

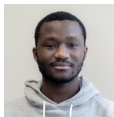
Multi-Objective Bandits Revisited

Emilie Kaufmann



based on collaborations with

Cyrille Koné, Laura Richert and Marc Jourdan



Statistics Seminar, ENSAE, May 2025

Bandits for adaptive clinical trials?



$B(p_1)$



$B(p_2)$



$B(p_3)$



$B(p_4)$



$B(p_5)$

For the t -th patient in a clinical trial,

- choose a **treatment (arm)** A_t
- observe its **efficacy (reward/response)**

$$X_t \in \{0, 1\} : \mathbb{P}(X_t = 1 | A_t = a) = p_a$$

Adaptive treatment allocation / sampling rule:

A_t can be chosen based on past outcomes $A_1, X_1, \dots, A_{t-1}, X_{t-1}$

Bandits for adaptive clinical trials?



$B(p_1)$



$B(p_2)$



$B(p_3)$



$B(p_4)$



$B(p_5)$

For the t -th patient in a clinical trial,

- choose a **treatment (arm)** A_t
- observe its **efficacy (reward/response)**

$$X_t \in \{0, 1\} : \mathbb{P}(X_t = 1 | A_t = a) = p_a$$

Adaptive treatment allocation / sampling rule:

A_t can be chosen based on past outcomes $A_1, X_1, \dots, A_{t-1}, X_{t-1}$

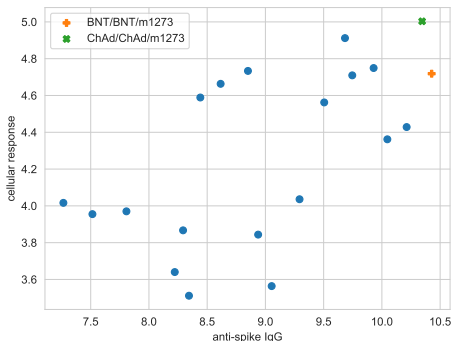
→ an idealized model for a *Phase III* trial

Specificities of early stage (*Phase I/II*) trials

Multiple responses are typically measured:

- side effects (toxicity)
- different indicators of biological efficacy (blood tests)

Vaccine design: different indicators of the immune response:



- binding antibodies
- neutralising antibodies for different variants
- cellular responses (T-cells ...)

$K = 20$ combinations of Covid vaccines (COVBOOST)

- 1 Pure Exploration in Multi-objective bandits
- 2 Best Arm Identification ($D = 1$)
- 3 Adaptive Pareto Exploration
- 4 Towards Optimal Algorithms

- 1 Pure Exploration in Multi-objective bandits
- 2 Best Arm Identification ($D = 1$)
- 3 Adaptive Pareto Exploration
- 4 Towards Optimal Algorithms

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.

Bandit (Pure Exploration) Algorithm

- (*sampling rule*) how is A_t selected based on past observation?
 - (*recommendation rule*) guess \hat{S}_t for a “good set of arms”
 - (*stopping rule*) decide whether to stop collecting observations
- Goal: make a confident guess with few samples

What is a good set of arms?

$$\mathcal{S}^* = \mathcal{S}^*(\mu_1, \dots, \mu_K) \subseteq [K]$$

- $k_* = \arg \max_k g(\mu_k)$ for some preference function $g : \mathbb{R}^D \rightarrow \mathbb{R}$, e.g. $g(\mu_k) = \sum_{d=1}^D w_d \mu_k^d$
- Feasible Set: all arms that satisfy some linear constraints [Katz-Samuels and Scott, 2018]
- Top Feasible Arm: a feasible arm maximizing one of the objectives [Katz-Samuels and Scott, 2019]
- All the arms that are not uniformly worse than the others
→ the **Pareto set** [Auer et al., 2016]

Pareto Set

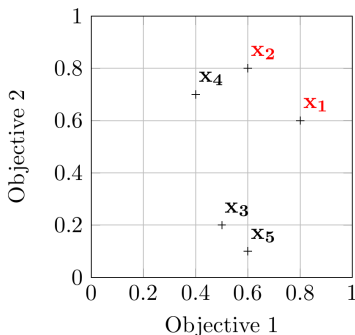
Let $\mathcal{X} \subset \mathbb{R}^D$ a set of vectors. Let $\mathbf{x}, \mathbf{y} \in \mathcal{X}$.

- \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

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

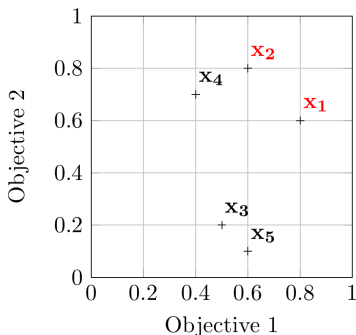
2 $\mathbf{x}_4 \prec \mathbf{x}_2$

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

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

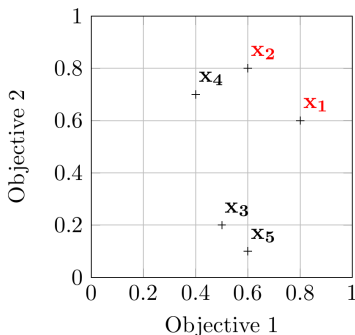
4 $\mathbf{x}_1 \not\prec \mathbf{x}_2$

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

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{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 $\boldsymbol{\nu} \in \mathcal{M}$,
 $\mathbb{P}_{\boldsymbol{\nu}}(\hat{S}_{\tau} \neq \mathcal{S}^*(\boldsymbol{\mu})) \leq \delta$.

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

- 1 Pure Exploration in Multi-objective bandits
- 2 Best Arm Identification ($D = 1$)
- 3 Adaptive Pareto Exploration
- 4 Towards Optimal Algorithms

Best Arm Identification with Fixed Confidence

$$\begin{aligned}\boldsymbol{\mu} &= (\mu_1, \dots, \mu_K) \in \mathbb{R}^K \\ i_*(\boldsymbol{\mu}) &= \arg \max_{k \in [K]} \mu_k\end{aligned}$$

Best Arm Identification algorithm:

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

Definition

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

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

3 approaches to Best Arm Identification

- Uniform sampling + **Eliminations**

Successive Eliminations [Even-Dar et al., 2006]

- Adaptive sampling based on **Confidence Intervals**

LUCB [Kalyanakrishnan et al., 2012], UGapE [Gabillon et al., 2012] ...

- **Lower Bound Inspired Algorithms**

e.g., [Garivier and Kaufmann, 2016, Degenne et al., 2019, Jourdan et al., 2022]

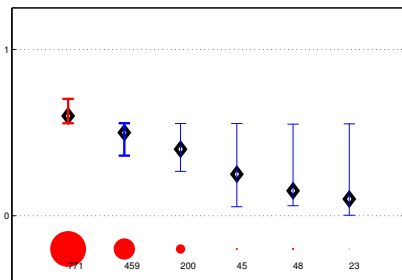
All algorithms rely on

$$N_k(t) := \sum_{s=1}^t \mathbb{1}(A_t = k), \quad \hat{\mu}_k(t) := \frac{1}{N_k(t)} \sum_{s=1}^t Y_{k,s}$$

where $(Y_{k,s})$ are the successive observations from arm k

LUCB: Lower and Upper Confidence Bounds

$$\mathcal{I}_k(t) = [\text{LCB}_k(t), \text{UCB}_k(t)].$$



- At round t , draw

$$B_t = \arg \max_{b \in [K]} \hat{\mu}_b(t)$$

$$C_t = \arg \max_{c \neq B_t} \text{UCB}_c(t)$$

- Stop at round t if

$$\text{LCB}_{B_t}(t) > \text{UCB}_{C_t}(t)$$

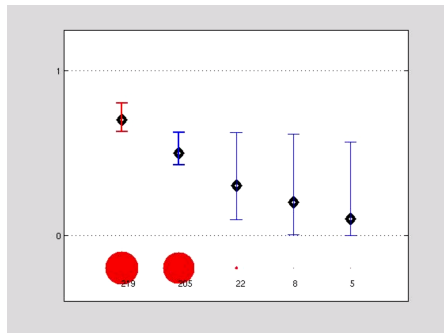
Theorem [Kalyanakrishnan et al., 2012]

For well-chosen confidence intervals, $\mathbb{P}_\nu(B_\tau = i_\star(\mu)) \geq 1 - \delta$ and

$$\mathbb{E}[\tau_\delta] = \mathcal{O} \left(\left[\sum_{a=1}^K \frac{1}{\Delta_a^2} \right] \ln \left(\frac{1}{\delta} \right) \right) \quad \Delta_k = \begin{cases} \mu_\star - \mu_k, & k \neq i_\star \\ \min_{i \neq i_\star} \Delta_i, & k = i_\star \end{cases}$$

LUCB: Lower and Upper Confidence Bounds

$$\mathcal{I}_k(t) = [\text{LCB}_k(t), \text{UCB}_k(t)].$$



- At round t , draw

$$B_t = \arg \max_{b \in [K]} \hat{\mu}_b(t)$$

$$C_t = \arg \max_{c \neq B_t} \text{UCB}_c(t)$$

- Stop at round t if

$$\text{LCB}_{B_t}(t) > \text{UCB}_{C_t}(t)$$

Theorem [Kalyanakrishnan et al., 2012]

For well-chosen confidence intervals, $\mathbb{P}_\nu(B_\tau = i_\star(\mu)) \geq 1 - \delta$ and

$$\mathbb{E}[\tau_\delta] = \mathcal{O} \left(\left[\sum_{a=1}^K \frac{1}{\Delta_a^2} \right] \ln \left(\frac{1}{\delta} \right) \right) \quad \Delta_k = \begin{cases} \mu_\star - \mu_k, & k \neq i_\star \\ \min_{i \neq i_\star} \Delta_i, & k = i_\star \end{cases}$$

A Sample Complexity Lower Bound

Lower Bound [Garivier and Kaufmann, 2016]

For δ -correct algorithms for Gaussian bandits of variance σ^2 ,

$$\mathbb{E}_{\mu}[\tau] \geq T_{\star}(\mu) \log \left(\frac{1}{3\delta} \right)$$

with

$$(T_{\star}(\mu))^{-1} = \sup_{\mathbf{w} \in \Delta_K} \inf_{\boldsymbol{\lambda} \in \text{Alt}(i_{\star}(\mu))} \sum_{a \in [K]} w_a \frac{(\mu_a - \lambda_a)^2}{2\sigma^2}$$

$$\Delta_K = \{\mathbf{w} \in [0, 1]^K : \sum_a w_a = 1\} \text{ and}$$
$$\text{Alt}(i) = \{\boldsymbol{\lambda} \in \mathbb{R}^K : i_{\star}(\boldsymbol{\lambda}) \neq i\}.$$

Proof. Information theoretic argument

For all ν' : $i_{\star}(\nu') \neq i_{\star}(\nu)$, for any δ -correct algorithm,

$$\sum_{a \in [K]} \mathbb{E}_{\nu}[N_a(\tau)] \text{KL}(\nu_a, \nu'_a) \geq \log \left(\frac{1}{3\delta} \right)$$

A Sample Complexity Lower Bound

The “minimal distance” has a closed form:

$$\inf_{\lambda \in \text{Alt}(i_*(\mu))} \sum_{a \in [K]} w_a \frac{(\mu_a - \lambda_a)^2}{2\sigma^2} = \min_{a \neq i_*} \frac{(\mu_a - \mu_{i_*})^2}{2\sigma^2 \left(\frac{1}{w_a} + \frac{1}{w_{i_*}} \right)}$$

but not the characteristic time

$$(T_*(\mu))^{-1} = \sup_{w \in \Delta_K} \min_{a \neq i_*} \frac{(\mu_a - \mu_{i_*})^2}{2\sigma^2 \left(\frac{1}{w_a} + \frac{1}{w_{i_*}} \right)}$$

Approximation of the characteristic time

$$\sum_{a=1}^K \frac{2\sigma^2}{\Delta_a^2} \leq T_*(\mu) \leq 2 \left(\sum_{a=1}^K \frac{2\sigma^2}{\Delta_a^2} \right)$$

→ Can we still match this (non-explicit) lower bound?

$$(T_{\star}(\mu))^{-1} = \sup_{w \in \Delta_K} \min_{a \neq i_{\star}} \frac{(\mu_a - \mu_{i_{\star}})^2}{2\sigma^2 \left(\frac{1}{w_a} + \frac{1}{w_{i_{\star}}} \right)}$$

Yes, with an appropriate stopping rule

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

where \hat{i}_t^* is the empirical best arm at time t

$$(T_{\star}(\boldsymbol{\mu}))^{-1} = \sup_{w \in \Delta_K} \min_{a \neq i_{\star}} \frac{(\mu_a - \mu_{i_{\star}})^2}{2\sigma^2 \left(\frac{1}{w_a} + \frac{1}{w_{i_{\star}}} \right)}$$

Yes, with an appropriate **GLR stopping rule**

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

where \hat{i}_t^{\star} is the empirical best arm at time t

→ **Generalized Likelihood Ratio** Statistic for testing

$$\mathcal{H}_0 : (i_{\star}(\boldsymbol{\mu}) \neq \hat{i}_t) \text{ against } \mathcal{H}_1 : (i_{\star}(\boldsymbol{\mu}) = \hat{i}_t)$$

$$(T_{\star}(\mu))^{-1} = \sup_{w \in \Delta_K} \min_{a \neq i_{\star}} \frac{(\mu_a - \mu_{i_{\star}})^2}{2\sigma^2 \left(\frac{1}{w_a} + \frac{1}{w_{i_{\star}}} \right)}$$

Yes, with an appropriate GLR stopping rule

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

where \hat{i}_t^* is the empirical best arm at time t
... and a sampling rule satisfying

$$\left(\frac{N_1(t)}{t}, \dots, \frac{N_K(t)}{t} \right) \rightarrow w^*(\mu)$$

where $w^*(\mu)$ is the maximizer in $w \in \Delta_K$

Tracking sampling rule: letting $U_t = \{a : N_a(t) < \sqrt{t}\}$,

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

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

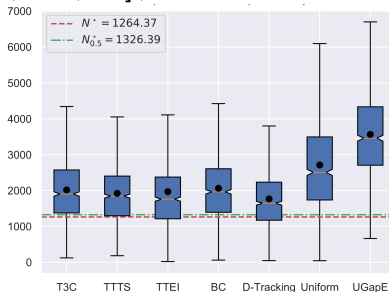
The Track-and-Stop strategy, that uses

- the **Tracking sampling rule**
- the **GLR stopping rule** with $\beta(t, \delta) \simeq \log \left(\frac{K \log(t)}{\delta} \right)$
- and recommends $\hat{i}_t = i_*(\hat{\mu}(t))$

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

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

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:

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

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

Limitation: Higher computational cost for TaS (w_\star)

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

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

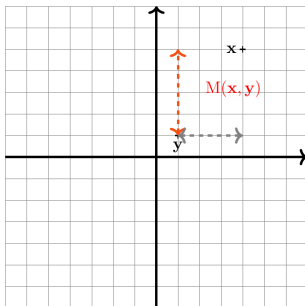
An algorithm is **δ -correct** if, for all $\boldsymbol{\nu}$, $\mathbb{P}_{\boldsymbol{\nu}}(\hat{\mathcal{S}}_{\tau} \neq \mathcal{S}^*(\boldsymbol{\mu})) \leq \delta$.

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

- 1 Pure Exploration in Multi-objective bandits
- 2 Best Arm Identification ($D = 1$)
- 3 Adaptive Pareto Exploration**
- 4 Towards Optimal Algorithms

A non-dominance measure

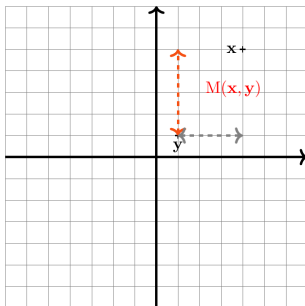
$$\begin{aligned}x \not\prec y &\Leftrightarrow \exists d, x^d \geq y^d, \\&\Leftrightarrow \exists d, x^d - y^d \geq 0, \\&\Leftrightarrow \underbrace{\max_{d \in [D]} (x^d - y^d)}_{:= M(x,y)} > 0,\end{aligned}$$



Interpretation: The larger $M(x,y)$ the “further” y is from dominating x

A non-dominance measure

$$\begin{aligned}x \not\prec y &\Leftrightarrow \exists d, x^d \geq y^d, \\&\Leftrightarrow \exists d, x^d - y^d \geq 0, \\&\Leftrightarrow \underbrace{\max_{d \in [D]} (x^d - y^d)}_{:= M(\mathbf{x}, \mathbf{y})} > 0,\end{aligned}$$



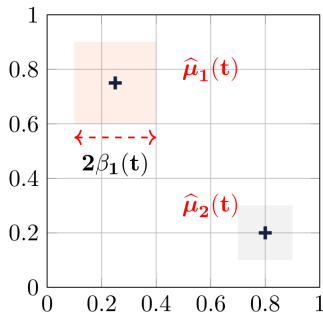
Interpretation: The larger $M(\mathbf{x}, \mathbf{y})$ the “further” \mathbf{y} is from dominating \mathbf{x}

$$M(i, j) := M(\mu_i, \mu_j)$$

Confidence Regions on $M(i, j)$

$\hat{\mu}_k(t) \in \mathbb{R}^D$ the empirical mean vector of arm k at time t

$$M(i, j; t) = M(\hat{\mu}_i(t), \hat{\mu}_j(t))$$



Confidence bonus for μ_k

$$\beta_k(t) \simeq \sqrt{2\sigma^2 \log\left(\frac{K \log(N_k(t))}{\delta}\right) \frac{1}{N_k(t)}}$$

and for $\mu_i - \mu_j$

$$\beta_{i,j}(t) \simeq \sqrt{2\sigma^2 \log\left(\frac{K^2 \log(N_k(t))}{\delta}\right) \left(\frac{1}{N_i(t)} + \frac{1}{N_j(t)}\right)}$$

Lemma

With probability $1 - \delta$, for all i, j, t ,

$$M(i, j) \geq M^-(i, j; t) := M(i, j; t) - \beta_{i,j}(t)$$

$$M(i, j) \leq M^+(i, j; t) := M(i, j; t) + \beta_{i,j}(t)$$

$$\text{OPT}(t) := \{i \in [K] : \forall j \in [K] \setminus \{i\}, M^-(i, j; t) > 0\}$$

Two interesting arms to explore:

- a potentially Pareto optimal arm

$$B_t = \arg \max_{i \in [K] \setminus \text{OPT}(t)} \min_{j \neq i} M^+(i, j; t)$$

- the arm that is the closest to potentially dominate it

$$C_t := \arg \min_{j \neq B_t} M^-(B_t, j; t)$$

Adaptive Pareto Exploration (APE)

selects the least sampled among these two candidate arms:

$$A_{t+1} = \arg \min_{a \in \{B_t, C_t\}} N_a(t)$$

Stopping rule

Letting $\hat{S}(t) = \mathcal{P}^*(\hat{\mu}_1(t), \dots, \hat{\mu}_K(t))$, the algorithm stops and recommends $\hat{S}_t = \hat{S}(t)$ when

- all arms in $\hat{S}(t)$ are confidently non-dominated:

$$Z_1(t) := \min_{i \in \hat{S}(t)} \min_{j \neq i} M^-(i, j; t) > 0$$

- all arms in $(\hat{S}(t))^c$ are confidently dominated:

$$Z_2(t) := \min_{i \notin \hat{S}(t)} \max_{j \neq i} [-M^+(i, j; t)] > 0$$

Stopping rule for (exact) PSI

$$\tau = \inf \left\{ t \in \mathbb{N} : Z_1(t) > 0, Z_2(t) > 0 \right\}$$

Stopping rule

Letting $\hat{S}(t) = \mathcal{P}^*(\hat{\mu}_1(t), \dots, \hat{\mu}_K(t))$, the algorithm stops and recommends $\hat{S}_t = \hat{S}(t)$ when

- all arms in $\hat{S}(t)$ are confidently non-dominated:

$$Z_1^\delta(t) := \min_{i \in \hat{S}(t)} \min_{j \neq i} M_\delta^-(i, j; t) > 0$$

- all arms in $(\hat{S}(t))^c$ are confidently dominated:

$$Z_2^\delta(t) := \min_{i \notin \hat{S}(t)} \max_{j \neq i} [-M_\delta^+(i, j; t)] > 0$$

Stopping rule for (exact) PSI

$$\tau_\delta = \inf \left\{ t \in \mathbb{N} : Z_1^\delta(t) > 0, Z_2^\delta(t) > 0 \right\}$$

Theorem [Kone et al., 2023]

Assume the observations are bounded in $[0, 1]^D$. Then, with probability larger than $1 - \delta$, APE with the stopping rule τ_δ outputs $\hat{S}_\tau = \mathcal{S}^*(\mu)$ and satisfies

$$\tau_\delta \leq \sum_{a=1}^K \frac{32}{\Delta_a^2} \log \left(\frac{2KD}{\delta} \log \left(\frac{32}{\Delta_a^2} \right) \right),$$

for an appropriate notion of “Pareto gap”.

- ➔ same scaling as the bound of [Auer et al., 2016] for an elimination-based algorithm, with better constants and a $\log \log(1/\Delta)$ versus $\log(1/\Delta)$

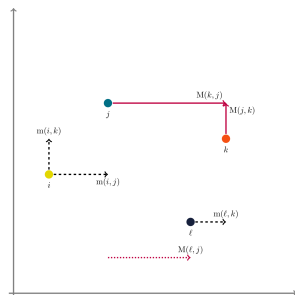
For sub-optimal arms $i \notin \mathcal{S}^*(\mu)$,

$$\Delta_i := \max_{j \in \mathcal{S}^*} m(i, j), \quad m(i, j) = -M(j, i)$$

while for optimal arms $i \in \mathcal{S}^*$, $\Delta_i = \min(\delta_i^+, \delta_i^-)$ where

$$\delta_i^+ := \min_{j \in \mathcal{S}^* \setminus \{i\}} \min(M(i, j), M(j, i))$$

$$\delta_i^- := \min_{j \in [K] \setminus \mathcal{S}^*} \{[M(j, i)]_+ + \Delta_j\}$$



APE can further be combined with different stopping rules to tackle different **relaxations** of PSI, e.g. $\min(\tau, \tau^k)$ where

$$\tau^k = \inf\{t \in \mathbb{N} : |\text{OPT}(t)| \geq k\}$$

to identify **at most k Pareto optimal arms**.

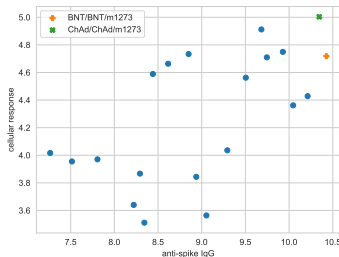
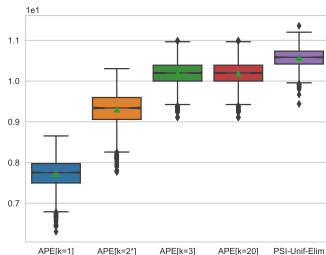
Theorem [Kone et al., 2023]

Assume the observations are bounded in $[0, 1]^D$. Then, with probability larger than $1 - \delta$, APE with the stopping rule $\tau_\delta \wedge \tau^k$ outputs $\hat{S}_\tau = \mathcal{P}^*(\mu)$ and satisfies

$$\tau_\delta \leq \sum_{a=1}^K \frac{32}{\widetilde{\Delta}_a^2} \log \left(\frac{2KD}{\delta} \log \left(\frac{32}{\widetilde{\Delta}_a^2} \right) \right),$$

for a relaxation $\widetilde{\Delta}_a = \max(\Delta_a, h_k)$.

Numerical results



(Log) Empirical sample complexity of APE (with a k -relaxation) compared to the algorithm of [Auer et al., 2016] on simulated CovBoost data [Munro et al., 2021]

- improved practical performance
- the k -relaxation (provably) reduces the sample complexity

- 1 Pure Exploration in Multi-objective bandits
- 2 Best Arm Identification ($D = 1$)
- 3 Adaptive Pareto Exploration
- 4 Towards Optimal Algorithms

For arms that are multi-variate Gaussian (known covariance Σ), could we further try to match the lower bound?

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

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

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

- ➔ The structure of the alternative is complex for PSI, making even the computation of “minimal distance” challenging...

For arms that are multi-variate Gaussian (known covariance Σ), could we further try to match the lower bound?

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

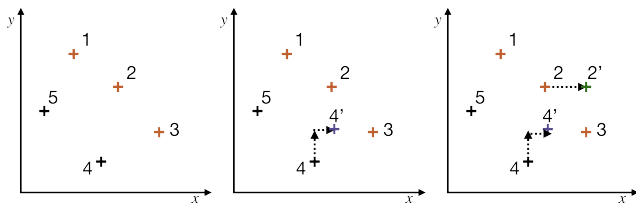
$$T^*(\mu)^{-1} = \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(\mathcal{S}^*(\mu))} \left(\sum_{k=1}^K w_k \frac{1}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right).$$

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

- ➔ The structure of the alternative is complex for PSI, making even the computation of “minimal distance” challenging...

Computing the Minimal Distance

- there are many ways to alter the Pareto set



- no closed-form is known for the minimal distance

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

- 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

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

is already 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 [Ménard, 2019, Degenne et al., 2019] also rely on **minimal distance** computation.

A Fully Sampling-Based Approach

Posterior Sampling for PSI (PSIPS)

(simplified)

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 return**
 $\hat{S}_t = \mathcal{S}^*(\hat{\mu}(t))$
- Else, take the first m such that $\mathcal{S}^*(\tilde{\theta}^m) \neq \mathcal{S}^*(\hat{\mu}(t))$
Update an online learning algorithm on Δ_K with the gain

$$g_t(w) = \sum_{a=1}^K w_a \frac{1}{2} \|\hat{\mu}_a(t) - \tilde{\theta}_a^m\|_{\Sigma^{-1}}^2$$

to get w_t . Select arm $A_t \sim (1 - \gamma_t)w_t + \gamma_t w_{\text{exp}}$

[Kone et al., 2025], inspired by PEPS [Li et al., 2024]

A Fully Sampling-Based Approach

Posterior Sampling for PSI (PSIPS)

(simplified)

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 return**
 $\hat{S}_t = \mathcal{S}^*(\hat{\mu}(t))$
- Else, take the first m such that $\mathcal{S}^*(\tilde{\theta}^m) \neq \mathcal{S}^*(\hat{\mu}(t))$
Update an online learning algorithm on Δ_K with the gain

$$g_t(w) = \sum_{a=1}^K w_a \frac{1}{2} \|\hat{\mu}_a(t) - \tilde{\theta}_a^m\|_{\Sigma^{-1}}^2$$

to get w_t . Select arm $A_t \sim (1 - \gamma_t)w_t + \gamma_t w_{\text{exp}}$

[Kone et al., 2025], inspired by PEPS [Li et al., 2024]

Sample complexity

Using budget M and inflation c such that

$$\limsup_{\delta \rightarrow 0} \frac{c(t, \delta) \log M(t, \delta)}{\log(1/\delta)} \leq 1,$$

PSIPS satisfies

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\text{PS}}]}{\log(1/\delta)} \leq T_{\star}(\mu)$$

Rationale. the truncated posterior density is close to

$$\begin{aligned} q_t(\lambda) &\propto \exp \left(- \sum_k N_{t,k} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right) \cdot \mathbb{1}_{\lambda \in \text{Alt}(S_t)} \\ &\propto q_{t-1}(\lambda) \cdot \exp \left(- \|\mu_{A_t} - \lambda_{A_t}\|_{\Sigma^{-1}}^2 \right) \end{aligned}$$

which mirrors the behavior of the **continuous Exponential Weights** algorithm under quadratic loss.

Sample complexity

Using budget M and inflation c such that

$$\limsup_{\delta \rightarrow 0} \frac{c(t, \delta) \log M(t, \delta)}{\log(1/\delta)} \leq 1,$$

PSIPS satisfies

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\text{PS}}]}{\log(1/\delta)} \leq T_{\star}(\mu)$$

Rationale. the truncated posterior density is close to

$$\begin{aligned} q_t(\lambda) &\propto \exp \left(- \sum_k N_{t,k} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right) \cdot \mathbb{1}_{\lambda \in \text{Alt}(S_t)} \\ &\propto q_{t-1}(\lambda) \cdot \exp \left(- \|\mu_{A_t} - \lambda_{A_t}\|_{\Sigma^{-1}}^2 \right) \end{aligned}$$

which mirrors the behavior of the **continuous Exponential Weights** algorithm under quadratic loss.

The calibration of the PS stopping rule is not as easy as the GLR: it requires a lower bound on

$$\Pi_t(\text{Alt}(S_t)^c) \quad \text{when } S_t \neq \mathcal{S}^*$$

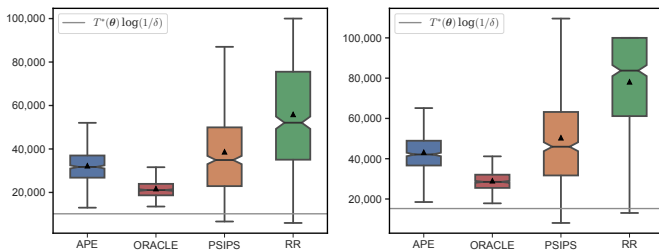
and thus some anti-concentration results.

Lemma

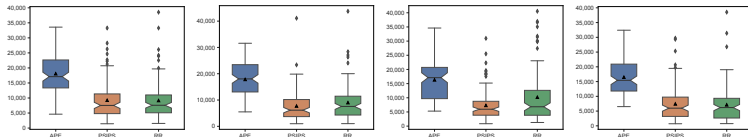
For PSIPS to be δ -correct we can choose

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

- CovBoost ($d = 3$) for $\delta = 0.1$ (left) and $\delta = 0.01$ (right)



- Random Gaussian instances with $K = 10$ for $d \in \{3, 4, 5, 6\}$



We proposed two approaches to Pareto Set Identification in the Fixed Confidence Setting:

- **Adaptive Pareto Exploration**: finite time bound, sub-optimal in the asymptotic regime $\delta \rightarrow 0$
 - **PSIPS**, a (tractable !) Lower Bound Inspired algorithm, optimal in the asymptotic regime
- which one should we use in practise?

The **sampling-based stopping rule** is an interesting alternative to the GLR stopping rule for any complex pure exploration problem

Perspective: multi-objective bandit algorithms always sample *all* the marginals of the chosen arm → can we also adaptively select which marginals to observe?



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.



Degenne, R., Koolen, W. M., and Ménard, P. (2019).
Non-asymptotic pure exploration by solving games.
In Advances in Neural Information Processing Systems (NeurIPS).



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.



Gabillon, V., Ghavamzadeh, M., and Lazaric, A. (2012).
Best Arm Identification: A Unified Approach to Fixed Budget and Fixed
Confidence.
In Advances in Neural Information Processing Systems.



Garivier, A. and Kaufmann, E. (2016).
Optimal best arm identification with fixed confidence.
In Proceedings of the 29th Conference On Learning Theory.



Jourdan, M., Degenne, R., Baudry, D., de Heide, R., and Kaufmann, E. (2022).
Top two algorithms revisited.
In Advances in Neural Information Processing Systems (NeurIPS).



Kalyanakrishnan, S., Tewari, A., Auer, P., and Stone, P. (2012).
PAC subset selection in stochastic multi-armed bandits.
In International Conference on Machine Learning (ICML).



Katz-Samuels, J. and Scott, C. (2018).
Feasible arm identification.
In International Conference on Machine Learning (ICML).



Katz-Samuels, J. and Scott, C. (2019).
Top feasible arm identification.
In International Conference on Artificial Intelligence and Statistics (AISTATS).



Kaufmann, E. and Koolen, W. (2021).
Mixture martingales revisited with applications to sequential tests and confidence intervals.
Journal of Machine Learning Research, 22(246).



Kone, C., Jourdan, M., and Kaufmann, E. (2025).
Pareto set identification with posterior sampling.
In AISTATS.



Kone, C., Kaufmann, E., and Richert, L. (2023).
Adaptive algorithms for relaxed pareto set identification.
In Advances in Neural Information Processing Systems (NeurIPS).



Li, Z., Jamieson, K., and Jain, L. (2024).
Optimal exploration is no harder than Thompson sampling.
In Proceedings of The 27th International Conference on Artificial Intelligence and Statistics. PMLR.



Ménard, P. (2019).

Gradient ascent for active exploration in bandit problems.

arXiv 1905.08165.



Munro, A.-P.-S., Janani, L., Cornelius, V., and et al. (2021).

Safety and immunogenicity of seven COVID-19 vaccines as a third dose (booster) following two doses of ChAdOx1 nCov-19 or BNT162b2 in the UK (COV-BOOST): a blinded, multicentre, randomised, controlled, phase 2 trial.

The Lancet, 398(10318):2258–2276.

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.

