(t))} \sum_{a=1}^K N_a(t) d(\hat{\mu}_a(t), \lambda_a)$$

for the sampling rule.

- 1 (Optimal) Pure Exploration: A General Recipe
- 2 Best Arm Identification
- 3 Pareto Set Identification

$$i_*(\mu) = a_*(\mu) = \arg \max_{a \in [K]} \mu_a$$

Using that $\text{Alt}(\mu) = \bigcup_{a \neq a_*(\mu)} \{\lambda : \lambda_a > \lambda_{a_*}\}$, the minimal distance can be computed in closed-form:

$$\begin{aligned} & \inf_{\lambda \in \text{Alt}(\mu)} \sum_{i=1}^K w_i d(\mu_i, \lambda_i) \\ &= \min_{a \neq a_*} \inf_{\lambda: \lambda_a > \lambda_{a_*}} \sum_{i=1}^K w_i d(\mu_i, \lambda_i) \\ &= \min_{a \neq a_*} \min_{\lambda \in (\mu_a, \mu_{a_*})} [w_{a_*} d(\mu_{a_*}, \lambda) + w_a d(\mu_a, \lambda)] \\ &= \min_{a \neq a_*} \left[w_{a_*} d \left(\mu_{a_*}, \frac{w_{a_*} \mu_{a_*} + w_a \mu_a}{w_{a_*} + w_a} \right) + w_a d \left(\mu_a, \frac{w_{a_*} \mu_{a_*} + w_a \mu_a}{w_{a_*} + w_a} \right) \right] \end{aligned}$$

for exponential families.

Example: Gaussian bandits

For Gaussian bandits with variance σ^2 :

$$\inf_{\lambda \in \text{Alt}(\mu)} \sum_{i=1}^K w_i \frac{(\mu_i - \lambda_i)^2}{2\sigma^2} = \min_{a \neq a_*} \frac{(\mu_* - \mu_a)^2}{2\sigma^2 \left(\frac{1}{w_{a_*}} + \frac{1}{w_a} \right)}$$

hence

$$\tau_\delta = \inf \left\{ t \in \mathbb{N} : \min_{a \neq \hat{a}_t^*} \frac{(\hat{\mu}_{\hat{a}_t^*}(t) - \hat{\mu}_a(t))^2}{2\sigma^2 \left(\frac{1}{N_{\hat{a}_t^*}(t)} + \frac{1}{N_a(t)} \right)} > \beta(t, \delta) \right\}$$

Example: Gaussian bandits

For Gaussian bandits with variance σ^2 :

$$\inf_{\lambda \in \text{Alt}(\mu)} \sum_{i=1}^K w_i \frac{(\mu_i - \lambda_i)^2}{2\sigma^2} = \min_{a \neq a_*} \frac{(\mu_* - \mu_a)^2}{2\sigma^2 \left(\frac{1}{w_{a_*}} + \frac{1}{w_a} \right)}$$

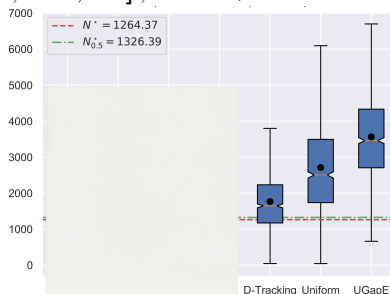
hence

$$\tau_\delta = \inf \left\{ t \in \mathbb{N} : \min_{a \neq \hat{a}_t^*} \frac{(\hat{\mu}_{\hat{a}_t^*}(t) - \hat{\mu}_a(t))^2}{2\sigma^2 \left(\frac{1}{N_{\hat{a}_t^*}(t)} + \frac{1}{N_a(t)} \right)} > \beta(t, \delta) \right\}$$

But $w_*(\mu)$ still doesn't have a closed form

- we propose an efficient approximation algorithm for exponential families in [Garivier and Kaufmann, 2016]

Empirical distribution of τ_δ for $\delta = 0.01$ for different algorithms on $\mu = [1, 0.8, 0.75, 0.7]$, $\sigma^2 = 1$, estimated on 1000 runs



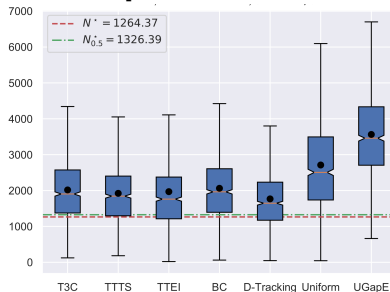
Using the right stopping rule makes a difference:

$$\text{UGapE : } \forall a \neq \hat{i}_t^* \quad , \quad \hat{\mu}_{\hat{i}_t^*}(t) - \hat{\mu}_a(t) > \sqrt{\frac{2\sigma^2\beta(t, \delta)}{N_{a,t}}} + \sqrt{\frac{2\sigma^2\beta(t, \delta)}{N_{a,t}}}$$

$$\text{GLR : } \forall a \neq \hat{i}_t^* \quad , \quad \hat{\mu}_{\hat{i}_t^*}(t) - \hat{\mu}_a(t) > \sqrt{2\sigma^2\beta(t, \delta) \left(\frac{1}{N_{a,t}} + \frac{1}{N_{a,t}} \right)}$$

Limitation: Computing w^* is costly

Empirical distribution of τ_δ for $\delta = 0.01$ for different algorithms on $\mu = [1, 0.8, 0.75, 0.7]$, $\sigma^2 = 1$, estimated on 1000 runs



Using the right stopping rule makes a difference:

$$\text{UGapE : } \forall a \neq \hat{i}_t^* \quad , \quad \hat{\mu}_{\hat{i}_t^*}(t) - \hat{\mu}_a(t) > \sqrt{\frac{2\sigma^2\beta(t, \delta)}{N_{a,t}}} + \sqrt{\frac{2\sigma^2\beta(t, \delta)}{N_{a,t}}}$$

$$\text{GLR : } \forall a \neq \hat{i}_t^* \quad , \quad \hat{\mu}_{\hat{i}_t^*}(t) - \hat{\mu}_a(t) > \sqrt{2\sigma^2\beta(t, \delta) \left(\frac{1}{N_{a,t}} + \frac{1}{N_{a,t}} \right)}$$

Efficient alternatives to Tracking exist, e.g. Top Two algorithms

- 1 (Optimal) Pure Exploration: A General Recepte
- 2 Best Arm Identification
- 3 Pareto Set Identification

Bandit model

- K arms ν_1, \dots, ν_K
- ν_k is a multi-variate distribution in \mathbb{R}^D with mean $\mu_k \in \mathbb{R}^D$
- Assumption: each marginal of ν_k is *sub-Gaussian*

In each round t , an agent selects an arm $A_t \in [K]$ and observes a response $\mathbf{X}_t \sim \nu_{A_t}$, independently from past observations.

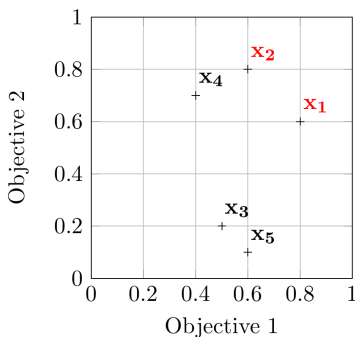
- What is a “good set of arms”?
a possibility: the **Pareto Set** of the mean vectors

$$\mathcal{S}^*(\mu) = \{ \text{“all arms that are not uniformly worse than any other arm”} \}$$

Pareto Set

Let $\mathcal{X} \subset \mathbb{R}^D$ a set of vectors. Let $\mathbf{x}, \mathbf{y} \in \mathcal{X}$.

- \mathbf{x} is (strictly) dominated by \mathbf{y} ($\mathbf{x} \prec \mathbf{y}$) if $\forall d \in [D], x^d < y^d$
- The Pareto Set is
$$\mathcal{P}(\mathcal{X}) := \{\mathbf{x} \in \mathcal{X} : \nexists \mathbf{y} \in \mathcal{X} \text{ such that } \mathbf{x} \prec \mathbf{y}\}$$
- A vector $\mathbf{x} \in \mathcal{P}(\mathcal{X})$ is called Pareto optimal



1 $\mathbf{x}_3 \prec \mathbf{x}_1$

2 $\mathbf{x}_4 \prec \mathbf{x}_2$

3 $\mathbf{x}_5 \prec \mathbf{x}_1$

4 $\mathbf{x}_1 \not\prec \mathbf{x}_2$

5 $\mathbf{x}_2 \not\prec \mathbf{x}_1$

$$\mathcal{P}(\mathcal{X}) = \{\mathbf{x}_1, \mathbf{x}_2\}$$

Pareto Set Identification with Fixed Confidence

$$\begin{aligned}\boldsymbol{\mu} &= (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K) \in (\mathbb{R}^D)^K \\ \mathcal{S}^*(\boldsymbol{\mu}) &= \{k \in [K] : \mu_k \in \mathcal{P}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)\}\end{aligned}$$

Pareto Set Identification algorithm:

- a **sampling rule** $A_t \in [K]$: what is the next arm to explore?
- get a new observation $\mathbf{X}_t \sim \nu_{A_t} \in \mathbb{R}^D$
- a **recommendation rule** $\hat{\mathcal{S}}_t$: a guess for $\mathcal{S}^*(\boldsymbol{\mu})$
- a **stopping rule** τ : when to stop the data collection?

Definition

An algorithm is **δ -correct** (on \mathcal{M}) if, for all $\boldsymbol{\nu} \in \mathcal{M}$,
 $\mathbb{P}_{\boldsymbol{\nu}}(\hat{\mathcal{S}}_{\tau} \neq \mathcal{S}^*(\boldsymbol{\mu})) \leq \delta$.

Goal: a δ -correct algorithm with small **sample complexity** $\mathbb{E}_{\boldsymbol{\nu}}[\tau]$

Theorem

For arms that are multi-variate Gaussian (known covariance Σ), any δ -correct algorithm for Pareto Set Identification satisfies, for all $\mu \in (\mathbb{R}^D)^K$,

$$\mathbb{E}_{\mu}[\tau_{\delta}] \geq T^*(\mu) \log \left(\frac{1}{3\delta} \right)$$

where

$$T^*(\mu)^{-1} = \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(\mathcal{S}^*(\mu))} \left(\sum_{k=1}^K w_k \text{KL}(\mathcal{N}(\mu_a, \Sigma), \mathcal{N}(\lambda_a, \Sigma)) \right).$$

with $\text{Alt}(\mathcal{S}) = \{\lambda \in (\mathbb{R}^D)^K : \mathcal{S}^*(\lambda) \neq \mathcal{S}\}$.

Theorem

For arms that are multi-variate Gaussian (known covariance Σ), any δ -correct algorithm for Pareto Set Identification satisfies, for all $\mu \in (\mathbb{R}^D)^K$,

$$\mathbb{E}_{\mu}[\tau_{\delta}] \geq T^*(\mu) \log \left(\frac{1}{3\delta} \right)$$

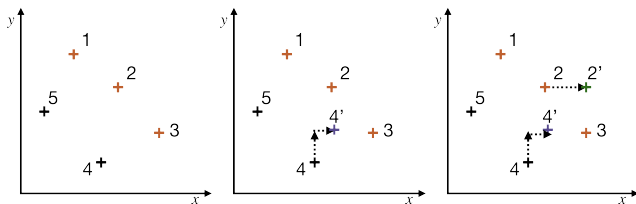
where

$$T^*(\mu)^{-1} = \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(\mathcal{S}^*(\mu))} \left(\sum_{k=1}^K w_k \frac{1}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right).$$

with $\text{Alt}(\mathcal{S}) = \{\lambda \in (\mathbb{R}^D)^K : \mathcal{S}^*(\lambda) \neq \mathcal{S}\}$.

Computing the Minimal Distance

- there are many ways to alter the Pareto set



- no closed-form is known for the minimal distance

$$(1) : w \mapsto \inf_{\lambda \in \text{Alt}(\mathcal{S})} \sum_k \frac{w_k}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2$$

- for $\Sigma = \sigma^2 I_d$, (1) can be computed by solving $O(K|S^*(\mu)|^d)$ separably convex problems [Crepon et al., 2024]

Track-And-Stop?

The GLR stopping rule

$$\tau = \inf \left\{ t \in \mathbb{N} : \inf_{\lambda \in \text{Alt}(\hat{S}_t^*)} \sum_{k=1}^K \frac{N_k(t)}{2} \|\hat{\mu}_k(t) - \lambda_k\|_{\Sigma^{-1}}^2 > \beta(t, \delta) \right\}$$

can be calibrated to attain correctness with

$$\beta(t, \delta) \simeq \log(1/\delta) + \log \log(1/\delta) + KD \log \log(Dt)$$

... but is computationally expansive due to the **minimal distance**.

The Tracking sampling rule is intractable as it further computes

$$w_\star(\mu) = \arg \max_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(S^\star(\mu))} \sum_k \frac{w_k}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2$$

→ existing alternative approaches based on online learning
e.g. [Ménard, 2019] also rely on **minimal distance** computation.

The Posterior (Re)Sampling Stopping Rule

PS Stopping rule

For all $m \leq M(t, \delta)$, sample $\tilde{\theta}^m = (\tilde{\theta}_1^m, \dots, \tilde{\theta}_K^m)$ with

$$\tilde{\theta}_a^m \sim \mathcal{N}\left(\hat{\mu}_a(t), \frac{c(t, \delta)}{N_a(t)} \Sigma\right)$$

If for all m , $\mathcal{S}^*(\tilde{\theta}^m) = \mathcal{S}^*(\hat{\mu}(t))$, **stop** and recommend $\mathcal{S}^*(\hat{\mu}(t))$

- inspired by the TS-Explore strategy for Combinatorial bandits [Wang and Zhu, 2022]
- analyzed in [Kone et al., 2025] for PSI, together with a tractable sampling rule giving asymptotic optimality

Let $\beta(t, \delta)$ be such that \mathcal{E}_δ holds w.p. at least $1 - \delta$:

$$\mathcal{E}_\delta = \bigcap_{t \geq 1} \underbrace{\left(\sum_k N_{t,k} \|\mu_k - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 < 2\beta(t, \delta) \right)}_{\mathcal{E}_\delta^t}$$

Then,

$$\begin{aligned} \mathbb{P}_\nu(\tau < \infty, \hat{\mathcal{S}}_\tau \neq \mathcal{S}^*) &\leq \delta/2 + \mathbb{P}_\nu(\tau < \infty \text{ and } \hat{\mathcal{S}}_\tau \neq \mathcal{S}^*, \mathcal{E}_{\delta/2}) \\ &\leq \delta/2 + \sum_{t \geq 1} \mathbb{P}_\nu(\tau = t \text{ and } \hat{\mathcal{S}}_t \neq \mathcal{S}^*, \mathcal{E}_{\delta/2}^t) \\ &= \delta/2 + \sum_{t \geq 1} \mathbb{E}_\nu \left[\mathbb{1}_{\hat{\mathcal{S}}_t \neq \mathcal{S}^*} \mathbb{1}_{\mathcal{E}_{\delta/2}^t} \mathbb{P}_\nu(\tau = t \mid \mathcal{H}_{t-1}) \right] \end{aligned}$$

$$\begin{aligned}\mathbb{P}_{\nu}(\tau = t \mid \mathcal{H}_{t-1}) &\leq \mathbb{P}_{\nu}(\forall m \leq M(t, \delta), \mathcal{S}^*(\tilde{\theta}_t^m) = \hat{\mathcal{S}}_t \mid \mathcal{H}_{t-1}) \\ &= (1 - \mathbb{P}_{\nu}(\mathcal{S}^*(\tilde{\theta}_t^1) \neq \hat{\mathcal{S}}_t \mid \mathcal{H}_{t-1}))^{M(t, \delta)}, \\ &= (1 - \Pi_t(\text{Alt}(\hat{\mathcal{S}}_t)))^{M(t, \delta)} \\ &\leq \exp\left(-\Pi_t(\text{Alt}(\hat{\mathcal{S}}_t))M(t, \delta)\right)\end{aligned}$$

hence the error probability is bounded by

$$\frac{\delta}{2} + \sum_{t \geq 1} \mathbb{E}_{\nu} \left[\mathbb{1}_{\hat{\mathcal{S}}_t \neq \mathcal{S}^*(\mu)} \mathbb{1}_{\mathcal{E}_{\delta/2}^t} \exp\left(-\Pi_t(\text{Alt}(\hat{\mathcal{S}}_t))M(t, \delta)\right) \right].$$

The tricky part of the proof is then to get a lower bound on $\Pi_t(\text{Alt}(\hat{\mathcal{S}}_t))$ (Gaussian anti-concentration)

$$\begin{aligned}\mathbb{P}_{\nu}(\tau = t \mid \mathcal{H}_{t-1}) &\leq \mathbb{P}_{\nu}(\forall m \leq M(t, \delta), \mathcal{S}^*(\tilde{\theta}_t^m) = \hat{\mathcal{S}}_t \mid \mathcal{H}_{t-1}) \\ &= (1 - \mathbb{P}_{\nu}(\mathcal{S}^*(\tilde{\theta}_t^1) \neq \hat{\mathcal{S}}_t \mid \mathcal{H}_{t-1}))^{M(t, \delta)}, \\ &= (1 - \Pi_t(\text{Alt}(\hat{\mathcal{S}}_t)))^{M(t, \delta)} \\ &\leq \exp\left(-\Pi_t(\text{Alt}(\hat{\mathcal{S}}_t))M(t, \delta)\right)\end{aligned}$$

hence the error probability is bounded by

$$\frac{\delta}{2} + \sum_{t \geq 1} \mathbb{E}_{\nu} \left[\mathbb{1}_{\hat{\mathcal{S}}_t \neq \mathcal{S}^*(\mu)} \mathbb{1}_{\mathcal{E}_{\delta/2}^t} \exp\left(-\Pi_t(\text{Alt}(\hat{\mathcal{S}}_t))M(t, \delta)\right) \right].$$

The tricky part of the proof is then to get a lower bound on $\Pi_t(\text{Alt}(\hat{\mathcal{S}}_t))$ (Gaussian anti-concentration)

Lemma [Kone et al., 2025]

The PS Stopping rule is δ correct for

$$c(t, \delta) \simeq \frac{\log(\log(t)/\delta)}{\log(1/\delta)} \text{ and } M(t, \delta) \simeq \frac{\log(t/\delta)}{\delta}$$

Two generic stopping rules for pure exploration tasks in bandits:

- the GLR stopping rule that is easy to calibrate
(for exponential families)
- the PS stopping rule that can be easier to compute
(but harder to calibrate)

In these approaches, ϵ -processes are hidden in the proofs... but the resulting calibration are a bit conservative in practice.

- ➔ Tighter calibrations?
- ➔ When is PS “better” than GLR?

- Aurélien Garivier, Emilie Kaufmann
Optimal Best Arm Identification with Fixed Confidence
(COLT 2016)
- Emilie Kaufmann, Wouter M. Koolen
Mixture Martingales Revisited with Applications to Sequential Tests and Confidence Intervals (JMLR 2021)



- Marc Jourdan, Cyrille Koné, Emilie Kaufmann
Pareto Set Identification with Posterior Sampling. (AISTATS 2025)



-  Auer, P., Chiang, C., Ortner, R., and Drugan, M. M. (2016). Pareto front identification from stochastic bandit feedback. In *AISTATS*.
-  Crepon, É., Garivier, A., and Koolen, W. M. (2024). Sequential learning of the pareto front for multi-objective bandits. In *AISTATS*.
-  Even-Dar, E., Mannor, S., and Mansour, Y. (2006). Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems. *Journal of Machine Learning Research*, 7:1079–1105.
-  Garivier, A. and Kaufmann, E. (2016). Optimal best arm identification with fixed confidence. In *Proceedings of the 29th Conference On Learning Theory*.
-  Kaufmann, E. and Koolen, W. (2021). Mixture martingales revisited with applications to sequential tests and confidence intervals. *Journal of Machine Learning Research*, 22(246).
-  Kaufmann, E., Koolen, W., and Garivier, A. (2018). Sequential test for the lowest mean: From Thompson to Murphy Sampling. In *Advances in Neural Information Processing Systems (NeurIPS)*.
-  Kone, C., Jourdan, M., and Kaufmann, E. (2025). Pareto set identification with posterior sampling. In *AISTATS*.
-  Locatelli, A., Gutzeit, M., and Carpentier, A. (2016).

An optimal algorithm for the thresholding bandit problem.

In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 1690–1698.



Ménard, P. (2019).

Gradient ascent for active exploration in bandit problems.

arXiv 1905.08165.



Robbins, H. (1970).

Statistical Methods Related to the law of the iterated logarithm.

Annals of Mathematical Statistics, 41(5):1397–1409.



Wang, S. and Zhu, J. (2022).

Thompson sampling for (combinatorial) pure exploration.

In *International Conference on Machine Learning*. PMLR.

On the effect of correlation

We evaluate the performance of PSIPS on a 5-arm, 2-dimensional Gaussian instance with **correlated objectives**.

- Covariance matrix: Σ_ρ with unit variances and correlation $\rho \in (-1, 1)$.
- $\rho = 0$: objectives are independent.
- $\rho \rightarrow +1$ (resp. $\rho \rightarrow -1$): strongly positively (resp. negatively) correlated objectives.

