# BANDIT PROBLEMS
## Part III - Bandits for Optimization

RLSS, Lille, July 2019

$K$ **arms** $\leftrightarrow K$ probability distributions : $\nu_a$ has mean $\mu_a$



$\nu_1 \qquad \nu_2 \qquad \nu_3 \qquad \nu_4 \qquad \nu_5$

At round $t$, an agent:

▶ chooses an arm $A_t$

▶ receives a reward $R_t = X_{A_t, t} \sim \nu_{A_t}$

Sequential sampling strategy (**bandit algorithm**):

$$A_{t+1} = F_t(A_1, R_1, \ldots, A_t, R_t).$$

**Goal:** Maximize $\mathbb{E}\left[\sum_{t=1}^{T} R_t\right]$

$K$ **arms** $\leftrightarrow$ $K$ probability distributions : $\nu_a$ has mean $\mu_a$



$\nu_1$      $\nu_2$      $\nu_3$      $\nu_4$      $\nu_5$

At round $t$, an agent:

- chooses an arm $A_t$
- receives a sample $X_t = X_{A_t,t} \sim \nu_{A_t}$

Sequential sampling strategy (**bandit algorithm**):

$$A_{t+1} = F_t(A_1, X_1, \ldots, A_t, X_t).$$

**Goal:** Maximize $\mathbb{E}\left[\sum_{t=1}^{T} X_t\right] \rightarrow$ not the only possible goal!

$\mathcal{B}(\mu_1)$ $\qquad$ $\mathcal{B}(\mu_2)$ $\qquad$ $\mathcal{B}(\mu_3)$ $\qquad$ $\mathcal{B}(\mu_4)$ $\qquad$ $\mathcal{B}(\mu_5)$
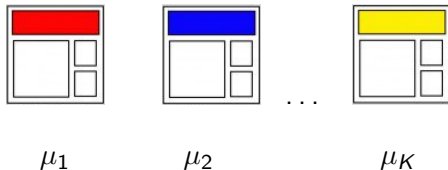
For the $t$-th patient in a clinical study,

▶ chooses a treatment $A_t$

▶ observes a response $X_t \in \{0, 1\}$: $\mathbb{P}(X_t = 1) = \mu_{A_t}$

**Maximize rewards** $\leftrightarrow$ cure as many patients as possible

$\mathcal{B}(\mu_1)$    $\mathcal{B}(\mu_2)$    $\mathcal{B}(\mu_3)$    $\mathcal{B}(\mu_4)$    $\mathcal{B}(\mu_5)$

For the $t$-th patient in a clinical study,

► chooses a treatment $A_t$
► observes a response $X_t \in \{0, 1\}$: $\mathbb{P}(X_t = 1) = \mu_{A_t}$

**Maximize rewards** $\leftrightarrow$ cure as many patients as possible

**Alternative goal:** identify as quickly as possible the best treatment (without trying to cure patients during the study)

## Should we maximize rewards?

Probability that some version of a website generates a conversion:



$\mu_1$        $\mu_2$        $\mu_K$

**Best version**: $a_\star = \underset{a=1,\ldots,K}{\operatorname{argmax}} \mu_a$
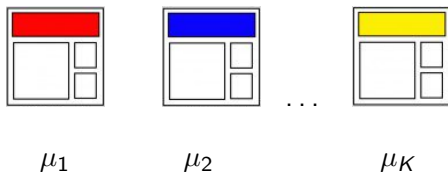
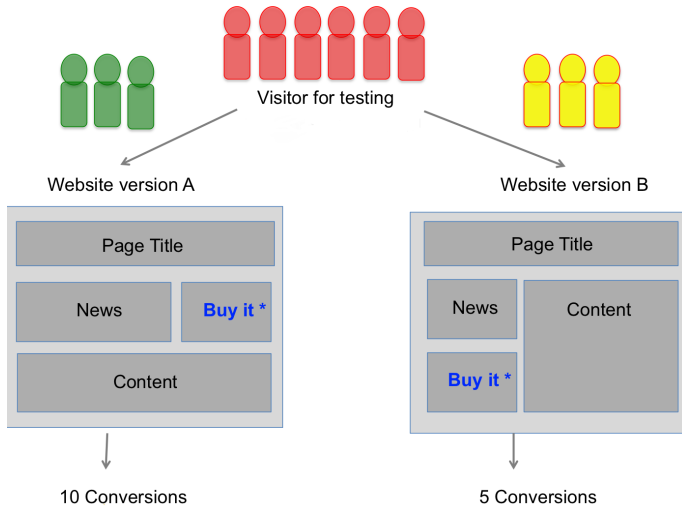**Sequential protocol**: for the $t$-th visitor:

- ▶ display version $A_t$
- ▶ observe conversion indicator $X_t \sim \mathcal{B}(\mu_{A_t})$.

**Maximize rewards** $\leftrightarrow$ maximize the number of conversions

# Should we maximize rewards?

Probability that some version of a website generates a conversion:



$$\mu_1 \qquad \mu_2 \qquad \qquad \mu_K$$

**Best version**: $a_\star = \underset{a=1,\dots,K}{\operatorname{argmax}} \mu_a$

**Sequential protocol**: for the $t$-th visitor:

▶ display version $A_t$

▶ observe conversion indicator $X_t \sim \mathcal{B}(\mu_{A_t})$.

**Maximize rewards** $\leftrightarrow$ maximize the number of conversions

**Alternative goal:** identify the best version
(without trying to maximize conversions during the test)

# A/B Testing ($K = 2$)



Visitor for testing

Website version A

Website version B

| Page Title |
| News | **Buy it \*** |
| Content |

| Page Title |
| News | Content |
| **Buy it \*** | |

10 Conversions

5 Conversions

**Goal:** find the best version.

A way to do A/B Testing:

- ▶ allocate $n_A$ users to page A and $n_B$ users to page B (decided in advance, often $n_A = n_B$)
- ▶ perform a statistical test of "A better than B"

A way to do A/B Testing:

► allocate $n_A$ users to page A and $n_B$ users to page B
   (decided in advance, often $n_A = n_B$)

► perform a statistical test of "A better than B"

Alternative: fully adaptive A/B Testing

► sequentially choose which version to allocate to each visitor

► (adaptively choose when to stop the experiment)

➜ best arm identification in a bandit model

# Finding the best arm in a bandit model

# Pure Exploration in Bandit Models

**Goal:** identify an arm with large mean as quickly and accurately as possible $\simeq$ identify

$$a_\star = \underset{a=1,\dots,K}{\mathrm{argmax}}\ \mu_a.$$

**Algorithm:** made of three components:

- ➜ sampling rule: $A_t$ (arm to explore)
- ➜ recommendation rule: $B_t$ (current guess for the best arm)
- ➜ stopping rule $\tau$ (when do we stop exploring?)

## Probability of error [Even-Dar et al., 2006, Audibert et al., 2010]

The probability of error after $n$ rounds is

$$p_\nu(n) = \mathbb{P}_\nu\left(B_n \neq a_\star\right).$$

**Goal:** identify an arm with large mean as quickly and accurately as possible $\simeq$ identify

$$a_\star = \underset{a=1,\ldots,K}{\mathrm{argmax}} \ \mu_a.$$

**Algorithm:** made of three components:

→ sampling rule: $A_t$ (arm to explore)

→ recommendation rule: $B_t$ (current guess for the best arm)

→ stopping rule $\tau$ (when do we stop exploring?)

## Simple regret [Bubeck et al., 2009]

The simple regret after $n$ rounds is

$$r_\nu(n) = \mu_\star - \mu_{B_n}.$$

**Goal:** identify an arm with large mean as quickly and accurately as possible $\simeq$ identify

$$a_\star = \underset{a=1,\dots,K}{\mathrm{argmax}} \ \mu_a.$$

**Algorithm:** made of three components:

➜ sampling rule: $A_t$ (arm to explore)

➜ recommendation rule: $B_t$ (current guess for the best arm)

➜ stopping rule $\tau$ (when do we stop exploring?)

## Simple regret [Bubeck et al., 2009]

The simple regret after $n$ rounds is

$$r_\nu(n) = \mu_\star - \mu_{B_n}.$$

$$\Delta_{\min} p_\nu(n) \le \mathbb{E}_\nu[r_\nu(n)] \le \Delta_{\max} p_\nu(n)$$

# Several objectives

**Algorithm:** made of three components:

➜ sampling rule: $A_t$ (arm to explore)

➜ recommendation rule: $B_t$ (current guess for the best arm)

➜ stopping rule $\tau$ (when do we stop exploring?)

▶ **Objectives studied in the literature:**

| Fixed-confidence setting | Fixed-budget setting | Anytime exploration |
|---|---|---|
| input: risk parameter $\delta$ (tolerance parameter $\epsilon$) | input: budget $T$ | no input |
| minimize $\mathbb{E}[\tau]$ $\mathbb{P}(B_\tau \neq a_\star) \leq \delta$ or $\mathbb{P}(r_\nu(\tau) < \epsilon) \leq \delta$ | $\tau = T$ minimize $\mathbb{P}(B_T \neq a_\star)$ or $\mathbb{E}[r_T(\nu)]$ | for all $t$, minimize $p_\nu(t)$ or $\mathbb{E}[r_\nu(t)]$ |
| [Even-Dar et al., 2006] | [Bubeck et al., 2009] [Audibert et al., 2010] | [Jun and Nowak, 2016] |

# Outline

**Context:** bounded rewards ($\nu_a$ supported in $[0, 1]$)

We know good algorithms to maximize rewards, for example UCB($\alpha$)

$$A_{t+1} = \underset{a=1,\ldots,K}{\operatorname{argmax}} \; \hat{\mu}_a(t) + \sqrt{\alpha \frac{\ln(t)}{N_a(t)}}$$

▶ How good is it for best arm identification?

# Can we use UCB?

**Context:** bounded rewards ($\nu_a$ supported in $[0, 1]$)

We know good algorithms to maximize rewards, for example UCB($\alpha$)

$$A_{t+1} = \operatorname*{argmax}_{a=1,\dots,K} \hat{\mu}_a(t) + \sqrt{\alpha \frac{\ln(t)}{N_a(t)}}$$

▶ How good is it for best arm identification?

**Possible recommendation rules**:

| | |
|---|---|
| Empirical Best Arm (EBA) | $B_t = \operatorname{argmax}_a \ \hat{\mu}_a(t)$ |
| Most Played Arm (MPA) | $B_t = \operatorname{argmax}_a \ N_a(t)$ |
| Empirical Distribution of Plays (EDP) | $B_t \sim p_t$, where $p_t = \left( \frac{N_1(t)}{t}, \dots, \frac{N_K(t)}{t} \right)$ |

[Bubeck et al., 2009]

▶ **UCB + Empirical Distribution of Plays**

$$
\begin{aligned}
\mathbb{E}\left[r_\nu(n)\right] &= \mathbb{E}\left[\mu_\star - \mu_{B_n}\right] = \mathbb{E}\left[\sum_{b=1}^{K}(\mu_\star - \mu_b)\mathbb{1}_{(B_n=b)}\right] \\
&= \mathbb{E}\left[\sum_{b=1}^{K}(\mu_\star - \mu_b)\mathbb{P}(B_n = b|\mathcal{F}_n)\right] \\
&= \mathbb{E}\left[\sum_{b=1}^{K}(\mu_\star - \mu_b)\frac{N_b(n)}{n}\right] \\
&= \frac{1}{n}\sum_{b=1}^{K}(\mu_\star - \mu_b)\mathbb{E}[N_b(n)] \\
&= \frac{\mathcal{R}_\nu(n)}{n}.
\end{aligned}
$$

➜ a conversion from cumulated regret to simple regret!

▶ **UCB + Empirical Distribution of Plays**

$$\mathbb{E}\left[r_\nu\left(\texttt{UCB}(\alpha), n\right)\right] \leq \frac{\mathcal{R}_\nu(\texttt{UCB}(\alpha), n)}{n} \leq \frac{C(\nu)\ln(n)}{n}$$

▶ **UCB + Empirical Distribution of Plays**

$$\mathbb{E}\left[r_\nu\left(\text{UCB}(\alpha), n\right)\right] \leq \frac{\mathcal{R}_\nu(\text{UCB}(\alpha), n)}{n} \leq C\sqrt{\frac{K\alpha \ln(n)}{n}}$$

# Can we use UCB?

► **UCB + Empirical Distribution of Plays**

$$\mathbb{E}\left[r_\nu\left(\text{UCB}(\alpha), n\right)\right] \leq \frac{\mathcal{R}_\nu(\text{UCB}(\alpha), n)}{n} \leq C\sqrt{\frac{K\alpha \ln(n)}{n}}$$

► **UCB + Most Played Arm**

### Theorem [Bubeck et al., 2009]

With the Most Play Armed as a recommendation rule, for $n$ large enough,

$$\mathbb{E}\left[r_\nu\left(\text{UCB}(\alpha), n\right)\right] \leq C\sqrt{\frac{K\alpha \ln(n)}{n}}.$$

(more precise problem-dependent analysis in [Bubeck et al., 2009])

# Are those results good?

$$\mathbb{E}_\nu \left[ r_\nu \left( \mathtt{UCB}(\alpha), n \right) \right] \simeq \sqrt{\frac{K\alpha \ln(n)}{n}}.$$

▶ the uniform allocation strategy can beat UCB!

**Theorem** [Bubeck et al., 2009]

For $n \geq 4K \ln(K)/\Delta_{\min}^2$, the simple regret decays exponentially :

$$\mathbb{E}_\nu \left[ r_\nu \left( \mathtt{Unif}, n \right) \right] \leq \Delta_{\max} \exp \left( -\frac{1}{8} \frac{n}{K} \Delta_{\min}^2 \right)$$

▶ the smaller the cumulative regret, the larger the simple regret

**Theorem** [Bubeck et al., 2009]

$$''\mathbb{E}[r_\nu(\mathcal{A}, n)] \geq \frac{\Delta_{\min}}{2} \exp \left( -C \times \mathcal{R}_\nu(\mathcal{A}, n) \right)''$$

**Answer:**

▶ UCB has to be coupled with an appropriate recommendation rule

▶ it is not guaranteed to perform better than uniform exploration...

**Answer:**

- ▶ UCB has to be coupled with an appropriate recommendation rule
- ▶ it is not guaranteed to perform better than uniform exploration...

**Variants of UCB?**

- ▶ UCB-E for the fixed-budget setting [Audibert et al., 2010]
- ▶ LIL-UCB for the fixed-confidence setting [Jamieson et al., 2014]

**Other algorithms?**

- ▶ many specific algorithm for best arm identification

# Fixed Budget: Sequential Halving

**Input:** total number of plays $T$

**Idea:** split the budget in $\log_2(K)$ phases of equal length, eliminate the worst half of the remaining arms after each phase.

---

**Initialisation:** $S_0 = \{1, \dots, K\}$;
**For** $r = 0$ **to** $\lceil \ln_2(K) \rceil - 1$, **do**
    sample each arm $a \in S_r$ $\quad t_r = \left\lfloor \frac{T}{|S_r| \lceil \log_2(K) \rceil} \right\rfloor$ times;
    let $\hat{\mu}_a^r$ be the empirical mean of arm $a$;
    let $S_{r+1}$ be the set of $\lceil |S_r|/2 \rceil$ arms with largest $\hat{\mu}_a^r$
**Output:** $B_T$ the unique arm in $S_{\lceil \log_2(K) \rceil}$

---

## Theorem [Karnin et al., 2013]

Letting $H_2(\nu) = \max_{a \neq a_\star} a\Delta_{[a]}^{-2}$, for any bounded bandit instance,

$$\mathbb{P}_\nu \left( B_T \neq a_\star \right) \leq 3 \log_2(K) \exp\left( -\frac{T}{8 \log_2(K) H_2(\nu)} \right).$$

# Fixed Confidence: Successive Elimination

**Input:** risk parameter $\delta \in (0, 1)$.
**Idea:** sample all remaining arm uniformly and perform eliminations of arms which look sub-optimal

**Initialization:** $\mathcal{S} = \{1, \dots, K\}$
**While** $|\mathcal{S}| > 1$
    Draw all arms in $\mathcal{S}$. $t \leftarrow t + |\mathcal{S}|$.
    $\mathcal{S} \leftarrow \mathcal{S} \backslash \{a\}$ if $\max_{i \in \mathcal{S}} \hat{\mu}_i(t) - \hat{\mu}_a(t) \geq 2 \sqrt{\frac{\ln(Kt^2/\delta)}{t}}$.
**Output:** the unique arm $B_\tau \in \mathcal{S}$.

## Theorem [Even-Dar et al., 2006]

Successive Elimination satisfies $\mathbb{P}_\nu \left( B_\tau = a_\star \right) \geq 1 - \delta$. Moreover,

$$\mathbb{P}_\nu \left( \tau_\delta = O \left( \sum_{a=2}^{K} \frac{1}{\Delta_a^2} \ln \left( \frac{K}{\delta \Delta_a} \right) \right) \right) \geq 1 - \delta.$$

# Fixed-confidence: LUCB

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



- At round $t$, draw

$$B_t = \underset{b}{\operatorname{argmax}} \ \hat{\mu}_b(t)$$
$$C_t = \underset{c \neq B_t}{\operatorname{argmax}} \ \mathrm{UCB}_c(t)$$

- Stop at round $t$ if

$$\mathrm{LCB}_{B_t}(t) > \mathrm{UCB}_{C_t}(t) - \epsilon$$

## Theorem [Kalyanakrishnan et al., 2012]

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

$$\mathbb{E}\left[\tau_\delta\right] = O\left(\left[\frac{1}{\Delta_2^2 \vee \epsilon^2} + \sum_{a=2}^{K} \frac{1}{\Delta_a^2 \vee \epsilon^2}\right] \ln\left(\frac{1}{\delta}\right)\right)$$

# Outline

# Minimizing the sample complexity

**Context:** Exponential family bandit model

$$\nu = (\nu_1, \ldots, \nu_K) \quad \leftrightarrow \quad \boldsymbol{\mu} = (\mu_1, \ldots, \mu_K)$$

*(Bernoulli, Gaussian with known variance, Poisson...)*

**Algorithm:** made of three components:

➜ sampling rule: $A_t$ (arm to explore)

➜ stopping rule $\tau$ (when do we stop exploring?)

➜ recommendation rule: $B_\tau$ (guess for the best arm when stopping)

## Objective

▶ a $\delta$-correct strategy: for all $\boldsymbol{\mu}$ with a unique optimal arm,

$$\mathbb{P}_{\boldsymbol{\mu}}\left(B_\tau = a_\star\right) \geq 1 - \delta.$$

▶ with a small sample complexity $\mathbb{E}_{\boldsymbol{\mu}}[\tau_\delta]$.

➜ minimal sample complexity?

# A first lower bound

**Divergence function:** $\mathrm{kl}(\mu, \mu') = \mathrm{KL}(\nu_\mu, \nu_{\mu'})$.

---

**Change of distribution lemma** [Kaufmann et al., 2016]

$\mu$ and $\lambda$ be such that $a_\star(\mu) \neq a_\star(\lambda)$. For any $\delta$-correct algorithm,
$$\sum_{a=1}^{K} \mathbb{E}_\mu[N_a(\tau)] \mathrm{kl}(\mu_a, \lambda_a) \geq \mathrm{kl}_{\mathsf{Ber}}(\delta, 1 - \delta).$$

---

▶ For any $a \in \{2, \ldots, K\}$, introducing $\lambda$:



$$\left\{ \begin{array}{rcl} \lambda_a & = & \mu_1 + \epsilon \\ \lambda_i & = & \mu_i, \ \ \text{if} \ \ i \neq a \end{array} \right.$$

$$\mathbb{E}_\mu[N_a(\tau)] \mathrm{kl}(\mu_a, \mu_1 + \epsilon) \geq \mathrm{kl}_{\mathsf{Ber}}(\delta, 1 - \delta)$$

$$\mathbb{E}_\mu[N_a(\tau)] \geq \frac{1}{\mathrm{kl}(\mu_a, \mu_1)} \ln\left(\frac{1}{3\delta}\right).$$

# A first lower bound

**Divergence function:** $\mathrm{kl}(\mu, \mu') = \frac{(\mu - \mu')^2}{2\sigma^2}$ (Gaussian distributions).

> **Change of distribution lemma** [Kaufmann et al., 2016]
>
> $\mu$ and $\lambda$ be such that $a_\star(\mu) \neq a_\star(\lambda)$. For any $\delta$-correct algorithm,
> $$\sum_{a=1}^{K} \mathbb{E}_{\mu}[N_a(\tau)]\mathrm{kl}(\mu_a, \lambda_a) \geq \mathrm{kl}_{\mathrm{Ber}}(\delta, 1 - \delta).$$

▶ For any $a \in \{2, \dots, K\}$, introducing $\lambda$:



$$\left\{ \begin{array}{rcl} \lambda_a & = & \mu_1 + \epsilon \\ \lambda_i & = & \mu_i, \ \ \text{if} \ \ i \neq a \end{array} \right.$$

$$\mathbb{E}_{\mu}[N_a(\tau)]\mathrm{kl}(\mu_a, \mu_1 + \epsilon) \geq \mathrm{kl}_{\mathrm{Ber}}(\delta, 1 - \delta)$$
$$\mathbb{E}_{\mu}[N_a(\tau)] \geq \frac{1}{\mathrm{kl}(\mu_a, \mu_1)} \ln\left(\frac{1}{3\delta}\right).$$

# A first lower bound

**Divergence function:** $\mathrm{kl}(\mu, \mu') = \mathrm{KL}(\nu_\mu, \nu_{\mu'})$.

### Change of distribution lemma [Kaufmann et al., 2016]

$\mu$ and $\lambda$ be such that $a_\star(\mu) \neq a_\star(\lambda)$. For any $\delta$-correct algorithm,
$$\sum_{a=1}^{K} \mathbb{E}_\mu[N_a(\tau)]\mathrm{kl}(\mu_a, \lambda_a) \geq \mathrm{kl}_{\mathrm{Ber}}(\delta, 1 - \delta).$$

▶ One obtains

### Theorem [Kaufmann et al., 2016]

For any $\delta$-correct algorithm,
$$\mathbb{E}_\mu[\tau] \geq \left( \frac{1}{\mathrm{kl}(\mu_1, \mu_2)} + \sum_{a=2}^{K} \frac{1}{\mathrm{kl}(\mu_a, \mu_1)} \right) \ln\left( \frac{1}{3\delta} \right).$$

# A first lower bound

**Divergence function:** $\mathrm{kl}(\mu, \mu') = \mathrm{KL}(\nu_\mu, \nu_{\mu'})$.

---

**Change of distribution lemma** [Kaufmann et al., 2016]

$\mu$ and $\lambda$ be such that $a_\star(\mu) \neq a_\star(\lambda)$. For any $\delta$-correct algorithm,
$$\sum_{a=1}^{K} \mathbb{E}_\mu[N_a(\tau)]\mathrm{kl}(\mu_a, \lambda_a) \geq \mathrm{kl}_{\mathsf{Ber}}(\delta, 1 - \delta).$$

---

▶ One obtains

---

**Theorem** [Kaufmann et al., 2016]

For any $\delta$-correct algorithm,
$$\mathbb{E}_\mu[\tau] \geq \left( \frac{1}{\mathrm{kl}(\mu_1, \mu_2)} + \sum_{a=2}^{K} \frac{1}{\mathrm{kl}(\mu_a, \mu_1)} \right) \ln \left( \frac{1}{3\delta} \right).$$

---

➜ not tight enough...

# The best possible lower bound

$\boldsymbol{\mu}$ and $\boldsymbol{\lambda}$ be such that $a_\star(\boldsymbol{\mu}) \neq a_\star(\boldsymbol{\lambda})$. For any $\delta$-correct algorithm,

$$\sum_{a=1}^{K} \mathbb{E}_{\boldsymbol{\mu}}[N_a(\tau)] \mathrm{kl}(\mu_a, \lambda_a) \geq \mathrm{kl}_{\mathrm{Ber}}(\delta, 1 - \delta).$$

▶ Let $\mathrm{Alt}(\boldsymbol{\mu}) = \{\boldsymbol{\lambda} : a_\star(\boldsymbol{\lambda}) \neq a_\star(\boldsymbol{\mu})\}$.

$$\inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\boldsymbol{\mu})} \sum_{a=1}^{K} \mathbb{E}_{\boldsymbol{\mu}}[N_a(\tau)] \mathrm{kl}(\mu_a, \lambda_a) \geq \mathrm{kl}_{\mathrm{Ber}}(\delta, 1 - \delta)$$

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \times \inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\boldsymbol{\mu})} \sum_{a=1}^{K} \frac{\mathbb{E}_{\boldsymbol{\mu}}[N_a(\tau)]}{\mathbb{E}_{\boldsymbol{\mu}}[\tau]} \mathrm{kl}(\mu_a, \lambda_a) \geq \ln\left(\frac{1}{3\delta}\right)$$

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \times \left( \sup_{w \in \Sigma_K} \inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\boldsymbol{\mu})} \sum_{a=1}^{K} w_a \mathrm{kl}(\mu_a, \lambda_a) \right) \geq \ln\left(\frac{1}{3\delta}\right)$$

# The best possible lower bound

**Theorem** [Garivier and Kaufmann, 2016]

For any $\delta$-PAC algorithm,

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

where

$$T_{\star}(\boldsymbol{\mu})^{-1} = \sup_{w \in \Sigma_K} \inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\boldsymbol{\mu})} \left(\sum_{a=1}^{K} w_a \mathrm{kl}(\mu_a, \lambda_a)\right).$$

Moreover, the vector of optimal proportions,

$$w_{\star}(\boldsymbol{\mu}) = \operatorname*{argmax}_{w \in \Sigma_K} \inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\boldsymbol{\mu})} \left(\sum_{a=1}^{K} w_a \mathrm{kl}(\mu_a, \lambda_a)\right)$$

is well-defined, and can be computed efficiently.

# Outline

$\hat{\boldsymbol{\mu}}(t) = (\hat{\mu}_1(t), \dots, \hat{\mu}_K(t))$: vector of empirical means

▶ Introducing

$$U_t = \left\{ a : N_a(t) < \sqrt{t} \right\},$$

one has

$$A_{t+1} \in \begin{cases} \underset{a \in U_t}{\operatorname{argmin}}\ N_a(t) \text{ if } U_t \neq \emptyset & (\textit{forced exploration}) \\ \underset{1 \leq a \leq K}{\operatorname{argmax}} \left[ (w_\star(\hat{\boldsymbol{\mu}}(t)))_a - \frac{N_a(t)}{t} \right] & (\textit{tracking}) \end{cases}$$

## Lemma

Under the Tracking sampling rule,

$$\mathbb{P}_{\boldsymbol{\mu}} \left( \lim_{t \to \infty} \frac{N_a(t)}{t} = (w_\star(\boldsymbol{\mu}))_a \right) = 1.$$

# How to match the lower bound? Stopping rule.

▶ a Generalized Likelihood Ratio:

$$\hat{Z}(t) = \ln \frac{\ell\left(X_1, \ldots, X_t; \hat{\boldsymbol{\mu}}(t)\right)}{\max\limits_{\boldsymbol{\lambda} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))} \ell(X_1, \ldots, X_t; \boldsymbol{\lambda})} = \inf\limits_{\boldsymbol{\lambda} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))} \sum_{a=1}^{K} N_a(t) \mathrm{kl}(\hat{\mu}_a(t), \lambda_a)$$

➜ high value of $\hat{Z}(t)$ rejects the hypothesis "$\boldsymbol{\mu} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))$".

## Stopping and recommendation rule

$$\tau_\delta = \inf\left\{ t \in \mathbb{N} : \hat{Z}(t) > \beta(t, \delta) \right\}$$

$$B_\tau = \underset{a=1,\ldots,K}{\mathrm{argmax}}\, \hat{\mu}_a(\tau).$$

(can be traced back to [Chernoff, 1959])

# How to match the lower bound? Stopping rule.

▶ a Generalized Likelihood Ratio:

$$\hat{Z}(t) = \ln \frac{\ell\left(X_1, \ldots, X_t; \hat{\boldsymbol{\mu}}(t)\right)}{\max_{\boldsymbol{\lambda} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))} \ell(X_1, \ldots, X_t; \boldsymbol{\lambda})} = \inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))} \sum_{a=1}^{K} N_a(t) \mathrm{kl}(\hat{\mu}_a(t), \lambda_a)$$

➔ high value of $\hat{Z}(t)$ rejects the hypothesis "$\boldsymbol{\mu} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))$".

## Stopping and recommendation rule

$$
\begin{aligned}
\tau_\delta &= \inf\left\{ t \in \mathbb{N} : \hat{Z}(t) > \beta(t, \delta) \right\} \\
B_\tau &= \operatorname*{argmax}_{a=1,\ldots,K} \hat{\mu}_a(\tau).
\end{aligned}
$$

(can be traced back to [Chernoff, 1959])

➔ How to pick the threshold $\beta(t, \delta)$?

# A $\delta$-correct stopping rule

$$
\begin{aligned}
\hat{Z}(t) &= \inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))} \sum_{a=1}^{K} N_a(t) \mathrm{kl}(\hat{\mu}_a(t), \lambda_a) \\
&= \min_{b \neq B_t} \inf_{\{\boldsymbol{\lambda} : \lambda_{B_t} \leq \lambda_b\}} \sum_{a=1}^{K} N_a(t) \mathrm{kl}(\hat{\mu}_a(t), \lambda_a)
\end{aligned}
$$

# A $\delta$-correct stopping rule

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

$$= \min_{b \neq B_t} \inf_{\{\lambda : \lambda_{B_t} \leq \lambda_b\}} \left[ N_{B_t}(t)\mathrm{kl}(\hat{\mu}_{B_t}(t), \lambda_{B_t}) + N_b(t)\mathrm{kl}(\hat{\mu}_b(t), \lambda_b) \right]$$

# A $\delta$-correct stopping rule

$$
\begin{aligned}
\hat{Z}(t) &= \inf_{\lambda \in \mathrm{Alt}(\hat{\mu}(t))} \sum_{a=1}^{K} N_a(t) \mathrm{kl}(\hat{\mu}_a(t), \lambda_a) \\
&= \min_{b \neq B_t} \inf_{\{\lambda : \lambda_{B_t} \leq \lambda_b\}} \left[ N_{B_t}(t) \mathrm{kl}(\hat{\mu}_{B_t}(t), \lambda_{B_t}) + N_b(t) \mathrm{kl}(\hat{\mu}_b(t), \lambda_b) \right]
\end{aligned}
$$

$$
\mathbb{P}\left( B_{\tau_\delta} \neq a_\star \right) \leq \mathbb{P}\left( \exists t \in \mathbb{N}_\star, \exists a \neq a_\star : B_t = a, \hat{Z}(t) > \beta(t, \delta) \right)
$$

$$
\leq \mathbb{P}\left( \exists t \in \mathbb{N}_\star, \exists a \neq a_\star : \inf_{\lambda_a \leq \lambda_{a_\star}} \sum_{i \in \{a, a_\star\}} N_i(t) \mathrm{kl}(\hat{\mu}_i(t), \lambda_i) > \beta(t, \delta) \right)
$$

$$
\leq \sum_{a \neq a_\star} \underbrace{\mathbb{P}\left( \exists t \in \mathbb{N}_\star : N_a(t) \mathrm{kl}(\hat{\mu}_a(t), \mu_a) + N_{a_\star}(t) \mathrm{kl}(\hat{\mu}_{a_\star}(t), \mu_{a_\star}) > \beta(t, \delta) \right)}_{\text{requires simultaneous deviations on the two arms}}
$$

Bernoulli case: [Garivier and Kaufmann, 2016]
Exponential families: [Kaufmann and Koolen, 2018]

# An asymptotically optimal algorithm

## Theorem

The Track-and-Stop strategy, that uses

- the Tracking sampling rule
- the GLRT stopping rule with

$$\beta(t, \delta) \simeq \ln\left(\frac{K-1}{\delta}\right) + 2\ln\ln\left(\frac{K-1}{\delta}\right) + 6\ln(\ln(t))$$

- and recommends $B_\tau = \underset{a=1\ldots K}{\operatorname{argmax}}\ \hat{\mu}_a(\tau)$

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

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau_\delta]}{\ln(1/\delta)} = T_\star(\boldsymbol{\mu}).$$

# An asymptotically optimal algorithm

## Theorem

The Track-and-Stop strategy, that uses

- the Tracking sampling rule
- the GLRT stopping rule with

$$\beta(t,\delta) \simeq \ln\left(\frac{K-1}{\delta}\right) + 2\ln\ln\left(\frac{K-1}{\delta}\right) + 6\ln(\ln(t))$$

- and recommends $B_\tau = \underset{a=1\dots K}{\text{argmax}}\ \hat{\mu}_a(\tau)$

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

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau_\delta]}{\ln(1/\delta)} = T_\star(\boldsymbol{\mu}).$$

**Why?**

$$\tau_\delta = \inf\left\{ t \in \mathbb{N}_\star : \inf_{\boldsymbol{\lambda} \in \text{Alt}(\hat{\mu}(t))} \sum_{a=1}^{K} N_a(t)\text{kl}\left(\hat{\mu}_a(t), \lambda_a\right) > \beta(t,\delta) \right\}$$

# An asymptotically optimal algorithm

## Theorem

The Track-and-Stop strategy, that uses

- the Tracking sampling rule
- the GLRT stopping rule with

$$\beta(t, \delta) \simeq \ln\left(\frac{K-1}{\delta}\right) + 2\ln\ln\left(\frac{K-1}{\delta}\right) + 6\ln(\ln(t))$$

- and recommends $B_\tau = \underset{a=1\ldots K}{\operatorname{argmax}}\ \hat{\mu}_a(\tau)$

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

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau_\delta]}{\ln(1/\delta)} = T_\star(\boldsymbol{\mu}).$$

**Why?**

$$\tau_\delta = \inf\left\{ t \in \mathbb{N}_\star : t \times \inf_{\boldsymbol{\lambda} \in \mathrm{Alt}(\hat{\boldsymbol{\mu}}(t))} \sum_{a=1}^{K} \frac{N_a(t)}{t}\mathrm{kl}\left(\hat{\mu}_a(t), \lambda_a\right) > \beta(t, \delta) \right\}$$

# An asymptotically optimal algorithm

## Theorem

The Track-and-Stop strategy, that uses

- the Tracking sampling rule
- the GLRT stopping rule with

$$\beta(t, \delta) \simeq \ln\left(\frac{K-1}{\delta}\right) + 2\ln\ln\left(\frac{K-1}{\delta}\right) + 6\ln(\ln(t))$$

- and recommends $B_\tau = \underset{a=1...K}{\operatorname{argmax}} \hat{\mu}_a(\tau)$

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

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau_\delta]}{\ln(1/\delta)} = T_\star(\boldsymbol{\mu}).$$

**Why?**

$$\tau_\delta \simeq \inf\left\{ t \in \mathbb{N}_\star : t \times \inf_{\boldsymbol{\lambda} \in \operatorname{Alt}(\boldsymbol{\mu})} \sum_{a=1}^K (w_\star(\boldsymbol{\mu}))_a \operatorname{kl}(\mu_a, \lambda_a) > \beta(t, \delta) \right\}$$

# An asymptotically optimal algorithm

**Theorem**

The Track-and-Stop strategy, that uses

▶ the Tracking sampling rule

▶ the GLRT stopping rule with

$$\beta(t,\delta) \simeq \ln\left(\frac{K-1}{\delta}\right) + 2\ln\ln\left(\frac{K-1}{\delta}\right) + 6\ln(\ln(t))$$

▶ and recommends $B_\tau = \underset{a=1\dots K}{\operatorname{argmax}} \hat{\mu}_a(\tau)$

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

$$\limsup_{\delta\to 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau_\delta]}{\ln(1/\delta)} = T_\star(\boldsymbol{\mu}).$$

**Why?**

$$\tau_\delta \simeq \inf\left\{t \in \mathbb{N}_\star : t \times T_\star^{-1}(\boldsymbol{\mu}) > \beta(t,\delta)\right\}$$

▶ $\boldsymbol{\mu}_1 = [0.5\ 0.45\ 0.43\ 0.4]$, such that

$$w_\star(\boldsymbol{\mu}_1) = [0.417\ 0.390\ 0.136\ 0.057]$$

▶ $\boldsymbol{\mu}_2 = [0.3\ 0.21\ 0.2\ 0.19\ 0.18]$, such that

$$w_\star(\boldsymbol{\mu}_2) = [0.336\ 0.251\ 0.177\ 0.132\ 0.104]$$

In practice, set the threshold to $\beta(t, \delta) = \ln\left(\frac{\ln(t)+1}{\delta}\right)$.

|  | Track-and-Stop | GLRT-SE | KL-LUCB | KL-SE |
|---|---|---|---|---|
| $\boldsymbol{\mu}_1$ | 4052 | 4516 | 8437 | 9590 |
| $\boldsymbol{\mu}_2$ | 1406 | 3078 | 2716 | 3334 |

Table: Expected number of draws $\mathbb{E}_{\boldsymbol{\mu}}[\tau_\delta]$ for $\delta = 0.1$, averaged over $N = 3000$ experiments.

- ▶ TaS: a lower-bound inspired algorithm
- ▶ KL-UCB versus TaS: very different sampling rules!



- ▶ Two recent improvements of the Tracking sampling rule:
  - → relax the need for computing the optimal weights at every round [Ménard, 2019]
  - → get rid of the forced exploration by using Upper Confidence Bounds [Degenne et al., 2019]

# Open problems

▶ Unlike previously mentioned strategies, the sampling rule of TaS is **anytime**, i.e. does not depend on a budget $T$ or a risk parameter $\delta$.
   - → is it also a good strategy in the fixed budget setting?
   - → can control $\mathbb{E}[r_\nu(\mathtt{TaS}, t)]$ for any $t$?

▶ Fixed-budget setting: no exactly matching upper and lower bound best lower bounds: [Carpentier and Locatelli, 2016]

▶ Top-Two Thompson Sampling [Russo, 2016]: a Bayesian (anytime) strategy that is optimal in a different (Bayesian, asymptotic) sense
   - → can we obtain frequentist guarantees for this algorithm?

# Beyond Best Arm Identification

# Outline

**Context**: exponential family bandit model

$$\nu \leftrightarrow \boldsymbol{\mu} = (\mu_1, \ldots, \mu_K) \in \mathcal{I}^K$$

**Goal:** Given $M$ regions of $\mathcal{I}^K$, $\mathcal{R}_1, \ldots, \mathcal{R}_M$, the goal is to identify one region to which $\boldsymbol{\mu}$ belongs.

**Formalization**: build a
- → sampling rule $(A_t)$
- → stopping rule $\tau$
- → recommendation rule $\hat{\imath}_\tau \in \{1, \ldots, M\}$

such that, for some risk parameter $\delta$,

$$\mathbb{P}_{\boldsymbol{\mu}} \left( \boldsymbol{\mu} \notin \mathcal{R}_{\hat{\imath}_\tau} \right) \leq \delta \quad \text{and} \quad \mathbb{E}_{\boldsymbol{\mu}}[\tau] \text{ is small.}$$

# Active Identification in Bandit Models

**Context**: exponential family bandit model

$$\nu \leftrightarrow \boldsymbol{\mu} = (\mu_1, \ldots, \mu_K) \in \mathcal{I}^K$$

**Goal:** Given $M$ regions of $\mathcal{I}^K$, $\mathcal{R}_1, \ldots, \mathcal{R}_M$, the goal is to identify one region to which $\boldsymbol{\mu}$ belongs.

**Two cases:**

▶ $\mathcal{R}_1, \ldots, \mathcal{R}_M$ form a partition : a lower-bound inspired sampling rule + a GLRT stopping rule essentially works
⚠ computing the optimal allocation may be difficult
[Juneja and Krishnasamy, 2019, Kaufmann and Koolen, 2018]

▶ $\mathcal{R}_1, \ldots, \mathcal{R}_M$ are overlapping regions : a Track-and-Stop approach does not always work [Degenne and Koolen, 2019]

# Outline

**Anomaly detection**: given a threshold $\theta$:

▶ find all the arms whose mean is below $\theta$ [Locatelli et al., 2016]

▶ find whether there is an arm with mean below $\theta$
[Kaufmann et al., 2018]

**Phase I clinical trial**: find the arm with mean closest to the threshold...
with increasing means. [Garivier et al., 2017]



$$\mathcal{R}_a = \left\{ \boldsymbol{\lambda} \in \mathtt{Inc} : a = \operatorname*{argmin}_i |\theta - \lambda_i| \right\}$$

# Bandit and games

Find the best move at the root of a game tree by actively sampling its leaves. $s_\star = \underset{s \in \mathcal{C}(s_0)}{\arg\max} \, V_s$.



$$V_s = \begin{cases} \mu_s & \text{if } s \in \mathcal{L}, \\ \max_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MAX node}, \\ \min_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MIN node}. \end{cases}$$

➜ more details later

# Bandit for Optimization in a Larger Space

# Outline

# Bandit problems from an optimization perspective

$$f : \{1, \dots, K\} \longrightarrow \mathbb{R} \qquad \max_{a=1,\dots,K} f(a) \ ?$$

**Sequential evaluations:** at time $t$, choose $A_t \in \{1, \dots, K\}$, observe

$$X_t \sim \nu_{A_t} \quad \text{where} \quad \nu_a \quad \text{has mean} \quad f(a).$$

After $T$ observations,

### Minimize the cumulative regret

$$\text{minimize} \ \ \mathbb{E}\left[\sum_{t=1}^{T}(f(a_\star) - f(A_t))\right]$$

### Minimize the simple regret (optimization error)

If $B_T$ is a guess of the $\mathrm{argmax}$

$$\text{minimize} \ \ \mathbb{E}\left[f(a_\star) - f(B_T)\right]$$

$$f : \{1, \ldots, K\} \longrightarrow \mathbb{R} \qquad \max_{a=1,\ldots,K} f(a) \; ?$$

**Sequential evaluations:** at time $t$, choose $A_t \in \{1, \ldots, K\}$, observe

$$X_t \sim \nu_{A_t} \quad \text{where} \quad \nu_a \quad \text{has mean} \quad f(a).$$



*sequential optimization of a discrete function*
*based on noisy observations*

$$f : \mathcal{X} \longrightarrow \mathbb{R} \qquad \max_{x \in \mathcal{X}} \ f(x) \ ?$$

**Sequential evaluations:** at time $t$, choose $x_t \in \mathcal{X}$, observe

$$y_t = f(x_t) + \epsilon_t$$



*sequential optimization of a black-box function*
*based on noisy observations*

$$f : \mathcal{X} \longrightarrow \mathbb{R} \qquad \max_{x \in \mathcal{X}} f(x) \text{ ?}$$

**Sequential evaluations:** at time $t$, choose $x_t \in \mathcal{X}$, observe

$$y_t = f(x_t) + \epsilon_t$$

After $T$ observations,

### Minimize the cumulative regret

$$\text{minimize } \mathbb{E}\left[\sum_{t=1}^{T}(f(x_\star) - f(x_t))\right]$$

### Minimize the simple regret (optimization error)

If $z_T$ is a guess of the $\mathrm{argmax}$

$$\text{minimize } \mathbb{E}\left[f(x_\star) - f(z_T)\right]$$

# Black Box Optimization

- learning based on (costly, noisy) function evaluations only!
- no access to gradients of $f$

$$x \longrightarrow \boxed{f} \longrightarrow y$$

## Examples of function $f$

- a costly PDE solver (numerical experiments)
- a simulator of the effect of a chemical compound (drug discovery)
- training and evaluation a neural network
  (hyper-parameter optimization)

How to choose the next querry?

How to choose the next querry?

# Outline

# Hierarchical partitioning

Bandit in metric spaces [Kleinberg et al., 2008]
$\mathcal{X}$-armed bandits [Bubeck et al., 2011]

**Idea:** Partition the space, and adaptively choose in which cell to sample



▶ For any **depth** $h$, $\mathcal{X}$ is partitioned in $K^h$ cells $(\mathcal{P}_{h,i})_{0 \leq K^h - 1}$.

▶ $K$-ary tree $\mathcal{T}$ where depth $h = 0$ is the whole $\mathcal{X}$.

# Hierarchical Optimistic Optimization (HOO)

**Assumptions (given a metric $\ell(x, y)$)**

- $f(x_\star) - f(y) \leq f(x_\star) - f(x) + \max\{f(x_\star) - f(x); \ell(x, y)\}$
- $\sup_{(x,y)\in\mathcal{P}_{h,i}} \ell(x, y) \leq \nu\rho^h$.

**Idea:** Use Upper-Confidence Bounds on the maximum values of the function in each cell to guide exploration



$$U_{(h,i)}(t) = \hat{\mu}_{(h,i)}(t) + \sqrt{\frac{2\ln(t)}{N_{(h,i)}(t)}} + \nu\rho^h$$

$$B_{(h,i)}(t) = \min\left\{U_{(h,i)}(t); B_{(h+1,2i)}(t); \right.$$
$$\left. B_{(h+1,2i+1)}(t)\right\}$$

[Bubeck et al., 2011]

# Results

## Cumulative regret of HOO

$$\mathbb{E}\left[\sum_{t=1}^{T}(f(x_\star) - f(x_t))\right] \leq C T^{\frac{d+1}{d+2}}(\ln(T))^{\frac{1}{d+2}}$$

for some near-optimality dimension $d$.

➜ how to turn this into an optimizer?

$z_T \sim \mathcal{U}(x_1, \ldots, x_T)$, then

$$\mathbb{E}[f(x_\star) - f(z_T)] = \frac{\mathcal{R}(\texttt{HOO},T)}{T} \leq C\left(\frac{\ln(T)}{T}\right)^{\frac{1}{d+2}}$$



*a tree built by HOO*

► **Many variants!**

DOO, SOO [Munos, 2011], StoSOO [Valko et al., 2013], POO
[Grill et al., 2015], GPO [Shang et al., 2019]...

# Outline

# Gaussian Process Regression

**Assumption.** the function $f$ is drawn from some Gaussian Process :

$$f \sim \mathcal{GP}(0, k(x, y)).$$

i.e. for any distinct points $x_1, \ldots, x_\ell$ in $\mathcal{X}$,

$$\begin{pmatrix} f(x_1) \\ f(x_2) \\ \ldots \\ f(x_\ell) \end{pmatrix} \sim \mathcal{N}(0, K) \quad \text{where} \quad K = (k(x_i, x_j))_{1 \leq i,j \leq \ell}$$

## Bayesian inference

Given some (possibly noisy) observations of $f$ in $x_1, \ldots, x_t$, the posterior on all the function value in any point is Gaussian

$$f(y)|x_1, \ldots, x_t \sim \mathcal{N}\left(\mu_t(y), \sigma_t(y)^2\right)$$

[Rasmussen and Williams, 2005]

➜ use the current GP posterior to pick the new point to select



**Example:** $[\mu_t(y) \pm \beta\sigma_t(y)]$ is a kind of confidence interval on $f(y)$.

# GP-UCB

**GP-UCB** [Srinivas et al., 2012] selects at round $t+1$

$$x_{t+1} = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \ \mu_t(x) + \sqrt{\beta(t,\delta)}\sigma_t(x)$$



➜ Bayesian and frequentist guarantees in terms of cumulative regret for different $\beta(t,\delta)$: $\mathbb{P}\left(\mathcal{R}_T(\texttt{GP-UCB},T) \leq C\sqrt{T\beta(T,\delta)\gamma_T}\right) \geq 1-\delta$.

# BO Algorithms

More generally, many Bayesian Optimization algorithms optimize an acquisition function that depends on the posterior and select

$$x_{t+1} = \underset{x \in \mathcal{X}}{\text{argmax}} \; \alpha_t(x)$$

Many other acquisition functions: [Shahriari et al., 2016]

▶ Expected improvement

▶ Probability of improvement

▶ Entropy Search ...

**Remark:** optimization the acquisition function is another (non-trivial) optimization problem!

▶ Thompson Sampling?

# Bandit Tools for Planning in Games

# Outline

**Goal:** decide for the next move based on evaluation of possible trajectories in the game, ending with a random evaluation.

**A famous bandit approach:** [Kocsis and Szepesvári, 2006]
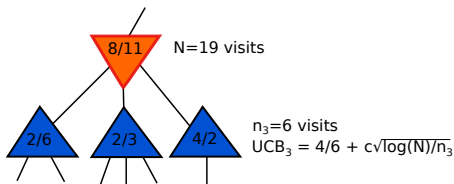→ use UCB in each node to decide the next children to explore
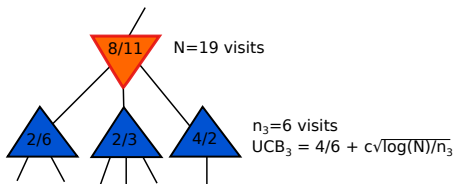
# The UCT algorithm

$N(s)$ : number of visits of node $s$
$S(s)$ : number of visits finishing ending with the root player winning

## UCT in a MaxMin Tree

In a MAX node $s$ ($=$ root player move), go towards the children

$$\underset{c \in \mathcal{C}(s)}{\mathrm{argmax}} \ \frac{S(c)}{N(c)} + c\sqrt{\frac{\ln(N(s))}{N(c)}}$$



8/11    N=19 visits

2/6    2/3    4/2

$n_3$=6 visits
$UCB_3 = 4/6 + c\sqrt{\log(N)/n_3}$
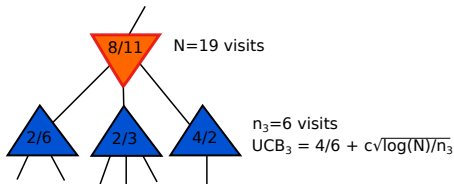
# The UCT algorithm

$N(s)$ : number of visits of node $s$
$S(s)$ : number of visits finishing ending with the root player winning

## UCT in a MaxMin Tree

In a MIN node $s$ (= adversary move), go towards the children

$$\operatorname*{argmin}_{c \in \mathcal{C}(s)} \frac{S(c)}{N(c)} - c\sqrt{\frac{\ln(N(s))}{N(c)}}$$



8/11    N=19 visits

2/6    2/3    4/2

$n_3 = 6$ visits
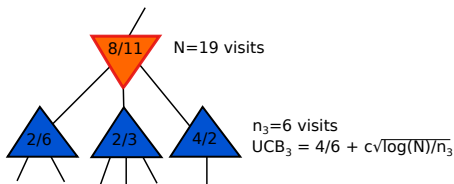$UCB_3 = 4/6 + c\sqrt{\log(N)/n_3}$

# The UCT algorithm

$N(s)$ : number of visits of node $s$
$S(s)$ : number of visits finishing ending with the root player winning

## UCT in a MaxMin Tree

In a MAX node $s$ ($=$ root player move), go towards the children

$$\underset{c \in \mathcal{C}(s)}{\mathrm{argmax}} \ \frac{S(c)}{N(c)} + c \sqrt{\frac{\ln(N(s))}{N(c)}}$$



8/11   N=19 visits

2/6   2/3   4/2

$n_3$=6 visits
$UCB_3 = 4/6 + c\sqrt{\log(N)/n_3}$

# The UCT algorithm

$N(s)$ : number of visits of node $s$

$S(s)$ : number of visits finishing ending with the root player winning

### UCT in a MaxMin Tree

In a MAX node $s$ (= root player move), go towards the children

$$\underset{c \in \mathcal{C}(s)}{\operatorname{argmax}} \ \frac{S(c)}{N(c)} + c\sqrt{\frac{\ln(N(s))}{N(c)}}$$



N=19 visits

$n_3$=6 visits

$UCB_3 = 4/6 + c\sqrt{\log(N)/n_3}$

$+$ first Go AI based on variants of UCT ($+$ heuristics)

# The UCT algorithm

$N(s)$ : number of visits of node $s$
$S(s)$ : number of visits finishing ending with the root player winning

## UCT in a MaxMin Tree

In a MAX node $s$ (= root player move), go towards the children

$$\operatorname*{argmax}_{c \in \mathcal{C}(s)} \frac{S(c)}{N(c)} + c\sqrt{\frac{\ln(N(s))}{N(c)}}$$



8/11  N=19 visits

2/6  2/3  4/2

$n_3$=6 visits
$UCB_3$ = 4/6 + $c\sqrt{\log(N)/n_3}$

— UCT is *not* based on statistically-valid confidence intervals
— no sample complexity guarantees
— should we really minimize rewards?

# Outline

A fixed MAXMIN game tree $\mathcal{T}$, with leaves $\mathcal{L}$.

▽ MAX node ($=$ root player move)

▲ MIN node ($=$ adversary move)

● Leaf $\ell$: stochastic oracle $\mathcal{O}_\ell$ that evaluates the position
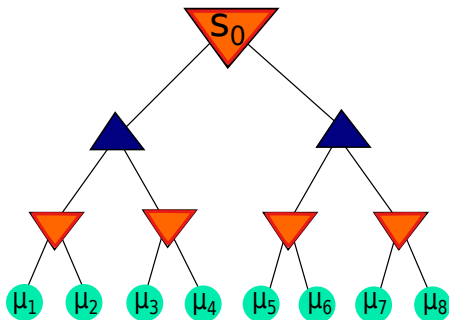
At round $t$ a **MCTS algorithm**:

- picks a path down to a leaf $L_t$
- get an evaluation of this leaf $X_t \sim \mathcal{O}_{L_t}$

Assumption: i.i.d. sucessive evaluations, $\mathbb{E}_{X \sim \mathcal{O}_\ell}[X] = \mu_\ell$

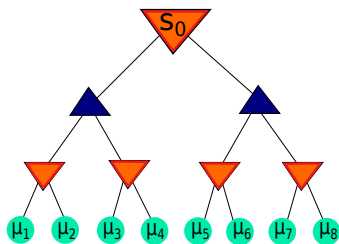A MCTS algorithm should find the best move at the root:

$$V_s = \begin{cases} \mu_s & \text{if } s \in \mathcal{L}, \\ \max_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MAX node,} \\ \min_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MIN node.} \end{cases}$$

$$s^* = \underset{s \in \mathcal{C}(s_0)}{\operatorname{argmax}} V_s$$

# A structured BAI problem

MCTS algorithm: $(L_t, \tau, \hat{s}_\tau)$, where

- $L_t$ is the sampling rule
- $\tau$ is the stopping rule
- $\hat{s}_\tau \in \mathcal{C}(s_0)$ is the recommendation rule



**Goal:** an $(\epsilon, \delta)$-PAC MCTS algorithm:

$$\mathbb{P}(\boldsymbol{V}_{\hat{s}_\tau} \geq \boldsymbol{V}_{s^*} - \boldsymbol{\epsilon}) \geq \boldsymbol{1} - \boldsymbol{\delta}$$

with a small sample complexity $\boldsymbol{\tau}$.                    [Teraoka et al., 2014]

MCTS algorithm: $(L_t, \tau, \hat{s}_\tau)$, where

- $L_t$ is the sampling rule
- $\tau$ is the stopping rule
- $\hat{s}_\tau \in \mathcal{C}(s_0)$ is the recommendation rule



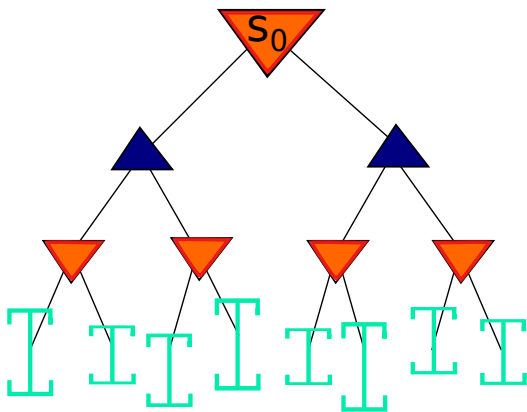**Idea**: use LUCB on the depth-one nodes

→ requires confidence intervals on the values $(V_s)_{s \in \mathcal{C}_0}$
→ requires to identify a leaf to sample starting from $s \in \mathcal{C}_0$

Using the samples collected for the leaves, one can build, for $\ell \in \mathcal{L}$,

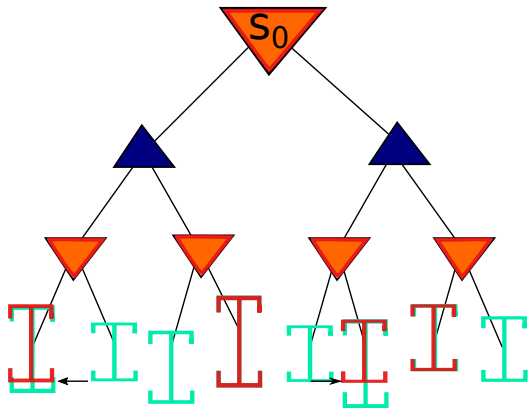$$[\text{LCB}_\ell(t), \text{UCB}_\ell(t)] \quad \text{a confidence interval on } \mu_\ell$$



**Idea:** Propagate these confidence intervals up in the tree

MAX node:

$$\mathrm{UCB}_s(t) = \max_{c \in \mathcal{C}(s)} \mathrm{UCB}_c(t) \quad \mathrm{LCB}_s(t) = \max_{c \in \mathcal{C}(s)} \mathrm{LCB}_c(t)$$



**Idea:** Propagate these confidence intervals up in the tree

MAX node:

$$\text{UCB}_s(t) = \max_{c \in \mathcal{C}(s)} \text{UCB}_c(t) \quad \text{LCB}_s(t) = \max_{c \in \mathcal{C}(s)} \text{LCB}_c(t)$$



**Idea:** Propagate these confidence intervals up in the tree

MIN node:

$$\text{UCB}_s(t) = \min_{c \in \mathcal{C}(s)} \text{UCB}_c(t) \quad \text{LCB}_s(t) = \min_{c \in \mathcal{C}(s)} \text{LCB}_c(t)$$



**Idea:** Propagate these confidence intervals up in the tree

$\ell_s(t)$: representative leaf of internal node $s \in \mathcal{T}$.



**Idea:** alternate optimistic/pessimistic moves starting from $s$

▶ run a BAI algorithm on the depth-on nodes

$$\rightarrow \text{selects} \quad R_t \in \mathcal{C}_0$$

▶ sample the representative leaf associated to that node:

$$L_t = \ell_{R_t}(t)$$

($\simeq$ from depth one, run UCT based on *statistically valid* CIs)

▶ update the confidence intervals

▶ stop when the BAI algorithm tell us to

▶ recommand the depth-one node chosen by the BAI algorithm

# Theoretical guarantees

For some exploration function $\beta$, define

$$\mathrm{LCB}_\ell(t) = \hat{\mu}_\ell(t) - \sqrt{\frac{\beta(N_\ell(t), \delta)}{2N_\ell(t)}},$$

$$\mathrm{UCB}_\ell(t) = \hat{\mu}_\ell(t) + \sqrt{\frac{\beta(N_\ell(t), \delta)}{2N_\ell(t)}}.$$

## Theorem [Kaufmann et al., 2018]

Choosing

$$\beta(s, \delta) \simeq \ln\left(\frac{|\mathcal{L}| \ln(s)}{\delta}\right),$$

LUCB-MCTS and UGapE-MCTS are $(\epsilon, \delta)$-PAC and

$$\mathbb{P}\left(\tau = O\left(H_\epsilon^*(\boldsymbol{\mu}) \ln\left(\frac{1}{\delta}\right)\right)\right) \geq 1 - \delta$$
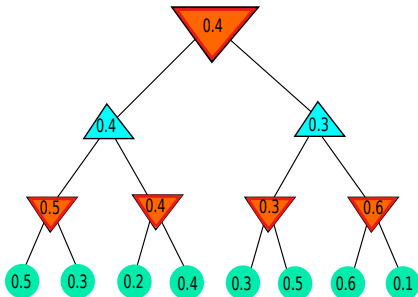
for UGapE-MCTS.

$$H_\epsilon^*(\boldsymbol{\mu}) := \sum_{\ell \in \mathcal{L}} \frac{1}{\Delta_\ell^2 \vee \Delta_*^2 \vee \epsilon^2}$$

where

$$\Delta_* := V(s^*) - V(s_2^*)$$

$$\Delta_\ell := \max_{s \in \mathtt{Ancestors}(\ell) \setminus \{s_0\}} \left| V_{\mathtt{Parent}(s)} - V_s \right|$$



*(slightly improved complexity in the work of [Huang et al., 2017])*

▶ Optimal and efficient algorithms for solving best action identification in a maxmin tree...

▶ ... and other generic active identification problems

▶ UCT versus BAI-MCTS on large-scale problem?

▶ Best arm identification and regret minimization are two different problems that require different sampling rules

▶ Upper and Lower Confidence Bounds are useful in both settings

▶ Optimal algorithms for BAI are inspired by the lower bounds (cf. structured bandits)

▶ Tools for BAI $\rightarrow$ more general Active Identification problems

▶ Bandit tools inspire methods for sequential optimization in large spaces (games trees or continuous spaces)

# That's all!



now you're ready to pull the right arm ;-)

Audibert, J.-Y., Bubeck, S., and Munos, R. (2010).
Best Arm Identification in Multi-armed Bandits.
In *Proceedings of the 23rd Conference on Learning Theory.*

Bubeck, S., Munos, R., and Stoltz, G. (2009).
Pure Exploration in multi-armed bandit problems.
In *Algorithmtic Learning.*

Bubeck, S., Munos, R., Stoltz, G., and Szepesvári, C. (2011).
X-armed bandits.
*Journal of Machine Learning Research*, 12:1587–1627.

Carpentier, A. and Locatelli, A. (2016).
Tight (lower) bounds for the fixed budget best arm identification bandit problem.
In *Conference on Learning Theory (COLT).*

Chernoff, H. (1959).
Sequential design of Experiments.
*The Annals of Mathematical Statistics*, 30(3):755–770.

Degenne, R., Koolen, W., and Ménard, P. (2019).
Non asymptotic pure exploration by solving games.
*arXiv:1906.10431.*

Degenne, R. and Koolen, W. M. (2019).
Pure Exploration with Multiple Correct Answers.
*arXiv:1902.03475.*

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 (to appear).*

Garivier, A., Ménard, P., and Rossi, L. (2017).
Thresholding bandit for dose-ranging: The impact of monotonicity.
*arXiv:1711.04454.*

Grill, J., Valko, M., and Munos, R. (2015).
Black-box optimization of noisy functions with unknown smoothness.
In *Advances in Neural Information Processing Systems.*

Huang, R., Ajallooeian, M. M., Szepesvári, C., and Müller, M. (2017).
Structured best arm identification with fixed confidence.
In *International Conference on Algorithmic Learning Theory (ALT).*

Jamieson, K., Malloy, M., Nowak, R., and Bubeck, S. (2014).
lil'UCB: an Optimal Exploration Algorithm for Multi-Armed Bandits.
In *Proceedings of the 27th Conference on Learning Theory*.

Jun, K.-W. and Nowak, R. (2016).
Anytime exploration for multi-armed bandits using confidence information.
In *International Conference on Machine Learning (ICML)*.

Juneja, S. and Krishnasamy, S. (2019).
Sample complexity of partition identification using multi-armed bandits.
In *Conference on Learning Theory (COLT)*.

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)*.

Karnin, Z., Koren, T., and Somekh, O. (2013).
Almost optimal Exploration in multi-armed bandits.
In *International Conference on Machine Learning (ICML)*.

Kaufmann, E., Cappé, O., and Garivier, A. (2016).
On the Complexity of Best Arm Identification in Multi-Armed Bandit Models.
*Journal of Machine Learning Research*, 17(1):1–42.

📄 Kaufmann, E. and Koolen, W. (2018).
Mixture martingales revisited with applications to sequential tests and confidence intervals.
*arXiv:1811.11419.*

📄 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).*

📄 Kleinberg, R., Slivkins, A., and Upfal, E. (2008).
Multi-armed bandits in metric spaces.
In *Proceedings of the 40th Annual ACM Symposium on Theory of Computing.*

📄 Kocsis, L. and Szepesvári, C. (2006).
Bandit based monte-carlo planning.
In *Proceedings of the 17th European Conference on Machine Learning*, ECML'06, pages 282–293, Berlin, Heidelberg. Springer-Verlag.

📄 Locatelli, A., Gutzeit, M., and Carpentier, A. (2016).
An optimal algorithm for the thresholding bandit problem.
In *International Conference on Machine Learning (ICML).*

📄 Ménard, P. (2019).

Gradient ascent for active exploration in bandit problems.
*arXiv:1905.08165*.

Munos, R. (2011).
Optimistic optimization of a deterministic function without the knowledge of its smoothness.
In *Advances in Neural Information Processing Systems*.

Rasmussen, C. E. and Williams, C. K. I. (2005).
*Gaussian Processes for Machine Learning*.
The MIT Press.

Russo, D. (2016).
Simple bayesian algorithms for best arm identification.
In *Proceedings of the 29th Conference On Learning Theory*.

Shahriari, B., Swersky, K., Wang, Z., Adams, R. P., and de Freitas, N. (2016).
Taking the human out of the loop: A review of bayesian optimization.
*Proceedings of the IEEE*, 104(1):148–175.

Shang, X., Kaufmann, E., and Valko, M. (2019).
General parallel optimization a without metric.
In *Algorithmic Learning Theory (ALT)*.

Srinivas, N., Krause, A., Kakade, S., and Seeger, M. (2012).
Information-Theoretic Regret Bounds for Gaussian Process Optimization in the Bandit Setting.
*IEEE Transactions on Information Theory*, 58(5):3250–3265.

Teraoka, K., Hatano, K., and Takimoto, E. (2014).
Efficient sampling method for Monte Carlo tree search problem.
*IEICE Transactions on Infomation and Systems*, pages 392–398.

Valko, M., Carpentier, A., and Munos, R. (2013).
Stochastic simultaneous optimistic optimization.
In *Proceedings of the 30th International Conference on Machine Learning (ICML)*.