

OPTIMAL BEST ARM IDENTIFICATION WITH FIXED CONFIDENCE IN RESTLESS BANDITS

P. N. KARTHIK

INSTITUTE OF DATA SCIENCE
NATIONAL UNIVERSITY OF SINGAPORE
NOVEMBER 20, 2023

JOINT WORK WITH



VINCENT TAN
NUS, SINGAPORE

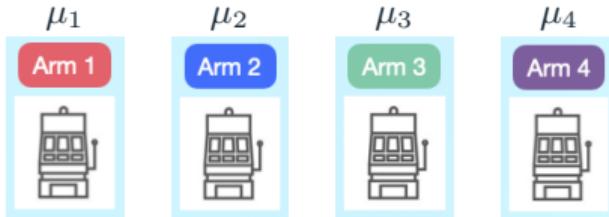


ARPAN MUKHERJEE
RPI, NEW YORK



ALI TAJER
RPI, NEW YORK

PRELIMINARIES: FIXED-CONFIDENCE BAI



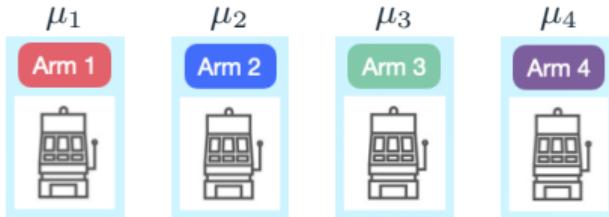
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n) Arm (A_n) Observation (\bar{X}_n)
0 1 $X_{0,1}$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI



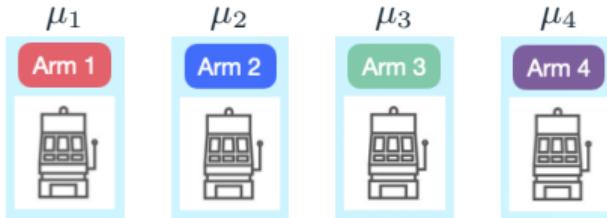
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI



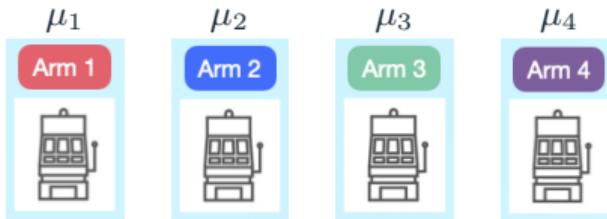
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI



$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$
3	4	$X_{3,4}$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

μ_1	μ_2	μ_3	μ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$
3	4	$X_{3,4}$
4	3	$X_{4,3}$
5	3	$X_{5,4}$
6	2	$X_{6,2}$
7	1	$X_{7,1}$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

μ_1	μ_2	μ_3	μ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$
3	4	$X_{3,4}$
4	3	$X_{4,3}$
5	3	$X_{5,4}$
6	2	$X_{6,2}$
7	1	$X_{7,1}$

$$A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

μ_1	μ_2	μ_3	μ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$
3	4	$X_{3,4}$
4	3	$X_{4,3}$
5	3	$X_{5,4}$
6	2	$X_{6,2}$
7	1	$X_{7,1}$

$$A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$$

τ – stopping time

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

μ_1	μ_2	μ_3	μ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$
3	4	$X_{3,4}$
4	3	$X_{4,3}$
5	3	$X_{5,4}$
6	2	$X_{6,2}$
7	1	$X_{7,1}$

$A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$
 $\tau - \text{stopping time}$
 $a - \text{guess of the best arm}$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

μ_1	μ_2	μ_3	μ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$
3	4	$X_{3,4}$
4	3	$X_{4,3}$
5	3	$X_{5,4}$
6	2	$X_{6,2}$
7	1	$X_{7,1}$

$A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$
 $\tau - \text{stopping time}$
 $a - \text{guess of the best arm}$

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_4)$$

$$a^*(\boldsymbol{\mu}) = \arg \max_a \mu_a$$

Given $\delta \in (0, 1)$:

$$\mathbb{P}_{\boldsymbol{\mu}}(\tau < +\infty) = 1, \quad \mathbb{P}_{\boldsymbol{\mu}}(a = a^*(\boldsymbol{\mu})) \geq 1 - \delta$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

μ'_1	μ'_2	μ'_3	μ'_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

Time (n)	Arm (A_n)	Observation (\bar{X}_n)
0	1	$X_{0,1}$
1	2	$X_{1,2}$
2	3	$X_{2,3}$
3	4	$X_{3,4}$
4	3	$X_{4,3}$
5	3	$X_{5,4}$
6	2	$X_{6,2}$
7	1	$X_{7,1}$

$$A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$$

τ – stopping time

a – guess of the best arm

$$\begin{aligned} \boldsymbol{\mu}' &= (\mu'_1, \dots, \mu'_4) \\ a^*(\boldsymbol{\mu}') &= \arg \max_a \mu'_a \end{aligned}$$

Given $\delta \in (0, 1)$:

$$\mathbb{P}_{\boldsymbol{\mu}'}(\tau < +\infty) = 1, \quad \mathbb{P}_{\boldsymbol{\mu}'}(a = a^*(\boldsymbol{\mu}')) \geq 1 - \delta$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

- At time n : π_n – a rule for choosing $A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$ or stopping

PRELIMINARIES: FIXED-CONFIDENCE BAI

- At time n : π_n – a rule for choosing $A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$ or stopping
- $\pi = \{\pi_n\}_{n \geq 0}$

PRELIMINARIES: FIXED-CONFIDENCE BAI

- At time n : π_n – a rule for choosing $A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$ or stopping
- $\pi = \{\pi_n\}_{n \geq 0}$
- $\tau(\pi)$ – stopping time

PRELIMINARIES: FIXED-CONFIDENCE BAI

- At time n : π_n – a rule for choosing $A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$ or stopping
- $\pi = \{\pi_n\}_{n \geq 0}$
- $\tau(\pi)$ – stopping time
- $a(\pi)$ – guess of the best arm

PRELIMINARIES: FIXED-CONFIDENCE BAI

- At time n : π_n – a rule for choosing $A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$ or stopping
- $\pi = \{\pi_n\}_{n \geq 0}$
- $\tau(\pi)$ – stopping time
- $a(\pi)$ – guess of the best arm
- Given $\delta \in (0, 1)$, our interest is in the class

$$\Pi(\delta) := \left\{ \pi : \mathbb{P}_{\mu'}(\tau(\pi) < +\infty) = 1, \quad \mathbb{P}_{\mu'}(a(\pi) = a^*(\mu')) \geq 1 - \delta \quad \forall \mu' \right\}$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

- At time n : π_n – a rule for choosing $A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$ or stopping
- $\pi = \{\pi_n\}_{n \geq 0}$
- $\tau(\pi)$ – stopping time
- $a(\pi)$ – guess of the best arm
- Given $\delta \in (0, 1)$, our interest is in the class

$$\Pi(\delta) := \left\{ \pi : \mathbb{P}_{\mu'}(\tau(\pi) < +\infty) = 1, \quad \mathbb{P}_{\mu'}(a(\pi) = a^*(\mu')) \geq 1 - \delta \quad \forall \mu' \right\}$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

- At time n : π_n – a rule for choosing $A_n \leftarrow (A_{0:n-1}, \bar{X}_{0:n-1})$ or stopping
- $\pi = \{\pi_n\}_{n \geq 0}$
- $\tau(\pi)$ – stopping time
- $a(\pi)$ – guess of the best arm
- Given $\delta \in (0, 1)$, our interest is in the class

$$\Pi(\delta) := \left\{ \pi : \mathbb{P}_{\mu'}(\tau(\pi) < +\infty) = 1, \quad \mathbb{P}_{\mu'}(a(\pi) = a^*(\mu')) \geq 1 - \delta \quad \forall \mu' \right\}$$

How does $\inf_{\pi \in \Pi(\delta)} \mathbb{E}_\mu[\tau(\pi)]$ grow as $\delta \downarrow 0$ for any given μ ?

PRELIMINARIES: FIXED-CONFIDENCE BAI

Theorem

Garivier and Kaufmann (2016)

Assuming there are K arms and $\mu = (\mu_1, \dots, \mu_K)$,

$$\lim_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}_{\mu}[\tau(\pi)]}{\log(1/\delta)} \geq \frac{1}{T^*(\mu)},$$

where the constant $T^*(\mu)$ is defined as

$$T^*(\mu) := \sup_{\nu \in \Sigma_K} \inf_{\lambda \in \text{ALT}(\mu)} \sum_{a=1}^K \nu(a) \text{KL}(\mu_a, \lambda_a)$$

$$\text{ALT}(\mu) := \left\{ \lambda = (\lambda_1, \dots, \lambda_K) : a^*(\lambda) \neq a^*(\mu) \right\}$$

PRELIMINARIES: FIXED-CONFIDENCE BA

$$T^*(\mu) = \sup_{\nu \in \Sigma_K} \inf_{\lambda \in \text{ALT}(\mu)} \sum_{a=1}^K \nu(a) \text{KL}(\mu_a, \lambda_a)$$

$$\psi(\nu, \mu) = \inf_{\lambda \in \text{ALT}(\mu)} \sum_{a=1}^K \nu(a) \text{KL}(\mu_a, \lambda_a), \quad \nu \in \Sigma_K, \mu \in \mathbb{R}^K, \quad \nu_\mu^* = \arg \sup_{\nu} \psi(\nu, \mu)$$

Parameter
Estimation

1

Optimal
Distribution

2

Certainty
Equivalence

3

$$\hat{\mu}(n) = [\hat{\mu}_1(n), \dots, \hat{\mu}_K(n)]^\top$$

$$\nu_{\hat{\mu}(n)}^* = \arg \sup_{\nu} \psi(\nu, \hat{\mu}(n))$$

$$A_{n+1} \sim \eta \nu^{\text{unif}} + (1-\eta) \nu_{\hat{\mu}(n)}^*$$

Key ideas: $\hat{\mu}(n) \xrightarrow{n \rightarrow \infty} \mu$, $\nu_{\hat{\mu}(n)}^*(a) \xrightarrow{n \rightarrow \infty} \nu_\mu^*(a)$, $\frac{N_a(n)}{n} \xrightarrow{n \rightarrow \infty} \nu_\mu^*(a) \text{ a.s. } \forall a$

PRELIMINARIES: FIXED-CONFIDENCE BAI

Parameter
Estimation

Optimal
Distribution

Certainty
Equivalence

1

2

3

$$\hat{\mu}(n) = [\hat{\mu}_1(n), \dots, \hat{\mu}_K(n)]^\top$$

$$\nu_{\hat{\mu}(n)}^* = \arg \sup_{\nu} \psi(\nu, \hat{\mu}(n))$$

$$A_{n+1} \sim \eta \nu^{\text{unif}} + (1-\eta) \nu_{\hat{\mu}(n)}^*$$

Theorem

For a threshold-based stopping rule using the GLLR along with the above arms selection rule,

$$\limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau]}{\log(1/\delta)} \leq \frac{1}{(1-\eta) T^*(\mu)}$$

PRELIMINARIES: FIXED-CONFIDENCE BAI

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

- Single-parameter exponential family (SPEF)

- $T^*(\boldsymbol{\theta}) = \sup_{\nu} \inf_{\lambda \in \text{ALT}(\boldsymbol{\theta})} \sum_{a=1}^K \nu(a) \text{KL}(\theta_a, \lambda_a)$

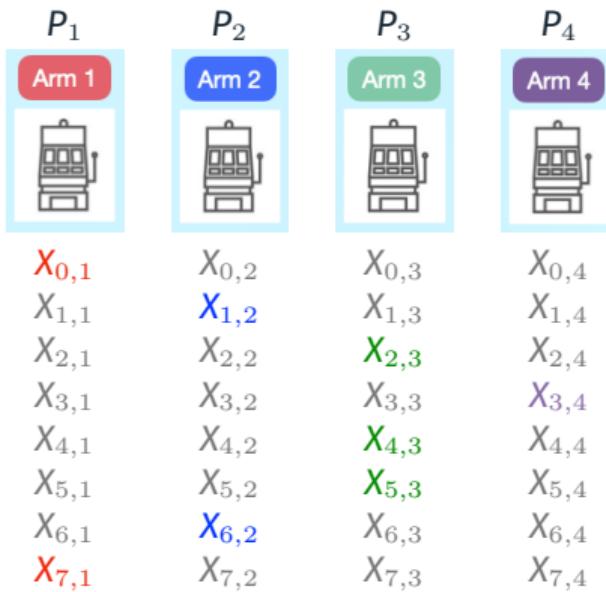
- $\theta \iff \mu_\theta = \mathbb{E}_\theta[X]$
(one-to-one correspondence)

This talk: extensions to **restless Markovian arms**

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_4)$$

$$a^*(\boldsymbol{\theta}) = \arg \max_a \theta_a$$

BAI WITH RESTLESS MARKOVIAN ARMS



- Common, finite state space
- P_a is ergodic for each a
- μ_a – stationary distribution of P_a

BAI WITH RESTLESS MARKOVIAN ARMS

P_1	P_2	P_3	P_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

- Common, finite state space
- P_a is ergodic for each a
- μ_a – stationary distribution of P_a
- How to define “mean” of arm a ?

BAI WITH RESTLESS MARKOVIAN ARMS

P_1	P_2	P_3	P_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

- Common, finite state space
- P_a is ergodic for each a
- μ_a – stationary distribution of P_a
- How to define “mean” of arm a ?
- Mean \iff TPM?

BAI WITH RESTLESS MARKOVIAN ARMS

P_1	P_2	P_3	P_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

- Common, finite state space
- P_a is ergodic for each a
- μ_a – stationary distribution of P_a
- How to define “mean” of arm a ?
- Mean \iff TPM?

Answer: SPEF for Markov chains!

MODEL: EXPONENTIAL FAMILY OF TPMs

$\Theta \subset \mathbb{R}$ - known parameter space

P – irreducible TPM on \mathcal{S} (**generator**)

Reward function $f : \mathcal{S} \rightarrow \mathbb{R}$

$$\theta \mapsto \tilde{P}_\theta(i, j) = P(i, j) e^{\theta f(j)}, \quad i, j \in \mathcal{S}$$

Property of \tilde{P}_θ : there exist vectors u_θ, v_θ such that

$$u_\theta(i) > 0 \quad \forall i, \quad v_\theta(i) > 0 \quad \forall i, \quad \tilde{P}_\theta v_\theta = \rho(\theta) v_\theta, \quad u_\theta^\top \tilde{P}_\theta = \rho(\theta) u_\theta^\top$$

$$\theta \mapsto P_\theta(i, j) = \frac{v_\theta(j)}{\rho(\theta) v_\theta(i)} \times \tilde{P}_\theta(i, j), \quad i, j \in \mathcal{S}$$

MODEL: EXPONENTIAL FAMILY OF TPMs

$\Theta \subset \mathbb{R}$ - known parameter space

P – irreducible TPM on \mathcal{S} (**generator**)

Reward function $f : \mathcal{S} \rightarrow \mathbb{R}$

$$\theta \mapsto \tilde{P}_\theta(i, j) = P(i, j) e^{\theta f(j)}, \quad i, j \in \mathcal{S}$$

Property of \tilde{P}_θ : there exist vectors u_θ, v_θ such that

$$u_\theta(i) > 0 \quad \forall i, \quad v_\theta(i) > 0 \quad \forall i, \quad \tilde{P}_\theta v_\theta = \rho(\theta) v_\theta, \quad u_\theta^\top \tilde{P}_\theta = \rho(\theta) u_\theta^\top$$

$$\theta \mapsto P_\theta(i, j) = \frac{v_\theta(j)}{\rho(\theta) v_\theta(i)} \times \tilde{P}_\theta(i, j), \quad i, j \in \mathcal{S}$$

Ergodic ($P_\theta \iff \mu_\theta = u_\theta \odot v_\theta$) ✓

$\theta \mapsto \eta_\theta = \langle f, \mu_\theta \rangle$ strictly increasing bijection ✓

ANALYSIS

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
			
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
<hr/> $X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
<hr/> $X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
<hr/>			
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9				

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9				

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
<hr/>			
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2			

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$			

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
<hr/>			
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1		

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$		

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
<hr/>			
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
<hr/>			
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
<hr/> $X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays				
n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10				

Last observed states				
n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10				

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
<hr/>	<hr/>	<hr/>	<hr/>
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3			

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$			

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
<hr/>	<hr/>	<hr/>	<hr/>
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2		

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$		

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
<hr/>	<hr/>	<hr/>	<hr/>
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
<hr/>	<hr/>	<hr/>	<hr/>
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
			
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

$$\mathbf{d}(n) = [d_a(n) : a \in [K]]^\top$$

$$\mathbf{i}(n) = [i_a(n) : a \in [K]]^\top$$

ANALYSIS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

$$\mathbf{d}(n) = [d_a(n) : a \in [K]]^\top$$

$$\mathbf{i}(n) = [i_a(n) : a \in [K]]^\top$$

$$(\mathbf{d}(n), \mathbf{i}(n)) \longrightarrow A_n \longrightarrow (\mathbf{d}(n+1), \mathbf{i}(n+1)) \longrightarrow A_{n+1} \longrightarrow (\mathbf{d}(n+2), \mathbf{i}(n+2)) \longrightarrow \dots$$

MARKOV DECISION PROCESS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

State space: $\mathbb{S} = \{(\mathbf{d}, \mathbf{i})\} \subset \mathbb{N}^K \times \mathcal{S}^K$ countably infinite

Arm delays				
n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

Last observed states				
n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

$$\mathbf{d}(n) = [d_a(n) : a \in [K]]^\top$$

$$\mathbf{i}(n) = [i_a(n) : a \in [K]]^\top$$

MARKOV DECISION PROCESS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Action space: $[K]$ finite

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

$$\mathbf{d}(n) = [d_a(n) : a \in [K]]^\top$$

$$\mathbf{i}(n) = [i_a(n) : a \in [K]]^\top$$

MARKOV DECISION PROCESS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays				
n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

Last observed states				
n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

$$\mathbf{d}(n) = [d_a(n) : a \in [K]]^\top$$

$$\mathbf{i}(n) = [i_a(n) : a \in [K]]^\top$$

Transition probabilities: $\Pr_{\theta}(\mathbf{d}(9), \mathbf{i}(9) | \mathbf{d}(8), \mathbf{i}(8), A_8 = 2) = \Pr_{\theta}(X_{8,2} | X_{6,2}) = P_{\theta_2}^2(X_{6,2}, X_{8,2})$

MARKOV DECISION PROCESS

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

Arm delays

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

Last observed states

n	$i_1(n)$	$i_2(n)$	$i_3(n)$	$i_4(n)$
8	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$
9	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$
10	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

$$\mathbf{d}(n) = [d_a(n) : a \in [K]]^\top$$

$$\mathbf{i}(n) = [i_a(n) : a \in [K]]^\top$$

Transition probabilities: $\Pr_{\theta}(\mathbf{d}(10), \mathbf{i}(10) | \mathbf{d}(9), \mathbf{i}(9), A_9 = 3) = P_{\theta_3}^4(X_{5,3}, X_{9,3})$

TPM powers

FLOW CONSTRAINT

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
<hr/>			
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

$$(\mathbf{d}(n), \mathbf{i}(n)) \xrightarrow{\text{entry}} A_n \xrightarrow{\text{exit}} (\mathbf{d}(n+1), \mathbf{i}(n+1)) \xrightarrow{\text{exit}} \dots$$

$$A_{n+1} \xrightarrow{\text{entry}} (\mathbf{d}(n+2), \mathbf{i}(n+2)) \xrightarrow{\text{exit}} A_{n+2} \xrightarrow{\text{entry}} \dots$$

$$\xrightarrow{\text{entry}} (\mathbf{d}, \mathbf{i}) \xrightarrow{\text{exit}}$$

$$N(n, \mathbf{d}, \mathbf{i}, a) = \sum_{t=K}^n \mathbf{1}_{\{\mathbf{d}(t)=\mathbf{d}, \mathbf{i}(t)=\mathbf{i}, A_t=a\}}$$

$$N(n, \mathbf{d}, \mathbf{i}) = \sum_{a=1}^K N(n, \mathbf{d}, \mathbf{i}, a)$$

FLOW CONSTRAINT

Under every policy π ,

$$\left| \mathbb{E}[N(n, \mathbf{d}', \mathbf{i}')] - \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}} \sum_{a=1}^K \mathbb{E}[N(n, \mathbf{d}, \mathbf{i}, a)] \Pr_{\theta}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) \right| \leq 1 \quad \forall (\mathbf{d}', \mathbf{i}') \in \mathbb{S}, \quad \forall n \geq K$$

CONSTRAINT ON MAXIMUM DELAY

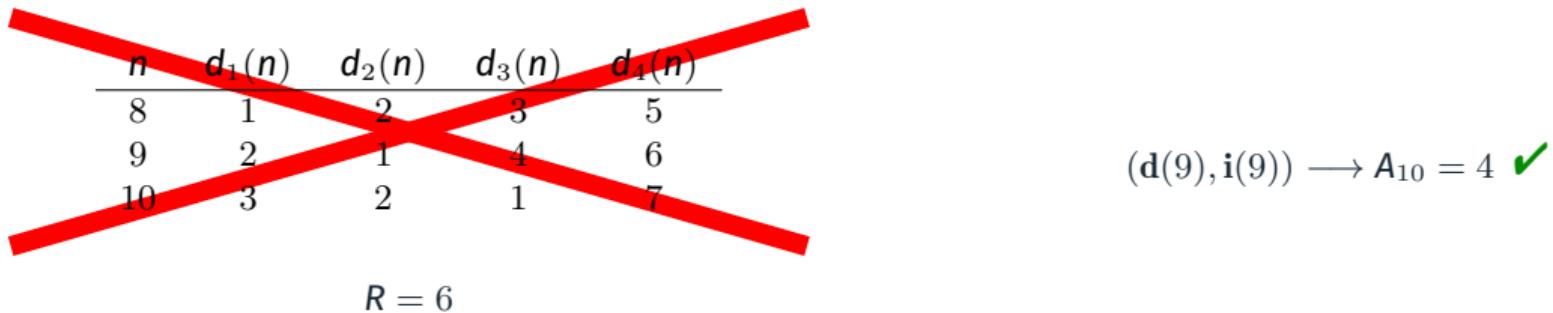
MAXIMUM DELAY CONSTRAINT

- Fix $R \in \mathbb{N}$ such that $R \gg K$
- **Forcefully pull** an arm if its delay equals R at any given time
- $\mathbb{S}_R = \{(\mathbf{d}, \mathbf{i}) \in \mathbb{S} : \max_a d_a \leq R\}$, $\mathbb{S}_{R,a} = \{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R : d_a = R\}$ finite state space
- $\Pr_{\theta} \longrightarrow \Pr_{\theta,R}$

n	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

MAXIMUM DELAY CONSTRAINT

- Fix $R \in \mathbb{N}$ such that $R \gg K$
- **Forcefully pull** an arm if its delay equals R at any given time
- $\mathbb{S}_R = \{(\mathbf{d}, \mathbf{i}) \in \mathbb{S} : \max_a d_a \leq R\}$, $\mathbb{S}_{R,a} = \{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R : d_a = R\}$ finite state space
- $\Pr_{\theta} \longrightarrow \Pr_{\theta,R}$



FLOW CONSTRAINT + R -MAX-DELAY CONSTRAINT

$$\left| \mathbb{E}[N(n, \mathbf{d}', \mathbf{i}')] - \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \mathbb{E}[N(n, \mathbf{d}, \mathbf{i}, a)] \Pr_{\theta, R}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) \right| \leq 1 \quad \forall (\mathbf{d}', \mathbf{i}') \in \mathbb{S}_R, \quad \forall n \geq K,$$

$$N(n, \mathbf{d}, \mathbf{i}, a) = N(n, \mathbf{d}, \mathbf{i}) \quad \text{if } d_a = R$$

OBJECTIVE RESTATED

θ_1	θ_2	θ_3	θ_4
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$
$X_{8,1}$	$X_{8,2}$	$X_{8,3}$	$X_{8,4}$
$X_{9,1}$	$X_{9,2}$	$X_{9,3}$	$X_{9,4}$
<hr/>			
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

- $\boldsymbol{\theta} = [\theta_a : a \in [K]]^\top$ unknown
- $a^*(\boldsymbol{\theta}) = \arg \max_a \eta_{\theta_a} = \arg \max_a \langle f, \mu_{\theta_a} \rangle$

Objective: $\inf_{\pi} \mathbb{E}[\tau_{\pi}]$

- ✓ $\pi \in \Pi(\delta)$
- ✓ π satisfies **R-max-delay constraint**

CONVERSE

CONVERSE

- Fix θ . Assume $a^*(\theta) = 1$
- Set of **alternative instances**:

$$\begin{aligned}\text{ALT}(\theta) &= \{\lambda \in \Theta^K : \exists a \neq 1 \text{ such that } \eta_{\lambda_a} > \eta_{\lambda_1}\} \\ &= \{\lambda \in \Theta^K : \exists a \neq 1 \text{ such that } \lambda_a > \lambda_1\} = \bigcup_{a=1}^K \{\lambda \in \Theta^K : \lambda_a > \lambda_1\}\end{aligned}$$

CONVERSE

- Fix θ . Assume $a^*(\theta) = 1$
- Set of **alternative instances**:

$$\begin{aligned}\text{ALT}(\theta) &= \{\lambda \in \Theta^K : \exists a \neq 1 \text{ such that } \eta_{\lambda_a} > \eta_{\lambda_1}\} \\ &= \{\lambda \in \Theta^K : \exists a \neq 1 \text{ such that } \lambda_a > \lambda_1\} = \bigcup_{a=1}^K \{\lambda \in \Theta^K : \lambda_a > \lambda_1\}\end{aligned}$$

Proposition

$$\inf_{\pi \in \Pi(\delta)} \mathbb{E}[\tau_\pi] \geq \frac{\delta \log \frac{\delta}{1-\delta} + (1-\delta) \log \frac{1-\delta}{\delta}}{T_R^*(\theta)},$$

$$T_R^*(\theta) = \sup_{\nu \in \Sigma_R(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{KL}(\Pr_{\theta, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda, R}(\cdot | \mathbf{d}, \mathbf{i}, a)),$$

$$\Sigma_R(\theta) = \left\{ \nu : \checkmark \text{prob. dist.}, \quad \checkmark \boxed{\text{flow constraint}}, \quad \checkmark \boxed{\text{R-max-delay constraint}} \right\}$$

CONVERSE – 1

$$\inf_{\pi \in \Pi(\delta)} \mathbb{E}[\tau_\pi] \geq \frac{\delta \log \frac{\delta}{1-\delta} + (1-\delta) \log \frac{1-\delta}{\delta}}{T_R^\star(\theta)}$$

CONVERSE – 1

$$\inf_{\pi \in \Pi(\delta)} \mathbb{E}[\tau_\pi] \geq \frac{\delta \log \frac{\delta}{1-\delta} + (1-\delta) \log \frac{1-\delta}{\delta}}{T_R^\star(\theta)} \sim \frac{\log \frac{1}{\delta}}{T_R^\star(\theta)}$$

CONVERSE – 1

$$\inf_{\pi \in \Pi(\delta)} \mathbb{E}[\tau_\pi] \geq \frac{\delta \log \frac{\delta}{1-\delta} + (1-\delta) \log \frac{1-\delta}{\delta}}{T_R^*(\theta)} \sim \frac{\log \frac{1}{\delta}}{T_R^*(\theta)}$$

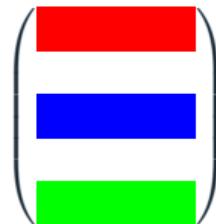
$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}[\tau_\pi]}{\log(1/\delta)} \geq \frac{1}{T_R^*(\theta)}$$

CONVERSE – 2

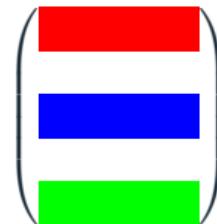
$$T_R^*(\boldsymbol{\theta}) = \sup_{\nu \in \Sigma_R(\boldsymbol{\theta})} \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

CONVERSE – 2

$$T_R^*(\boldsymbol{\theta}) = \sup_{\nu \in \Sigma_R(\boldsymbol{\theta})} \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$



$$\Pr_{\boldsymbol{\theta}, R}, \quad \boldsymbol{\theta} = [\theta_a : a \in [K]]^\top$$



$$\Pr_{\boldsymbol{\lambda}, R}, \quad \boldsymbol{\lambda} = [\lambda_a : a \in [K]]^\top$$

$$\Pr_{\boldsymbol{\theta}, R}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) = \begin{cases} P_{\theta_a}^{d_a}(i_a, i'_a), & (\mathbf{d}, \mathbf{i}) \xrightarrow{a} (\mathbf{d}', \mathbf{i}'), \\ 0, & \text{otherwise} \end{cases}$$

CONVERSE – 3

$$\Sigma_R(\theta) = \left\{ \nu : \checkmark \text{prob. dist.}, \quad \checkmark \text{flow constraint}, \quad \checkmark \text{R-max-delay constraint} \right\}$$

$$\begin{aligned} \Sigma_R(\theta) = & \left\{ \nu : \quad \nu(\mathbf{d}, \mathbf{i}, a) \geq 0 \quad \forall (\mathbf{d}, \mathbf{i}, a), \right. \\ & \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) = 1, \\ & \sum_{a=1}^K \nu(\mathbf{d}', \mathbf{i}', a) = \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) \Pr_{\theta, R}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) \quad \forall (\mathbf{d}', \mathbf{i}') \in \mathbb{S}_R, \\ & \left. \nu(\mathbf{d}, \mathbf{i}, a) = \sum_{a'=1}^K \nu(\mathbf{d}, \mathbf{i}, a') \quad \text{if } d_a = R, \ a \in [K] \right\} \end{aligned}$$

CONVERSE – 4

$$T_R^*(\boldsymbol{\theta}) = \sup_{\nu \in \Sigma_R(\boldsymbol{\theta})} \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

$$\psi(\nu, \boldsymbol{\theta}) = \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a)), \quad \nu \in \Sigma_R(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta^K$$

Lemma

CONVERSE – 4

$$T_R^*(\boldsymbol{\theta}) = \sup_{\nu \in \Sigma_R(\boldsymbol{\theta})} \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

$$\psi(\nu, \boldsymbol{\theta}) = \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a)), \quad \nu \in \Sigma_R(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta^K$$

Lemma

- ψ is continuous

Continuity of ψ – Berge's maximum theorem for non-compact sets ([Feinberg et al., 2014](#))

CONVERSE – 4

$$T_R^*(\boldsymbol{\theta}) = \sup_{\nu \in \Sigma_R(\boldsymbol{\theta})} \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

$$\psi(\nu, \boldsymbol{\theta}) = \inf_{\boldsymbol{\lambda} \in \text{ALT}(\boldsymbol{\theta})} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\boldsymbol{\theta}, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a)), \quad \nu \in \Sigma_R(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta^K$$

Lemma

- ψ is continuous
- $\boldsymbol{\theta} \mapsto \mathcal{W}^*(\boldsymbol{\theta}) = \{\nu \in \Sigma_R(\boldsymbol{\theta}) : \psi(\nu, \boldsymbol{\theta}) = T_R^*(\boldsymbol{\theta})\}$ is upper-semicontinuous

Continuity of ψ – Berge's maximum theorem for non-compact sets ([Feinberg et al., 2014](#))

CONVERSE – 4

$$T_R^*(\theta) = \sup_{\nu \in \Sigma_R(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\theta, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

$$\psi(\nu, \theta) = \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\theta, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda, R}(\cdot | \mathbf{d}, \mathbf{i}, a)), \quad \nu \in \Sigma_R(\theta), \theta \in \Theta^K$$

Lemma

- ψ is continuous
- $\theta \mapsto \mathcal{W}^*(\theta) = \{\nu \in \Sigma_R(\theta) : \psi(\nu, \theta) = T_R^*(\theta)\}$ is upper-semicontinuous

Continuity of ψ – Berge's maximum theorem for non-compact sets ([Feinberg et al., 2014](#))

Key idea for achievability: $\left[\frac{N(n, \mathbf{d}, \mathbf{i}, a)}{n} : (\mathbf{d}, \mathbf{i}, a) \in \mathbb{S}_R \times [K] \right]^\top \longrightarrow \mathcal{W}^*(\theta)$

COMPARISON

Restless arms

$$T_R^*(\theta) = \sup_{\nu \in \Sigma_R(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\theta, R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

Independent observations from arms
Garivier and Kaufmann (2016)

$$T^*(\theta) = \sup_{\nu} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{a=1}^K \nu(a) \text{KL}(\theta_a, \lambda_a)$$

ACHIEVABILITY

UNIFORM POLICY AND ERGODICITY

$$\pi^{\text{unif}}(a|\mathbf{d}, \mathbf{i}) = \begin{cases} \frac{1}{K}, & (\mathbf{d}, \mathbf{i}) \in \bigcup_{a'=1}^K \mathbb{S}_{R,a'}, \\ 1, & (\mathbf{d}, \mathbf{i}) \in \mathbb{S}_{R,a}, \\ 0, & (\mathbf{d}, \mathbf{i}) \in \bigcup_{a' \neq a} \mathbb{S}_{R,a'}. \end{cases}$$

UNIFORM POLICY AND ERGODICITY

$$\pi^{\text{unif}}(a|\mathbf{d}, \mathbf{i}) = \begin{cases} \frac{1}{K}, & (\mathbf{d}, \mathbf{i}) \in \bigcup_{a'=1}^K \mathbb{S}_{R,a'}, \\ 1, & (\mathbf{d}, \mathbf{i}) \in \mathbb{S}_{R,a}, \\ 0, & (\mathbf{d}, \mathbf{i}) \in \bigcup_{a' \neq a} \mathbb{S}_{R,a'}. \end{cases}$$

Transition kernel under $\theta \in \Theta^K$ and π^{unif} :

$$\Pr_{\theta, \pi^{\text{unif}}}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}) = \sum_{a=1}^K \Pr_{\theta, R}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) \cdot \pi^{\text{unif}}(a | \mathbf{d}, \mathbf{i}) \quad \forall (\mathbf{d}, \mathbf{i}), (\mathbf{d}', \mathbf{i}') \in \mathbb{S}_R$$

UNIFORM POLICY AND ERGODICITY

$$\pi^{\text{unif}}(a|\mathbf{d}, \mathbf{i}) = \begin{cases} \frac{1}{K}, & (\mathbf{d}, \mathbf{i}) \in \bigcup_{a'=1}^K \mathbb{S}_{R,a'}, \\ 1, & (\mathbf{d}, \mathbf{i}) \in \mathbb{S}_{R,a}, \\ 0, & (\mathbf{d}, \mathbf{i}) \in \bigcup_{a' \neq a} \mathbb{S}_{R,a'}. \end{cases}$$

Transition kernel under $\theta \in \Theta^K$ and π^{unif} :

$$\Pr_{\theta, \pi^{\text{unif}}}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}) = \sum_{a=1}^K \Pr_{\theta, R}(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) \cdot \pi^{\text{unif}}(a | \mathbf{d}, \mathbf{i}) \quad \forall (\mathbf{d}, \mathbf{i}), (\mathbf{d}', \mathbf{i}') \in \mathbb{S}_R$$

Lemma

$\Pr_{\theta, \pi^{\text{unif}}}$ is ergodic for all $\theta \in \Theta^K$

Stationary distribution of $\Pr_{\theta, \pi^{\text{unif}}} = \mu_{\theta}^{\text{unif}}$,

$$\nu_{\theta}^{\text{unif}}(\mathbf{d}, \mathbf{i}, a) = \mu_{\theta}^{\text{unif}}(\mathbf{d}, \mathbf{i}) \cdot \pi^{\text{unif}}(a | \mathbf{d}, \mathbf{i}) \quad \forall (\mathbf{d}, \mathbf{i}, a) \in \mathbb{S}_R \times [K]$$

ACHIEVABILITY – ARM SELECTION RULE

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

ARM SELECTION RULE – 1

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$$\hat{\theta}(n)$$



$$N_a(n) = \sum_{t=0}^n \mathbf{1}_{\{A_t=a\}}$$

$$\hat{\eta}^a(n) = \begin{cases} 0, & N_a(n) = 0, \\ \frac{1}{N_a(n)} \sum_{t=0}^n \mathbf{1}_{\{A_t=a\}} f(\bar{X}_t), & N_a(n) > 0 \end{cases}$$

$$\hat{\eta}_a(n) \iff \hat{\theta}_a(n)$$

$$\hat{\theta}(n) = [\hat{\theta}_a(n) : a \in [K]]$$

ARM SELECTION RULE – 2

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$\hat{\theta}(n)$

$\mathcal{W}^*(\hat{\theta}(n))$

— — — — — — — — — —

$$\mathcal{W}^*(\hat{\theta}(n)) = \arg \sup_{\nu \in \Sigma_R(\hat{\theta}(n))} \inf_{\lambda \in \text{ALT}(\hat{\theta}(n))} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\hat{\theta}(n), R}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

ARM SELECTION RULE – 3

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$$\hat{\theta}(n)$$

$$\mathcal{W}^*(\hat{\theta}(n))$$

$$\pi_{\hat{\theta}(n)}^\eta, \pi_n$$

Fix $\eta \in (0, 1)$ and $\varepsilon_n \rightarrow 0$.
Pick $\nu_n^* \in \mathcal{W}^*(\hat{\theta}(n))$

Wrapper 1

$$\pi_{\hat{\theta}(n)}^\eta(a|\mathbf{d}, \mathbf{i}) = \frac{\eta \nu_{\hat{\theta}(n)}^{\text{unif}}(\mathbf{d}, \mathbf{i}, a) + (1 - \eta) \nu_n^*(\mathbf{d}, \mathbf{i}, a)}{\eta \mu_{\hat{\theta}(n)}^{\text{unif}}(\mathbf{d}, \mathbf{i}) + (1 - \eta) \sum_{a'=1}^K \nu_n^*(\mathbf{d}, \mathbf{i}, a')}$$

Wrapper 2

$$\pi_n = \varepsilon_n \pi^{\text{unif}} + (1 - \varepsilon_n) \pi_{\hat{\theta}(n-1)}^\eta \quad \forall n \geq K$$

ARM SELECTION RULE – 4

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$$\hat{\theta}(n)$$

$$\mathcal{W}^*(\hat{\theta}(n))$$

$$\pi_{\hat{\theta}(n)}^\eta, \pi_n$$

$$A_n \sim \pi_n(\cdot | \mathbf{d}(n), \mathbf{i}(n))$$

ACHIEVABILITY – ARM SELECTION RULE

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$$\hat{\eta}(n) \iff \hat{\theta}(n)$$

$$\mathcal{W}^*(\hat{\theta}(n))$$

η, ε_n -mixtures

$$A_n \sim \pi_n(\cdot \mid \mathbf{d}(n), \mathbf{i}(n))$$

Lemma

ACHIEVABILITY – ARM SELECTION RULE

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$$\hat{\eta}(n) \iff \hat{\theta}(n)$$

$$\mathcal{W}^*(\hat{\theta}(n))$$

η, ε_n -mixtures

$$A_n \sim \pi_n(\cdot \mid d(n), i(n))$$

Lemma

For $S_R = |\mathbb{S}_R|$ and $\varepsilon_n = n^{-\frac{1}{2(1+S_R)}}$,

$$\lim_{n \rightarrow \infty} d_\infty \left(\left[\frac{N(n, d, i, a)}{n} : (d, i, a) \in \mathbb{S}_R \times [K] \right]^\top, \mathcal{W}_\eta^*(\theta) \right) = 0 \quad a.s.,$$

where $\mathcal{W}_\eta^*(\theta) = \{\eta \nu_\theta^{unif} + (1 - \eta) \nu : \nu \in \mathcal{W}^*(\theta)\}$

ACHIEVABILITY – STOPPING RULE

Empirical transition matrix:

$$\widehat{Q}_n(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) = \begin{cases} \frac{1}{N(n, \mathbf{d}, \mathbf{i}, a)} \sum_{t=K}^n \mathbf{1}_{\{(d(t), i(t)) = (d, i), A_t = a, (d(t+1), i(t+1)) = (d', i')\}}, & N(n, \mathbf{d}, \mathbf{i}, a) > 0, \\ \frac{1}{S_R}, & N(n, \mathbf{d}, \mathbf{i}, a) = 0 \end{cases}$$

ACHIEVABILITY – STOPPING RULE

Empirical transition matrix:

$$\hat{Q}_n(\mathbf{d}', \mathbf{i}' | \mathbf{d}, \mathbf{i}, a) = \begin{cases} \frac{1}{N(n, \mathbf{d}, \mathbf{i}, a)} \sum_{t=K}^n \mathbf{1}_{\{(d(t), i(t)) = (d, i), A_t = a, (d(t+1), i(t+1)) = (d', i')\}}, & N(n, \mathbf{d}, \mathbf{i}, a) > 0, \\ \frac{1}{S_R}, & N(n, \mathbf{d}, \mathbf{i}, a) = 0 \end{cases}$$

Test statistic

$$Z(n) = \inf_{\boldsymbol{\lambda} \in \text{ALT}(\hat{\theta}(n))} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K N(n, \mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\hat{Q}_n(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\boldsymbol{\lambda}, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

ACHIEVABILITY – STOPPING RULE

$$Z(n) = \inf_{\lambda \in \text{Alt}(\hat{\theta}(n))} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K N(n, \mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\hat{Q}_n(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda, R}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

Threshold:

$$\zeta(n, \delta) = \log\left(\frac{1}{\delta}\right) + (S_R - 1) \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}_R} \sum_{a=1}^K \log\left(e\left[1 + \frac{N(n, \mathbf{d}, \mathbf{i}, a)}{S_R - 1}\right]\right)$$

Stopping rule:

$$\tau_\delta = \inf\{n \geq K : Z(n) \geq \zeta(n, \delta)\}$$

Recommendation rule:

$$\hat{a} = \arg \max_a \hat{\eta}_a(n)$$

PERFORMANCE

■ Error prob. $\leq \delta$

■ Almost sure upper bound

$$\limsup_{\delta \downarrow 0} \frac{\tau_\delta}{\log(1/\delta)} \leq \frac{1}{(1-\eta) T_R^*(\theta)} \quad \text{a.s.}$$

■ Upper bound in expectation

$$\limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \frac{1}{(1-\eta) T_R^*(\theta)}$$

$$\limsup_{\eta \downarrow 0} \limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \frac{1}{T_R^*(\theta)}$$

Algorithm 1 Restless Tracking

Input: $K, R, \eta, \delta > 0$

Output: $\hat{a} \in [K]$ (best arm).

Initialization: $n = 0, N_a(n) = 0, \hat{\eta}_a(n) = 0$ for all $a \in [K]$,
 $N(n, d, i, a) = 0$ for all $(d, i, a) \in \mathbb{S}_R \times [K]$, $\text{stop} = 0$

- 1: **for** $n < K$ **do**
- 2: Select arm $A_n = n + 1$.
- 3: **end for**
- 4: **while** $\text{stop} == 0$ **do**
- 5: Update $(d(n), i(n))$; update $\hat{\eta}_a(n)$ for each $a \in [K]$
- 6: $\hat{\eta}(n) \iff \hat{\theta}(n)$
- 7: Evaluate $Z(n)$
- 8: **if** $Z(n) \geq \zeta(n, \delta)$ **then**
- 9: $\text{stop} = 1$
- 10: $\tau_\delta = n$
- 11: $\hat{a} = \arg \max_a \hat{\eta}_a(n)$
- 12: **else**
- 13: Select $A_n \sim \pi_n(\cdot | d(n), i(n))$
- 14: $n \leftarrow n + 1$.
- 15: **end if**
- 16: **end while**
- 17: **return** \hat{a} .

MAIN RESULT

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$$\hat{\eta}(n) \iff \hat{\theta}(n)$$

$$\mathcal{W}^*(\hat{\theta}(n))$$

η, ε_n -mixtures

$$A_n \sim \pi_n(\cdot \mid d(n), i(n))$$

Theorem

$$\frac{1}{T_R^*(\theta)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}[\tau_\pi]}{\log(1/\delta)} \leq \limsup_{\eta \downarrow 0} \limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \frac{1}{T_R^*(\theta)}$$

REMOVING THE MAX-DELAY CONSTRAINT

REMOVING THE MAX-DELAY CONSTRAINT

Lemma (Monotonicity)

$$T_R^*(\theta) \leq T_{R+1}^*(\theta) \text{ for all } R$$

REMOVING THE MAX-DELAY CONSTRAINT

Lemma (Monotonicity)

$$T_R^*(\theta) \leq T_{R+1}^*(\theta) \text{ for all } R$$

$\lim_{R \rightarrow \infty} T_R^*(\theta)$ exists, $\lim_{R \rightarrow \infty} \Sigma_R(\theta) = \Sigma(\theta) = \left\{ \nu : \checkmark \text{ prob. dist.}, \quad \checkmark \text{ flow constraint} \right\}$

REMOVING THE MAX-DELAY CONSTRAINT

Lemma (Monotonicity)

$$T_R^*(\theta) \leq T_{R+1}^*(\theta) \text{ for all } R$$

$\lim_{R \rightarrow \infty} T_R^*(\theta)$ exists, $\lim_{R \rightarrow \infty} \Sigma_R(\theta) = \Sigma(\theta) = \left\{ \nu : \checkmark \text{ prob. dist., } \checkmark \text{ flow constraint} \right\}$

$$T^*(\theta) = \sup_{\nu \in \Sigma(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\theta}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

REMOVING THE MAX-DELAY CONSTRAINT

Lemma (Monotonicity)

$$T_R^*(\theta) \leq T_{R+1}^*(\theta) \text{ for all } R$$

$\lim_{R \rightarrow \infty} T_R^*(\theta)$ exists, $\lim_{R \rightarrow \infty} \Sigma_R(\theta) = \Sigma(\theta) = \left\{ \nu : \checkmark \text{ prob. dist., } \checkmark \text{ flow constraint} \right\}$

$$T^*(\theta) = \sup_{\nu \in \Sigma(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\theta}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

$$\lim_{R \rightarrow \infty} T_R^*(\theta) = T^*(\theta) ??$$

REMOVING THE MAX-DELAY CONSTRAINT

Lemma (Monotonicity)

$$T_R^*(\theta) \leq T_{R+1}^*(\theta) \text{ for all } R$$

$\lim_{R \rightarrow \infty} T_R^*(\theta)$ exists, $\lim_{R \rightarrow \infty} \Sigma_R(\theta) = \Sigma(\theta) = \left\{ \nu : \checkmark \text{ prob. dist.}, \quad \checkmark \text{ flow constraint} \right\}$

$$T^*(\theta) = \sup_{\nu \in \Sigma(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\theta}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

$$\lim_{R \rightarrow \infty} T_R^*(\theta) = T^*(\theta) ??$$

Independent observations from arms ✓

Restless arms – OPEN

REMOVING THE MAX-DELAY CONSTRAINT

Lemma (Monotonicity)

$$T_R^*(\theta) \leq T_{R+1}^*(\theta) \text{ for all } R$$

$\lim_{R \rightarrow \infty} T_R^*(\theta)$ exists, $\lim_{R \rightarrow \infty} \Sigma_R(\theta) = \Sigma(\theta) = \left\{ \nu : \checkmark \text{ prob. dist.}, \quad \checkmark \text{ flow constraint} \right\}$

$$T^*(\theta) = \sup_{\nu \in \Sigma(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\mathbf{d}, \mathbf{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\mathbf{d}, \mathbf{i}, a) D_{\text{KL}}(\Pr_{\theta}(\cdot | \mathbf{d}, \mathbf{i}, a) \| \Pr_{\lambda}(\cdot | \mathbf{d}, \mathbf{i}, a))$$

$$\lim_{R \rightarrow \infty} T_R^*(\theta) = T^*(\theta)??$$

Independent observations from arms ✓

Restless arms – OPEN

$$\frac{1}{T^*(\theta)} \leq \frac{1}{T_R^*(\theta)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}[\tau_{\pi}]}{\log(1/\delta)} \leq \limsup_{\eta \downarrow 0} \limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau_{\delta}]}{\log(1/\delta)} \leq \frac{1}{T_R^*(\theta)} \stackrel{??}{\leq} \frac{1}{T^*(\theta)}$$

MAIN RESULT

Parameter
Estimation

Certainty
Equivalence

Wrapper 1
Wrapper 2

Conditional
Sampling

1

2

3

4

$$\hat{\eta}(n) \iff \hat{\theta}(n)$$

$$\mathcal{W}^*(\hat{\theta}(n))$$

η, ε_n -mixtures

$$A_n \sim \pi_n(\cdot \mid d(n), i(n))$$

Theorem

$$\frac{1}{T_R^*(\theta)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}[\tau_\pi]}{\log(1/\delta)} \leq \limsup_{\eta \downarrow 0} \limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \frac{1}{T_R^*(\theta)}$$

REFERENCES

- Feinberg, E. A., Kasyanov, P. O., and Voorneveld, M. (2014). Berge's maximum theorem for noncompact image sets. *Journal of Mathematical Analysis and Applications*, 413(2):1040–1046.
- Garivier, A. and Kaufmann, E. (2016). Optimal best arm identification with fixed confidence. In *Conference on Learning Theory*, pages 998–1027. PMLR.