

Contributions to the Optimal Solution of Several Bandit Problems

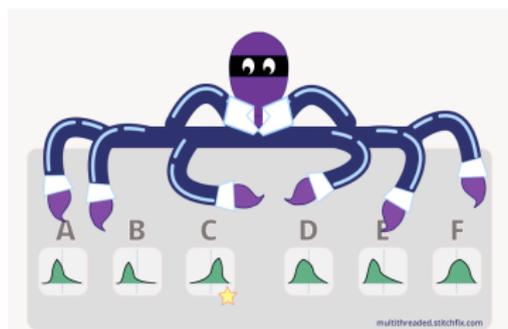
Emilie Kaufmann



HDR defense
November 13th, 2020

The stochastic MAB model

Arms = probability distributions an agent can choose from:



In each round t , the agent

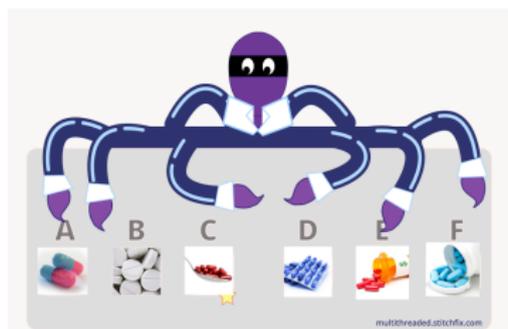
- selects arm $A_t \in \{1, \dots, A\}$
- observes a sample $X_t \sim \nu_{A_t}$
independent from past data

sequential protocol:

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

The stochastic MAB model

Arms = probability distributions an agent can choose from:



In each round t , the agent

- selects arm $A_t \in \{1, \dots, A\}$
- observes a sample $X_t \sim \mathcal{B}(\mu_{A_t})$
independent from past data

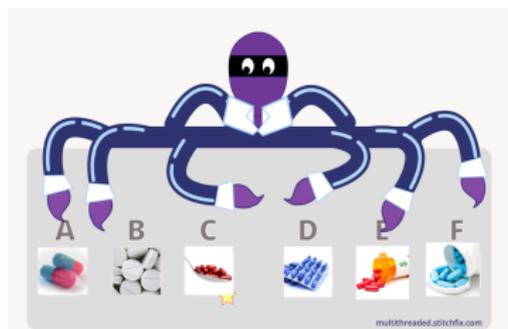
sequential protocol:

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

Assumption (in this work): arms are simple distribution parameterized by their means (e.g. **Bernoulli**, exponential families)

The stochastic MAB model

Arms = probability distributions an agent can choose from:



In each round t , the agent

- selects arm $A_t \in [A]$
- observes a sample $X_t \sim \mathcal{B}(\mu_{A_t})$
independent from past data

sequential protocol:

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

Assumption (in this work): arms are simple distribution parameterized by their means (e.g. **Bernoulli**, exponential families)

Notation: $\nu_a = \nu_{\mu_a}$, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_A) \in \mathcal{I}^A$.

rewards maximization... with a twist

- feedback \neq reward [Ch. 1]
- structured bandits [Ch. 1]
- multi-player bandits [Ch. 2]

pure exploration

- a generic stopping rule for active identification [Ch. 3]
- the complexity of best arm identification [Ch. 4]
- two MCTS-related examples [Ch. 5]

Common emphasis on designing *optimal* algorithms

rewards maximization... with a twist

- feedback \neq reward 
- structured bandits 
- multi-player bandits 

pure exploration

- a generic stopping rule for active identification
- the complexity of best arm identification
- two MCTS-related examples



Common emphasis on designing *optimal* algorithms

Research *motivated by* some applications

rewards maximization... with a twist

- feedback \neq reward 
- structured bandits 
- multi-player bandits 

pure exploration

- a generic stopping rule for active identification
- the complexity of best arm identification
- two MCTS-related examples



Common emphasis on designing *optimal* algorithms

Research *motivated by* some applications

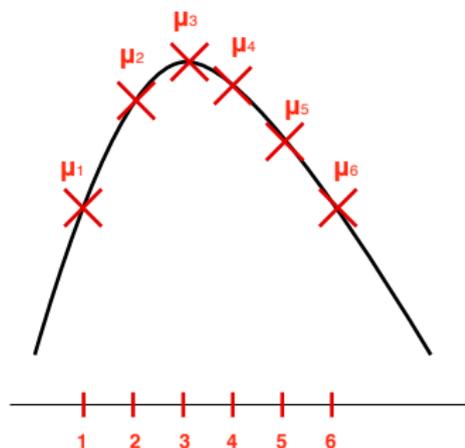
- Lower bounds...
and how they inspire algorithms
- Mixture martingales for new deviation inequalities
- Recent tools for the analysis of Thompson Sampling
[Agrawal and Goyal, 2013, Russo, 2016]

- 1 Thompson Sampling for a Structured Bandit Problem
- 2 The Complexity of Pure Exploration
- 3 Thompson Sampling for Pure Exploration?

- 1 Thompson Sampling for a Structured Bandit Problem
- 2 The Complexity of Pure Exploration
- 3 Thompson Sampling for Pure Exploration?

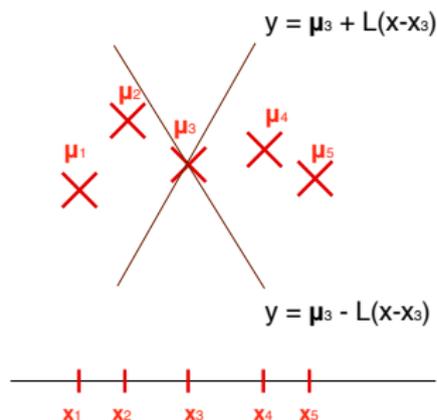
Structured bandits

- Classical bandits: $\mu = (\mu_1, \dots, \mu_A) \in \mathcal{I}^A$
 - Structured bandits: $\mu = (\mu_1, \dots, \mu_A) \in \mathcal{S} \subset \mathcal{I}^A$
- can we exploit the knowledge of \mathcal{S} to gain more reward?



unimodal bandit

[Combes and Proutière, 2014]



Lipschitz bandit

[Magureanu et al., 2014]

Lower Bounds can help

In each round t , the agent

- selects arm $A_t \in [A]$, observes a reward $X_t \sim \nu_{A_t}$

Goal: maximize the expected total reward \leftrightarrow minimize the regret

$$\begin{aligned}\mathcal{R}_\mu(\mathcal{A}, T) &= \mu_* T - \mathbb{E}_\mu \left[\sum_{t=1}^T X_t \right] \\ &= \sum_{a \in [A]} (\mu_* - \mu_a) \mathbb{E}_\mu [N_a(T)]\end{aligned}$$

$N_a(T)$: number of selections of arm a up to round T .

Theorem [Graves and Lai, 1997] (Theorem 1.8 in the HDR document)

Let \mathcal{A} be such that $\forall \mu \in \mathcal{S}, \forall \alpha \in (0, 1], \mathcal{R}_\mu(\mathcal{A}, T) = o(T^\alpha)$.

$$\forall \mu \in \mathcal{S}, \lim_{T \rightarrow \infty} \frac{\mathcal{R}_\mu(\mathcal{A}, T)}{\log(T)} \geq C_{\mathcal{S}}(\mu).$$

$\rightarrow \mathcal{A}$ is **asymptotically optimal** if $\mathcal{R}_\mu(\mathcal{A}, T) = C_{\mathcal{S}}(\mu) \log(T) + o(\log(T))$

Lower Bounds can help

In each round t , the agent

- selects arm $A_t \in [A]$, observes a reward $X_t \sim \nu_{A_t}$

Goal: maximize the expected total reward \leftrightarrow minimize the regret

$$\begin{aligned}\mathcal{R}_\mu(\mathcal{A}, T) &= \mu_{a_*} T - \mathbb{E}_\mu \left[\sum_{t=1}^T X_t \right] \\ &= \sum_{a \in [A]} (\mu_{a_*} - \mu_a) \mathbb{E}_\mu [N_a(T)]\end{aligned}$$

$N_a(T)$: number of selections of arm a up to round T .

Theorem [Graves and Lai, 1997] (Theorem 1.8 in the HDR document)

Let \mathcal{A} be such that $\forall \mu \in \mathcal{S}, \forall \alpha \in (0, 1], \mathcal{R}_\mu(\mathcal{A}, T) = o(T^\alpha)$.

$$\forall \mu \in \mathcal{S}, \lim_{T \rightarrow \infty} \frac{\mathcal{R}_\mu(\mathcal{A}, T)}{\log(T)} \geq C_{\mathcal{S}}(\mu).$$

$\rightarrow \mathcal{A}$ is **asymptotically optimal** if $\mathcal{R}_\mu(\mathcal{A}, T) = C_{\mathcal{S}}(\mu) \log(T) + o(\log(T))$

Lower bounds can help

$C_S(\boldsymbol{\mu})$ features the Kullback-Leibler divergence $d(\boldsymbol{\mu}, \boldsymbol{\mu}') := \text{KL}(\nu_{\boldsymbol{\mu}}, \nu_{\boldsymbol{\mu}'})$

- $\mathcal{S} = \mathcal{I}^A$, $C_S(\boldsymbol{\mu}) = \sum_{a=1}^A \frac{\mu_{\star} - \mu_a}{d(\mu_a, \mu_{\star})}$ [Lai and Robbins, 1985]
- in general, $C_S(\boldsymbol{\mu})$ has **no closed-form expression** (solution of a complex optimization problem)

Special case [Combes and Proutière, 2014]

$\boldsymbol{\mu}$ is **unimodal with respect to a graph** $G = ([A], E)$: for all $a \in [A]$ there exists an increasing path to the optimal arm a_{\star} :

$$(a_1 = a, \dots, a_{m_a} = a_{\star}) : (a_i, a_{i+1}) \in E \text{ and } \mu_{a_i} < \mu_{a_{i+1}}.$$

For graphical unimodal bandits,

$$C_S(\boldsymbol{\mu}) = \sum_{a \in \mathcal{N}_G(a_{\star})} \frac{\mu_{\star} - \mu_a}{d(\mu_a, \mu_{\star})} \quad \mathcal{N}_G(a_{\star}) = \{a : (a, a_{\star}) \in E\}$$

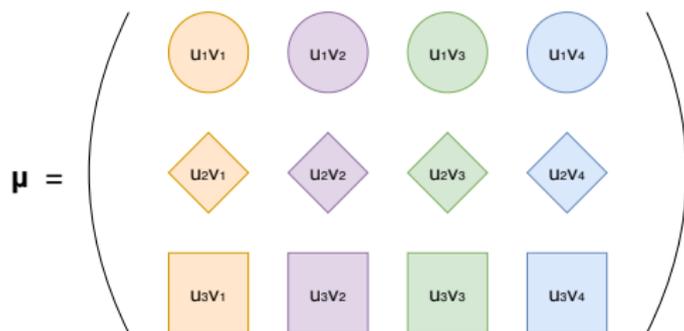
→ an optimal algorithm focusses on **neighbors of the optimal arm**

Solving Rank-One Bandits

$$\mathcal{S}_{R1} = \left\{ \boldsymbol{\mu} = (\mu_{(k,\ell)})_{\substack{1 \leq k \leq K \\ 1 \leq \ell \leq L}} \mid \exists \mathbf{u} \in [0, 1]^K, \mathbf{v} \in [0, 1]^L : \mu_{(k,\ell)} = u_k v_\ell \right\}$$

[Katariya et al., 2017]

Example: content optimization with two independent factors



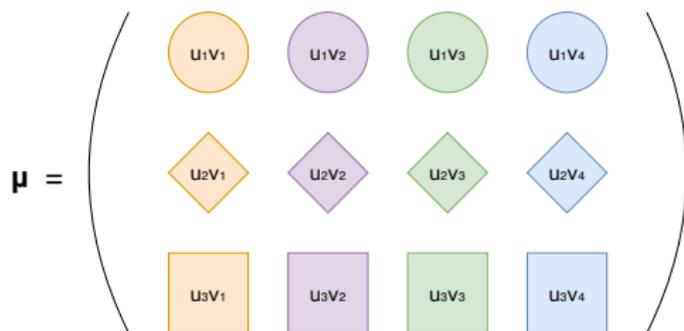
click probability $\mu_{(k,\ell)} = u_k \times v_\ell$

Solving Rank-One Bandits

$$\mathcal{S}_{R1} = \left\{ \boldsymbol{\mu} = (\mu_{(k,\ell)})_{\substack{1 \leq k \leq K \\ 1 \leq \ell \leq L}} \mid \exists \mathbf{u} \in [0, 1]^K, \mathbf{v} \in [0, 1]^L : \mu_{(k,\ell)} = u_k v_\ell \right\}$$

[Katariya et al., 2017]

Example: content optimization with two independent factors



click probability $\mu_{(k,\ell)} = u_k \times v_\ell$

Key observation

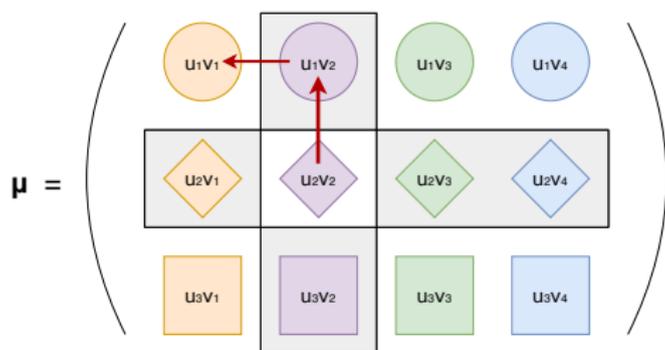
$\boldsymbol{\mu}$ is unimodal with respect to the graph $G_1 = ([K] \times [L], E)$
 $((i, j), (k, \ell)) \in E$ if $(i = k \text{ (x) or } j = \ell)$

Solving Rank-One Bandits

$$\mathcal{S}_{R1} = \left\{ \boldsymbol{\mu} = (\mu_{(k,\ell)})_{\substack{1 \leq k \leq K \\ 1 \leq \ell \leq L}} \mid \exists \mathbf{u} \in [0, 1]^K, \mathbf{v} \in [0, 1]^L : \mu_{(k,\ell)} = u_k v_\ell \right\}$$

[Katariya et al., 2017]

Example: content optimization with two independent factors



click probability $\mu_{(k,\ell)} = u_k \times v_\ell$

Key observation

$\boldsymbol{\mu}$ is unimodal with respect to the graph $G_1 = ([K] \times [L], E)$

$((i, j), (k, \ell)) \in E$ if $(i = k \text{ (x) or } j = \ell)$

Unimodal Thompson Sampling for Rank-One Bandits

Idea: use an optimal algorithm for graphical unimodal bandits

- **Unimodal Thompson Sampling** [Paladino et al., 2017]

UTS with parameter $\gamma \in \{2, 3, \dots\}$ for Bernoulli bandits

In each round $t + 1$:

- compute the empirical leader $B_{t+1} = \operatorname{argmax}_{a \in [A]} \hat{\mu}_a(t)$
- if $\ell_{B_{t+1}}(t + 1) = 0[\gamma]$, select $A_{t+1} = B_{t+1}$ (leader exploration)
- else, draw **posterior samples** for arms in $\mathcal{N}_G(B_{t+1}) \cup \{B_{t+1}\}$:

$$\theta_a(t) \sim \text{Beta}(S_a(t) + 1, N_a(t) - S_a(t) + 1)$$

$$\text{and } A_{t+1} = \operatorname{argmax}_{a \in \mathcal{N}_G(B_{t+1}) \cup \{B_{t+1}\}} \theta_a(t) \quad (\text{TS around the leader})$$

$S_a(t) = \sum_{s=1}^t X_s \mathbb{1}(A_s = a)$: sum of rewards from arm a

$\hat{\mu}_a(t) = S_a(t)/N_a(t)$: empirical mean of arm a

$\ell_b(t) = \sum_{s=1}^t \mathbb{1}(B_s = b)$: number of times arm b has been leader

Unimodal Thompson Sampling for Rank-One Bandits

Idea: use an optimal algorithm for graphical unimodal bandits

- **Unimodal Thompson Sampling** [Paladino et al., 2017]

UTS with parameter $\gamma \in \{2, 3, \dots\}$ for Bernoulli bandits

In each round $t + 1$:

- compute the empirical leader $B_{t+1} = \operatorname{argmax}_{(k,\ell) \in [K] \times [L]} \hat{\mu}_{(k,\ell)}(t)$
- if $\ell_{B_{t+1}}(t + 1) = 0[\gamma]$, select $A_{t+1} = B_{t+1}$ (leader exploration)
- else, draw **posterior samples** for arms in $\mathcal{N}_G(B_{t+1}) \cup \{B_{t+1}\}$:

$$\theta_{(k,\ell)}(t) \sim \text{Beta}(S_{(k,\ell)}(t) + 1, N_{(k,\ell)}(t) - S_{(k,\ell)}(t) + 1)$$

$$\text{and } A_{t+1} = \operatorname{argmax}_{(k,\ell) \in \{(k', B_t^2)\} \cup \{(B_t^1, \ell')\}} \theta_{(k,\ell)}(t) \quad (\text{TS around the leader})$$

$S_a(t) = \sum_{s=1}^t X_s \mathbb{1}(A_s = a)$: sum of rewards from arm a

$\hat{\mu}_a(t) = S_a(t)/N_a(t)$: empirical mean of arm a

$\ell_b(t) = \sum_{s=1}^t \mathbb{1}(B_s = b)$: number of times arm b has been leader

Theorem [Trinh, K., Vernade, Combes, ALT 2020]

Let μ be a unimodal bandit instance with respect to a graph G , with Bernoulli rewards. For all $\gamma \geq 2$, UTS with parameter γ satisfies, for every $\varepsilon > 0$,

$$\mathcal{R}_\mu(\text{UTS}(\gamma), T) \leq (1 + \varepsilon) \sum_{a \in \mathcal{N}_G(a_\star)} \frac{(\mu_\star - \mu_a)}{d(\mu_a, \mu_\star)} \log(T) + C(\mu, \gamma, \varepsilon).$$

- a novel analysis, valid for any leader exploration parameter γ , with $\gamma = 2$ being the best choice in practice
- UTS(γ) is asymptotically optimal for Rank-One bandits (matching the existing lower bound of [Katariya et al., 2017])
- ... and greatly outperforms the previous state-of-the-art

- 1 Thompson Sampling for a Structured Bandit Problem
- 2 The Complexity of Pure Exploration**
- 3 Thompson Sampling for Pure Exploration?

Active Identification in a bandit model

Goal: answer *some question* about the unknown mean vector $\mu = (\mu_1, \dots, \mu_A)$ by adaptively sampling the arms

Input:

- $\mathcal{R} \subseteq \mathcal{I}^A$ a subset that contains μ
- I regions $\mathcal{R}_1, \dots, \mathcal{R}_I$ such that $\mathcal{R} \subseteq \bigcup_{i=1}^I \mathcal{R}_i$

Output: one region \mathcal{R}_i that contains μ .

Active Identification with fixed-confidence

Given a risk parameter $\delta \in (0, 1)$, the goal is to build a

- sampling rule (A_t)
- stopping rule τ
- recommendation rule $\hat{i}_\tau \in [I]$

such that $\mathbb{P}_\mu(\mu \notin \mathcal{R}_{\hat{i}_\tau}) \leq \delta$ and the sample complexity τ is small.

Best Arm Identification

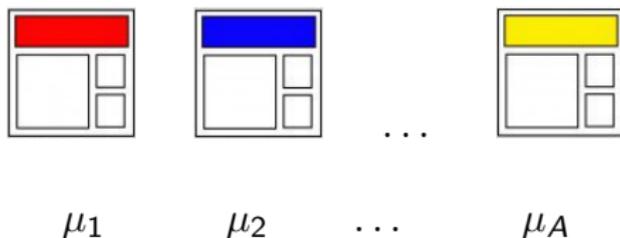
→ Identify the arm with largest mean:

$$\mathcal{R} = \left\{ \boldsymbol{\mu} \in \mathcal{I}^A : \exists a \in [A] : \mu_a > \max_{b \neq a} \mu_b \right\}$$

and $\mathcal{R}_i = \left\{ \boldsymbol{\mu} \in \mathcal{I}^A : \mu_i > \max_{b \neq i} \mu_b \right\}$ for $i \in [A]$

[Even-Dar et al., 2006]

Example: identify the version of a webpage with the largest conversion probability (A/B/C testing)

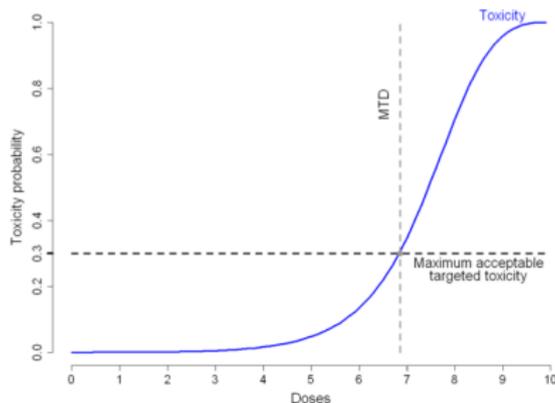


→ Identify the arm whose mean is the closest to some threshold:

$$\mathcal{R}_i = \left\{ \mu \in \mathcal{R} : |\mu_i - \theta| = \min_a |\mu_a - \theta| \right\}$$

[Garivier et al., 2019a] [Aziz, K., Rivière, JMLR 2021]

Motivation: identify the Maximum Tolerated Dose in a dose-finding clinical trial



Designing a good stopping rule

Let us fix some sampling rule $(A_t)_{t \in \mathbb{N}}$, giving a data stream

$$A_1, X_1, A_2, X_2, \dots, A_t, X_t, \dots \quad \text{where} \quad X_t \sim \nu_{\mu_{A_t}}$$

Goal: construct a *sequential test* (τ, \hat{i}_τ) for the hypotheses

$$\mathcal{H}_1 : (\mu \in \mathcal{R}_1) \quad \mathcal{H}_2 : (\mu \in \mathcal{R}_2) \quad \dots \quad \mathcal{H}_I : (\mu \in \mathcal{R}_I)$$

→ multiple, composite hypotheses (possibly overlapping)

Definition

A **δ -correct sequential test** is a pair (τ, \hat{i}_τ) where

- τ is a stopping time with respect to $\mathcal{F}_t = \sigma(X_1, \dots, X_t)$
- $\hat{i}_\tau \in [I]$ is \mathcal{F}_τ -measurable

such that $\forall \mu \in \mathcal{R}, \mathbb{P}_\mu(\tau < \infty, \mu \notin \mathcal{R}_{\hat{i}_\tau}) \leq \delta$.

Idea: run I statistical tests of

$$\tilde{\mathcal{H}}_0 : (\boldsymbol{\mu} \in \mathcal{R} \setminus \mathcal{R}_i) \text{ against } \tilde{\mathcal{H}}_1 : (\boldsymbol{\mu} \in \mathcal{R}_i)$$

in **parallel** until one of them rejects $\tilde{\mathcal{H}}_0$.

Individual test: a GLR Test rejects $\tilde{\mathcal{H}}_0$ for large values of the **Generalized Likelihood Ratio**

$$\frac{\sup_{\boldsymbol{\lambda} \in \mathcal{R}} \ell(X_1, \dots, X_t; \boldsymbol{\lambda})}{\sup_{\boldsymbol{\lambda} \in \mathcal{R} \setminus \mathcal{R}_i} \ell(X_1, \dots, X_t; \boldsymbol{\lambda})} = \inf_{\boldsymbol{\lambda} \in \mathcal{R} \setminus \mathcal{R}_i} \frac{\ell(X_1, \dots, X_t; \hat{\boldsymbol{\mu}}(t))}{\ell(X_1, \dots, X_t; \boldsymbol{\lambda})}$$

where $\ell(X_1, \dots, X_t; \boldsymbol{\lambda})$ is the likelihood of the observations under a bandit model with means $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_A)$.

[Wilks, 1938]

$\hat{\boldsymbol{\mu}}(t) = (\hat{\mu}_1(t), \dots, \hat{\mu}_A(t))$, Maximum Likelihood Estimator.

Parallel GLRT

Given some threshold function $\beta(t, \delta)$,

$$\tau_\delta = \inf \left\{ t \in \mathbb{N} : \max_{i \in [I]} \inf_{\lambda \in \mathcal{R} \setminus \mathcal{R}_i} \log \frac{\ell(X_1, \dots, X_t; \hat{\mu}(t))}{\ell(X_1, \dots, X_t; \lambda)} > \beta(t, \delta) \right\}$$
$$\hat{i}_{\tau_\delta} \in \arg \max_{i \in [I]} \inf_{\lambda \in \mathcal{R} \setminus \mathcal{R}_i} \log \frac{\ell(X_1, \dots, X_{\tau_\delta}; \hat{\mu}(\tau_\delta))}{\ell(X_1, \dots, X_{\tau_\delta}; \lambda)}$$

In an exponential family bandit model,

$$\log \frac{\ell(X_1, \dots, X_t; \hat{\mu}(t))}{\ell(X_1, \dots, X_t; \lambda)} = \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \lambda_a)$$

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

(rewards in a one-parameter exponential family: Bernoulli, Gaussian, Poisson...)

Parallel GLRT

Given some **threshold function** $\beta(t, \delta)$,

$$\tau_\delta = \inf \left\{ t \in \mathbb{N} : \max_{i \in [I]} \inf_{\lambda \in \mathcal{R} \setminus \mathcal{R}_i} \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right\}$$

$$\hat{i}_{\tau_\delta} \in \arg \max_{i \in [I]} \inf_{\lambda \in \mathcal{R} \setminus \mathcal{R}_i} \sum_{a \in [A]} N_a(\tau_\delta) d(\hat{\mu}_a(\tau_\delta), \lambda_a)$$

In an exponential family bandit model,

$$\log \frac{\ell(X_1, \dots, X_t; \hat{\mu}(t))}{\ell(X_1, \dots, X_t; \lambda)} = \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \lambda_a)$$

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

(rewards in a one-parameter exponential family: Bernoulli, Gaussian, Poisson...)

Upper bound on the error probability

For any sampling rule, under the GLRT stopping rule,

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}} \left(\tau_{\delta} < \infty, \boldsymbol{\mu} \notin \mathcal{R}_{\hat{\boldsymbol{\mu}}_{\tau_{\delta}}} \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \notin \mathcal{R}_i, \inf_{\boldsymbol{\lambda} \in \mathcal{R} \setminus \mathcal{R}_i} \sum_{a \in [A]} N_a(t) d(\hat{\boldsymbol{\mu}}_a(t), \boldsymbol{\lambda}_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \in \mathcal{R} \setminus \mathcal{R}_i, \sum_{a \in [A]} N_a(t) d(\hat{\boldsymbol{\mu}}_a(t), \boldsymbol{\mu}_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \sum_{a \in [A]} N_a(t) d(\hat{\boldsymbol{\mu}}_a(t), \boldsymbol{\mu}_a) > \beta(t, \delta) \right) \end{aligned}$$

Upper bound on the error probability

For any sampling rule, under the GLRT stopping rule,

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}} \left(\tau_{\delta} < \infty, \boldsymbol{\mu} \notin \mathcal{R}_{\hat{\tau}_{\delta}} \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \notin \mathcal{R}_i, \inf_{\lambda \in \mathcal{R} \setminus \mathcal{R}_i} \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \in \mathcal{R} \setminus \mathcal{R}_i, \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \end{aligned}$$

Wanted: a deviation inequality in which

→ deviations are measured with **KL-divergence**

Upper bound on the error probability

For any sampling rule, under the GLRT stopping rule,

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}} \left(\tau_{\delta} < \infty, \boldsymbol{\mu} \notin \mathcal{R}_{\hat{\tau}_{\delta}} \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \notin \mathcal{R}_i, \inf_{\lambda \in \mathcal{R} \setminus \mathcal{R}_i} \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \in \mathcal{R} \setminus \mathcal{R}_i, \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \\ & \leq \mathbb{P} \left(\exists t \in \mathbb{N}^*, \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \end{aligned}$$

Wanted: a deviation inequality in which

- deviations are measured with KL-divergence
- deviations are **uniform over time** (*martingales...*)

Upper bound on the error probability

For any sampling rule, under the GLRT stopping rule,

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\mu}} \left(\tau_{\delta} < \infty, \boldsymbol{\mu} \notin \mathcal{R}_{\hat{\tau}_{\delta}} \right) \\ \leq & \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \notin \mathcal{R}_i, \inf_{\lambda \in \mathcal{R} \setminus \mathcal{R}_i} \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right) \\ \leq & \mathbb{P} \left(\exists t \in \mathbb{N}^*, \exists i : \boldsymbol{\mu} \in \mathcal{R} \setminus \mathcal{R}_i, \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \\ \leq & \mathbb{P} \left(\exists t \in \mathbb{N}^*, \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \mu_a) > \beta(t, \delta) \right) \end{aligned}$$

Wanted: a deviation inequality in which

- deviations are measured with KL-divergence
- deviations are uniform over time (*martingales...*)
- deviations take into account **multiple arms** (*..products*)

A universal δ -correct stopping rule

Theorem [K. and Koolen, 2018, under review]

Let μ be an exponential family bandit model. There exists a threshold function $\mathcal{T}(x) \simeq x + \log(x)$ such that, for any subset $\mathcal{S} \subseteq [A]$, for all $x > 0$,

$$\mathbb{P}_{\mu} \left(\exists t \in \mathbb{N}^* : \sum_{a \in \mathcal{S}} N_a(t) d(\hat{\mu}_a(t), \mu_a) \geq 3 \sum_{a \in \mathcal{S}} \log(1 + \log N_a(t)) + |\mathcal{S}| \mathcal{T} \left(\frac{x}{|\mathcal{S}|} \right) \right) \leq e^{-x}.$$

Consequence: the Parallel GLRT stopping rule with threshold

$$\beta(t, \delta) = 3A \log(1 + \log t) + A \mathcal{T} \left(\frac{\log(1/\delta)}{A} \right)$$

is δ -correct

- for any active identification problem
- regardless of the [sampling rule](#)

A universal δ -correct stopping rule

Theorem [K. and Koolen, 2018, under review]

Let μ be an exponential family bandit model. There exists a threshold function $\mathcal{T}(x) \simeq x + \log(x)$ such that, for any subset $\mathcal{S} \subseteq [A]$, for all $x > 0$,

$$\mathbb{P}_{\mu} \left(\exists t \in \mathbb{N}^* : \sum_{a \in \mathcal{S}} N_a(t) d(\hat{\mu}_a(t), \mu_a) \geq 3 \sum_{a \in \mathcal{S}} \log(1 + \log N_a(t)) + |\mathcal{S}| \mathcal{T} \left(\frac{x}{|\mathcal{S}|} \right) \right) \leq e^{-x}.$$

Consequence: the Parallel GLRT stopping rule with threshold

$$\beta(t, \delta) \simeq \log(1/\delta) + A \log \log(1/\delta) + 3A \log \log(t)$$

is δ -correct

- for any active identification problem
- regardless of the [sampling rule](#)

A universal δ -correct stopping rule

Theorem [K. and Koolen, 2018, under review]

Let μ be an exponential family bandit model. There exists a threshold function $\mathcal{T}(x) \simeq x + \log(x)$ such that, for any subset $\mathcal{S} \subseteq [A]$, for all $x > 0$,

$$\mathbb{P}_{\mu} \left(\exists t \in \mathbb{N}^* : \sum_{a \in \mathcal{S}} N_a(t) d(\hat{\mu}_a(t), \mu_a) \geq 3 \sum_{a \in \mathcal{S}} \log(1 + \log N_a(t)) + |\mathcal{S}| \mathcal{T} \left(\frac{x}{|\mathcal{S}|} \right) \right) \leq e^{-x}.$$

Consequence: the Parallel GLRT stopping rule with threshold

$$\beta(t, \delta) \simeq \log(1/\delta) + A \log \log(1/\delta) + 3A \log \log(t)$$

is δ -correct

- for any active identification problem
- regardless of the [sampling rule](#)

The [sample complexity](#) τ_{δ} crucially depends on the sampling rule!

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\text{KL} \left(\mathbb{P}_{\mu}^{X_1, \dots, X_{\tau}}, \mathbb{P}_{\lambda}^{X_1, \dots, X_{\tau}} \right) \geq \text{kl} \left(\mathbb{P}_{\mu}(\hat{i}_{\tau} = i_*(\lambda)), \mathbb{P}_{\lambda}(\hat{i}_{\tau} = i_*(\lambda)) \right)$$

with $\text{kl}(x, y) = \text{KL}(\mathcal{B}(x), \mathcal{B}(y))$.

[Garivier et al., 2019b]

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\text{KL} \left(\mathbb{P}_{\mu}^{X_1, \dots, X_{\tau}}, \mathbb{P}_{\lambda}^{X_1, \dots, X_{\tau}} \right) \geq \underbrace{\text{kl} \left(\mathbb{P}_{\mu}(\hat{i}_{\tau} = i_*(\lambda)), \mathbb{P}_{\mu}(\hat{i}_{\tau} = i_*(\mu)) \right)}_{\leq \delta} + \underbrace{\text{kl} \left(\mathbb{P}_{\lambda}(\hat{i}_{\tau} = i_*(\mu)), \mathbb{P}_{\lambda}(\hat{i}_{\tau} = i_*(\lambda)) \right)}_{\geq 1 - \delta}$$

with $\text{kl}(x, y) = \text{KL}(\mathcal{B}(x), \mathcal{B}(y))$.

[Garivier et al., 2019b]

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\text{KL} \left(\mathbb{P}_{\mu}^{X_1, \dots, X_{\tau}}, \mathbb{P}_{\lambda}^{X_1, \dots, X_{\tau}} \right) \geq \text{kl}(\delta, 1 - \delta)$$

with $\text{kl}(x, y) = \text{KL}(\mathcal{B}(x), \mathcal{B}(y))$.

[Garivier et al., 2019b]

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\text{KL} \left(\mathbb{P}_{\mu}^{X_1, \dots, X_{\tau}}, \mathbb{P}_{\lambda}^{X_1, \dots, X_{\tau}} \right) \geq \log(1/(3\delta))$$

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\sum_{a \in [A]} \mathbb{E}_{\mu}[N_a(\tau)] d(\mu_a, \lambda_a) \geq \log(1/(3\delta))$$

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} \mathbb{E}_\mu[N_a(\tau)] d(\mu_a, \lambda_a) \geq \log(1/(3\delta))$$

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\mathbb{E}_{\mu}[\tau] \times \left[\inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} \underbrace{\frac{\mathbb{E}_{\mu}[N_a(\tau)]}{\mathbb{E}_{\mu}[\tau]}}_{w_a} d(\mu_a, \lambda_a) \right] \geq \log(1/(3\delta))$$

Best achievable sample complexity

$\mathcal{R} = \bigcup_{i=1}^I \mathcal{R}_i$ forms a partition

$i_*(\mu)$: unique region that contains μ .

Theorem [K. and Garivier, COLT 2016]

Any δ -correct algorithm satisfies, for all $\mu \in \mathcal{R}$,

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

with

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

$$\Sigma_A = \{w \in [0, 1]^A : \sum_{a \in [A]} w_a = 1\} \quad \text{Alt}(\mu) = \{\lambda : i_*(\lambda) \neq i_*(\mu)\}$$

Proof. change of distribution between μ and $\lambda : i_*(\lambda) \neq i_*(\mu)$

$$\mathbb{E}_{\mu}[\tau] \times \left[\sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a) \right] \geq \log(1/(3\delta))$$

An algorithm matching the lower bound should satisfy

$$\forall a \in [A], \frac{\mathbb{E}_{\mu}[N_a(\tau)]}{\mathbb{E}_{\mu}[\tau]} \simeq w_a^*(\mu)$$

for a vector of **optimal proportions**

$$\mathbf{w}^*(\mu) \in \operatorname{argmax}_{\mathbf{w} \in \Sigma_A} \inf_{\lambda \in \operatorname{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a).$$

Remark: in general $\mathbf{w}^*(\mu)$

- may be non unique
- may be hard to compute

A lower-bound-inspired sampling rule for BAI

Optimal proportions

For the Best Arm Identification (BAI) problem, we propose an efficient algorithm to compute $w^*(\mu)$ for any μ .

The Tracking sampling rule:

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

with $U_t = \{a : N_a(t) < \sqrt{t}\}$.

Lemma

Under the Tracking sampling rule,

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

Optimal Best Arm Identification

The Parallel GLRT for BAI:

$$\tau_\delta = \inf \left\{ t \in \mathbb{N}^* : \inf_{\lambda \in \text{Alt}(\hat{\mu}(t))} \sum_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \lambda_a) > \beta(t, \delta) \right\}$$

Characteristic time:

$$(T^*(\mu))^{-1} = \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a)$$

Theorem [K. and Garivier, COLT 2016]

The Track-and-Stop algorithm which uses

- the **Tracking** sampling rule
- the **Parallel GLRT** stopping rule τ_δ
- recommends the **empirical best arm** $\hat{a}_{\tau_\delta} = \arg \max_a \hat{\mu}_a(\tau_\delta)$

satisfies $\mathbb{P}_\mu(\hat{a}_{\tau_\delta} \neq a_*(\mu)) \leq \delta$ and $\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau_\delta]}{\log(1/\delta)} \leq T^*(\mu)$.

→ an **asymptotically optimal** algorithm for fixed-confidence BAI!

- 1 Thompson Sampling for a Structured Bandit Problem
- 2 The Complexity of Pure Exploration
- 3 Thompson Sampling for Pure Exploration?

Thompson Sampling for BAI

Track-and-Stop can be a bit computationally heavy due to the computation of $w^*(\hat{\mu}(t))$ in every round

→ more efficient Thompson Sampling based alternatives?

Top-Two Thompson Sampling [Russo, 2016]

Input: parameter $\beta \in (0, 1)$. In round $t + 1$:

- draw a posterior sample $\theta \sim \Pi_t$, $a_*(\theta) = \arg \max_a \theta_a$
- with probability β , select $A_{t+1} = a_*(\theta)$
- with probability $1 - \beta$, re-sample the posterior $\theta' \sim \Pi_t$ until $a_*(\theta') \neq a_*(\theta)$, select $A_{t+1} = a_*(\theta')$

[Russo, 2016] performs a Bayesian analysis of TTTS:

$$\Pi_t(\text{Alt}(\mu)) \lesssim C \exp(-t/T_\beta^*(\mu)) \quad \text{a.s.}$$

- New fixed-confidence guarantees for Gaussian bandits

Theorem [Shang, De Heide, K., Ménard, Valko, AISTATS 2020]

Using the TTTS sampling rule and the Parallel GLRT yields a δ -correct BAI algorithm satisfying

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

where

$$(T_{\beta}^*(\mu))^{-1} = \sup_{\substack{w \in \Sigma_A \\ w_{a^*}(\mu) = \beta}} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a)$$

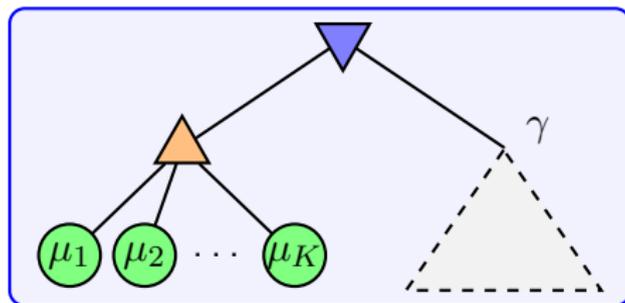
→ oracle tuning $\beta = w_{a^*}^*(\mu)$ needed for asymptotic optimality...

Comparing the Smallest Mean to a Threshold

Fix threshold γ , let $\mu_{\min} = \min_a \mu_a$. Does μ belong to

$$\mathcal{R}_{<} = \{\mu \in \mathcal{I}^A : \mu_{\min} < \gamma\}$$

or to $\mathcal{R}_{>} = \{\mu \in \mathcal{I}^A : \mu_{\min} > \gamma\}$?



Algorithm:

- sampling rule A_t
- stopping rule τ
- recommendation rule $\hat{m}_\tau \in \{<, >\}$.

Goal: $\mathbb{P}_\mu(\hat{m}_\tau \neq m^*(\mu)) \leq \delta$, small sample complexity τ .

Optimal allocation for this problem

For any δ -correct strategy,

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

Oracle allocation: $w^{\star}(\mu) = \operatorname{argmax}_{w \in \Sigma_A} \inf_{\lambda \in \operatorname{Alt}(\mu)} \sum_{a=1}^A w_a d(\mu_a, \lambda_a)$.

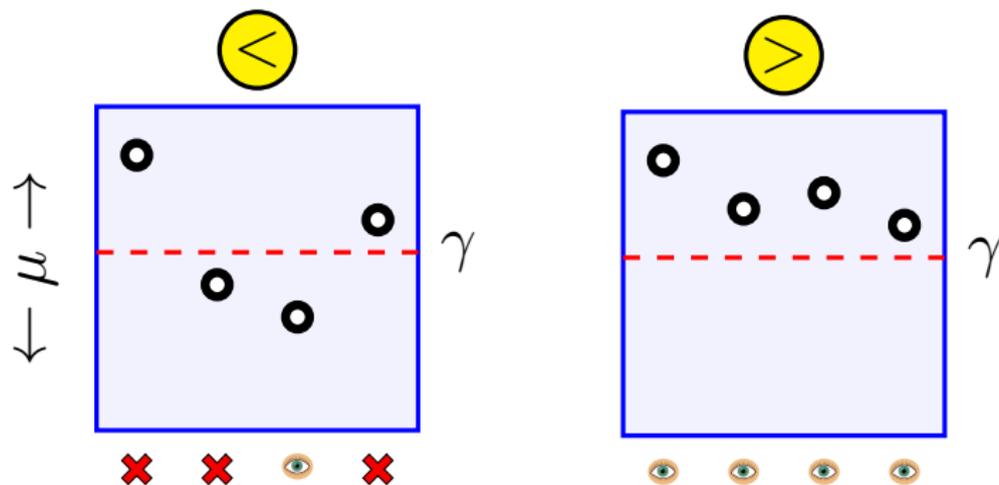
Closed-form expression for the **optimal allocation** :

$$w_a^{\star}(\mu) = \begin{cases} 1_{(a=a_{\min})} & \text{if } \mu \in \mathcal{R}_{<} \\ \frac{\frac{1}{d(\mu_a, \gamma)}}{\sum_j \frac{1}{d(\mu_j, \gamma)}} & \text{if } \mu \in \mathcal{R}_{>} \end{cases}$$

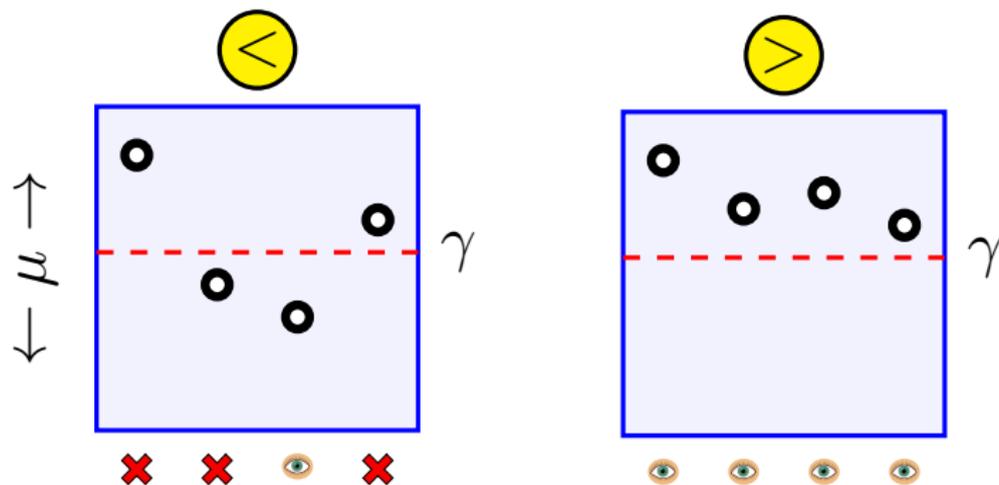
and the characteristic time

$$T_{\star}(\mu) = \begin{cases} \frac{1}{d(\mu_{\min}, \gamma)} & \text{if } \mu \in \mathcal{R}_{<} \\ \sum_a \frac{1}{d(\mu_a, \gamma)} & \text{if } \mu \in \mathcal{R}_{>} \end{cases}$$

Dichotomous Oracle Behaviour!



Dichotomous Oracle Behaviour!



Two different ideas to converge to those sampling profiles:

- **Thompson Sampling**

Sample $\theta(t) \sim \Pi_t$

Select $A_{t+1} = \arg \min_a \theta_a(t)$

(Π_t : posterior after t rounds)

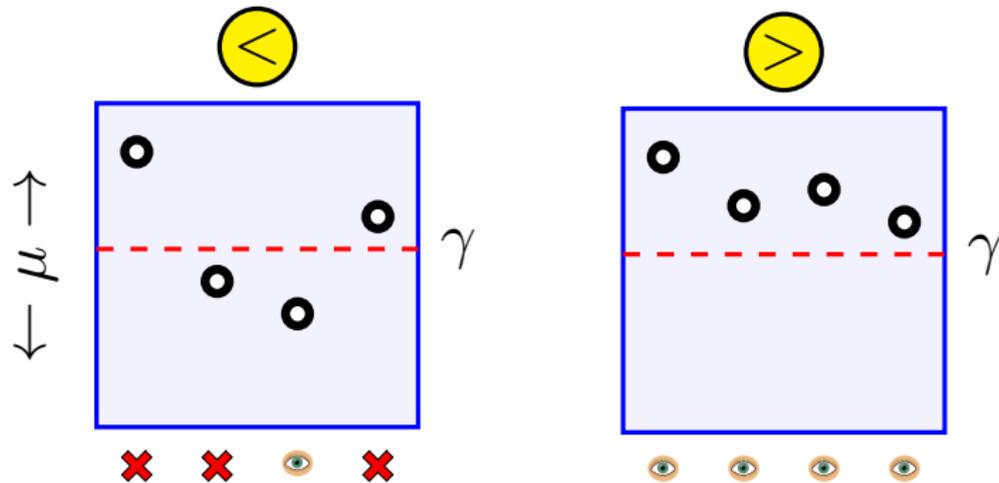
- **a LCB algorithm**

Compute a LCB on μ_a

Select $A_{t+1} = \arg \min_a \text{LCB}_a(t)$

(Lower Confidence Bound on μ_a)

A Solution: Murphy Sampling!



Murphy Sampling

Sample $\theta(t) \sim \Pi_t(\cdot | \min_a \theta_a < \gamma)$

Select $A_{t+1} = \arg \min_a \theta_a(t)$.

Idea: condition on *low* minimum mean

Properties of Murphy Sampling

Theorem [K., Koolen and Garivier, NeurIPS 2018]

For all exponential family bandit model μ , Murphy Sampling satisfies, for all a ,

$$\frac{N_a(t)}{t} \rightarrow w_a^*(\mu).$$

Sampling rule:

Thompson Sampling

Lower Confidence Bound

Murphy Sampling



Corollary [K., Koolen and Garivier, NeurIPS 2018]

Murphy Sampling combined with a “good” stopping rule satisfies

$$\limsup_{\delta \rightarrow 0} \frac{\tau_\delta}{\log \frac{1}{\delta}} \leq T_*(\mu), \text{ a.s.}$$

For both regret minimization and pure exploration:

- lower bounds are crucial to validate the (asymptotic) optimality of an algorithm
- ... and can also guide the design of optimal algorithms
- variants of Thompson Sampling provide efficient algorithms in different contexts

- Solving best arm identification in the fixed-budget setting
- Towards universal, optimal and efficient lower-bound inspired algorithms
- ... based on Thompson Sampling?
- Beyond “simple parameteric distributions”:
the power of re-sampling / sub-sampling based approaches?
- Beyond bandits:
pure exploration done right in reinforcement learning
- Sequential methods for drug design?

-  Agrawal, S. and Goyal, N. (2013).
Further Optimal Regret Bounds for Thompson Sampling.
In Proceedings of the 16th Conference on Artificial Intelligence and Statistics.
-  Aziz, M., Kaufmann, E., and Riviere, M. (2018).
On multi-armed bandit designs for dose-finding clinical trials.
arXiv:1903.07082.
-  Combes, R. and Proutière, A. (2014).
Unimodal bandits: Regret lower bounds and optimal algorithms.
In International Conference on Machine Learning (ICML).
-  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.
-  Garivier, A., Ménard, P., and Rossi, L. (2019a).
Thresholding bandit for dose-ranging: The impact of monotonicity.
In International Conference on Machine Learning, Artificial Intelligence and Applications.
-  Garivier, A., Ménard, P., and Stoltz, G. (2019b).
Explore first, exploit next: The true shape of regret in bandit problems.
Mathematics of Operation Research, 44(2):377–399.

-  Graves, T. and Lai, T. (1997).
Asymptotically Efficient adaptive choice of control laws in controlled markov chains.
SIAM Journal on Control and Optimization, 35(3):715–743.
-  Katariya, S., Kveton, B., Szepesvári, C., Vernade, C., and Wen, Z. (2017).
Bernoulli rank-1 bandits for click feedback.
In *IJCAI*.
-  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)*.
-  Lai, T. and Robbins, H. (1985).
Asymptotically efficient adaptive allocation rules.
Advances in Applied Mathematics, 6(1):4–22.
-  Magureanu, S., Combes, R., and Proutière, A. (2014).
Lipschitz Bandits: Regret lower bounds and optimal algorithms.
In *Proceedings on the 27th Conference On Learning Theory*.
-  Paladino, S., Trovò, F., Restelli, M., and Gatti, N. (2017).
Unimodal thompson sampling for graph-structured arms.
In *AAAI*.
-  Russo, D. (2016).
Simple Bayesian algorithms for best arm identification.
In *Proceedings of the 29th Conference on Learning Theory (COLT)*.
-  Shang, X., de Heide, R., Kaufmann, E., Ménard, P., and Valko, M. (2020).

Fixed-confidence guarantees for bayesian best-arm identification.

In International Conference on Artificial Intelligence and Statistics (AISTATS).



Trinh, C., Kaufmann, E., Vernade, C., and Combes, R. (2020).

Solving bernoulli rank-one bandits with unimodal thompson sampling.

In Algorithmic Learning Theory (ALT).



Wilks, S. (1938).

The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses.

The Annals of Mathematical Statistics, 9(1):60–62.

Characteristic time: (for $a_*(\mu) = 1$)

$$\begin{aligned}(T^*(\mu))^{-1} &= \sup_{w \in \Sigma_A} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a \in [A]} w_a d(\mu_a, \lambda_a) \\ &= \sup_{w \in \Sigma_A} \min_{a \neq 1} \left[w_1 d \left(\mu_1, \frac{w_1 \mu_1 + w_a \mu_a}{w_1 + w_a} \right) + w_a d \left(\mu_a, \frac{w_1 \mu_1 + w_a \mu_a}{w_1 + w_a} \right) \right]\end{aligned}$$

Parallel GLRT:

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

with

$$\begin{aligned}\hat{Z}(t) &= \inf_{\lambda \in \text{Alt}(\hat{\mu}(t))} \sum_{a \in [A]} w_a d(\hat{\mu}_a(t), \lambda_a) \\ &= \min_{a \neq \hat{a}_*(t)} \left[N_{\hat{a}_*(t)}(t) d(\hat{\mu}_{\hat{a}_*(t)}(t), \hat{\mu}_{\hat{a}_*(t), a}(t)) + N_a(t) d(\hat{\mu}_a(t), \hat{\mu}_{\hat{a}_*, a}(t)) \right],\end{aligned}$$

letting $\hat{\mu}_{a,b}(t) = \frac{N_a(t)\hat{\mu}_a(t) + N_b(t)\hat{\mu}_b(t)}{N_a(t) + N_b(t)}$.

Practical impact of Track-and-Stop

Using the right stopping rule can make a big difference in practice!

- $\mu_1 = [0.5 \ 0.45 \ 0.43 \ 0.4]$, such that

$$w_*(\mu_1) = [0.417 \ 0.390 \ 0.136 \ 0.057]$$

- $\mu_2 = [0.3 \ 0.21 \ 0.2 \ 0.19 \ 0.18]$, such that

$$w_*(\mu_2) = [0.336 \ 0.251 \ 0.177 \ 0.132 \ 0.104]$$

NB. GLRT with “stylized” threshold set to $\log\left(\frac{\log(t)+1}{\delta}\right)$.

	Track-and-Stop	GLRT-SE*	KL-LUCB	KL-SE*
μ_1	4052	4516	8437	9590
μ_2	1406	3078	2716	3334

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

* Successive Elimination

How to prove

$$\mathbb{P}_{\mu} \left(\exists t \in \mathbb{N}^* : \sum_{a \in \mathcal{S}} N_a(t) d(\hat{\mu}_a(t), \mu_a) \geq 3 \sum_{a \in \mathcal{S}} \log(1 + \log N_a(t)) + |\mathcal{S}| \mathcal{T} \left(\frac{x}{|\mathcal{S}|} \right) \right) \leq e^{-x} ?$$

Letting $X_a(t) = N_a(t) d(\hat{\mu}_a(t), \mu_a) - 3 \log(1 + \log N_a(t))$, find a martingale $M_a^\lambda(t)$ and a function $g : \Lambda \rightarrow \mathbb{R}$ such that

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

and such that $\prod_{a \in \mathcal{S}} M_a^\lambda(t)$ is still a martingale.

→ Cramer-Chernoff method + Doob inequality easily yields

$$\forall \lambda \in \Lambda, \mathbb{P} \left(\exists t \in \mathbb{N} : \sum_{a \in \mathcal{S}} X_a(t) > \frac{|\mathcal{S}| g(\lambda) + x}{\lambda} \right) \leq e^{-x}$$

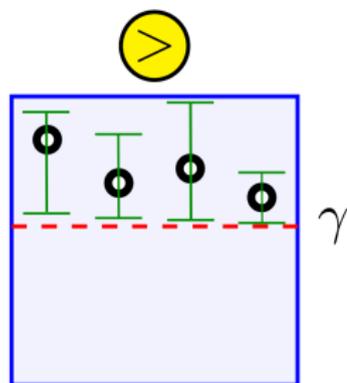
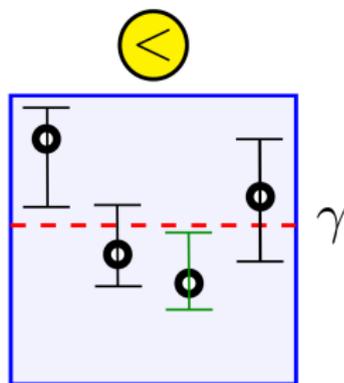
Building the martingale(s):

$$\tilde{Z}_a^\pi(t) = \int \exp \left(\eta S_a(t) - \phi_{\mu_a}(\eta) N_a(t) \right) d\pi(\eta)$$

for a well chosen continuous mixture of discrete priors.

Good stopping rules for the Smallest Minimum

Sufficient for asymptotic guarantees: a simple stopping rule based on **individual confidence intervals** $\tau^{\text{Box}} := \min(\tau_{<}; \tau_{>})$ where

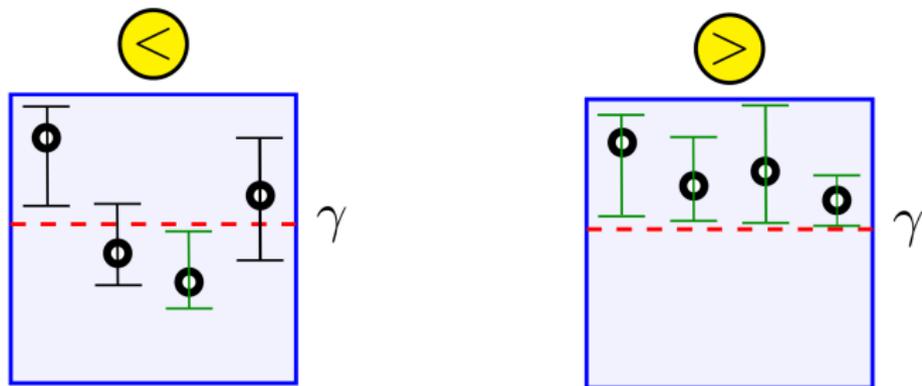


$$\tau_{<} = \inf\{t : \exists a : \text{UCB}_a(t) < \gamma\}$$

$$\tau_{>} = \inf\{t : \forall a, \text{LCB}_a(t) > \gamma\}$$

Good stopping rules for the Smallest Minimum

Sufficient for asymptotic guarantees: a simple stopping rule based on **individual confidence intervals** $\tau^{\text{Box}} := \min(\tau_{<}; \tau_{>})$ where



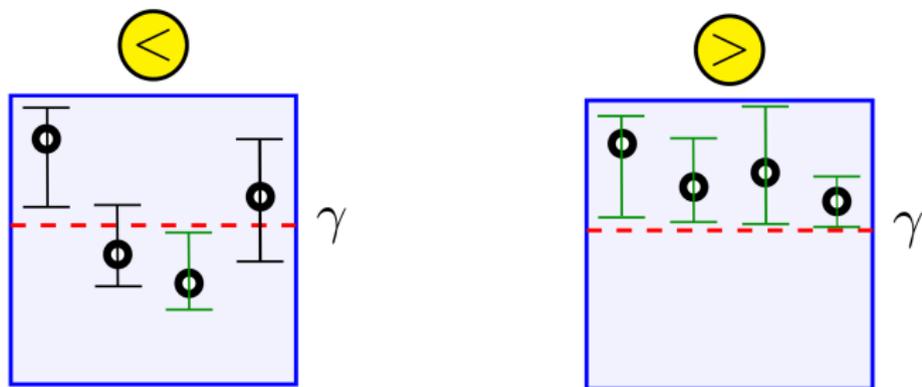
$$\tau_{<} = \inf\{t : \exists a : \text{UCB}_a(t) < \gamma\} \quad \tau_{>} = \inf\{t : \forall a, \text{LCB}_a(t) > \gamma\}$$

The Parallel GLRT?

$$\tau_{>}^{\text{GLRT}} = \inf \left\{ t \in \mathbb{N}^* : \min_{a \in [A]} N_a(t) d(\hat{\mu}_a(t), \gamma) \mathbb{1}(\hat{\mu}_a(t) \geq \gamma) > \beta(t, \delta) \right\}$$

Good stopping rules for the Smallest Minimum

Sufficient for asymptotic guarantees: a simple stopping rule based on **individual confidence intervals** $\tau^{\text{Box}} := \min(\tau_{<}; \tau_{>})$ where



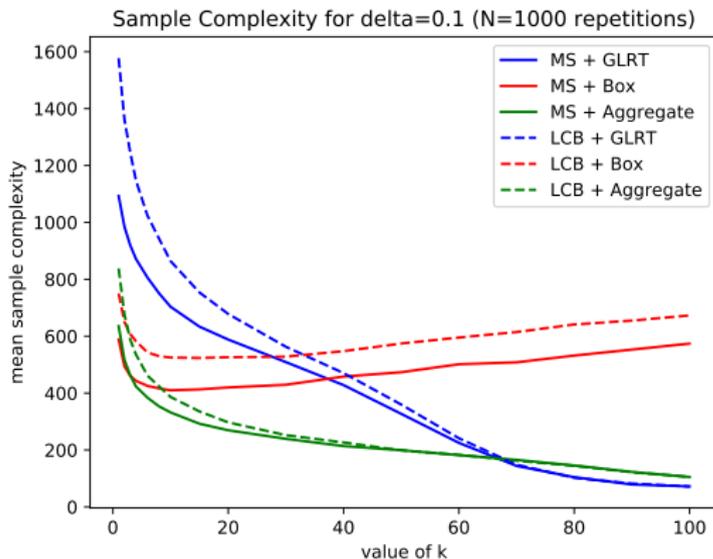
$$\tau_{<} = \inf\{t : \exists a : \text{UCB}_a(t) < \gamma\} \quad \tau_{>} = \inf\{t : \forall a, \text{LCB}_a(t) > \gamma\}$$

The Parallel GLRT?

$$\tau_{<}^{\text{GLRT}} = \inf \left\{ t \in \mathbb{N}^* : \sum_{a: \hat{\mu}_a(t) < \gamma} N_a(t) d(\hat{\mu}_a(t), \gamma_a) > \beta(t, \delta) \right\}$$

Practical performance of Murphy Sampling

Empirical sample complexity for a Gaussian instance with $\mu_a \in \{-1, 0\}$ and $\gamma = 0$ as a function of the number k of low arms

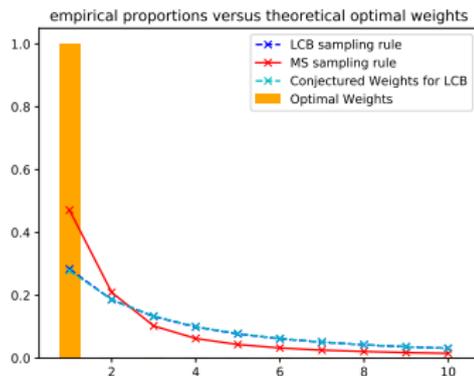


$$(\mu \in \mathcal{R}_<)$$

Convergence of Murphy Sampling

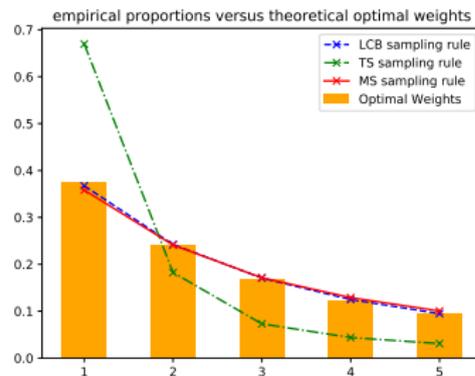
$$\mu = \text{linspace}(-1, 1, 10) \in \mathcal{R}_<$$

$$\gamma = 0$$



$$\mu = \text{linspace}(1/2, 1, 5) \in \mathcal{R}_>$$

$$\gamma = 0$$



Sampling proportions vs oracle, $\delta = e^{-23}$. Sampling proportions vs oracle, $\delta = e^{-7}$.