# Best Arm Identification with Limited Precision Sampling

P. N. Karthik

Institute of Data Science
National University of Singapore
March 23, 2023

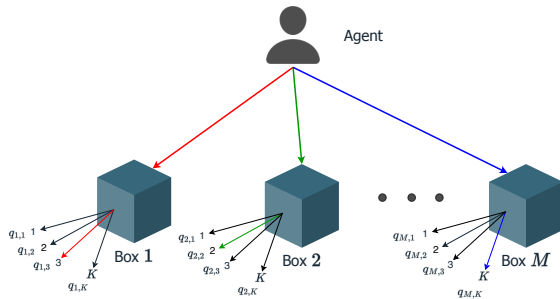# Joint Work With



Kota Srinivas Reddy
IIT Chennai



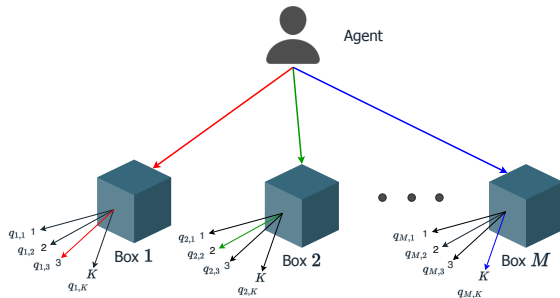Nikhil Karamchandani
IIT Mumbai



Jayakrishnan Nair
IIT Mumbai

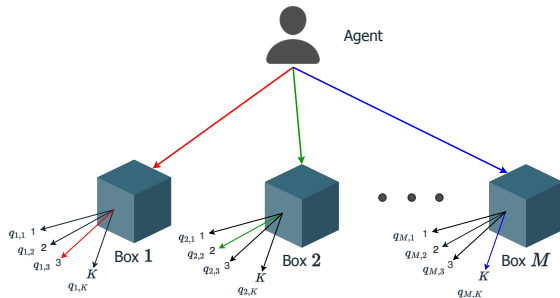# PROBLEM SETUP

- ■ *M* boxes, *K* arms per box

# PROBLEM SETUP



- *M* boxes, *K* arms per box
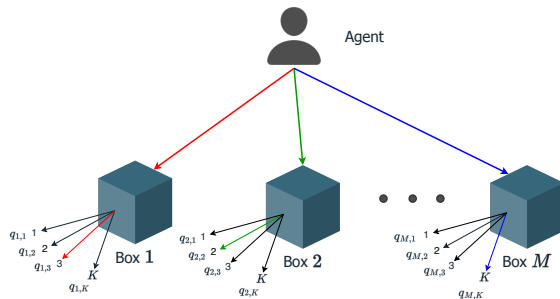- Agent can only pick boxes

# PROBLEM SETUP



- *M* boxes, *K* arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$

# PROBLEM SETUP



- *M* boxes, *K* arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$

# PROBLEM SETUP



- $M$ boxes, $K$ arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
- $\nu_k = \mathcal{N}(\mu_k, 1), \ k \in [K]$

# PROBLEM SETUP



- $M$ boxes, $K$ arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
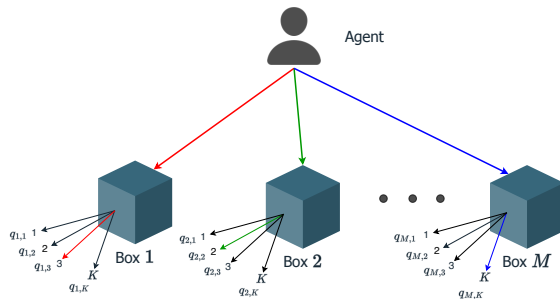- $\nu_k = \mathcal{N}(\mu_k, 1), \; k \in [K]$
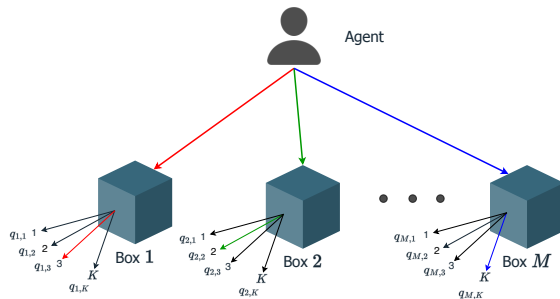- Unknowns:

# PROBLEM SETUP



- $M$ boxes, $K$ arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
- $\nu_k = \mathcal{N}(\mu_k, 1)$, $k \in [K]$
- Unknowns:
  - $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$
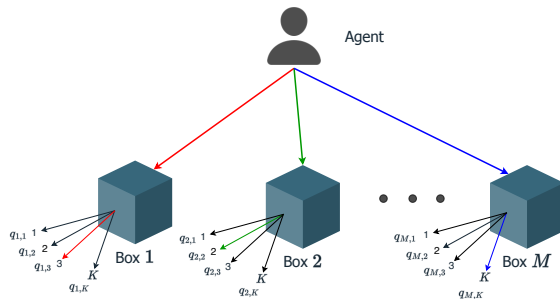
# PROBLEM SETUP



- $M$ boxes, $K$ arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
- $\nu_k = \mathcal{N}(\mu_k, 1), \ k \in [K]$
- Unknowns:
  - $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$
  - $\boldsymbol{\mu} = \{\mu_k\}_k$
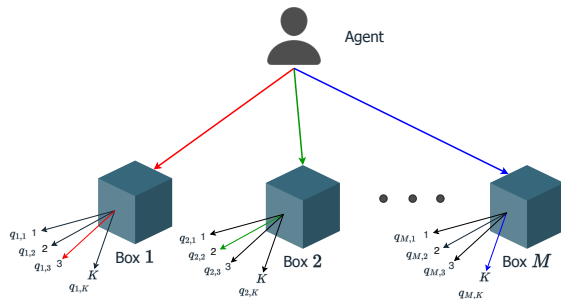
# PROBLEM SETUP



- $M$ boxes, $K$ arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
- $\nu_k = \mathcal{N}(\mu_k, 1), \ k \in [K]$
- Unknowns:
  - $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$
  - $\boldsymbol{\mu} = \{\mu_k\}_k$
- Best arm: $k^* = \arg\max_k \mu_k$
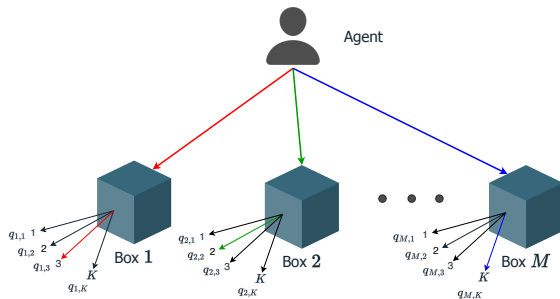
# PROBLEM SETUP



- $M$ boxes, $K$ arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
- $\nu_k = \mathcal{N}(\mu_k, 1),\ k \in [K]$
- Unknowns:
  - $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$
  - $\boldsymbol{\mu} = \{\mu_k\}_k$
- Best arm: $k^* = \arg\max_k \mu_k$
- Goal: Fixed-confidence BAI
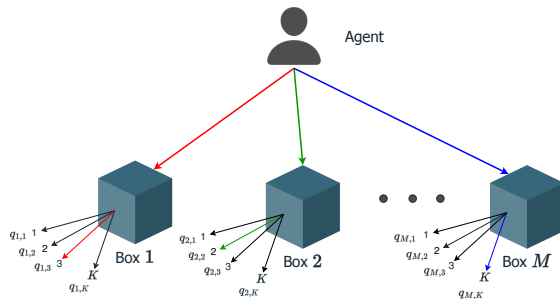
# PROBLEM SETUP



- $M$ boxes, $K$ arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
- $\nu_k = \mathcal{N}(\mu_k, 1)$, $k \in [K]$
- Unknowns:
  - $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$
  - $\boldsymbol{\mu} = \{\mu_k\}_k$
- Best arm: $k^* = \arg\max_k \mu_k$
- Goal: Fixed-confidence BAI

# PROBLEM SETUP



- *M* boxes, *K* arms per box
- Agent can only pick boxes
- $P(\text{arm} = k | \text{box} = m) = q_{m,k}$
- Instance: $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_K)$
- $\nu_k = \mathcal{N}(\mu_k, 1), \ k \in [K]$
- Unknowns:
  - $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$
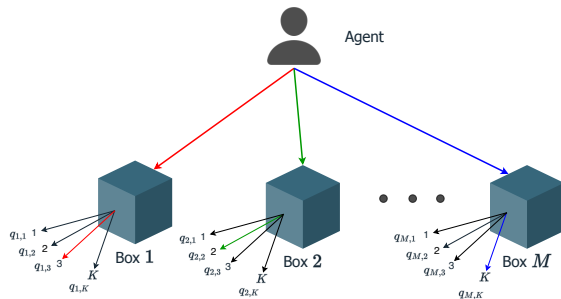  - $\boldsymbol{\mu} = \{\mu_k\}_k$
- Best arm: $k^* = \arg\max_k \mu_k$
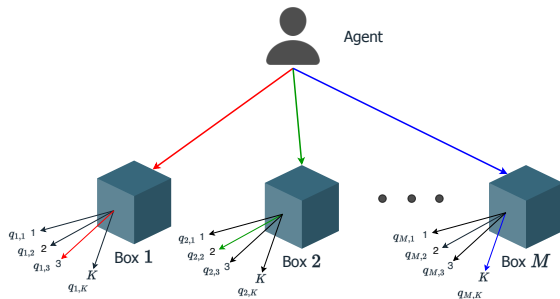- Goal: Fixed-confidence BAI

The key: determining the optimal box weight(s)

# Preliminaries

The key: determining the optimal box weight(s)



- Agent sees the pulled arm and its reward at each time

# PRELIMINARIES

The key: determining the optimal box weight(s)



- Agent sees the pulled arm and its reward at each time
- SE-type or LUCB-type algorithms cannot be applied verbatim

# PRELIMINARIES

The key: determining the optimal box weight(s)



- Agent sees the pulled arm and its reward at each time
- SE-type or LUCB-type algorithms cannot be applied verbatim
- To pull a certain arm $k$:

# PRELIMINARIES

> The key: determining the optimal box weight(s)



- Agent sees the pulled arm and its reward at each time
- SE-type or LUCB-type algorithms cannot be applied verbatim
- To pull a certain arm $k$:
  - If $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$ is known, then choose box $m^* \in \arg\max_m q_{m,k}$

# PRELIMINARIES

> The key: determining the optimal box weight(s)



- Agent sees the pulled arm and its reward at each time
- SE-type or LUCB-type algorithms cannot be applied verbatim
- To pull a certain arm $k$:
  - If $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$ is known, then choose box $m^* \in \arg\max_m q_{m,k}$
  - If $\boldsymbol{q} = \{q_{m,k}\}_{m,k}$ is unknown, then? If $\{q_{m,k}\}_k = \{q_{m'k}\}_k$ for all $m, m'$, then every box weight is optimal

# Outline

# ASYMPTOTIC ANALYSIS

# CONVERSE

Fix a problem instance $\boldsymbol{q}_0 = \{q_{m,k}^0\}_{m,k}$, $\boldsymbol{\mu}_0 = \{\mu_k^0\}_k$

# Converse

Fix a problem instance $\boldsymbol{q}_0 = \{q_{m,k}^0\}_{m,k}$, $\boldsymbol{\mu}_0 = \{\mu_k^0\}_k$

## Theorem

$$\liminf_{\delta \downarrow 0} \inf_{\pi \ \delta\text{-PC}} \frac{\mathbb{E}[\tau_\pi]}{\log(1/\delta)} \geq \frac{1}{T^*(\boldsymbol{q}_0, \boldsymbol{\mu}_0)},$$

where $T^*(\boldsymbol{q}_0, \boldsymbol{\mu}_0)$ is given by

$$T^*(\boldsymbol{q}_0, \boldsymbol{\mu}_0) = \sup_{\mathbf{w} \in \Sigma_M} \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\mu}_0)} \sum_{m=1}^{M} \sum_{k=1}^{K} w_m \ q_{m,k}^0 \ \frac{(\mu_k^0 - \lambda_k)^2}{2}.$$

The supremum is over $\Sigma_M = \{\mathbf{w} = (w_1, \ldots, w_M) : w_m \geq 0 \ \ \forall m, \ \ \sum_{m=1}^{M} w_m = 1\}$.

- From transportation Lemma 1 of Kaufmann et al. [2016],

$$\pi \ \delta\text{-PC} \implies \inf_{\boldsymbol{\lambda} \in \mathsf{ALT}(\boldsymbol{\mu}_0)} \sum_{k=1}^{K} \underbrace{\mathbb{E}[N_k(\tau_\pi)]}_{\#\text{ pulls of arm } k} \frac{(\mu_k^0 - \lambda_k)^2}{2} \geq D_{\mathsf{KL}}(\mathsf{Ber}(\delta) \| \mathsf{Ber}(1-\delta)).$$



$$T^*(\boldsymbol{q}_0, \boldsymbol{\mu}_0) = \sup_{\mathbf{w} \in \Sigma_M} \inf_{\boldsymbol{\lambda} \in \mathsf{ALT}(\boldsymbol{\mu}_0)} \sum_{m=1}^{M} \sum_{k=1}^{K} w_m \ q_{m,k}^0 \ \frac{(\mu_k^0 - \lambda_k)^2}{2}$$

- From transportation Lemma 1 of Kaufmann et al. [2016],

$$\pi \ \delta\text{-PC} \implies \inf_{\boldsymbol{\lambda} \in \mathsf{ALT}(\boldsymbol{\mu}_0)} \sum_{k=1}^{K} \underbrace{\mathbb{E}[N_k(\tau_\pi)]}_{\text{\# pulls of arm } k} \ \frac{(\mu_k^0 - \lambda_k)^2}{2} \geq D_{\mathsf{KL}}(\mathsf{Ber}(\delta) \| \mathsf{Ber}(1 - \delta)).$$

- For each $k \in [K]$,

$$\mathbb{E}[N_k(\tau_\pi)] = \sum_{m=1}^{M} q_{m,k}^0 \ \underbrace{\mathbb{E}[N(\tau_\pi, m)]}_{\text{\# selections of box } m} \ .$$



Agent

Box 1   Box 2   $\cdots$   Box $M$

$$T^*(\boldsymbol{q}_0, \boldsymbol{\mu}_0) \ = \ \sup_{\mathbf{w} \in \Sigma_M} \ \inf_{\boldsymbol{\lambda} \in \mathsf{ALT}(\boldsymbol{\mu}_0)} \ \sum_{m=1}^{M} \ \sum_{k=1}^{K} \ w_m \ q_{m,k}^0 \ \frac{(\mu_k^0 - \lambda_k)^2}{2}$$

# NON-UNIQUENESS OF OPTIMAL BOX WEIGHTS

### Example (1)

$\{q^0_{m,k}\}_k$ independent of $m$, i.e., $\{q^0_{m,k}\}_k = \{q^0_{m'k}\}_k = \{\alpha_k\}_k$ for all $m, m'$. In this case,

$$\sum_{m=1}^{M} w_m \, q^0_{m,k} = \sum_{m=1}^{M} w_m \, \alpha_k = \alpha_k \quad \forall k \in [K], \; \mathbf{w} \in \Sigma_M.$$

### Example (2)

$M = 2$, $K = 4$, $\boldsymbol{\mu}_0 = \{0.5, 0.4, 0.3, 0.3\}$, $\boldsymbol{q}_0 = \begin{pmatrix} 0.3 & 0.3 & 0.3 & 0.1 \\ 0.3 & 0.3 & 0.1 & 0.3 \end{pmatrix}$

For the above examples, every $\mathbf{w} \in \Sigma_M$ is optimal

# ACHIEVABILITY: PRELIMINARIES

- Set of optimal box weights under $(\boldsymbol{q}, \boldsymbol{\mu})$:

$$\mathcal{W}^{\star}(\boldsymbol{q}, \boldsymbol{\mu}) = \arg \sup_{w \in \Sigma_M} \inf_{\boldsymbol{\lambda} \in \mathsf{ALT}(\boldsymbol{\mu})} \sum_{m=1}^{M} \sum_{k=1}^{K} w_m \, q_{m,k} \, \frac{(\mu_k - \lambda_k)^2}{2}$$

# ACHIEVABILITY: PRELIMINARIES

- Set of optimal box weights under $(\boldsymbol{q}, \boldsymbol{\mu})$:

$$\mathcal{W}^\star(\boldsymbol{q}, \boldsymbol{\mu}) = \arg \sup_{w \in \Sigma_M} \inf_{\boldsymbol{\lambda} \in \mathsf{ALT}(\boldsymbol{\mu})} \sum_{m=1}^{M} \sum_{k=1}^{K} w_m \, q_{m,k} \, \frac{(\mu_k - \lambda_k)^2}{2}$$

- $(\boldsymbol{q}, \boldsymbol{\mu}) \mapsto \mathcal{W}^\star(\boldsymbol{q}, \boldsymbol{\mu})$ is compact-valued and upper hemicontinuous

# ACHIEVABILITY: PRELIMINARIES

- Set of optimal box weights under $(\boldsymbol{q}, \boldsymbol{\mu})$:

$$\mathcal{W}^{\star}(\boldsymbol{q}, \boldsymbol{\mu}) = \arg \sup_{w \in \Sigma_M} \inf_{\boldsymbol{\lambda} \in \mathsf{ALT}(\boldsymbol{\mu})} \sum_{m=1}^{M} \sum_{k=1}^{K} w_m \, q_{m,k} \, \frac{(\mu_k - \lambda_k)^2}{2}$$

- $(\boldsymbol{q}, \boldsymbol{\mu}) \mapsto \mathcal{W}^{\star}(\boldsymbol{q}, \boldsymbol{\mu})$ is compact-valued and upper hemicontinuous
- $\mathcal{W}^{\star}(\boldsymbol{q}, \boldsymbol{\mu})$ is convex for each $(\boldsymbol{q}, \boldsymbol{\mu})$

# ACHIEVABILITY: MODIFIED D-TRACKING (1)

- Parameter estimates at time $t$:

$$\hat{q}_{m,k}(t) = \frac{\text{\# times box } m \text{ selected and arm } k \text{ pulled}}{N(t, m)}, \quad \hat{\mu}_k(t) = \frac{1}{N_k(t)} \sum_{s=1}^{t} \mathbf{1}_{\{A_s=k\}} X_s$$

# ACHIEVABILITY: MODIFIED D-TRACKING (1)

- Parameter estimates at time $t$:

$$\hat{q}_{m,k}(t) = \frac{\text{\# times box } m \text{ selected and arm } k \text{ pulled}}{N(t,m)}, \quad \hat{\mu}_k(t) = \frac{1}{N_k(t)} \sum_{s=1}^{t} \mathbf{1}_{\{A_s=k\}} X_s$$

- $f(t) = \frac{\sqrt{t}}{\sqrt{M}}, \quad N(0,m) = 0$ for all $m$

# ACHIEVABILITY: MODIFIED D-TRACKING (1)

- Parameter estimates at time $t$:

$$\hat{q}_{m,k}(t) = \frac{\text{\# times box } m \text{ selected and arm } k \text{ pulled}}{N(t,m)}, \quad \hat{\mu}_k(t) = \frac{1}{N_k(t)} \sum_{s=1}^{t} \mathbf{1}_{\{A_s = k\}} X_s$$

- $f(t) = \frac{\sqrt{t}}{\sqrt{M}}, \quad N(0,m) = 0$ for all $m$

- $\underbrace{i_0 = 0, \quad i_{t+1} = (i_t \bmod M) + \mathbf{1}_{\left\{ \min_{m \in [M]} N(t,m) < f(t) \right\}}}_{\text{round-robin box selection counter}} \quad$ for all $t \geq 0$

# ACHIEVABILITY: MODIFIED D-TRACKING (1)

- Parameter estimates at time $t$:

$$\hat{q}_{m,k}(t) = \frac{\text{\# times box } m \text{ selected and arm } k \text{ pulled}}{N(t,m)}, \quad \hat{\mu}_k(t) = \frac{1}{N_k(t)} \sum_{s=1}^{t} \mathbf{1}_{\{A_s=k\}} X_s$$

- $f(t) = \frac{\sqrt{t}}{\sqrt{M}}, \quad N(0,m) = 0$ for all $m$

- $\underbrace{i_0 = 0, \quad i_{t+1} = (i_t \bmod M) + 1_{\{\min_{m \in [M]} N(t,m) < f(t)\}} \quad \text{for all } t \geq 0}_{\text{round-robin box selection counter}}$

- Let $\{\mathbf{w}(t) : t \geq 1\}$ be such that $\mathbf{w}(t+1) \in \mathcal{W}^\star(\hat{\mathbf{q}}(t), \hat{\boldsymbol{\mu}}(t))$ for all $t$

# ACHIEVABILITY: MODIFIED D-TRACKING (2)

■ The modified D-Tracking rule:

$$B_{t+1} = \begin{cases} i_t, & \min_{m \in [M]} N(t, m) < f(t), \\ b_t, & \text{otherwise}, \end{cases}$$

where $\{b_t : t \geq 1\}$ is specified by

$$b_t = \arg \min_{m \in \text{supp}(\sum_{s=1}^{t} w(s))} N(t, m) - \sum_{s=1}^{t} w_m(s).$$

# ACHIEVABILITY: MODIFIED D-TRACKING (2)

- The modified D-Tracking rule:

$$B_{t+1} = \begin{cases} i_t, & \min_{m\in[M]} N(t,m) < f(t), \\ b_t, & \text{otherwise}, \end{cases}$$

where $\{b_t : t \geq 1\}$ is specified by

$$b_t = \arg \min_{m\in\text{supp}(\sum_{s=1}^t w(s))} N(t,m) - \sum_{s=1}^t w_m(s).$$

- Inspired from Jedra and Proutiere [2020]

# Tracking the Optimal Set

- Define $d_\infty(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i|, \quad d_\infty(\mathbf{x}, C) = \min_{y \in C} d_\infty(x, y)$

# Tracking the Optimal Set

- Define $d_\infty(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i|, \quad d_\infty(\mathbf{x}, C) = \min_{y \in C} d_\infty(x, y)$

# Tracking the Optimal Set

- Define $d_\infty(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i|, \quad d_\infty(\mathbf{x}, C) = \min_{y \in C} d_\infty(x, y)$

### Lemma

*Under the modified D-tracking rule,*

$$\lim_{t \to \infty} d_\infty((N(t, m)/t)_{m \in [M]}, \mathcal{W}^\star(\mathbf{q}_0, \boldsymbol{\mu}_0)) = 0 \quad a.s..$$

# TRACKING THE OPTIMAL SET

- Define $d_\infty(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i|, \quad d_\infty(\mathbf{x}, C) = \min_{y \in C} d_\infty(x, y)$

### Lemma

*Under the modified D-tracking rule,*

$$\lim_{t \to \infty} d_\infty((N(t, m)/t)_{m \in [M]}, \mathcal{W}^\star(\boldsymbol{q}_0, \boldsymbol{\mu}_0)) = 0 \quad a.s..$$

- Inspired by Degenne and Koolen [2019], the key idea in the proof is to track the behaviour of $\bar{\mathbf{w}}(t) = \frac{1}{t} \sum_{s=1}^{t} \mathbf{w}(s) \in \mathcal{W}^\star(\hat{\boldsymbol{q}}(t-1), \hat{\boldsymbol{\mu}}(t-1))$

# Tracking the Optimal Set

- Define $d_\infty(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i|, \quad d_\infty(\mathbf{x}, C) = \min_{y \in C} d_\infty(x, y)$

### Lemma

*Under the modified D-tracking rule,*

$$\lim_{t \to \infty} d_\infty((N(t,m)/t)_{m \in [M]}, \mathcal{W}^\star(\boldsymbol{q}_0, \boldsymbol{\mu}_0)) = 0 \quad \text{a.s..}$$

- When $\mathcal{W}^\star(\boldsymbol{q}_0, \boldsymbol{\mu}_0) = \{\mathbf{w}^\star\}$, we recover the classical tracking result

$$\frac{N(t,m)}{t} \overset{t \to \infty}{\longrightarrow} w_m^\star \quad \forall m, \text{ a.s..}$$

# STOPPING & RECOMMENDATION RULES

■ The GLLR statistic between arms $a, b \in [K]$ at time $t$ is

$$Z_{a,b}(t) = \begin{cases} N_a(t)\frac{\left(\hat{\mu}_a(t) - \hat{\mu}_{a,b}(t)\right)^2}{2} + N_b(t)\frac{\left(\hat{\mu}_b(t) - \hat{\mu}_{a,b}(t)\right)^2}{2}, & \hat{\mu}_a(t) \geq \hat{\mu}_b(t), \\ -Z_{b,a}(t), & \text{otherwise}, \end{cases}$$

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

# STOPPING & RECOMMENDATION RULES

- The GLLR statistic between arms $a, b \in [K]$ at time $t$ is

$$Z_{a,b}(t) = \begin{cases} N_a(t) \frac{\left(\hat{\mu}_a(t) - \hat{\mu}_{a,b}(t)\right)^2}{2} + N_b(t) \frac{\left(\hat{\mu}_b(t) - \hat{\mu}_{a,b}(t)\right)^2}{2}, & \hat{\mu}_a(t) \geq \hat{\mu}_b(t), \\ -Z_{b,a}(t), & \text{otherwise}, \end{cases}$$

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

- Let $Z(t) = \max_a \min_{b \neq a} Z_{a,b}(t)$

# Stopping & Recommendation Rules

- The GLLR statistic between arms $a, b \in [K]$ at time $t$ is

$$Z_{a,b}(t) = \begin{cases} N_a(t)\frac{\left(\hat{\mu}_a(t) - \hat{\mu}_{a,b}(t)\right)^2}{2} + N_b(t)\frac{\left(\hat{\mu}_b(t) - \hat{\mu}_{a,b}(t)\right)^2}{2}, & \hat{\mu}_a(t) \geq \hat{\mu}_b(t), \\ -Z_{b,a}(t), & \text{otherwise}, \end{cases}$$

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

- Let $Z(t) = \max_a \min_{b \neq a} Z_{a,b}(t)$

- Given $\delta \in (0, 1)$, let $\beta(t, \delta, \rho) = \log \frac{C t^{1+\rho}}{\delta}$, where $C$ is a predetermined constant

# STOPPING & RECOMMENDATION RULES

- The GLLR statistic between arms $a, b \in [K]$ at time $t$ is

$$Z_{a,b}(t) = \begin{cases} N_a(t)\frac{\left(\hat{\mu}_a(t) - \hat{\mu}_{a,b}(t)\right)^2}{2} + N_b(t)\frac{\left(\hat{\mu}_b(t) - \hat{\mu}_{a,b}(t)\right)^2}{2}, & \hat{\mu}_a(t) \geq \hat{\mu}_b(t), \\ -Z_{b,a}(t), & \text{otherwise}, \end{cases}$$

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

- Let $Z(t) = \max_a \min_{b \neq a} Z_{a,b}(t)$

- Given $\delta \in (0, 1)$, let $\beta(t, \delta, \rho) = \log\frac{C\,t^{1+\rho}}{\delta}$, where $C$ is a predetermined constant

- Stopping rule:   $\tau_{\delta,\rho} = \min\{t \geq 1 : Z(t) \geq \beta(t, \delta, \rho) \text{ and } \min_{k \in [K]} N_k(t) > 0\}$

# STOPPING & RECOMMENDATION RULES

- The GLLR statistic between arms $a, b \in [K]$ at time $t$ is

$$Z_{a,b}(t) = \begin{cases} N_a(t) \frac{\left(\hat{\mu}_a(t) - \hat{\mu}_{a,b}(t)\right)^2}{2} + N_b(t) \frac{\left(\hat{\mu}_b(t) - \hat{\mu}_{a,b}(t)\right)^2}{2}, & \hat{\mu}_a(t) \geq \hat{\mu}_b(t), \\ -Z_{b,a}(t), & \text{otherwise}, \end{cases}$$

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

- Let $Z(t) = \max_a \min_{b \neq a} Z_{a,b}(t)$

- Given $\delta \in (0,1)$, let $\beta(t, \delta, \rho) = \log \frac{C t^{1+\rho}}{\delta}$, where $C$ is a predetermined constant

- Stopping rule:  $\tau_{\delta, \rho} = \min\{t \geq 1 : Z(t) \geq \beta(t, \delta, \rho) \text{ and } \min_{k \in [K]} N_k(t) > 0\}$

- Recommendation rule:  $\hat{k} = \arg\max_k \hat{\mu}_k(\tau_{\delta, \rho})$

# Results

Under the box sampling, stopping, and recommendation rules stated before:

## Theorem

- $P(\tau_{\delta,\rho} < \infty, \hat{k} \neq k^*) \leq \delta$

# Results

Under the box sampling, stopping, and recommendation rules stated before:

## Theorem

- $P(\tau_{\delta,\rho} < \infty, \hat{k} \neq k^*) \leq \delta$

- $\tau_{\delta,\rho}$ *satisfies*

$$\tau_{\delta,\rho} \leq \frac{1+\rho}{T^\star(\boldsymbol{q}_0, \boldsymbol{\mu_0})} \, \log(1/\delta) + o(\log(1/\delta)) \quad a.s..$$

*Hence,* $\quad P(\tau_{\delta,\rho} < \infty) = 1.$

# Results

Under the box sampling, stopping, and recommendation rules stated before:

## Theorem

- $P(\tau_{\delta,\rho} < \infty, \hat{k} \neq k^*) \leq \delta$

- $\tau_{\delta,\rho}$ *satisfies*

$$\tau_{\delta,\rho} \leq \frac{1+\rho}{T^\star(\boldsymbol{q}_0, \boldsymbol{\mu_0})} \, \log(1/\delta) + \boldsymbol{o}(\log(1/\delta)) \quad a.s..$$

*Hence,* $\quad P(\tau_{\delta,\rho} < \infty) = 1.$

- *Asymptotic upper bound on* $\mathbb{E}[\tau_{\delta,\rho}]$:

$$\limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau_{\delta,\rho}]}{\log(1/\delta)} \leq \frac{1+\rho}{T^*(\boldsymbol{q}_0, \boldsymbol{\mu}_0)}.$$

# NON-ASYMPTOTIC ANALYSIS:

# Preliminaries

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly

# PRELIMINARIES

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly

- In LUCB-type algorithms, two arms are pulled at each time instant

# Preliminaries

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly

- In LUCB-type algorithms, two arms are pulled at each time instant

- In our setup, the learner cannot pull arms directly

# PRELIMINARIES

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly

- In LUCB-type algorithms, two arms are pulled at each time instant

- In our setup, the learner cannot pull arms directly

- To maximise the chances of pulling a given arm $k$:

# PRELIMINARIES

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly

- In LUCB-type algorithms, two arms are pulled at each time instant

- In our setup, the learner cannot pull arms directly

- To maximise the chances of pulling a given arm $k$:

  - When $\boldsymbol{q}_0$ is known: select box $m^* \in \arg\max_m q_{m,k}^0$

# Preliminaries

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly

- In LUCB-type algorithms, two arms are pulled at each time instant

- In our setup, the learner cannot pull arms directly

- To maximise the chances of pulling a given arm $k$:
    - When $\boldsymbol{q}_0$ is known: select box $m^* \in \arg\max_m q^0_{m,k}$
    - When $\boldsymbol{q}_0$ is unknown: select $m^* \in \arg\max_m \hat{q}_{m,k}(t)$ at time $t$

# Preliminaries

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly
- In LUCB-type algorithms, two arms are pulled at each time instant
- In our setup, the learner cannot pull arms directly
- To maximise the chances of pulling a given arm $k$:
  - When $\boldsymbol{q}_0$ is known: select box $m^* \in \arg\max_m q_{m,k}^0$
  - When $\boldsymbol{q}_0$ is unknown: select $m^* \in \arg\max_m \hat{q}_{m,k}(t)$ at time $t$
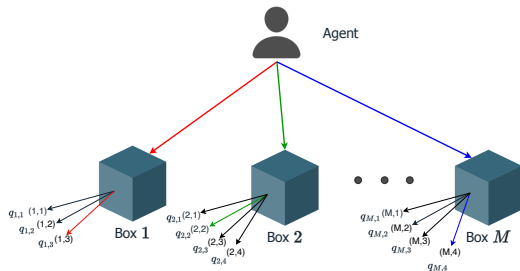- Open question: Are the above strategies optimal?

# PRELIMINARIES

- In SE-type algorithms, multiple arms are pulled at each time instant; sub-optimal arms are eliminated on-the-fly
- In LUCB-type algorithms, two arms are pulled at each time instant
- In our setup, the learner cannot pull arms directly
- To maximise the chances of pulling a given arm $k$:
  - When $\boldsymbol{q}_0$ is known: select box $m^* \in \arg\max_m q_{m,k}^0$
  - When $\boldsymbol{q}_0$ is unknown: select $m^* \in \arg\max_m \hat{q}_{m,k}(t)$ at time $t$
- Open question: Are the above strategies optimal?
- Simplified setting: arms partitioned across boxes

# SIMPLIFIED PROBLEM SETUP: PARTITION



Goal: fixed-confidence BAI

- Arms partitioned across boxes
- Arm $k$ of box $m$ indexed as $A_{m,k}$ or simply as $(m, k)$
- $\mathcal{A}_m$: set of arms in box $m$
- $\sum_{m=1}^{M} |\mathcal{A}_m| = K$
- Agent knows $\mathcal{A}_1, \ldots, \mathcal{A}_M$
- Unknowns:
  - $\boldsymbol{q}_0 = \{q_{m,k}^0\}_{m,k}$
  - $\boldsymbol{\mu}_0 = \{\mu_{m,k}^0\}_{m,k}$
- Best arm: $(m^*, k^*) = \arg\max_{m,k} \mu_{m,k}^0$

# Converse

- WLOG, let $(1,1)$ be the best arm
- Let $\Delta_{m,k} = \mu_{1,1}^0 - \mu_{m,k}^0$ for all $(m,k) \neq (1,1)$, and $\Delta_{1,1} = \min_{(m,k)\neq(1,1)} \Delta_{m,k}$

# Converse

- WLOG, let $(1,1)$ be the best arm
- Let $\Delta_{m,k} = \mu_{1,1}^0 - \mu_{m,k}^0$ for all $(m,k) \neq (1,1)$, and $\Delta_{1,1} = \min_{(m,k) \neq (1,1)} \Delta_{m,k}$

### Theorem

*Under any $\delta$-PC algorithm,*

$$\mathbb{E}\left[\tau_\pi\right] \geq \log\left(\frac{1}{2.4\,\delta}\right) \cdot \sum_{m=1}^{M} \max_{k \in \mathcal{A}_m} \frac{1}{q_{m,k}^0 \, \Delta_{m,k}^2}.$$

- Technique: change-of-measure arguments of Garivier and Kaufmann [2016]

# Achievability: Successive Elimination

Notations:

- $t_{m,k}(n)$: # pulls of arm $A_{m,k}$ up to round $n$

- $\alpha_\delta(x) = \sqrt{\frac{2 \ln (8 K x^2/\delta)}{x}}$

- $\text{UCB}_{m,k}(n) = \hat{\mu}_{m,k}(n) + \alpha_\delta(t_{m,k}(n))$

- $\text{LCB}_{m,k}(n) = \hat{\mu}_{m,k}(n) - \alpha_\delta(t_{m,k}(n))$

> Select box until each active arm is pulled $n$ times in round $n$

**Algorithm 1** Successive Elimination

**Input:** $K, M, \delta > 0, \mathcal{A}_m$ for $m \in [M]$
**Output:** $\hat{a} \in [K]$ (best arm).
　　**Initialization:** $S = [K], B = [M], S_m = [a_m], n = 0$,
　　$\hat{\mu}_{m,k}(n) = 0 \; \forall k, m, \; S_m = \mathcal{A}_m \; \forall m, t = 0$.

1: **while** $|S| > 1$ **do**
2: 　$n \leftarrow n + 1$
3: 　For each $m \in B$, select box $m$ until every active arm $A_{m,k}$ in box
　　$m$ is pulled at least $n$ times.
4: 　For every box selection, increment $t$ by $1$.
5: 　Update $t_{m,k}(n), \hat{\mu}_{m,k}(n), \text{UCB}_{m,k}(n)$ and $\text{LCB}_{m,k}(n)$ for all
　　the active arms.
6: 　**if** $\exists A_{m',k'} \in S$ such that $\text{UCB}_{m,k}(n) < \text{LCB}_{m',k'}(n)$ **then**
7: 　　$S_m \leftarrow S_m \backslash A_{m,k}, \quad S \leftarrow \bigcup_{m \in [M]} S_m$,
8: 　　$B \leftarrow \{m : S_m \neq \emptyset\}$.
9: 　**end if**
10: 　**if** $|S| = 1$ **then**
11: 　　$\hat{a} \leftarrow a \in S, \quad S \leftarrow \emptyset, \quad B \leftarrow \emptyset$.
12: 　**end if**
13: **end while**
14: **return** $\hat{a}$.

# Results

### Theorem

*Fix $\delta \in (0, 1)$. With probability greater than $1 - \delta$:*

- *The SE algorithm outputs the correct best arm*
- *The SE algorithm stops at time $\leq \sum_{m=1}^{M} U_m$, where $U_m$ is a random variable with*

$$P\left(U_m = \max_{k \in \mathcal{A}_m} O\left(\frac{\ln\left(\frac{K}{\delta \Delta_{m,k}}\right)}{q_{m,k}^0 \, \Delta_{m,k}^2}\right)\right) \geq 1 - \frac{\delta |\mathcal{A}_m|}{K}.$$

# Results

## Theorem

*Fix $\delta \in (0, 1)$. With probability greater than $1 - \delta$:*

- *The SE algorithm outputs the correct best arm*
- *The SE algorithm stops at time $\leq \sum_{m=1}^{M} U_m$, where $U_m$ is a random variable with*

$$P\left( U_m = \max_{k \in \mathcal{A}_m} O\left( \frac{\ln\left(\frac{K}{\delta \Delta_{m,k}}\right)}{q_{m,k}^0 \, \Delta_{m,k}^2} \right) \right) \geq 1 - \frac{\delta |\mathcal{A}_m|}{K}.$$

- Lower bound $= \Omega\left( \sum_{m=1}^{M} \max_{k \in \mathcal{A}_m} \frac{1}{q_{m,k}^0 \Delta_{m,k}^2} \right)$ (order-wise matching in problem unknowns)

# In Summary

- Problem studied: BAI with limited precision sampling
- Modified D-tracking algorithm to handle non-unique optimal box weights
- Partition setting: SE algorithm that selects each box until each active arm is pulled $n$ times in round $n$
- Non-partition setting: SE/LUCB-type algorithm design is an open question

Thank You!

Questions? Hit me up!
Email: karthik@nus.edu.sg
Web: https://karthikpn.com

# References

Rémy Degenne and Wouter M Koolen. Pure exploration with multiple correct answers. *Advances in Neural Information Processing Systems*, 32, 2019.

Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In *Conference on Learning Theory*, pages 998–1027. PMLR, 2016.

Yassir Jedra and Alexandre Proutiere. Optimal best-arm identification in linear bandits. *Advances in Neural Information Processing Systems*, 33:10007–10017, 2020.

Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.