

# ALMOST-OPTIMAL BEST RESTLESS MARKOV ARM IDENTIFICATION WITH FIXED CONFIDENCE

---

P. N. KARTHIK

INSTITUTE OF DATA SCIENCE  
NATIONAL UNIVERSITY OF SINGAPORE  
AUGUST 17, 2023

## JOINT WORK WITH



ARPAN MUKHERJEE  
RPI, NEW YORK

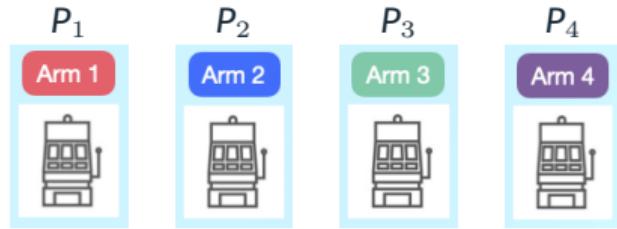


VINCENT TAN  
NUS, SINGAPORE



ALI TAJER  
RPI, NEW YORK

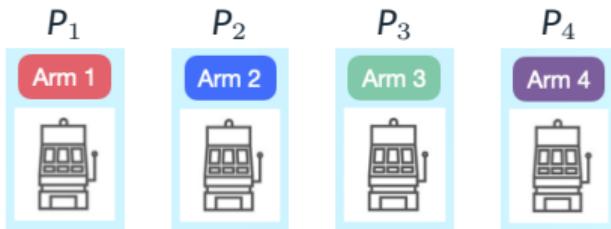
## PROBLEM SETUP



# PROBLEM SETUP

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

# PROBLEM SETUP



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

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

# PROBLEM SETUP

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

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

# PROBLEM SETUP

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

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

# PROBLEM SETUP

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

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	$X_{3,4}$

# PROBLEM SETUP

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

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	$X_{7,1}$

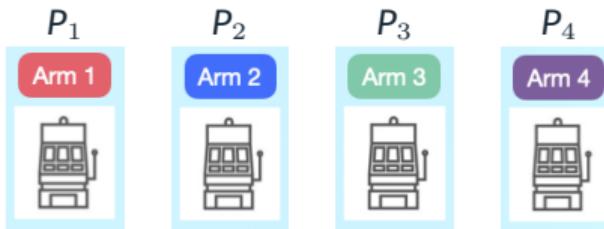
# PROBLEM SETUP

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

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	$X_{7,1}$

Goal: fixed-confidence best arm identification

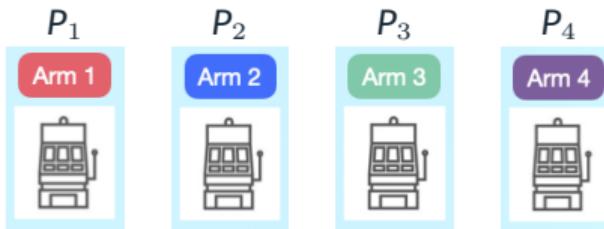
# PROBLEM SETUP



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

■ K-armed bandit,  $K \geq 2$

# PROBLEM SETUP



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

- $K$ -armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite
- Restless arms

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite
- Restless arms
- $P_a$  : TPM of arm  $a$ , ergodic

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite
- Restless arms
- $P_a$  : TPM of arm  $a$ , ergodic
- $\mu_a$  : stationary distribution of arm  $a$

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite
- Restless arms
- $P_a$  : TPM of arm  $a$ , ergodic
- $\mu_a$  : stationary distribution of arm  $a$
- $\{P_a, \mu_a : a \in [K]\}$  unknown

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite
- Restless arms
- $P_a$  : TPM of arm  $a$ , ergodic
- $\mu_a$  : stationary distribution of arm  $a$
- $\{P_a, \mu_a : a \in [K]\}$  unknown
- Known reward function  $f : \mathcal{S} \rightarrow \mathbb{R}$

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite
- Restless arms
- $P_a$  : TPM of arm  $a$ , ergodic
- $\mu_a$  : stationary distribution of arm  $a$
- $\{P_a, \mu_a : a \in [K]\}$  unknown
- Known reward function  $f : \mathcal{S} \rightarrow \mathbb{R}$
- $\eta_a = \langle f, \mu_a \rangle = \sum_{i \in \mathcal{S}} f(i) \mu_a(i)$

# PROBLEM SETUP

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

- K-armed bandit,  $K \geq 2$
- Arm  $a$ : homogeneous DTMC  $\{X_{t,a}\}_{t \geq 0}$
- State space of each arm:  $\mathcal{S}$ , finite
- Restless arms
- $P_a$  : TPM of arm  $a$ , ergodic
- $\mu_a$  : stationary distribution of arm  $a$
- $\{P_a, \mu_a : a \in [K]\}$  unknown
- Known reward function  $f : \mathcal{S} \rightarrow \mathbb{R}$
- $\eta_a = \langle f, \mu_a \rangle = \sum_{i \in \mathcal{S}} f(i) \mu_a(i)$
- Best arm  $a^* = \arg \max_a \eta_a$

# OBJECTIVE

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

$\pi$  - policy for best arm identification  
 $\tau_\pi$  - stopping time of policy  $\pi$

# OBJECTIVE

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

$\pi$  - policy for best arm identification  
 $\tau_\pi$  - stopping time of policy  $\pi$

# OBJECTIVE

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

$\pi$  - policy for best arm identification

$\tau_\pi$  - stopping time of policy  $\pi$

1.  $\Pr(\tau_\pi < +\infty) = 1$

# OBJECTIVE

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

$\pi$  - policy for best arm identification

$\tau_\pi$  - stopping time of policy  $\pi$

1.  $\Pr(\tau_\pi < +\infty) = 1$
2. **Fixed-confidence setting:** given any  $\delta \in (0, 1)$ ,

Instance	$P_1$	$P_2$	$P_3$	$P_4$	Arm 1	Arm 4	Other arms
$P_1$					$\geq 1 - \delta$	$\leq \delta$	$\leq \delta$
$Q_1$	$Q_2$	$Q_3$		$Q_4$	$\leq \delta$	$\geq 1 - \delta$	$\leq \delta$
$\Pr(a_\pi \neq a^*) \leq \delta \quad \forall \text{instances}$							

# OBJECTIVE

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

$\pi$  - policy for best arm identification

$\tau_\pi$  - stopping time of policy  $\pi$

1.  $\Pr(\tau_\pi < +\infty) = 1$
2. **Fixed-confidence setting:** given any  $\delta \in (0, 1)$ ,

Instance	$P_1$	$P_2$	$P_3$	$P_4$	Arm 1	Arm 4	Other arms
$P_1$	$P_1$	$P_2$	$P_3$	$P_4$	$\geq 1 - \delta$	$\leq \delta$	$\leq \delta$
$Q_1$	$Q_1$	$Q_2$	$Q_3$	$Q_4$	$\leq \delta$	$\geq 1 - \delta$	$\leq \delta$
$\Pr(a_\pi \neq a^*) \leq \delta \quad \forall \text{instances}$							

3.  $\Pi(\delta) = \{\pi \text{ satisfying 1. and 2.}\}$

# OBJECTIVE

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

$\pi$  - policy for best arm identification

$\tau_\pi$  - stopping time of policy  $\pi$

1.  $\Pr(\tau_\pi < +\infty) = 1$

2. **Fixed-confidence setting:** given any  $\delta \in (0, 1)$ ,

Instance	$P_1$	$P_2$	$P_3$	$P_4$	Arm 1	Arm 4	Other arms
$P_1$					$\geq 1 - \delta$	$\leq \delta$	$\leq \delta$
$Q_1$	$Q_2$	$Q_3$		$Q_4$	$\leq \delta$	$\geq 1 - \delta$	$\leq \delta$
$\Pr(a_\pi \neq a^*) \leq \delta \quad \forall \text{instances}$							

3.  $\Pi(\delta) = \{\pi \text{ satisfying 1. and 2.}\}$

Fixed-confidence BAI:  $\inf_{\pi \in \Pi(\delta)} \mathbb{E}[\tau_\pi]$

## MODEL: EXPONENTIAL FAMILY OF TPMs

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

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

Known 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  $\tilde{P}_\theta$ :** there exist vectors  $u_\theta, v_\theta$  such that

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

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

## MODEL: EXPONENTIAL FAMILY OF TPMs

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

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

Known 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  $\tilde{P}_\theta$ :** there exist vectors  $u_\theta, v_\theta$  such that

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

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

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

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

# OBJECTIVE

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

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

# OBJECTIVE

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
Arm 1	Arm 2	Arm 3	Arm 4
$X_{0,1}$	$X_{0,2}$	$X_{0,3}$	$X_{0,4}$
$X_{1,1}$	$X_{1,2}$	$X_{1,3}$	$X_{1,4}$
$X_{2,1}$	$X_{2,2}$	$X_{2,3}$	$X_{2,4}$
$X_{3,1}$	$X_{3,2}$	$X_{3,3}$	$X_{3,4}$
$X_{4,1}$	$X_{4,2}$	$X_{4,3}$	$X_{4,4}$
$X_{5,1}$	$X_{5,2}$	$X_{5,3}$	$X_{5,4}$
$X_{6,1}$	$X_{6,2}$	$X_{6,3}$	$X_{6,4}$
$X_{7,1}$	$X_{7,2}$	$X_{7,3}$	$X_{7,4}$

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

Fixed-confidence BAI:  $\inf_{\pi \in \Pi(\delta)} \mathbb{E}[\tau_\pi]$

# MARKOV DECISION PROCESS

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{6,2}$	$X_{5,3}$	$X_{3,4}$

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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				

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$			

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$		

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$	$X_{5,3}$	$X_{3,4}$

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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				

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	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	$X_{7,1}$			

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	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	$X_{7,1}$	$X_{8,2}$		

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

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

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

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

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

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

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

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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:  $\mathbb{S} = \{(\underline{d}, \underline{i})\} \subset \mathbb{N}^K \times \mathcal{S}^K$  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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

$$\begin{aligned}\underline{d}(n) &= [d_a(n) : a \in [K]]^\top \\ \underline{i}(n) &= [i_a(n) : a \in [K]]^\top\end{aligned}$$

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

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

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

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

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

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

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

# MARKOV DECISION PROCESS

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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	$X_{7,1}$	$X_{8,2}$	$X_{9,3}$	$X_{3,4}$

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

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

Transition probabilities:  $\Pr(\underline{d}(10), \underline{i}(10) | \underline{d}(9), \underline{i}(9), A_9 = 3) = P_{\theta_3}^4(X_{5,3}, X_{9,3})$  TPM powers

# FLOW CONSTRAINT

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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/>			
$X_{10,1}$	$X_{10,2}$	$X_{10,3}$	$X_{10,4}$

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

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

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

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

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

## FLOW CONSTRAINT

Under every policy  $\pi$ ,

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

## 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
- $\mathbb{S}_R = \{(\underline{d}, \underline{i}) \in \mathbb{S} : \max_a d_a \leq R\}, \quad \mathbb{S}_{R,a} = \{(\underline{d}, \underline{i}) \in \mathbb{S}_R : d_a = R\}$  finite state space
- $\xrightarrow{\text{entry}} (\underline{d}(n), \underline{i}(n)) \in \mathbb{S}_{R,a} \xrightarrow{\text{exit}} A_n = a$
- $Q_\theta \longrightarrow Q_{\theta,R}$

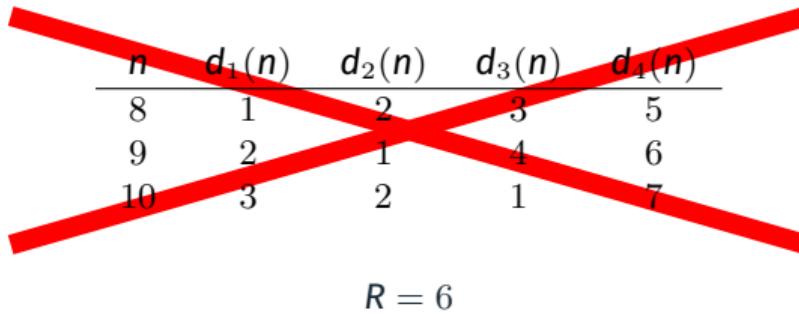
$n$	$d_1(n)$	$d_2(n)$	$d_3(n)$	$d_4(n)$
8	1	2	3	5
9	2	1	4	6
10	3	2	1	7

## 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
- $\mathbb{S}_R = \{(\underline{d}, \underline{i}) \in \mathbb{S} : \max_a d_a \leq R\}, \quad \mathbb{S}_{R,a} = \{(\underline{d}, \underline{i}) \in \mathbb{S}_R : d_a = R\}$  finite state space
- $\xrightarrow{\text{entry}} (\underline{d}(n), \underline{i}(n)) \in \mathbb{S}_{R,a} \xrightarrow{\text{exit}} A_n = a$
- $Q_\theta \longrightarrow Q_{\theta,R}$

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

$R = 6$



$(\underline{d}(9), \underline{i}(9)) \longrightarrow A_{10} = 4 \checkmark$

## FLOW CONSTRAINT + $R$ -MAX-DELAY CONSTRAINT

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

$$N(n, \underline{d}, i, a) = N(n, \underline{d}, i) \quad \forall (\underline{d}, i) \in \mathbb{S}_{R,a}, \quad a \in [K]$$

# OBJECTIVE RESTATED

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
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:  $\inf_{\pi} \mathbb{E}[\tau_{\pi}]$

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

# CONVERSE

## CONVERSE

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

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

## CONVERSE

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

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

### Proposition

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

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

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

## 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

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

## CONVERSE – 1

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

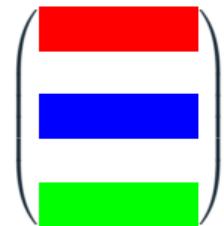
$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}[\tau_\pi]}{\log(1/\delta)} \geq \frac{1}{T_R^*(\theta)}$$

## CONVERSE – 2

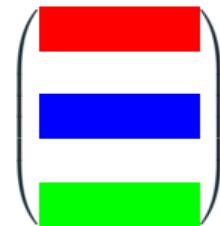
$$T_R^*(\theta) = \sup_{\nu \in \Sigma_R(\theta)} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) D_{\text{KL}}(Q_{\theta, R}(\cdot | \underline{d}, \underline{i}, a) \| Q_{\lambda, R}(\cdot | \underline{d},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

## CONVERSE – 2

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



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



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

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

## CONVERSE – 3

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

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

## CONVERSE – 4

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

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

Lemma

## CONVERSE – 4

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

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

### Lemma

- $\psi$  is continuous

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

## CONVERSE – 4

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

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

### Lemma

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

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

## CONVERSE – 4

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

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

### Lemma

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

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

Key idea for achievability:

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

## LOWER BOUNDS FROM PRIOR WORKS

Restless arms

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

Rested arms Moulos (2019)

$$T^*(\theta) = \sup_{\nu} \inf_{\lambda \in \text{ALT}(\theta)} \sum_{i \in \mathcal{S}} \sum_{a=1}^K \nu(a) \mu_{\theta_a}(i) D_{\text{KL}}(P_{\theta_a}(i, \cdot) \| P_{\lambda_a}(i, \cdot))$$

Independent observations from arms Garivier and Kaufmann (2016)

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

# ACHIEVABILITY

## UNIFORM POLICY AND ERGODICITY

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

## UNIFORM POLICY AND ERGODICITY

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

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

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

## UNIFORM POLICY AND ERGODICITY

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

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

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

### Lemma

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

Stationary distribution of  $Q_{\theta, \pi^{\text{unif}}} = \mu_{\theta}^{\text{unif}} > 0$ ,

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

## 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

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

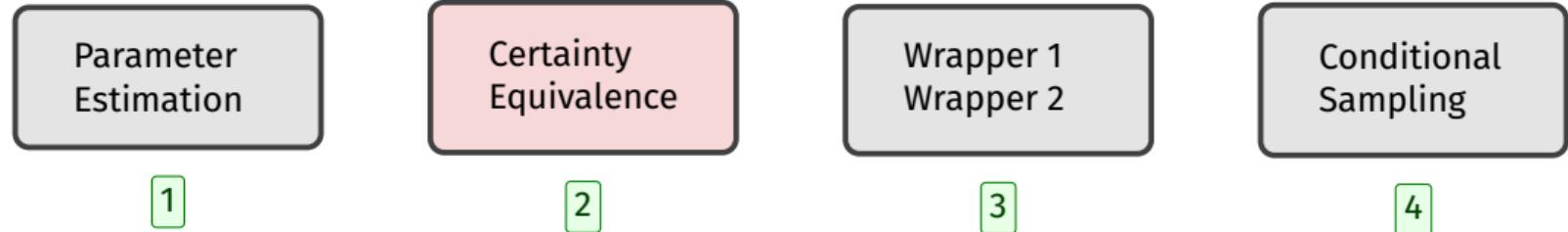
$$N_a(n) = \sum_{t=0}^n \mathbf{1}_{\{A_t=a\}}$$

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

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

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

## ARM SELECTION RULE – 2



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

## 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**

$$\pi_{\hat{\theta}(n)}^\eta(a|\underline{d}, \underline{i}) = \frac{\eta \nu_{\hat{\theta}(n)}^{\text{unif}}(\underline{d}, \underline{i}, a) + (1 - \eta) \nu_n^*(\underline{d}, \underline{i}, a)}{\eta \mu_{\hat{\theta}(n)}^{\text{unif}}(\underline{d}, \underline{i}) + (1 - \eta) \sum_{a'=1}^K \nu_n^*(\underline{d}, \underline{i}, a')}, \quad (\underline{d}, \underline{i}, a) \in \mathbb{S}_R \times [K]$$

**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

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

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

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

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

## ACHIEVABILITY – ARM SELECTION RULE

Parameter  
Estimation

Certainty  
Equivalence

Wrapper 1  
Wrapper 2

Conditional  
Sampling

1

2

3

4

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

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

$\eta, \varepsilon_n$ -mixtures

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

Lemma

## ACHIEVABILITY – ARM SELECTION RULE

Parameter  
Estimation

Certainty  
Equivalence

Wrapper 1  
Wrapper 2

Conditional  
Sampling

1

2

3

4

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

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

$\eta, \varepsilon_n$ -mixtures

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

Lemma

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

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

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

## ACHIEVABILITY – STOPPING RULE

Empirical transition matrix:

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

## ACHIEVABILITY – STOPPING RULE

Empirical transition matrix:

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

Test statistic

$$Z(n) = \inf_{\lambda \in \text{ALT}(\hat{\theta}(n))} \sum_{(\underline{d}, \underline{i}) \in S_R} \sum_{a=1}^K N(n, \underline{d}, \underline{i}, a) D_{\text{KL}}(\hat{Q}_n(\cdot | \underline{d}, \underline{i}, a) \| Q_{\lambda, R}(\cdot | \underline{d}, \underline{i}, a))$$

## ACHIEVABILITY – STOPPING RULE

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

Threshold:

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

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  $a \in [K]$ ,  
 $N(\underline{d}, \underline{i}, a) = 0$  for all  $(\underline{d}, \underline{i}, a) \in \mathbb{S}_R \times [K]$ ,  $\text{stop} = 0$

- 1: **for**  $n < K$  **do**
- 2:     Select arm  $A_n = n + 1$ .
- 3:     **end for**
- 4:     **while**  $\text{stop} == 0$  **do**
- 5:         Update  $(\underline{d}(n), \underline{i}(n))$ ; 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  $A_n \sim \pi_n(\cdot | \underline{d}(n), \underline{i}(n))$
- 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

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

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

$\eta, \varepsilon_n$ -mixtures

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

Theorem

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

# MORE MUSINGS

## MORE MUSINGS

Lemma (Monotonicity)

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

## MORE MUSINGS

### Lemma (Monotonicity)

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

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

## MORE MUSINGS

### Lemma (Montonicity)

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

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

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

## MORE MUSINGS

### Lemma (Montonicity)

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

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

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

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

## MORE MUSINGS

### Lemma (Monotonicity)

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

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

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

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

Independent observations from arms ✓

Restless arms – OPEN

## MORE MUSINGS

### Lemma (Monotonicity)

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

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

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

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

Independent observations from arms ✓

Restless arms – OPEN

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

## 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.
- Moulos, V. (2019). Optimal best Markovian arm identification with fixed confidence. *Advances in Neural Information Processing Systems*, 32.

# MAIN RESULT

Parameter  
Estimation

Certainty  
Equivalence

Wrapper 1  
Wrapper 2

Conditional  
Sampling

1

2

3

4

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

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

$\eta, \varepsilon_n$ -mixtures

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

Theorem

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