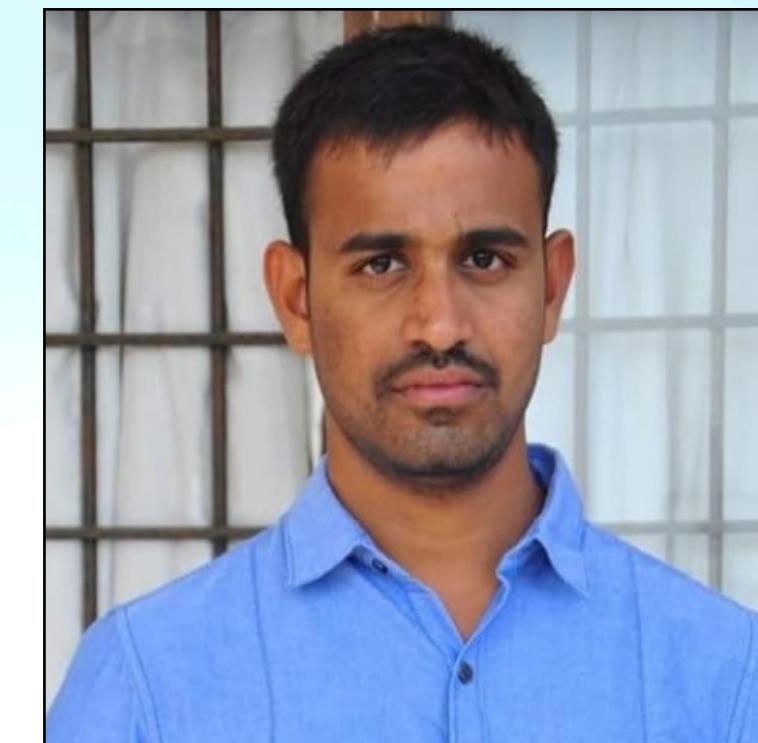


BEST RESTLESS MARKOV ARM IDENTIFICATION

2022 IEEE Information Theory Workshop, Mumbai, India



P. N. Karthik
karthik@nus.edu.sg



Srinivas Reddy Kota
ksreddy@nus.edu.sg



Vincent Y. F. Tan
vtan@nus.edu.sg

National University of Singapore
08 November 2022

PROBLEM SETUP & OBJECTIVE

PROBLEM SETUP & OBJECTIVE

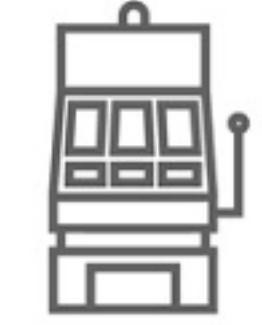
Arm 1



Arm 2



Arm 3



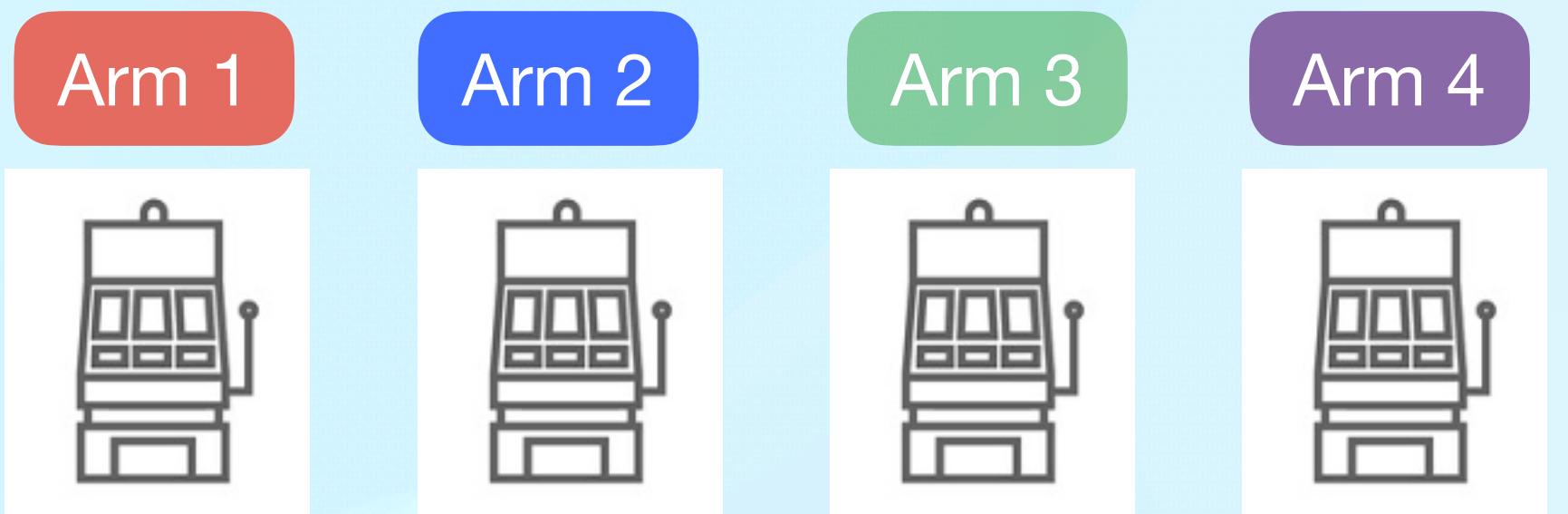
Arm 4



* $\text{TPM}(s)$: transition probability matrix(ces)

PROBLEM SETUP & OBJECTIVE

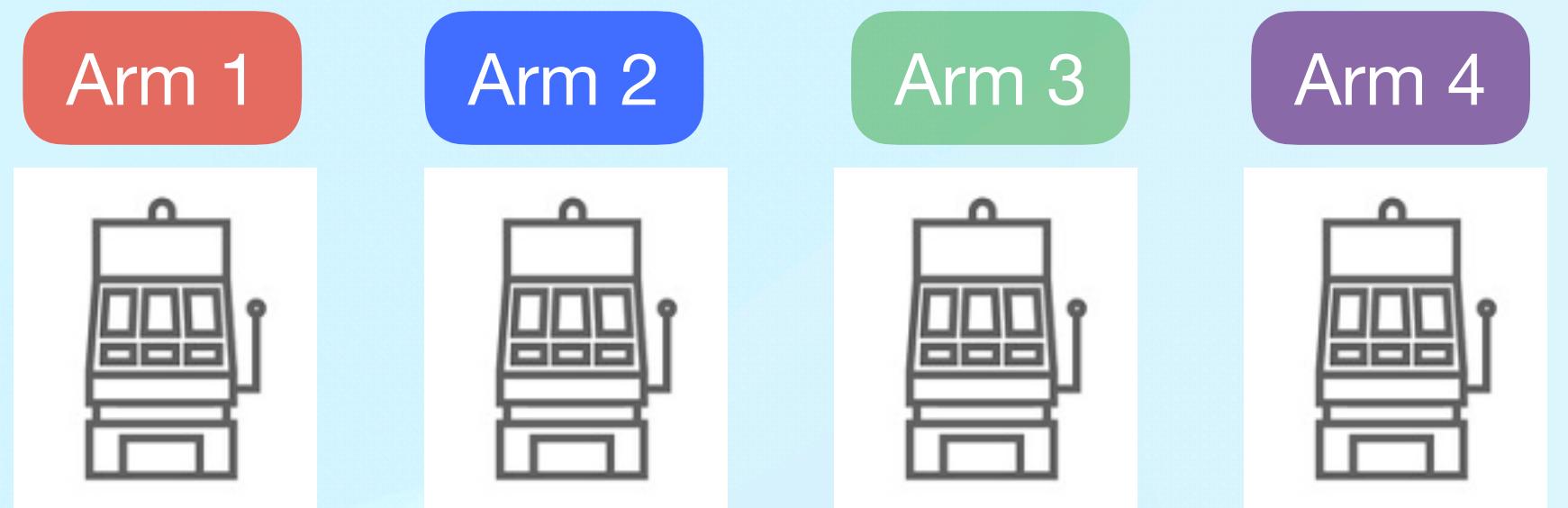
- A multi-armed bandit with $K \geq 2$ arms



* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

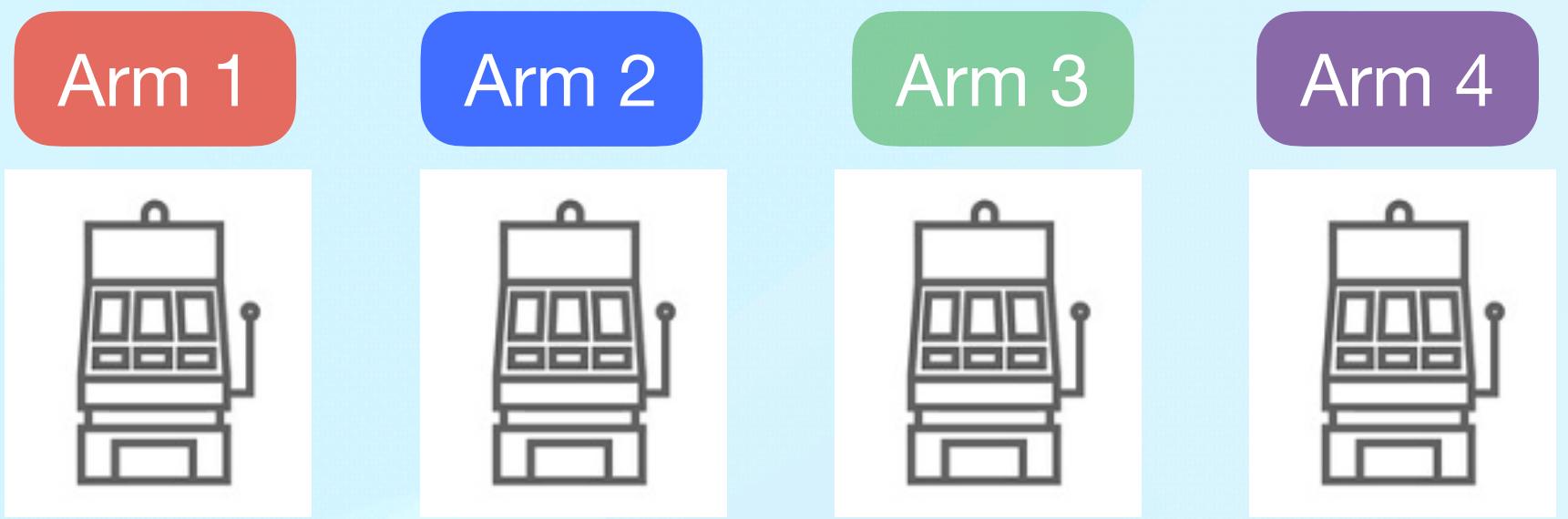
- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}



* TPM(s): transition probability matrix(ces)

PROBLEM SETUP & OBJECTIVE

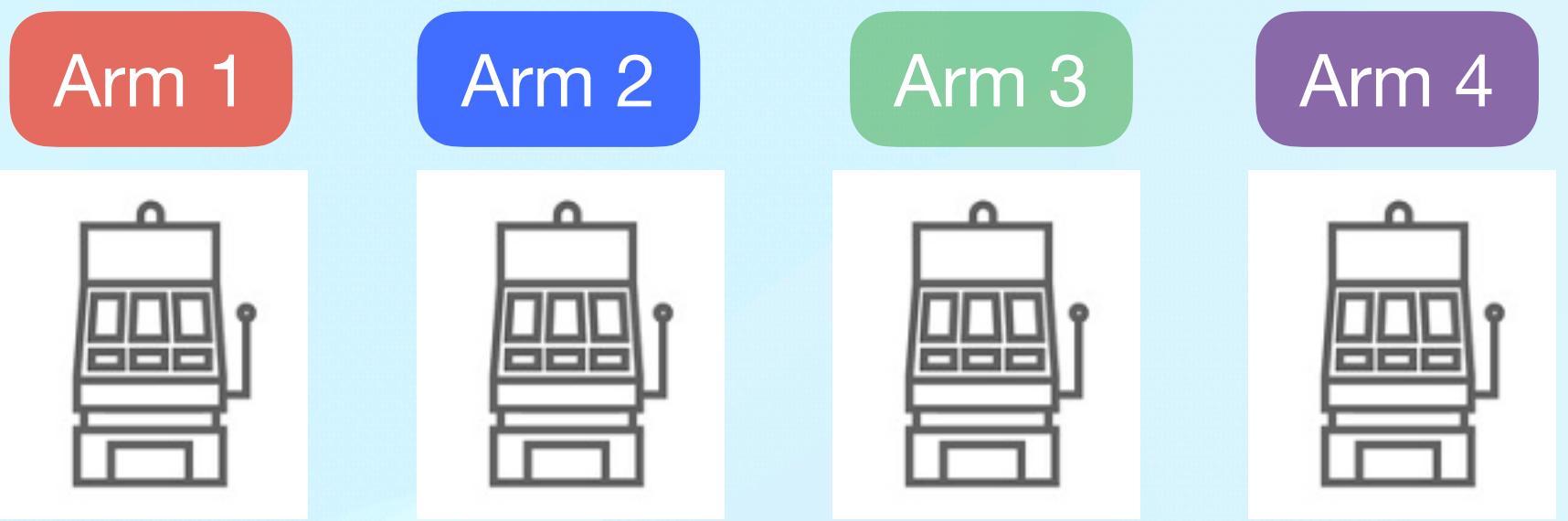
- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**



* TPM(s): transition probability matrix(ces)

PROBLEM SETUP & OBJECTIVE

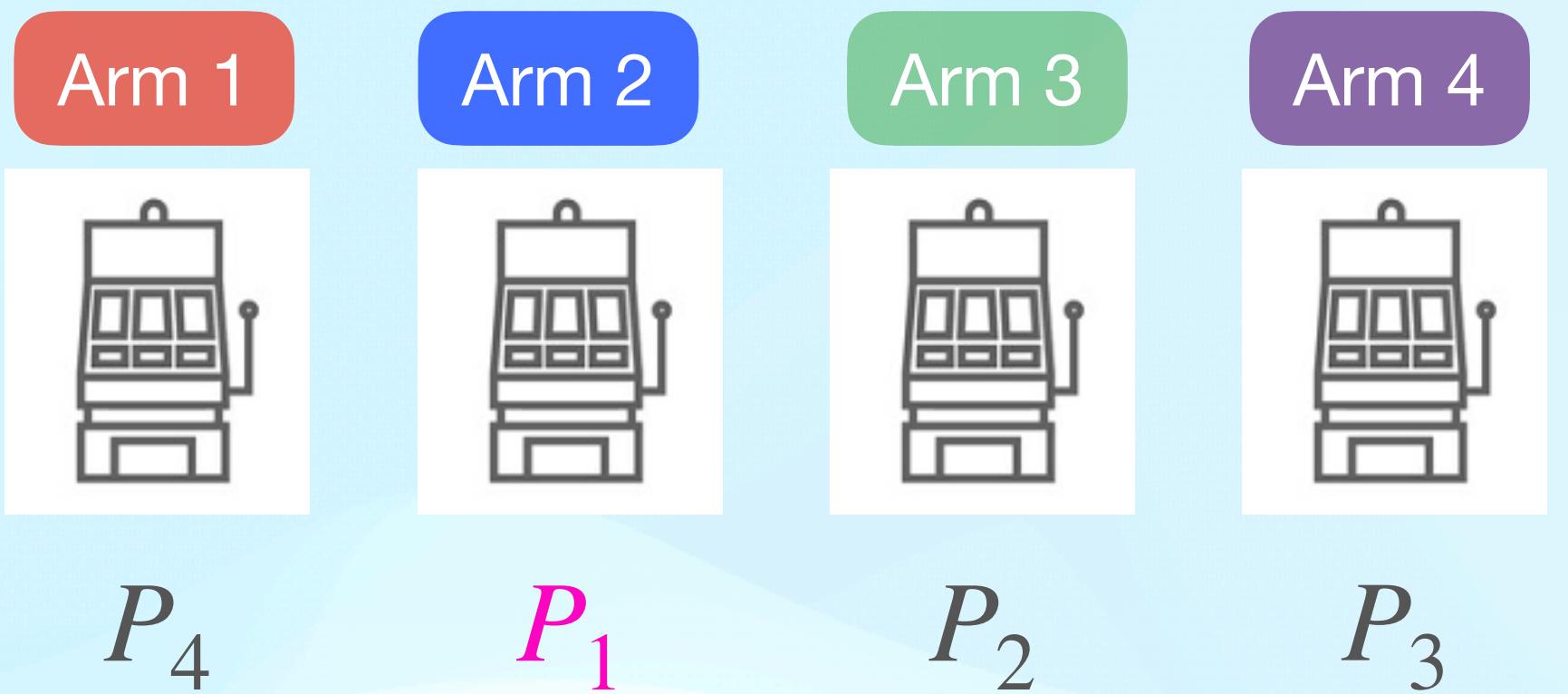
- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs* P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$



* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

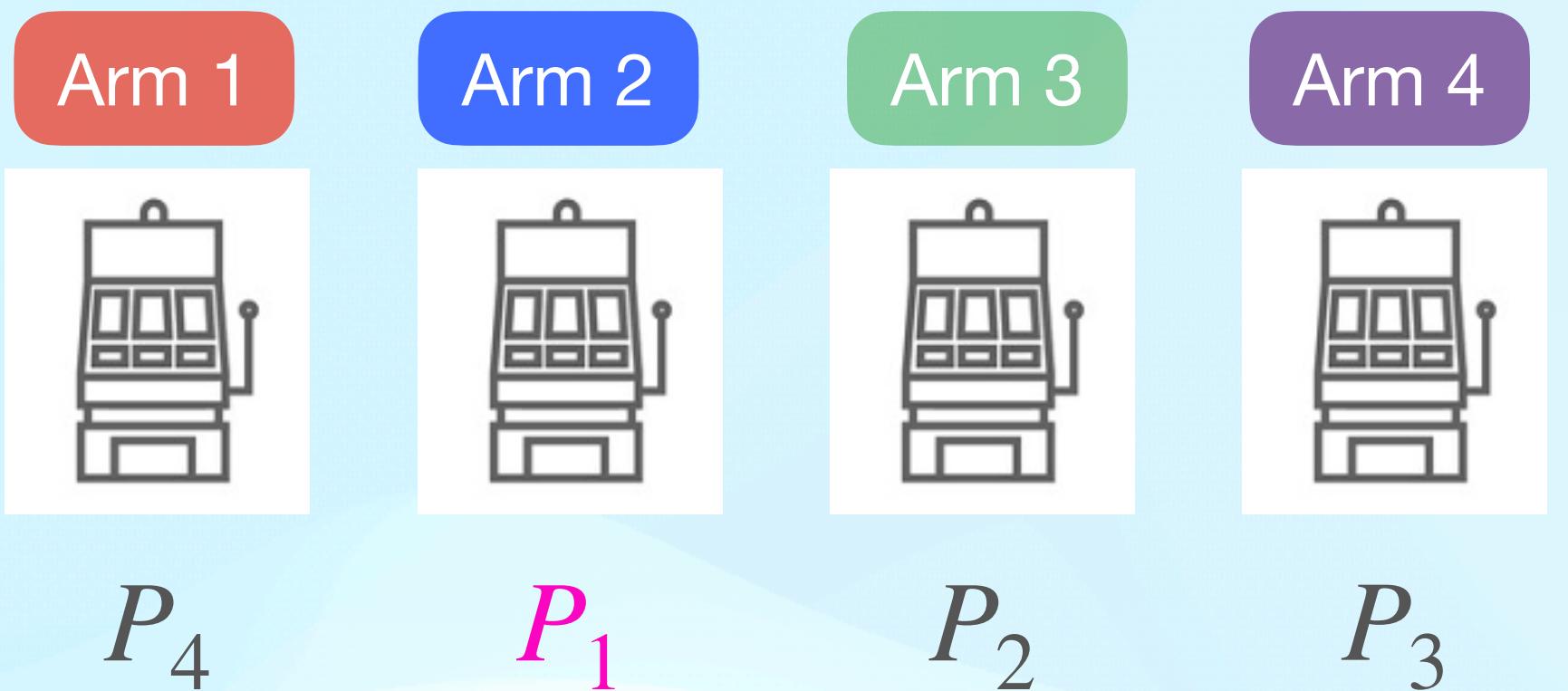
- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs* P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$



* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

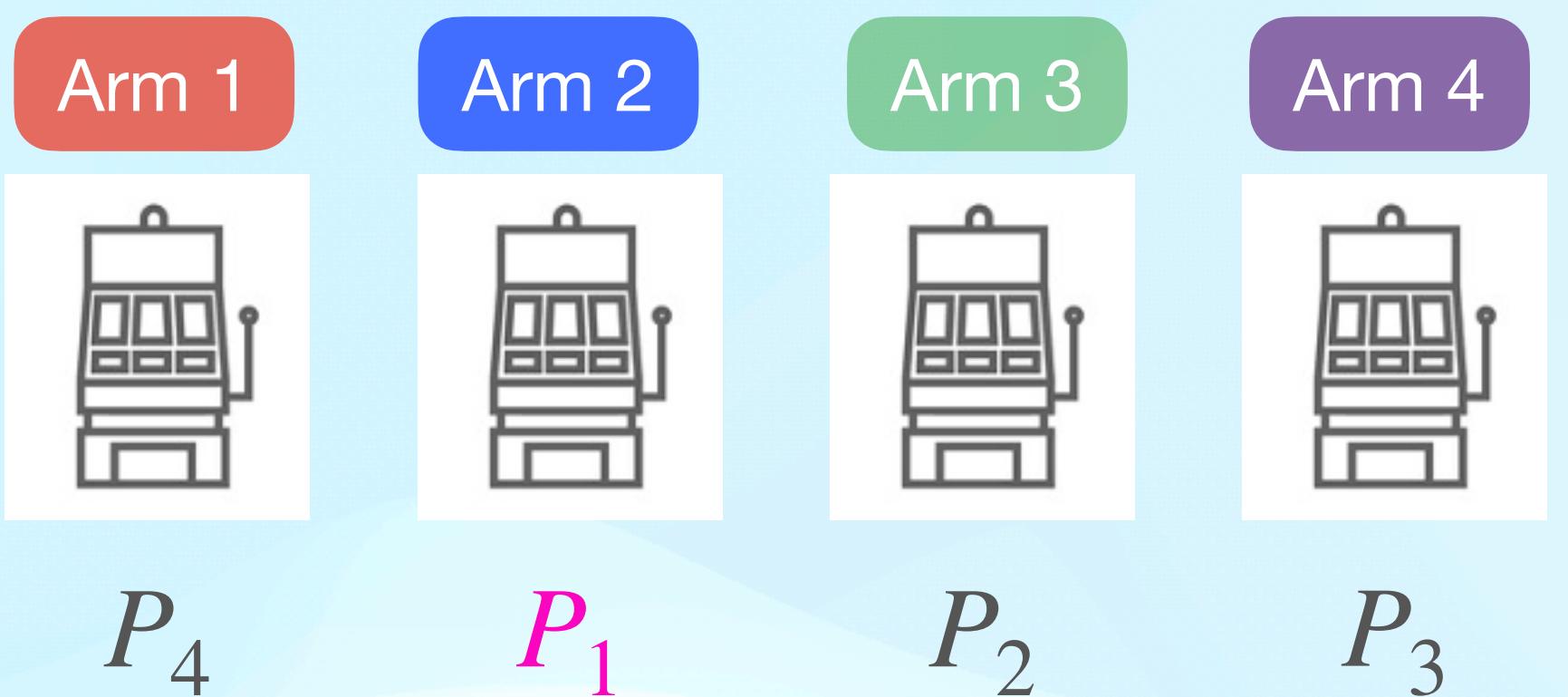
- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs^{*} P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- **TPMs are known, σ is unknown**



* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs^{*} P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- TPMs are **known**, σ is **unknown**
- Stationary distribution of TPM P_k is μ_k

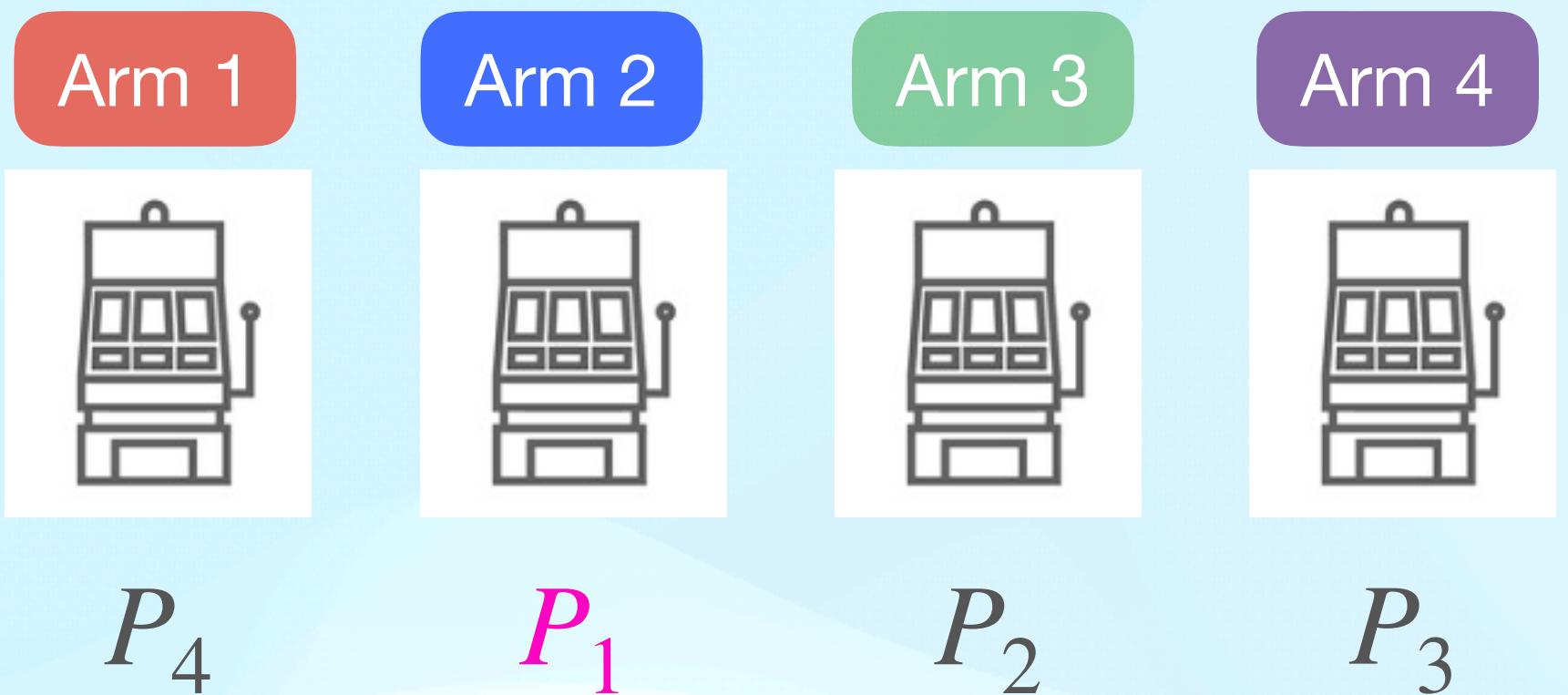


* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs* P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- **TPMs are known, σ is unknown**
- Stationary distribution of TPM P_k is μ_k
- Given $f : \mathcal{S} \rightarrow \mathbb{R}$, the **best arm**

$$a^\star := \arg \max_a \langle f, \mu_{\sigma(a)} \rangle = \arg \max_a \sum_{i \in \mathcal{S}} f(i) \mu_{\sigma(a)}(i)$$

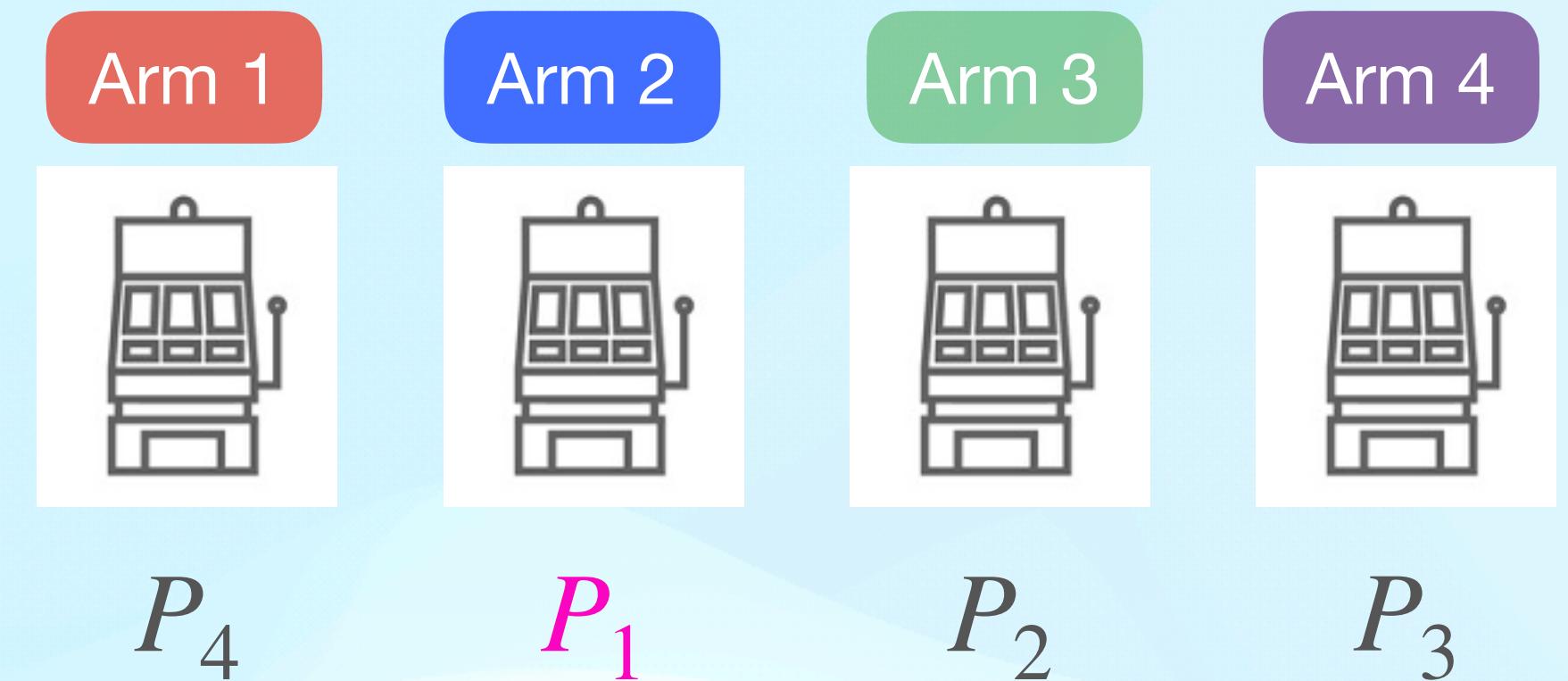


* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs* P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- TPMs are **known**, σ is **unknown**
- Stationary distribution of TPM P_k is μ_k
- Given $f : \mathcal{S} \rightarrow \mathbb{R}$, the **best arm**

$$a^\star := \arg \max_a \langle f, \mu_{\sigma(a)} \rangle = \arg \max_a \sum_{i \in \mathcal{S}} f(i) \mu_{\sigma(a)}(i)$$



goal

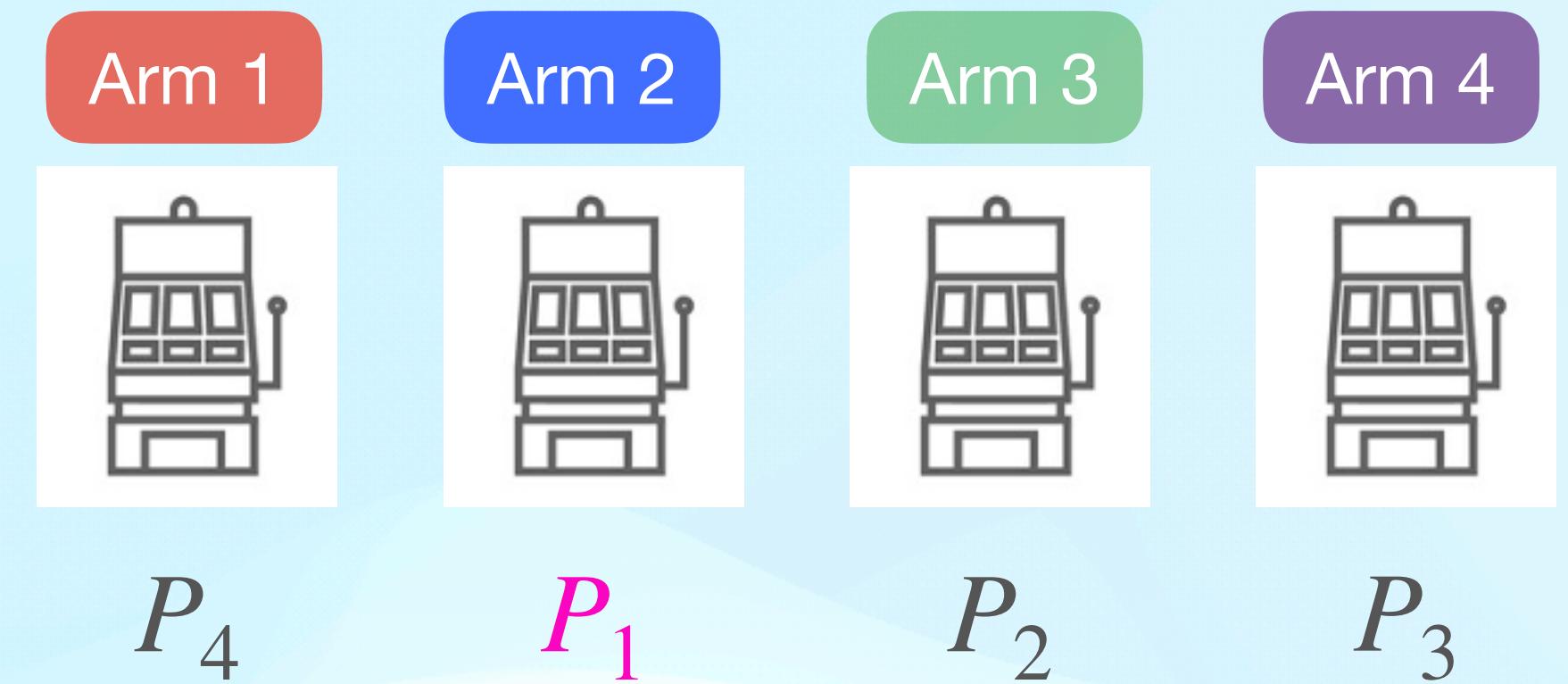
To determine the best arm quickly and with high confidence

* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs^{*} P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- TPMs are **known**, σ is **unknown**
- Stationary distribution of TPM P_k is μ_k
- Given $f : \mathcal{S} \rightarrow \mathbb{R}$, the **best arm**

$$a^\star := \arg \max_a \langle f, \mu_{\sigma(a)} \rangle = \arg \max_a \sum_{i \in \mathcal{S}} f(i) \mu_{\sigma(a)}(i)$$



accurately

$$\Pi(\delta) = \{\pi : P^\pi(\hat{a} \neq a^\star) \leq \delta\}$$

goal

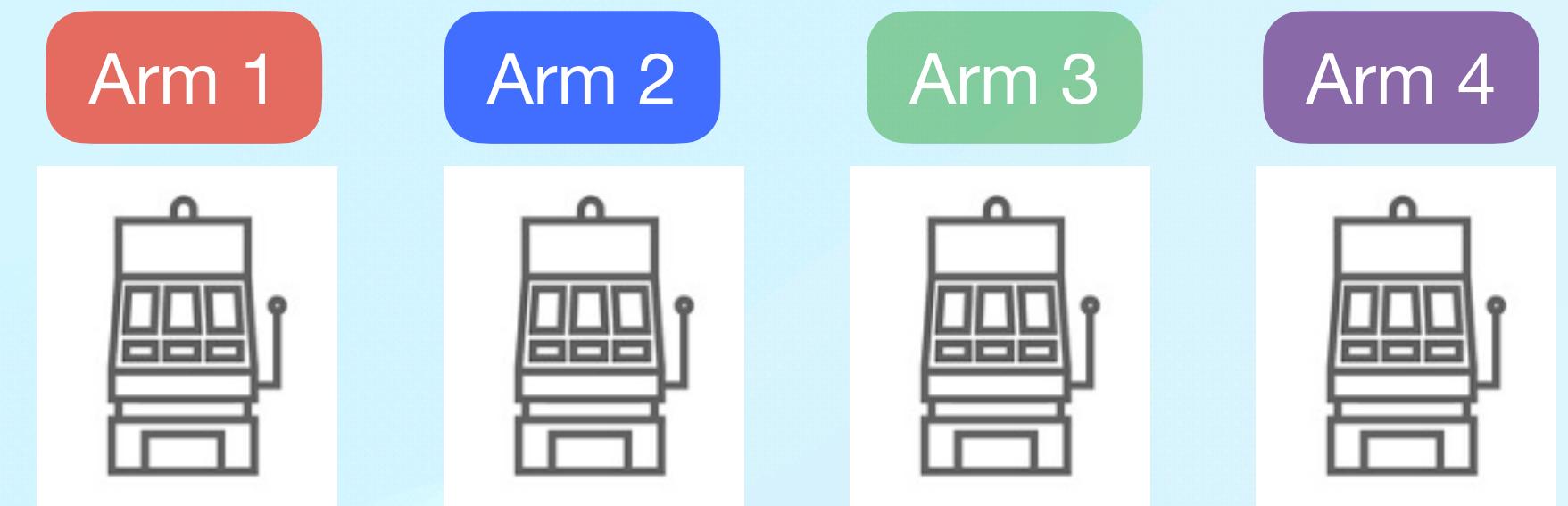
To determine the best arm quickly and with high confidence

* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs* P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- TPMs are **known**, σ is **unknown**
- Stationary distribution of TPM P_k is μ_k
- Given $f : \mathcal{S} \rightarrow \mathbb{R}$, the **best arm**

$$a^\star := \arg \max_a \langle f, \mu_{\sigma(a)} \rangle = \arg \max_a \sum_{i \in \mathcal{S}} f(i) \mu_{\sigma(a)}(i)$$



P_4

P_1

P_2

P_3

accurately

$$\Pi(\delta) = \{\pi : P^\pi(\hat{a} \neq a^\star) \leq \delta\}$$

quickly

$$\inf_{\pi \in \Pi(\delta)} \mathbb{E}^\pi[\text{stopping time under } \pi]$$

goal

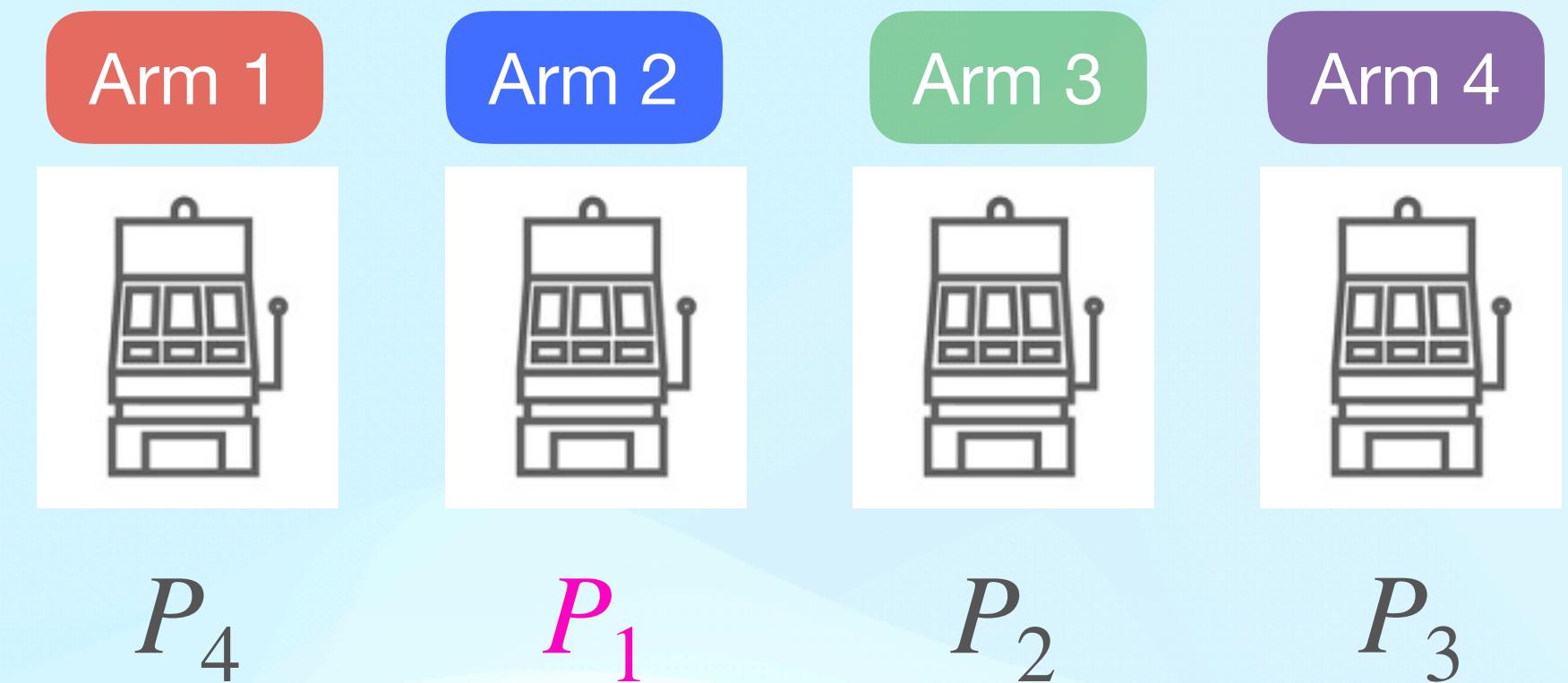
To determine the best arm quickly and with high confidence

* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs* P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- TPMs are **known**, σ is **unknown**
- Stationary distribution of TPM P_k is μ_k
- Given $f : \mathcal{S} \rightarrow \mathbb{R}$, the **best arm**

$$a^\star := \arg \max_a \langle f, \mu_{\sigma(a)} \rangle = \arg \max_a \sum_{i \in \mathcal{S}} f(i) \mu_{\sigma(a)}(i)$$



accurately $\Pi(\delta) = \{\pi : P^\pi(\hat{a} \neq a^\star) \leq \delta\}$

quickly $\inf_{\pi \in \Pi(\delta)} \mathbb{E}^\pi[\text{stopping time under } \pi]$

$\inf_{\pi \in \Pi(\delta)} \mathbb{E}^\pi[\text{stopping time under } \pi] \sim \Theta(\log(1/\delta))$

goal

To determine the best arm quickly and with high confidence

* TPM(s): transition probability matrix(es)

PROBLEM SETUP & OBJECTIVE

- A multi-armed bandit with $K \geq 2$ arms
- Arm: ergodic discrete-time **Markov** chain on **finite** state space \mathcal{S}
- Arms are **restless**
- Given TPMs* P_1, \dots, P_K and a permutation $\sigma : [K] \rightarrow [K]$,
the TPM of arm a is $P_{\sigma(a)}$
- TPMs are **known**, σ is **unknown**
- Stationary distribution of TPM P_k is μ_k
- Given $f : \mathcal{S} \rightarrow \mathbb{R}$, the **best arm**

$$a^\star := \arg \max_a \langle f, \mu_{\sigma(a)} \rangle = \arg \max_a \sum_{i \in \mathcal{S}} f(i) \mu_{\sigma(a)}(i)$$

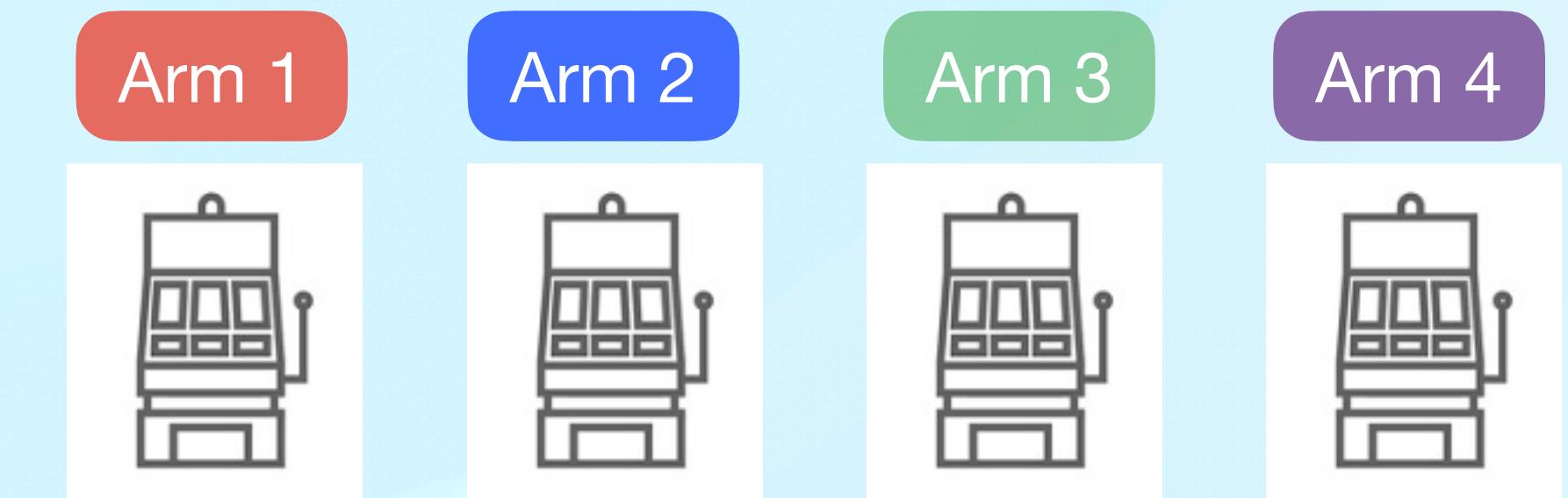
characterise or bound

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)}$$

goal

To determine the best arm quickly and with high confidence

* TPM(s): transition probability matrix(es)



P_4

P_1

P_2

P_3

accurately

$$\Pi(\delta) = \{\pi : P^\pi(\hat{a} \neq a^\star) \leq \delta\}$$

quickly

$$\inf_{\pi \in \Pi(\delta)} \mathbb{E}^\pi[\text{stopping time under } \pi]$$

$\inf_{\pi \in \Pi(\delta)}$

$$\mathbb{E}^\pi[\text{stopping time under } \pi] \sim \Theta(\log(1/\delta))$$

PRELIMINARIES

Arm 1



Arm 2



Arm 3



Arm 4



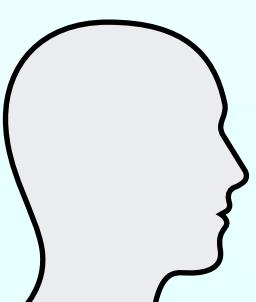
P_4

P_1

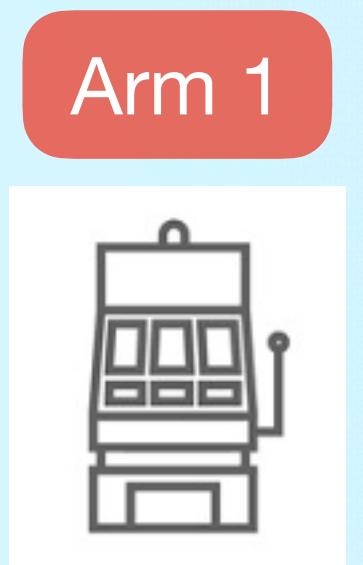
P_2

P_3

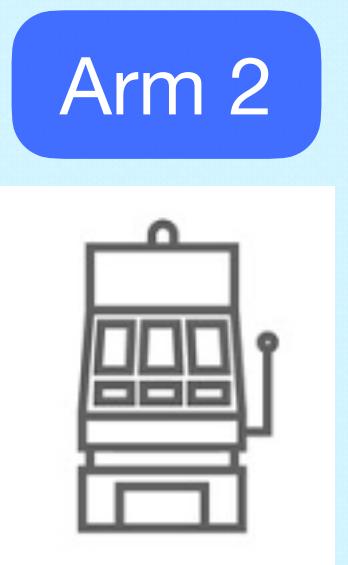
agent



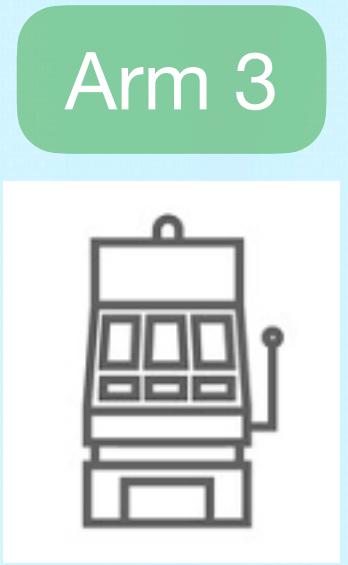
t
 a
 X_{ta}



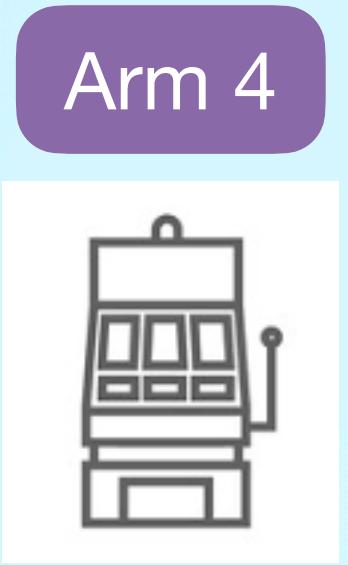
Arm 1



Arm 2



Arm 3



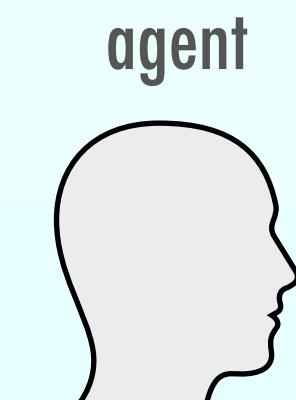
Arm 4

P₄

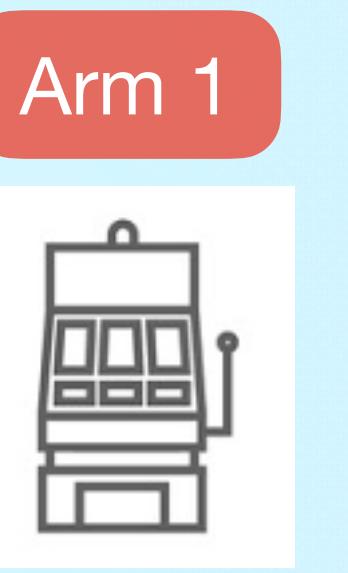
P_1

P_2

P_3



<i>f</i>							
<i>a</i>							
<i>Xia</i>							



Arm 1



Arm 2



Arm 3



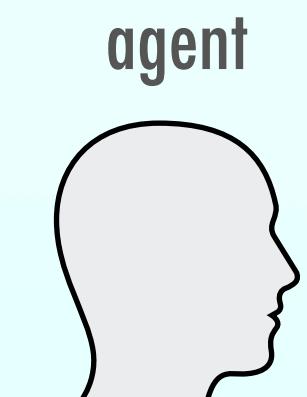
Arm 4

P_4

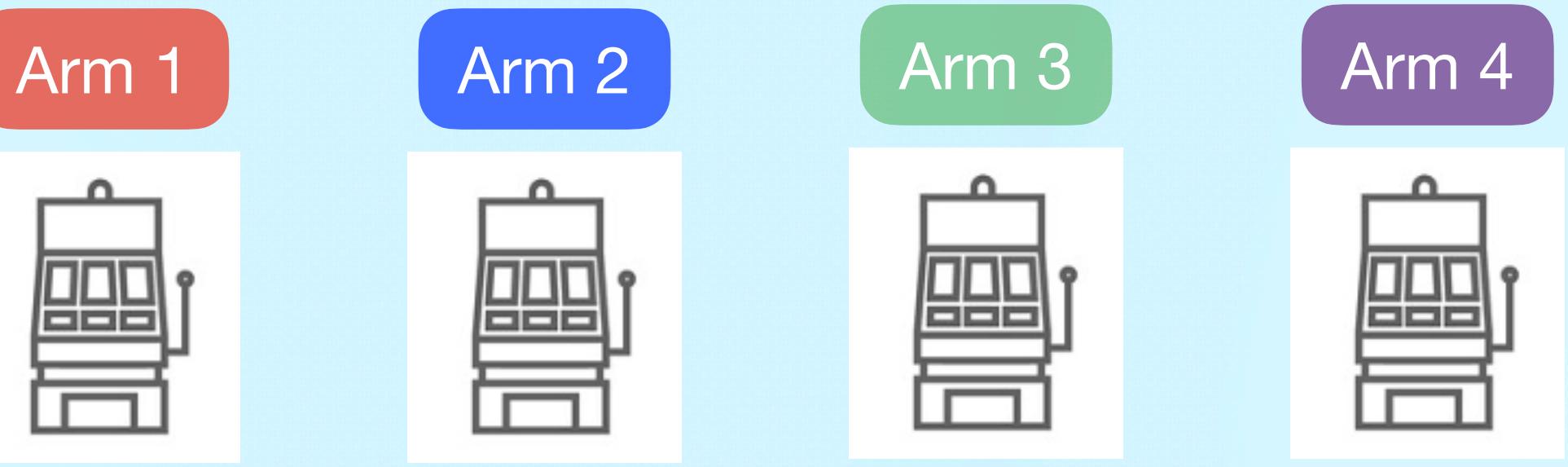
P_1

P_2

P_3



t							
a							
$X_{t,a}$							

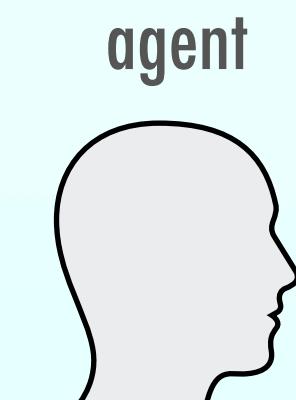


P_4

P_1

P_2

P_3



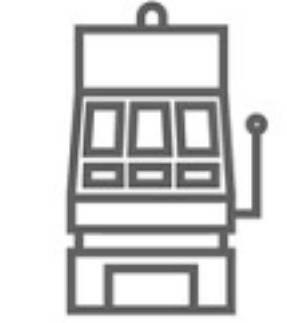
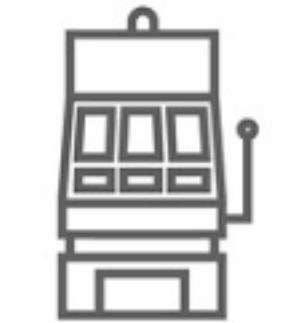
t	0						
a							
$X_{t,a}$							

Arm 1

Arm 2

Arm 3

Arm 4



P_4

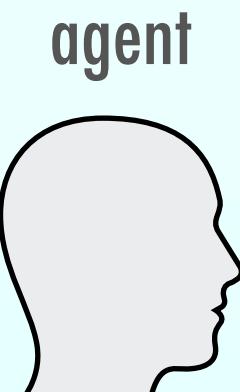
P_1

P_2

P_3

$t = 0$

$X_{0,1}$



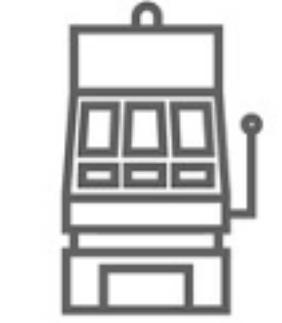
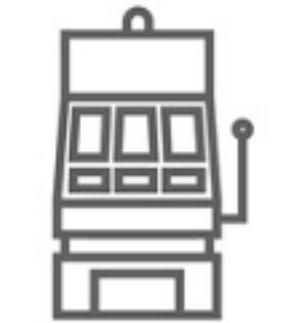
agent								
t	0							
a								
$X_{t,a}$	$X_{0,1}$							

Arm 1

Arm 2

Arm 3

Arm 4



P_4

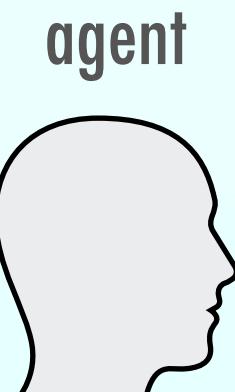
P_1

P_2

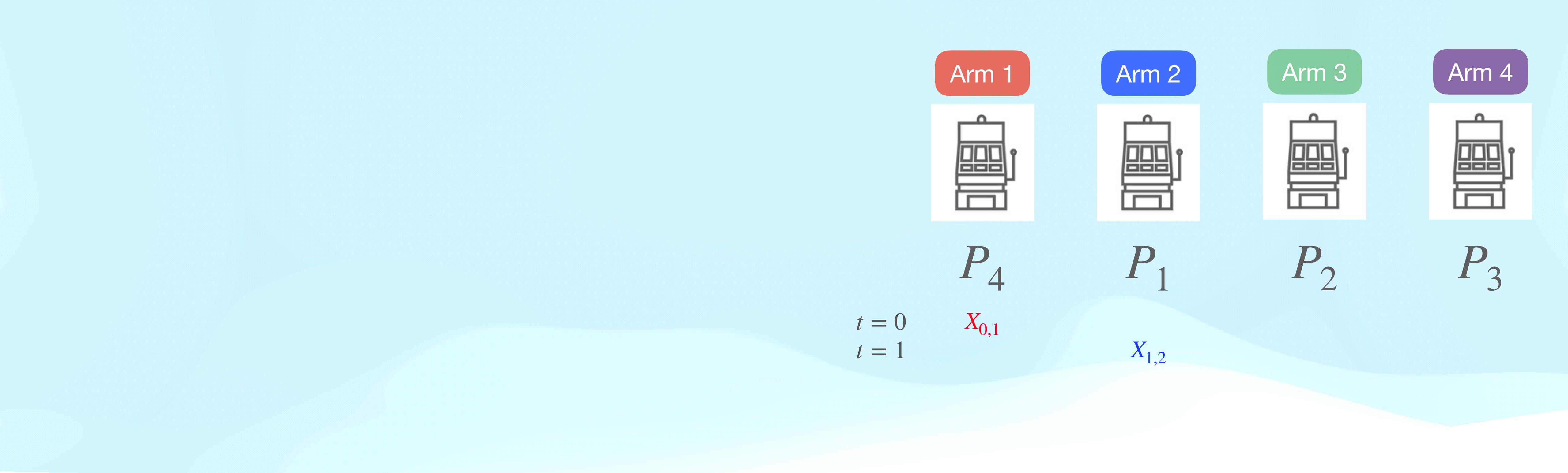
P_3

$t = 0$

$X_{0,1}$



agent									
t	0	1							
a	1	2							
$X_{t,a}$		$X_{0,1}$							



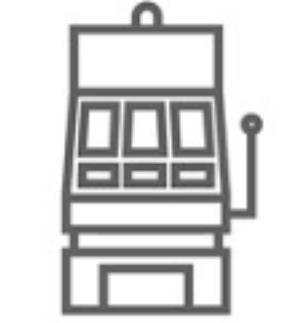
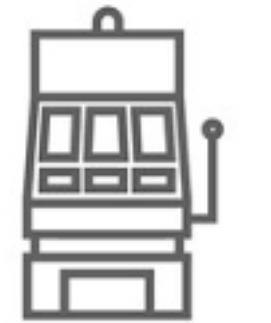
agent								
t	0	1						
a	1	2						
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$						

Arm 1

Arm 2

Arm 3

Arm 4



P_4

P_1

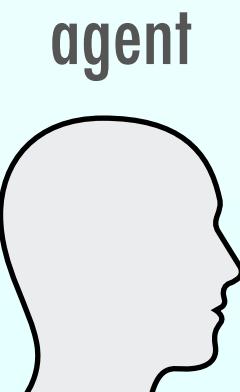
P_2

P_3

$t = 0$
 $t = 1$

$X_{0,1}$

$X_{1,2}$



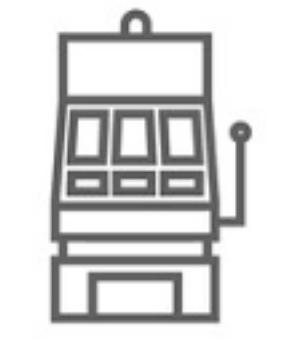
agent								
t	0	1	2					
a	1	2	3					
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$						

Arm 1

Arm 2

Arm 3

Arm 4



P_4

P_1

P_2

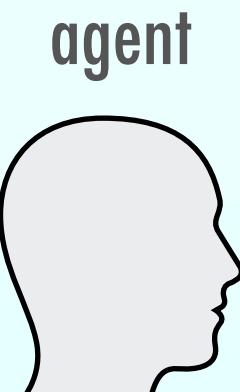
P_3

$t = 0$
 $t = 1$
 $t = 2$

$X_{0,1}$

$X_{1,2}$

$X_{2,3}$



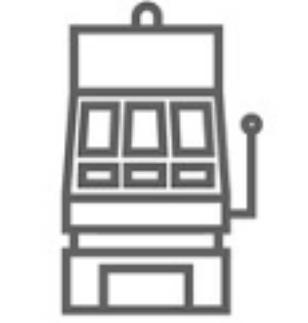
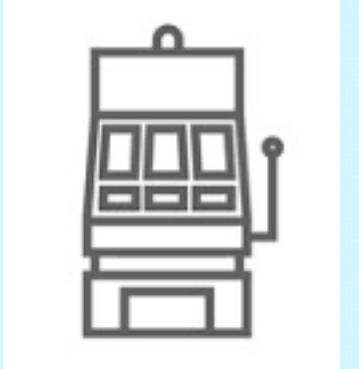
agent								
t	0	1	2					
a	1	2	3					
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$					

Arm 1

Arm 2

Arm 3

Arm 4



P_4

P_1

P_2

P_3

$t = 0$

$X_{0,1}$

$t = 1$

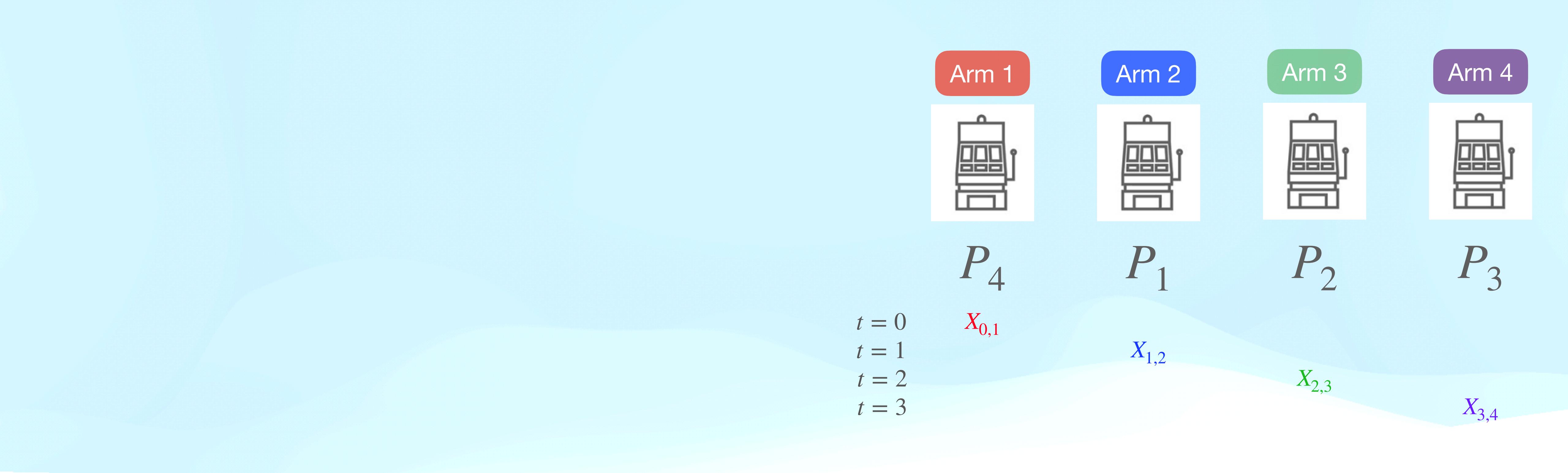
$X_{1,2}$

$t = 2$

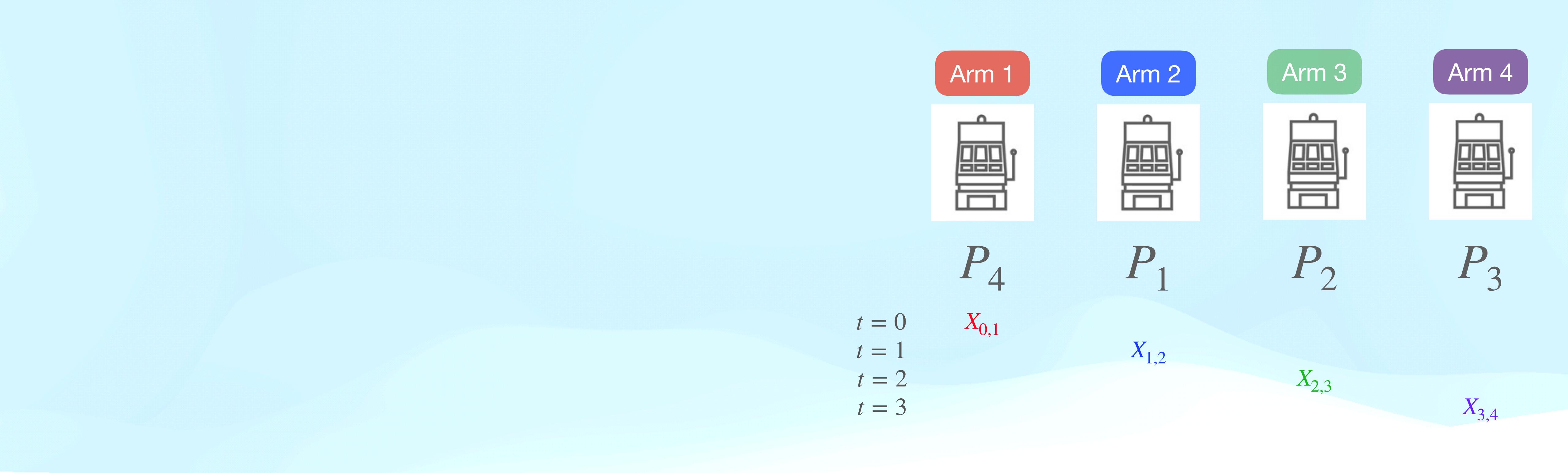
$X_{2,3}$



agent	0	1	2	3	4				
t	0	1	2	3	4				
a	1	2	3	4					
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$						



agent									
	t	0	1	2	3				
a	1	2	3	4					
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$					



agent							
	t	0	1	2	3	4	
	a	1	2	3	4	3	
	$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$		



P_3



Arm 3



Arm 2



Arm 1

$$\begin{aligned}t &= 0 \\ t &= 1 \\ t &= 2 \\ t &= 3 \\ t &= 4\end{aligned}$$

$$X_{0,1}$$

P_1

X_{1,2}

P_2

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

t	0	1	2	3	4		
a	1	2	3	4	3		
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,3}$		



P_3



Arm 3



Arm 2



Arm 1

$$\begin{aligned}t &= 0 \\ t &= 1 \\ t &= 2 \\ t &= 3 \\ t &= 4\end{aligned}$$

$$X_{0,1}$$

P_1

$$X_{1,2}$$

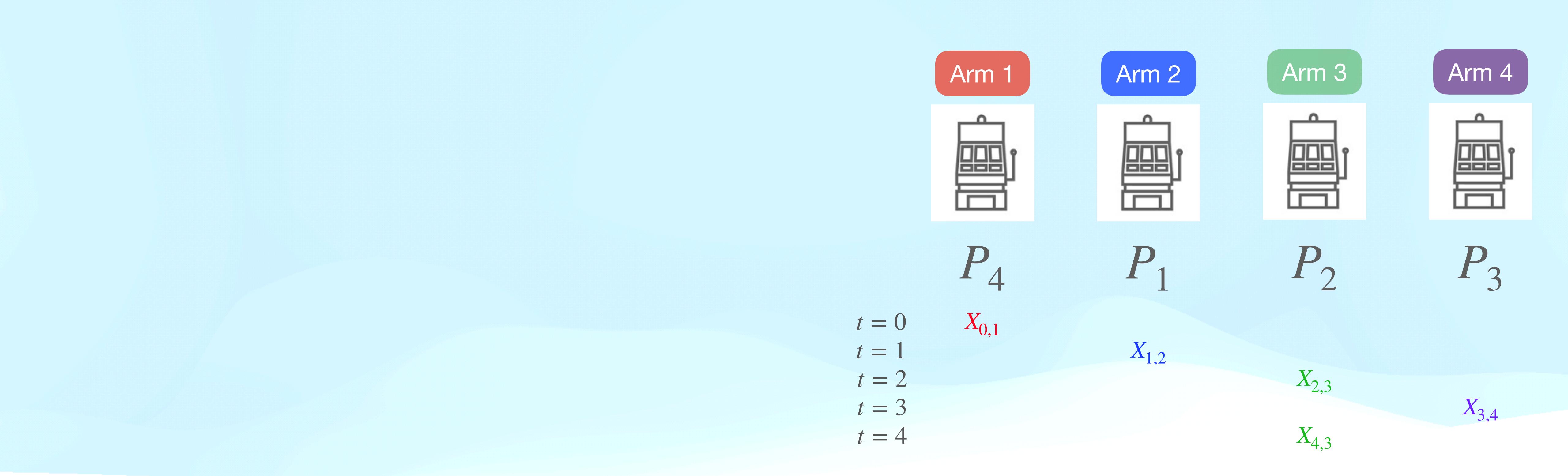
$$P_2$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

t	0	1	2	3	4	5		
a	1	2	3	4	3	3		
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,3}$			



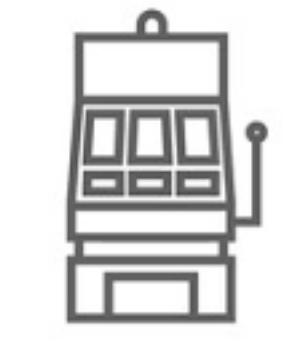
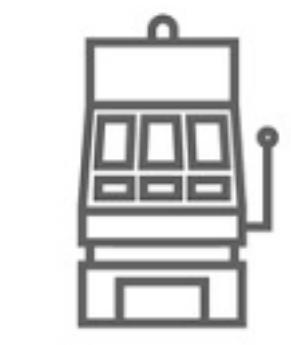
agent									
	t	0	1	2	3	4	5	6	
a	1	2	3	4	3	3	2		
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,3}$				

Arm 1

Arm 2

Arm 3

Arm 4



P_4

P_1

P_2

P_3

$t = 0$

$X_{0,1}$

$t = 1$

$X_{1,2}$

$t = 2$

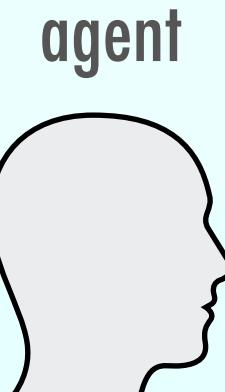
$X_{2,3}$

$t = 3$

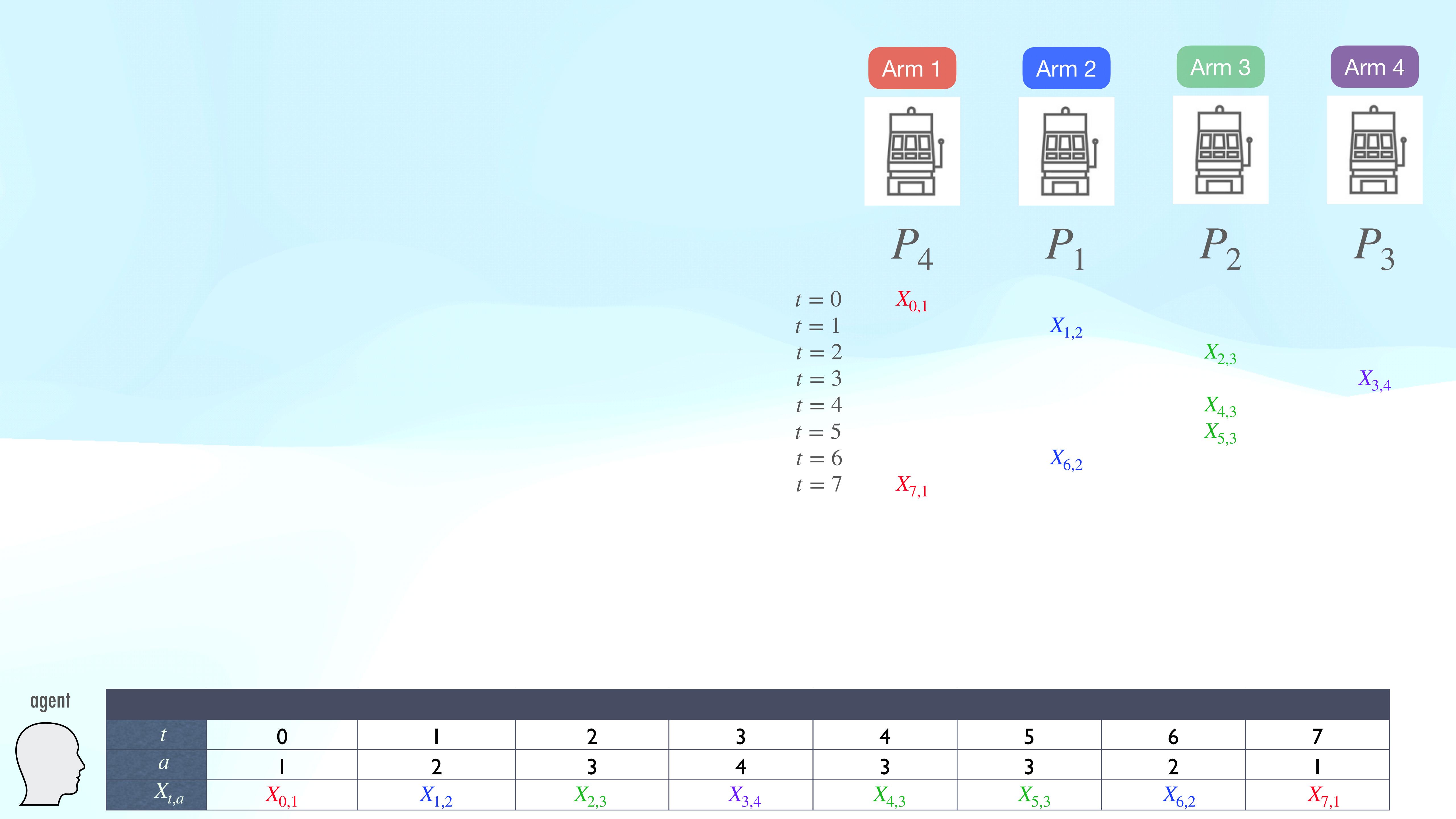
$X_{3,4}$

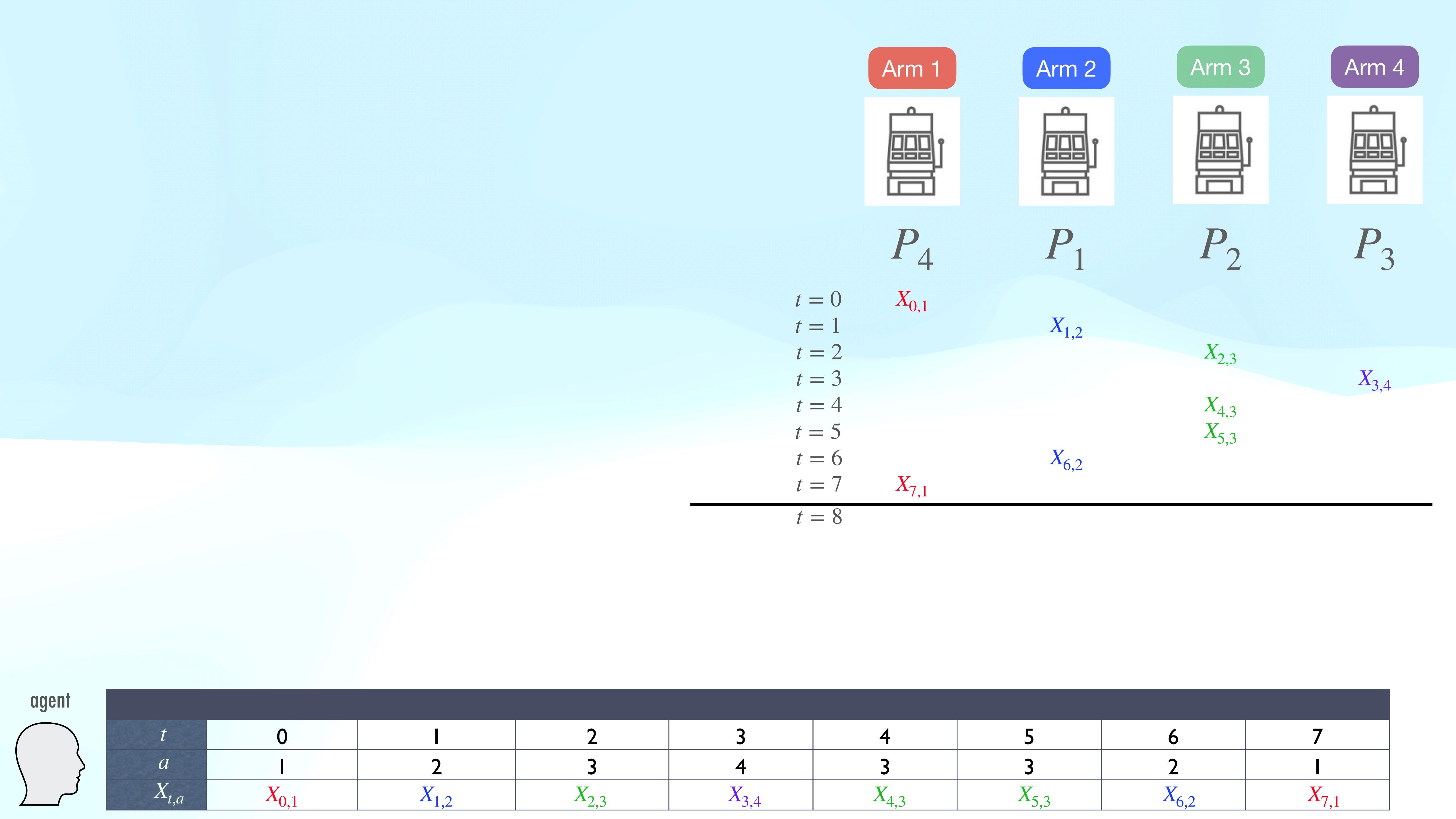
$t = 4$

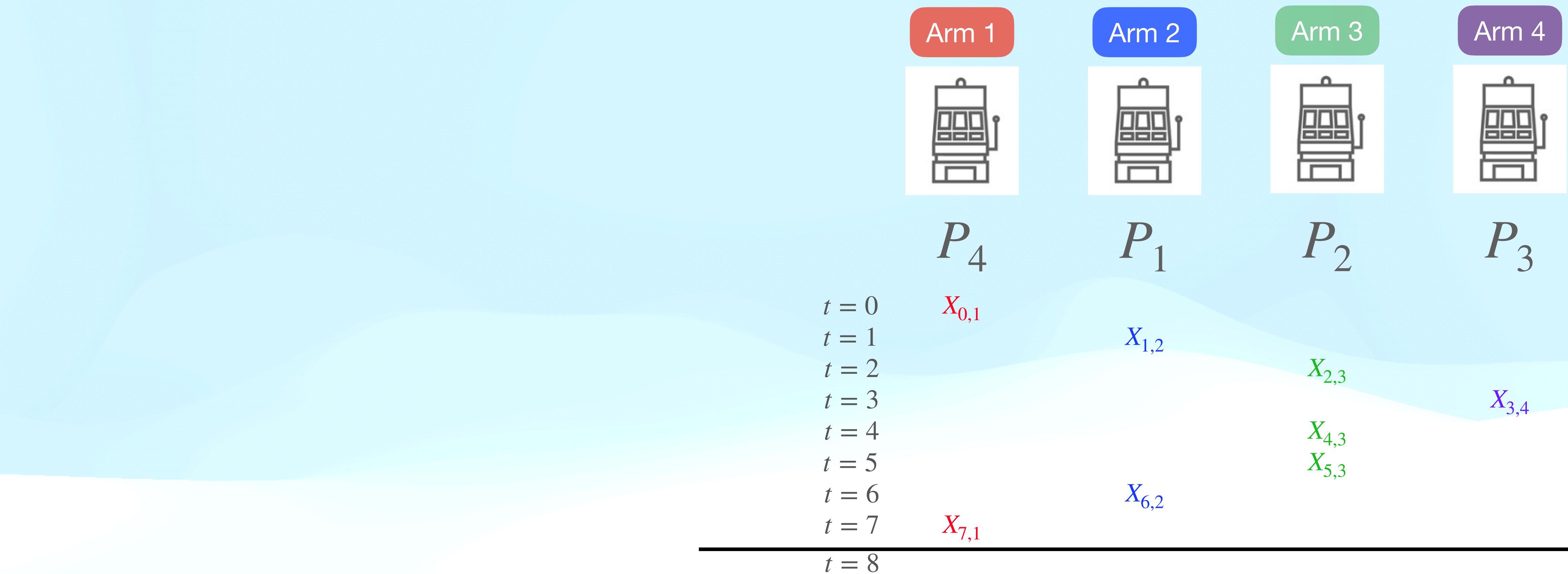
$X_{4,3}$



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







$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

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

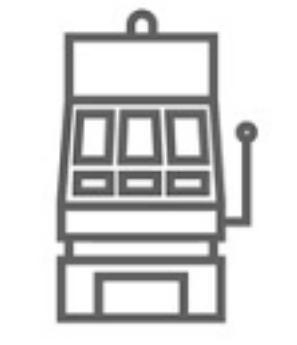
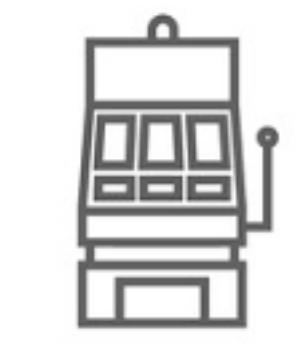
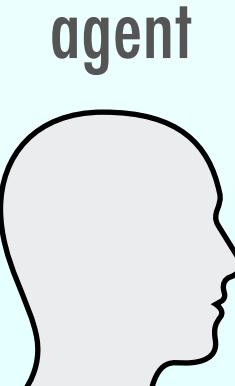
$t = 8$

Arm 1

Arm 2

Arm 3

Arm 4

 P_4 P_1 P_2 P_3 $t = 0$ $X_{0,1}$ $t = 1$ $X_{1,2}$ $t = 2$ $X_{2,3}$ $t = 3$ $X_{3,4}$ $t = 4$ $X_{4,3}$ $t = 5$ $X_{5,3}$ $t = 6$ $X_{6,2}$ $t = 7$ $X_{7,1}$ $t = 8$ $d_a(t) : \# \text{ time instants ago arm } a \text{ was previously observed w.r.t. time } t$ 

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

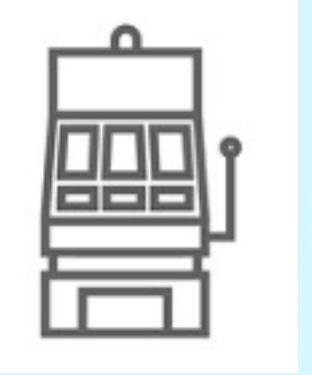
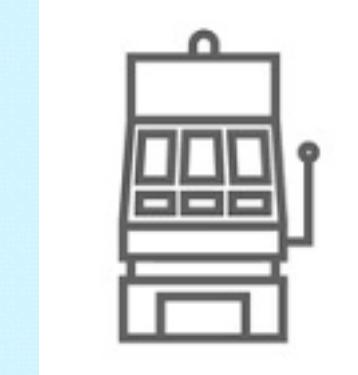
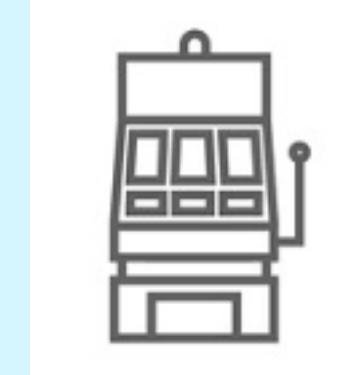
$$t = 8 \quad d_1(t) = 1$$

Arm 1

Arm 2

Arm 3

Arm 4



P_4

P_1

P_2

P_3

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

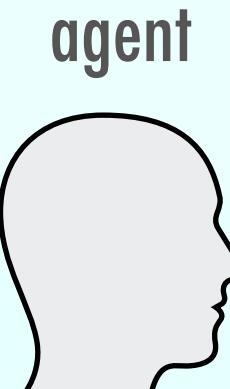
$$X_{6,2}$$

$$t = 7$$

$$X_{7,1}$$

$$t = 8$$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t



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

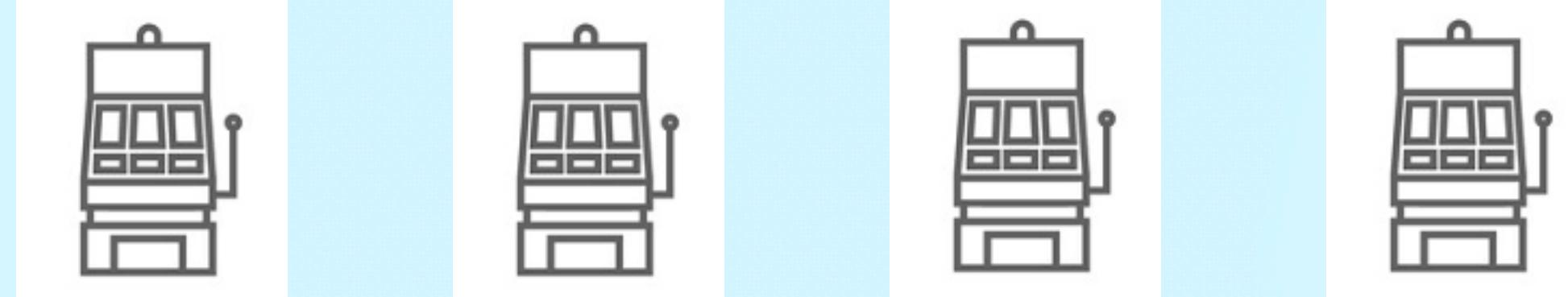
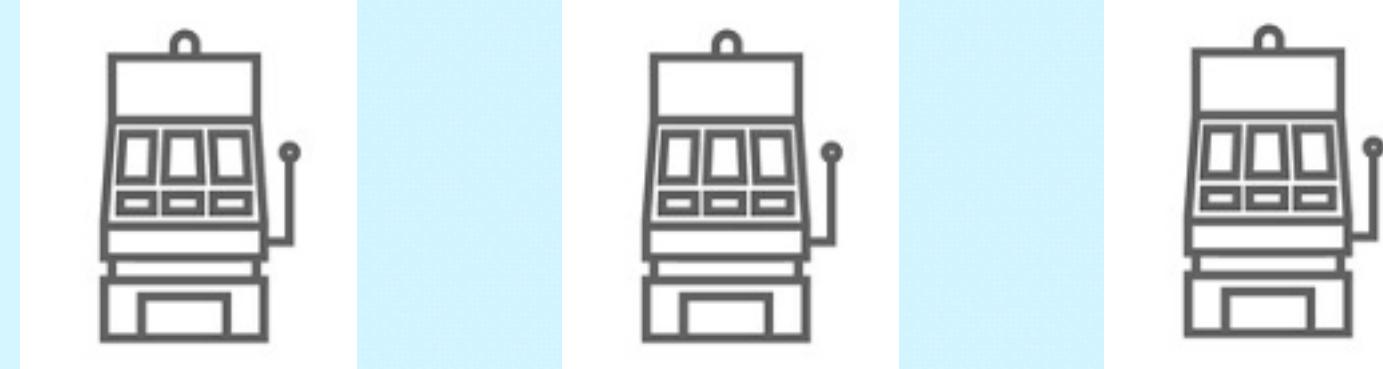
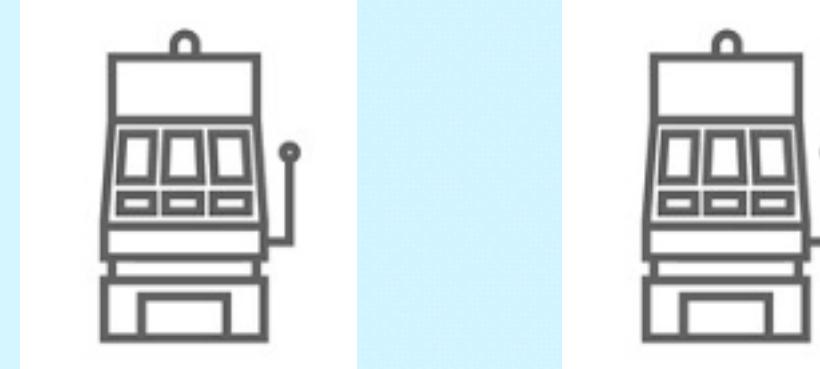
$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2$$

Arm 1

Arm 2

Arm 3

Arm 4



P_4

P_1

P_2

P_3

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

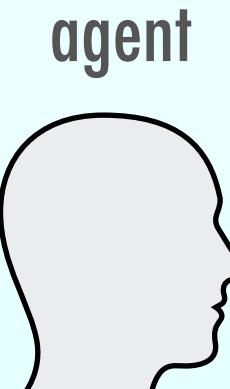
$$X_{6,2}$$

$$t = 7$$

$$X_{7,1}$$

$$t = 8$$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t



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

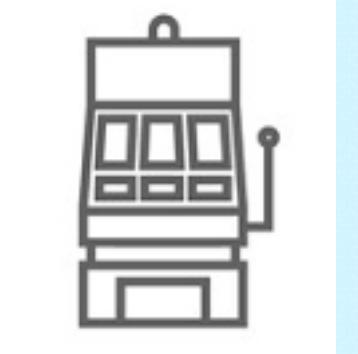
$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3$$

Arm 1

Arm 2

Arm 3

Arm 4



P_4

P_1

P_2

P_3

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

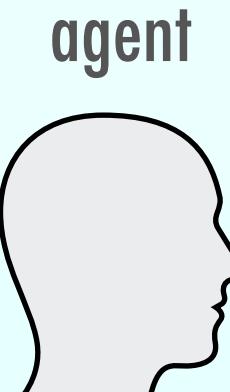
$$X_{6,2}$$

$$t = 7$$

$$X_{7,1}$$

$$t = 8$$

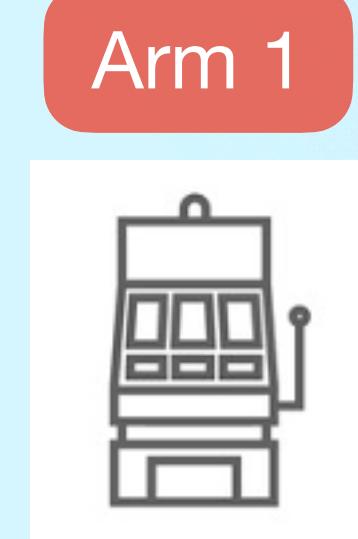
$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t



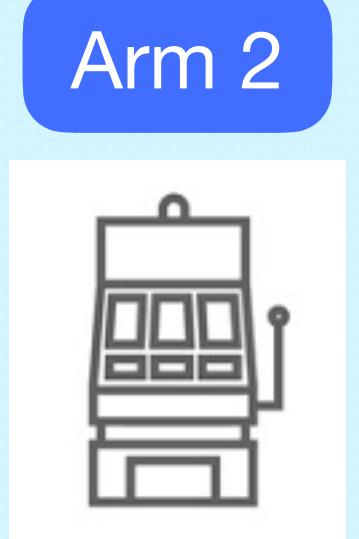
agent	t	0	1	2	3	4	5	6	7
a	1	2	3	4	3	3	2	1	
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,3}$	$X_{5,3}$	$X_{6,2}$	$X_{7,1}$	

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

Arm 1



Arm 2



Arm 3



Arm 4



P_4

$$t = 0$$

$X_{0,1}$

$$t = 1$$

$X_{1,2}$

$$t = 2$$

$X_{2,3}$

$$t = 3$$

$X_{3,4}$

$$t = 4$$

$X_{4,3}$

$$t = 5$$

$X_{5,3}$

$$t = 6$$

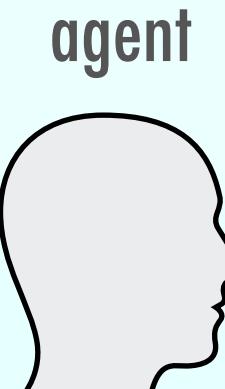
$X_{6,2}$

$$t = 7$$

$X_{7,1}$

$$t = 8$$

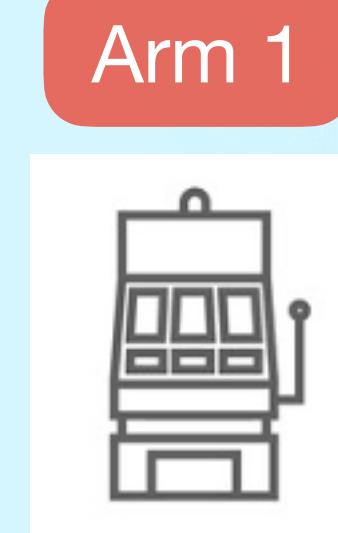
$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t



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

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

Arm 1



Arm 2



Arm 3



Arm 4



P_4

$$t = 0$$

$X_{0,1}$

$$t = 1$$

$X_{1,2}$

$$t = 2$$

$X_{2,3}$

$$t = 3$$

$X_{3,4}$

$$t = 4$$

$X_{4,3}$

$$t = 5$$

$X_{5,3}$

$$t = 6$$

$X_{6,2}$

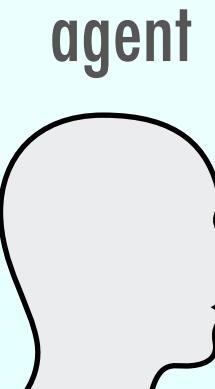
$$t = 7$$

$X_{7,1}$

$$t = 8$$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

$i_a(t)$: last observed state of arm a w.r.t. time t

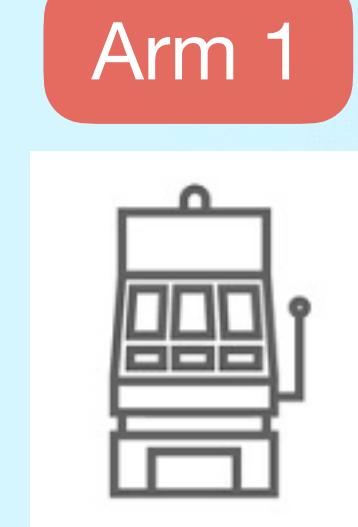


agent	t	0	1	2	3	4	5	6	7
a	1	2	3	4	3	3	2	1	
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,3}$	$X_{5,3}$	$X_{6,2}$	$X_{7,1}$	

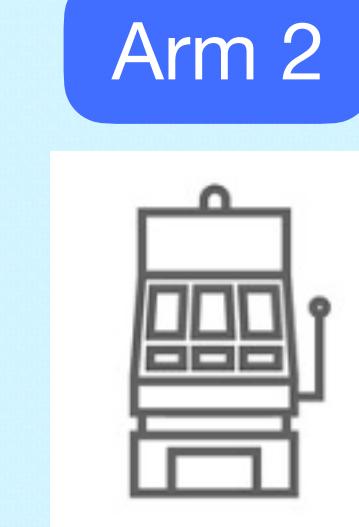
$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1}$$

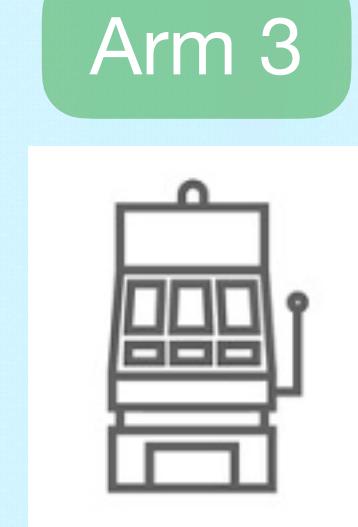
Arm 1



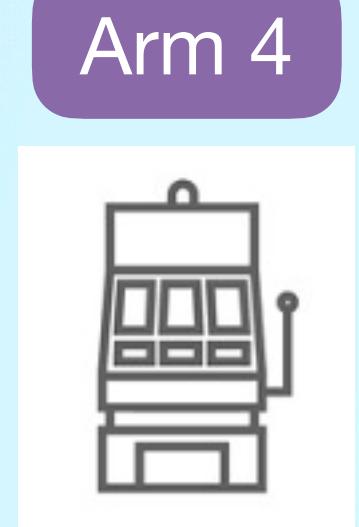
Arm 2



Arm 3



Arm 4



P_4

$$t = 0$$

$X_{0,1}$

$$t = 1$$

$X_{1,2}$

$$t = 2$$

$X_{2,3}$

$$t = 3$$

$X_{3,4}$

$$t = 4$$

$X_{4,3}$

$$t = 5$$

$X_{5,3}$

$$t = 6$$

$X_{6,2}$

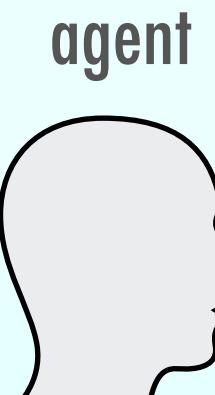
$$t = 7$$

$X_{7,1}$

$$t = 8$$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

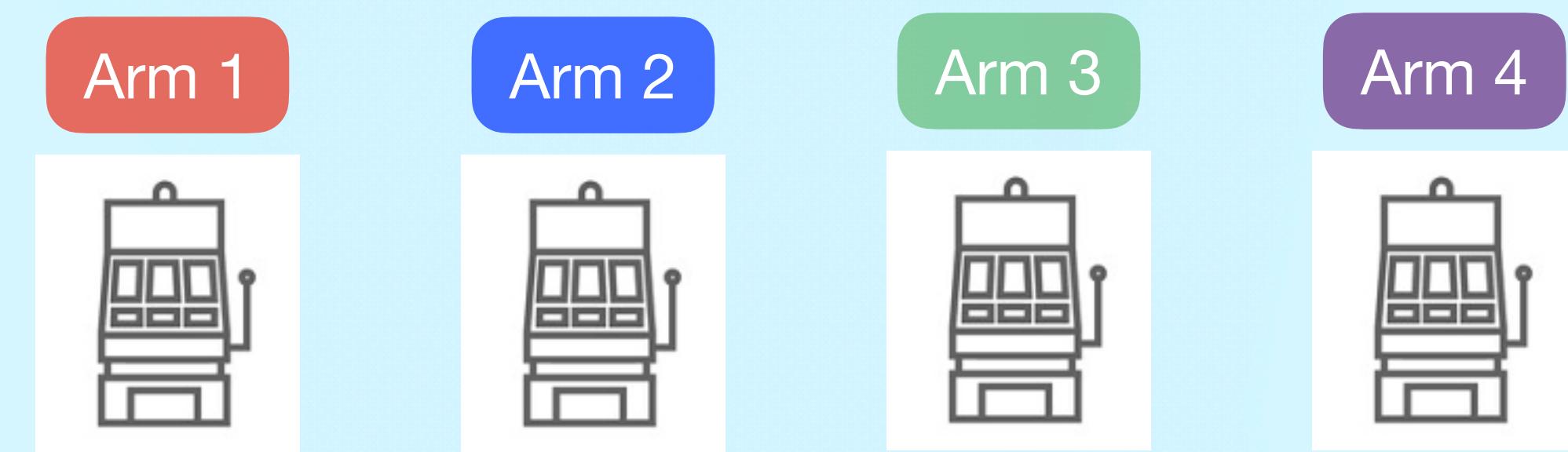
$i_a(t)$: last observed state of arm a w.r.t. time t



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

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2}$$



P_4 P_1 P_2 P_3

$t = 0 \quad X_{0,1}$
 $t = 1 \quad X_{1,2}$
 $t = 2$
 $t = 3$
 $t = 4$
 $t = 5$
 $t = 6 \quad X_{6,2}$
 $t = 7 \quad X_{7,1}$
 $t = 8$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t
 $i_a(t)$: last observed state of arm a w.r.t. time t

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

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3}$$

Arm 1



Arm 2



Arm 3



Arm 4



P_4

$$\begin{aligned} t &= 0 & X_{0,1} \\ t &= 1 \\ t &= 2 \\ t &= 3 \\ t &= 4 \\ t &= 5 \\ t &= 6 \\ t &= 7 & X_{7,1} \\ t &= 8 \end{aligned}$$

P_1

$X_{1,2}$

P_2

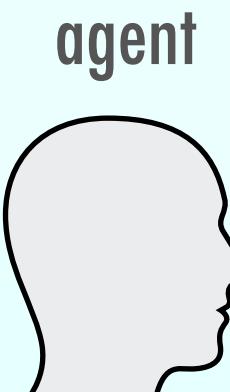
$X_{2,3}$

P_3

$X_{3,4}$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

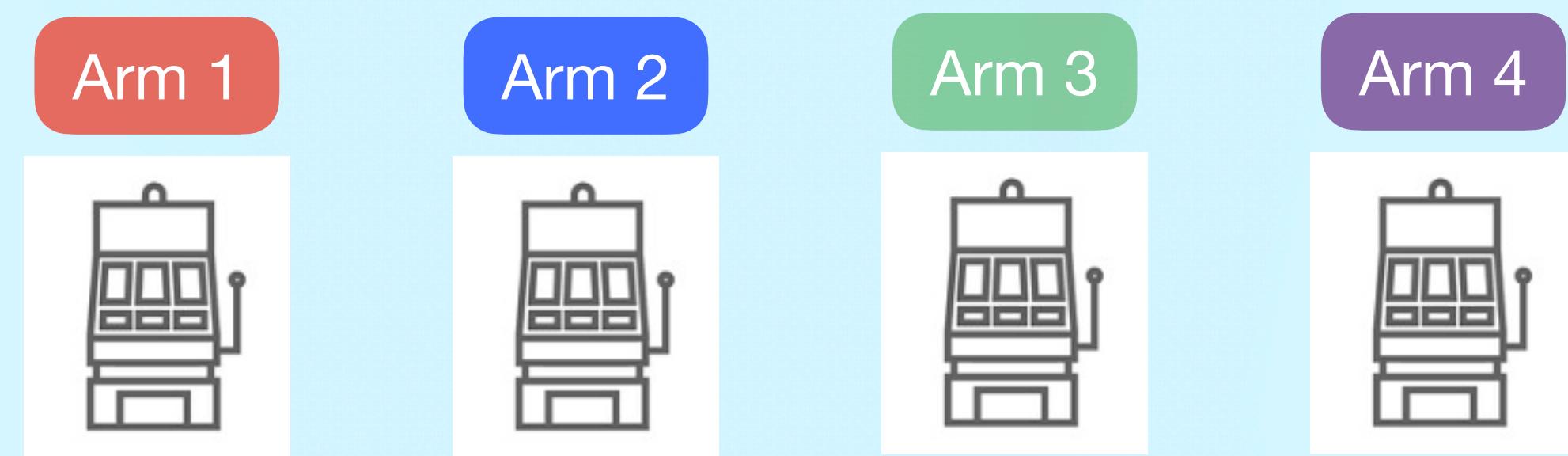
$i_a(t)$: last observed state of arm a w.r.t. time t



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

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3} \quad i_4(t) = X_{3,4}$$



P_4 P_1 P_2 P_3

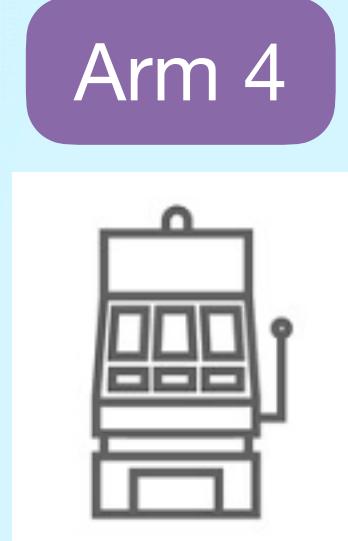
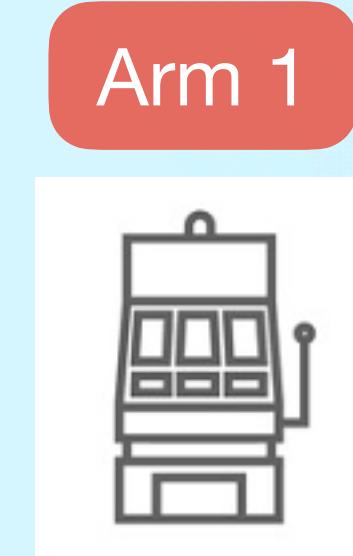
$t = 0 \quad X_{0,1}$
 $t = 1 \quad X_{1,2}$
 $t = 2$
 $t = 3$
 $t = 4$
 $t = 5$
 $t = 6 \quad X_{6,2}$
 $t = 7 \quad X_{7,1}$
 $t = 8$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t
 $i_a(t)$: last observed state of arm a w.r.t. time t

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

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3} \quad i_4(t) = X_{3,4}$$



P_4

P_1

P_2

P_3

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

$$X_{6,2}$$

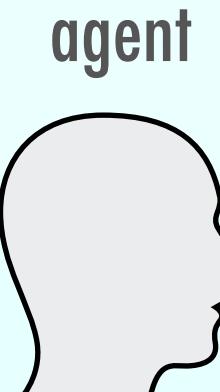
$$t = 7$$

$$X_{7,1}$$

$$t = 8$$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

$i_a(t)$: last observed state of arm a w.r.t. time t



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

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3} \quad i_4(t) = X_{3,4}$$

$$t = 5$$



P_4

P_1

P_2

P_3

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

$$X_{6,2}$$

$$t = 7$$

$$X_{7,1}$$

$$t = 8$$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

$i_a(t)$: last observed state of arm a w.r.t. time t

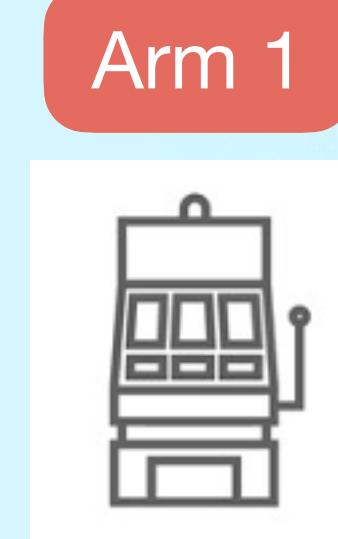
agent									
t	0	1	2	3	4	5	6	7	
a	1	2	3	4	3	3	2	1	
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,3}$	$X_{5,3}$	$X_{6,2}$	$X_{7,1}$	

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3} \quad i_4(t) = X_{3,4}$$

$$t = 5 \quad d_1(t) = 5 \quad d_2(t) = 4 \quad d_3(t) = 1 \quad d_4(t) = 2$$

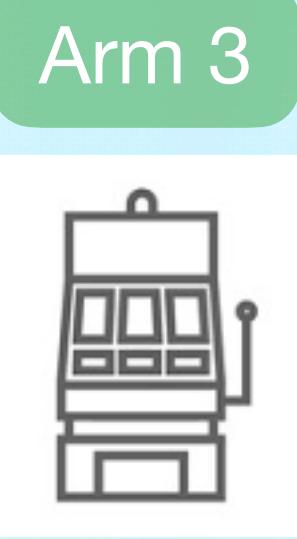
Arm 1



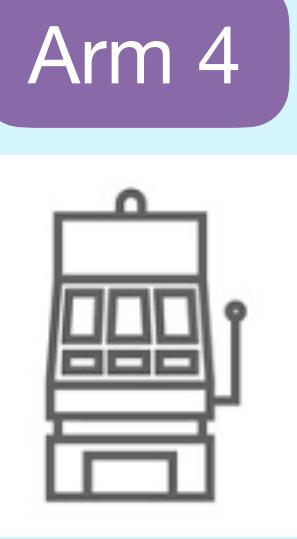
Arm 2



Arm 3



Arm 4



P_4

P_1

P_2

P_3

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

$$X_{6,2}$$

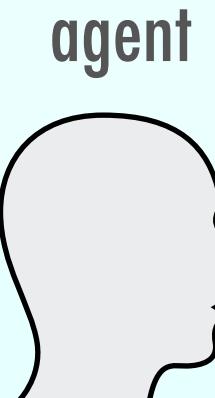
$$t = 7$$

$$X_{7,1}$$

$$t = 8$$

$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

$i_a(t)$: last observed state of arm a w.r.t. time t



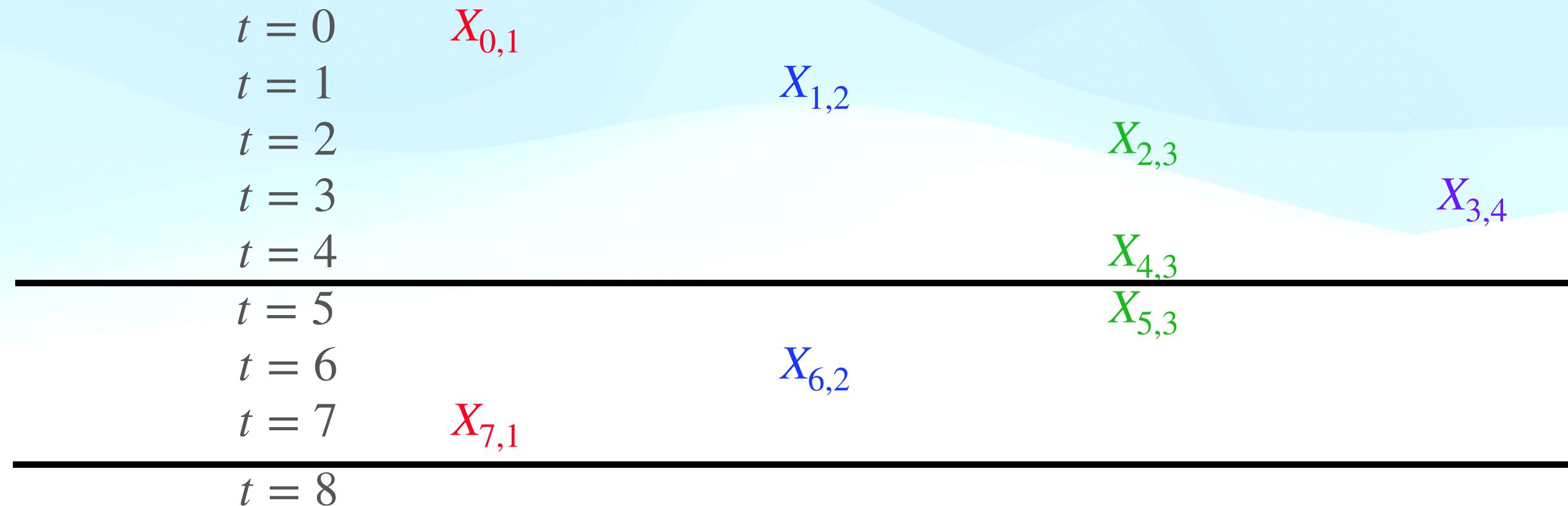
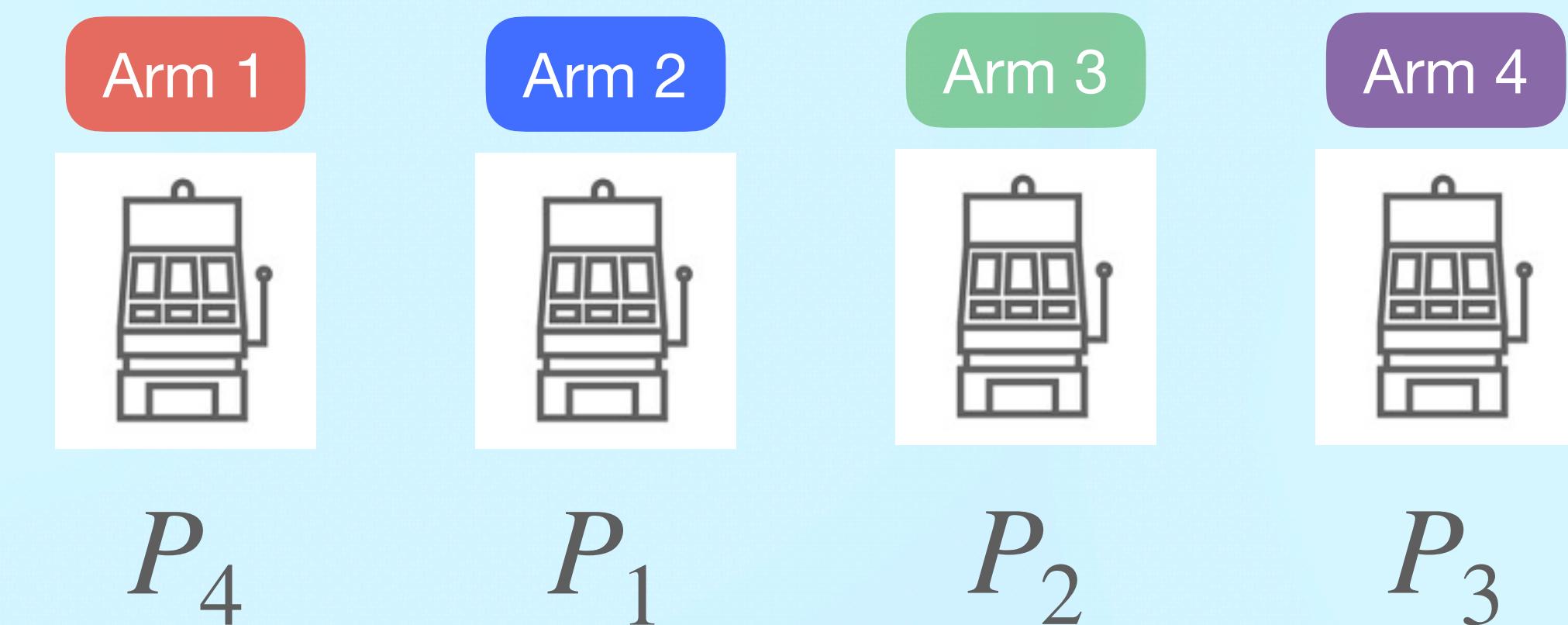
agent	t	0	1	2	3	4	5	6	7
a	1	2	3	4	3	3	2	1	
$X_{t,a}$	$X_{0,1}$	$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,3}$	$X_{5,3}$	$X_{6,2}$		$X_{7,1}$

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3} \quad i_4(t) = X_{3,4}$$

$$t = 5 \quad d_1(t) = 5 \quad d_2(t) = 4 \quad d_3(t) = 1 \quad d_4(t) = 2$$

$$i_1(t) = X_{0,1} \quad i_2(t) = X_{1,2} \quad i_3(t) = X_{4,3} \quad i_4(t) = X_{3,4}$$



$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

$i_a(t)$: last observed state of arm a w.r.t. time t

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

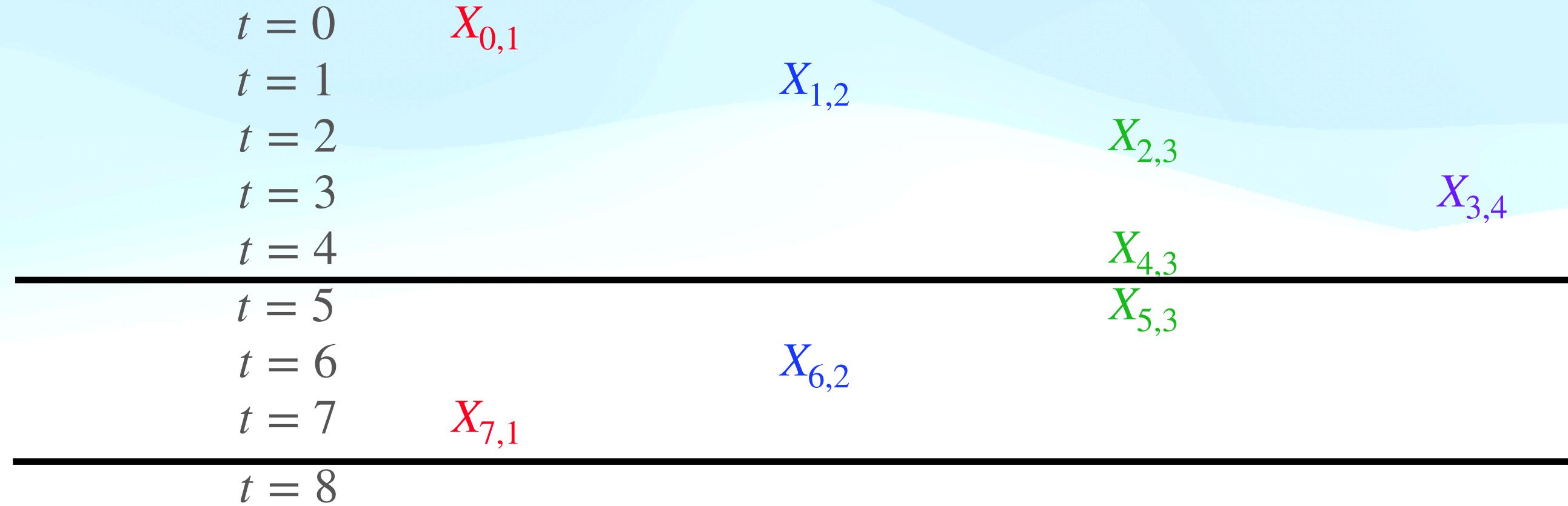
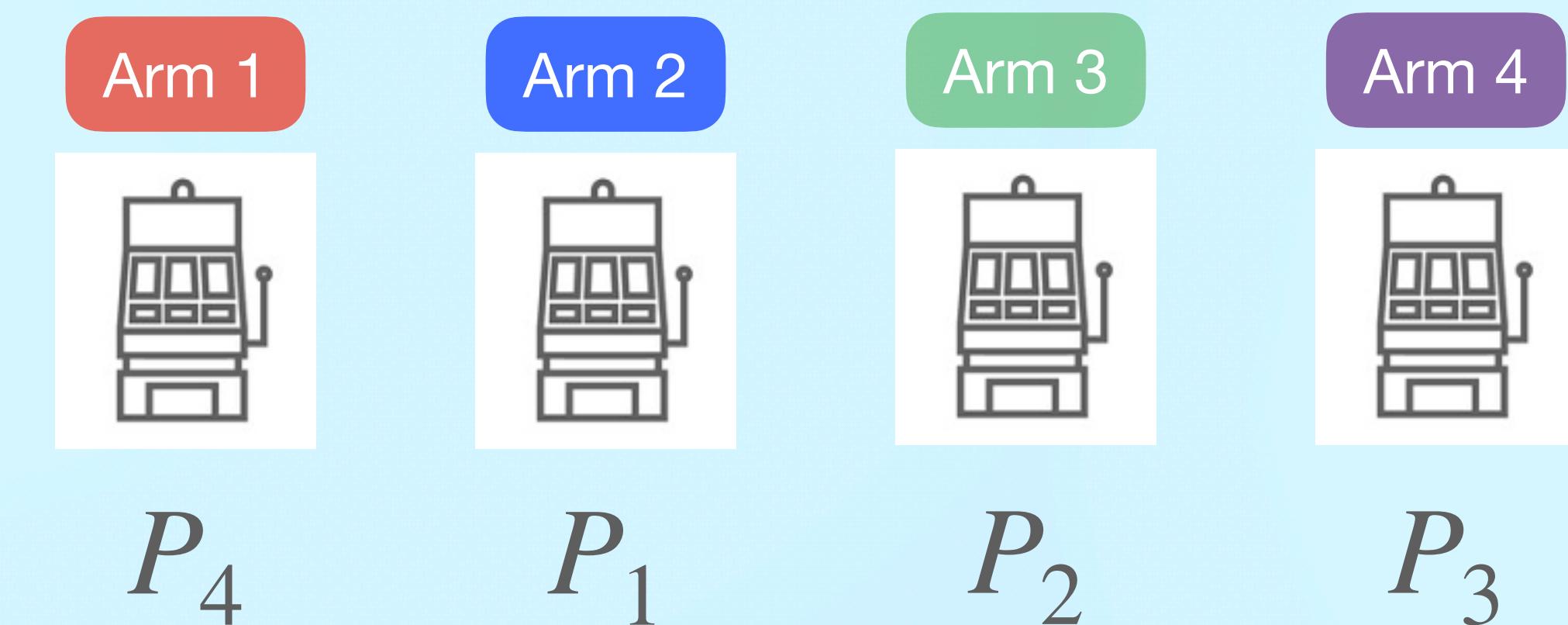
$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3} \quad i_4(t) = X_{3,4}$$

$$t = 5 \quad d_1(t) = 5 \quad d_2(t) = 4 \quad d_3(t) = 1 \quad d_4(t) = 2$$

$$i_1(t) = X_{0,1} \quad i_2(t) = X_{1,2} \quad i_3(t) = X_{4,3} \quad i_4(t) = X_{3,4}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$



$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

$i_a(t)$: last observed state of arm a w.r.t. time t

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

$$t = 8 \quad d_1(t) = 1 \quad d_2(t) = 2 \quad d_3(t) = 3 \quad d_4(t) = 5$$

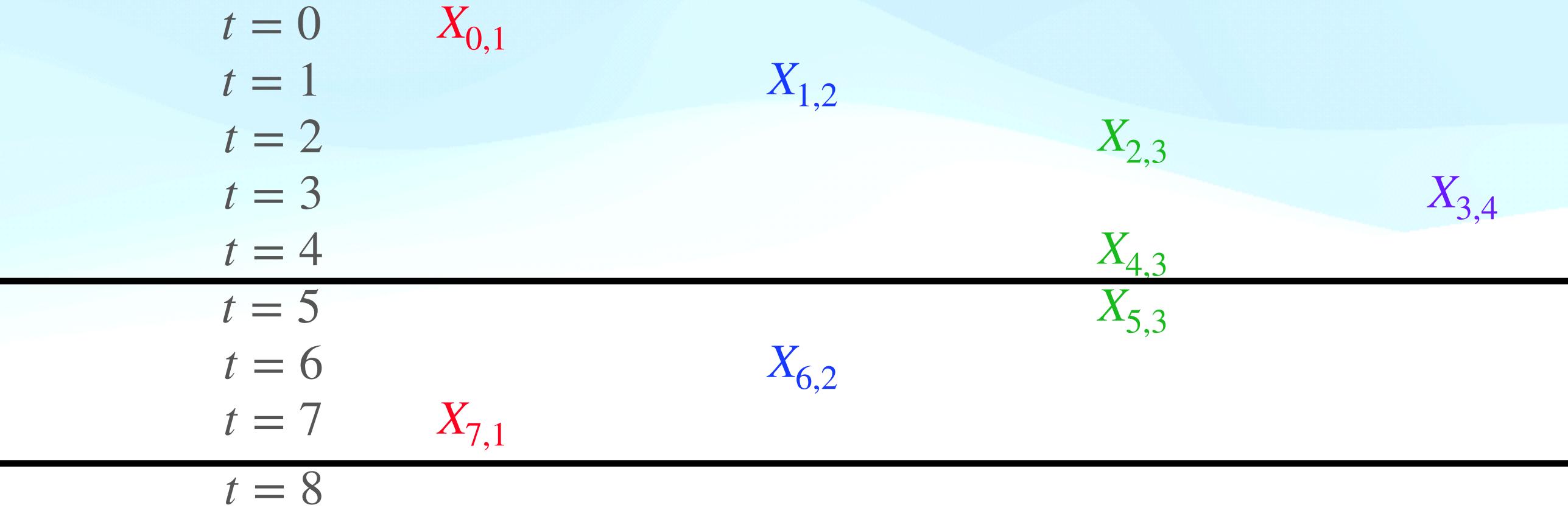
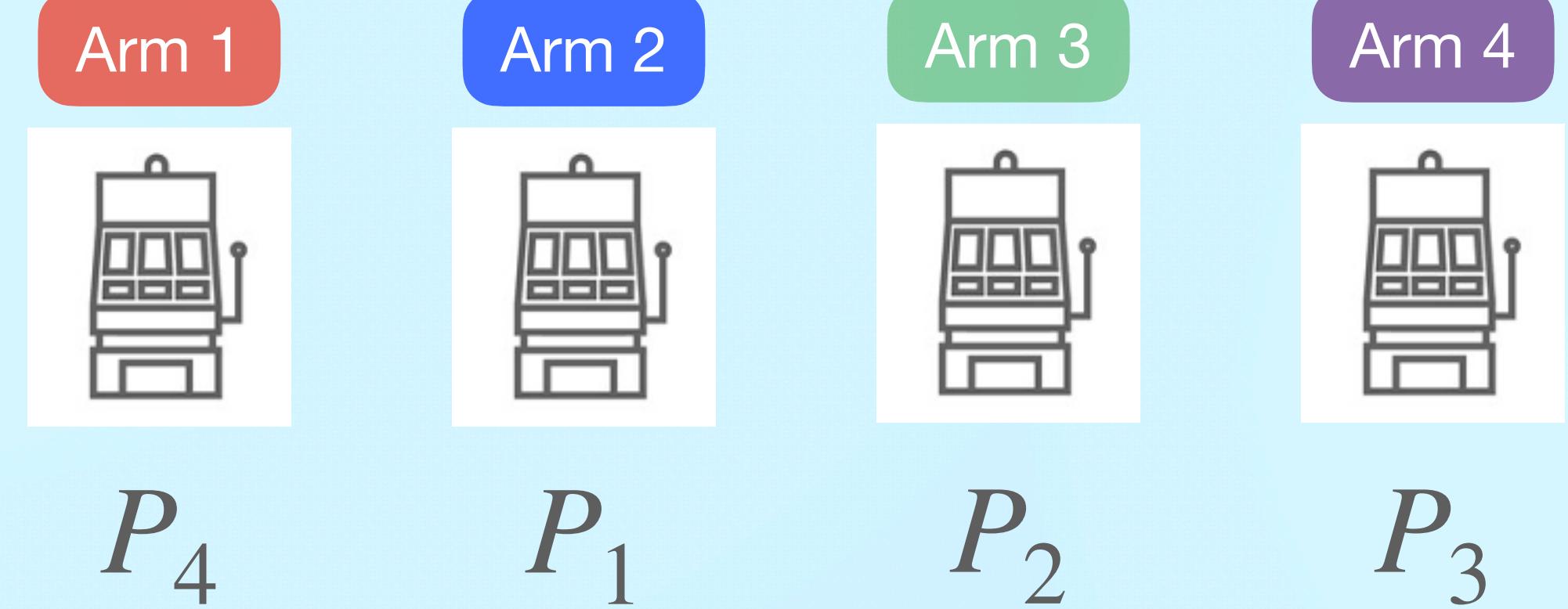
$$i_1(t) = X_{7,1} \quad i_2(t) = X_{6,2} \quad i_3(t) = X_{5,3} \quad i_4(t) = X_{3,4}$$

$$t = 5 \quad d_1(t) = 5 \quad d_2(t) = 4 \quad d_3(t) = 1 \quad d_4(t) = 2$$

$$i_1(t) = X_{0,1} \quad i_2(t) = X_{1,2} \quad i_3(t) = X_{4,3} \quad i_4(t) = X_{3,4}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$



$d_a(t)$: # time instants ago arm a was previously observed w.r.t. time t

$i_a(t)$: last observed state of arm a w.r.t. time t

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

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

AN MDP



$$P_4$$



$$P_1$$



$$P_2$$



$$P_3$$

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

$$X_{6,2}$$

$$t = 7$$

$$X_{7,1}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

AN MDP

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$



$$P_4$$



$$P_1$$



$$P_2$$



$$P_3$$

$$t = 0$$

$$X_{0,1}$$

$$t = 1$$

$$X_{1,2}$$

$$t = 2$$

$$X_{2,3}$$

$$t = 3$$

$$X_{3,4}$$

$$t = 4$$

$$X_{4,3}$$

$$t = 5$$

$$X_{5,3}$$

$$t = 6$$

$$X_{6,2}$$

$$t = 7$$

$$X_{7,1}$$

• $\underline{d}(t)$ is a vector of dimensions K .
• A_t is a function mapping $\underline{d}(t)$ to a set of actions.
• $\underline{i}(t)$ is a vector of dimensions K .

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

AN MDP



$$P_4$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$



$$P_1$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$



$$P_2$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$



$$P_3$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$

ARM PULLED

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

MDP

AN MDP



P_4



P_1



P_2



P_3

$$t = 0$$

$$t = 1$$

$$t = 2$$

$$t = 3$$

$$t = 4$$

$$t = 5$$

$$t = 6$$

$$t = 7$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

state space

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

MDP

AN MDP



$$P_4$$



$$P_1$$



$$P_2$$



$$P_3$$

$$t = 0$$

$$t = 1$$

$$t = 2$$

$$t = 3$$

$$t = 4$$

$$t = 5$$

$$t = 6$$

$$t = 7$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

state space

action space

MDP

AN MDP

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

$$\{1, \dots, K\}$$



$$P_4$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$



$$P_1$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

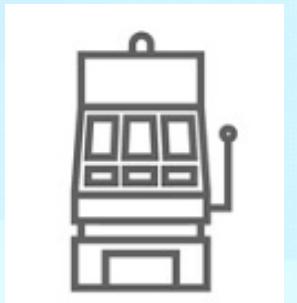
$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$



$$P_2$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$



$$P_3$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$



$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

state space

action space

MDP

state at time t

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

$$\{1, \dots, K\}$$

$$(\underline{d}(t), \underline{i}(t))$$



$$P_4$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$



$$P_1$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$



$$P_2$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{4,3}$$



$$P_3$$

$$X_{3,4}$$

$$\begin{aligned} t &= 0 \\ t &= 1 \\ t &= 2 \\ t &= 3 \\ t &= 4 \\ t &= 5 \\ t &= 6 \\ t &= 7 \end{aligned}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

state space

action space

MDP

state at time t

action at time t

AN MDP

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

$$\{1, \dots, K\}$$

$$(\underline{d}(t), \underline{i}(t))$$

$$A_t$$



$$P_4$$



$$P_1$$



$$P_2$$



$$P_3$$

$$\begin{aligned} t &= 0 \\ t &= 1 \\ t &= 2 \\ t &= 3 \\ t &= 4 \\ t &= 5 \\ t &= 6 \\ t &= 7 \end{aligned}$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

state space

action space

MDP

state at time t

action at time t

transition probabilities

AN MDP

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

$$\{1, \dots, K\}$$

$$(\underline{d}(t), \underline{i}(t))$$

$$A_t$$

$$\mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) = (P_{\sigma(a)})^{d_a}(i'_a | i_a)$$



$$P_4$$

$$X_{0,1}$$



$$P_1$$

$$X_{1,2}$$



$$P_2$$

$$X_{2,3}$$



$$P_3$$

$$X_{3,4}$$

$$\begin{aligned} t &= 0 \\ t &= 1 \\ t &= 2 \\ t &= 3 \\ t &= 4 \\ t &= 5 \\ t &= 6 \\ t &= 7 \end{aligned}$$

$$\begin{aligned} X_{0,1} \\ X_{1,2} \\ X_{2,3} \\ X_{3,4} \\ X_{4,3} \\ X_{5,3} \\ X_{6,2} \\ X_{7,1} \end{aligned}$$

$$\underline{d}(t) = (d_1(t), \dots, d_K(t))$$

$$\underline{i}(t) = (i_1(t), \dots, i_K(t))$$

$$(\underline{d}(t), \underline{i}(t)) \longrightarrow A_t \longrightarrow (\underline{d}(t+1), \underline{i}(t+1)) \longrightarrow A_{t+1} \longrightarrow \dots \dots$$

arm pulled at time t

state space

action space

MDP

state at time t

action at time t

transition probabilities

AN MDP

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

$$\{1, \dots, K\}$$

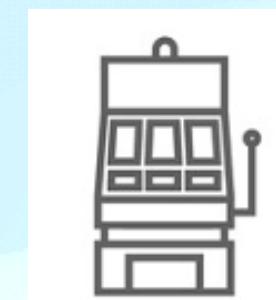
$$(\underline{d}(t), \underline{i}(t))$$

$$A_t$$

characterise or bound

$$\mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) = (P_{\sigma(a)})^{d_a}(i'_a | i_a)$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^{\pi}[\text{stopping time under } \pi]}{\log(1/\delta)}$$



$$P_4$$

$$X_{0,1}$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$

$$X_{6,2}$$

$$X_{7,1}$$



$$P_1$$

$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$



$$P_2$$

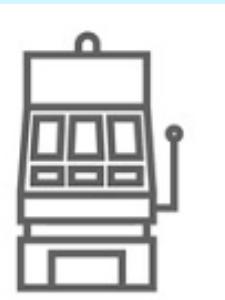
$$X_{1,2}$$

$$X_{2,3}$$

$$X_{3,4}$$

$$X_{4,3}$$

$$X_{5,3}$$



$$P_3$$

CONVERSE: LOWER BOUND

LOWER BOUND

characterise or bound

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)}$$

$$C = (P_{\sigma(1)}, \dots, P_{\sigma(K)})$$

LOWER BOUND

characterise or bound

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)}$$

$$C = (P_{\sigma(1)}, \dots, P_{\sigma(K)})$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^*(C)}$$

LOWER BOUND

characterise or bound

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)}$$

$$C = (P_{\sigma(1)}, \dots, P_{\sigma(K)})$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^\star(C)}$$

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

LOWER BOUND

characterise or bound

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)}$$

$$C = (P_{\sigma(1)}, \dots, P_{\sigma(K)})$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^\star(C)}$$

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

LOWER BOUND

characterise or bound

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)}$$

$$C = (P_{\sigma(1)}, \dots, P_{\sigma(K)})$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^*(C)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \text{ for all } (\underline{d}, \underline{i}, a)$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1$$

LOWER BOUND

characterise or bound

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)}$$

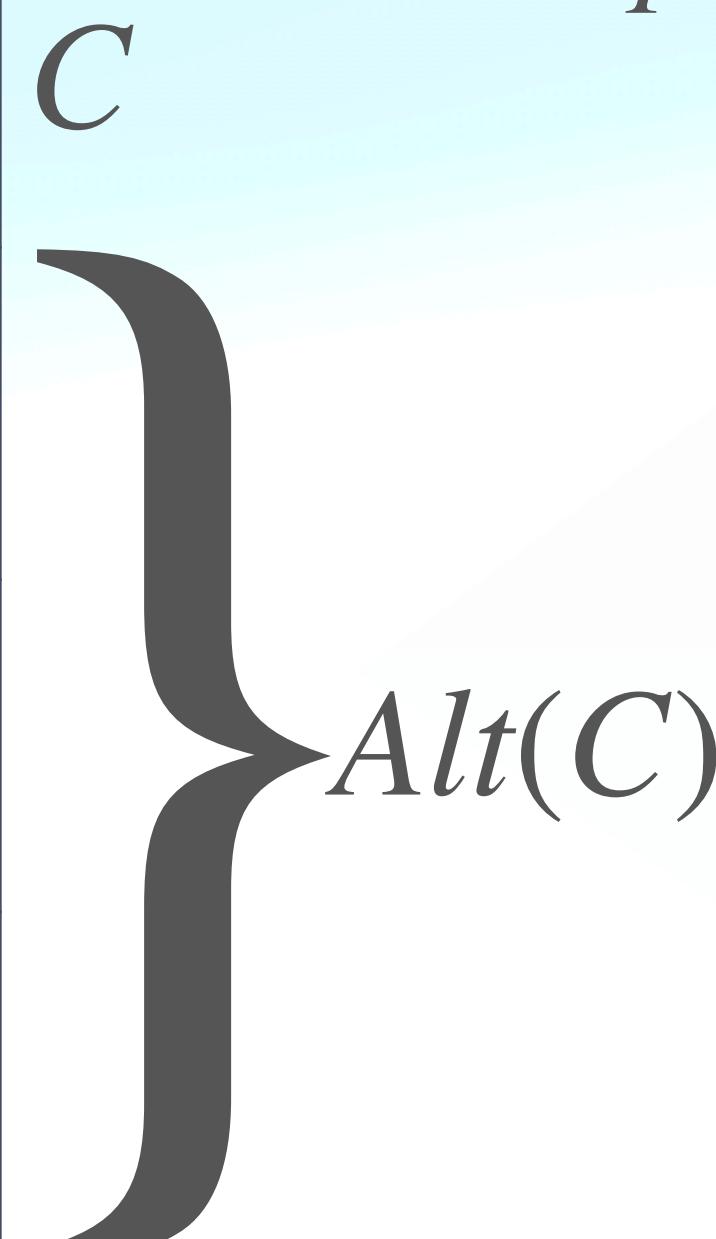
best arm = 1

1,2,3,4	1,2,4,3	1,3,2,4
1,3,4,2	1,4,2,3	1,4,3,2
2,1,3,4	2,1,4,3	3,1,2,4
3,1,4,2	4,1,2,3	4,1,3,2
2,3,1,4	2,4,1,3	3,2,1,4
3,4,1,2	4,2,1,3	4,3,1,2
2,3,4,1	2,4,3,1	3,2,4,1
3,4,2,1	4,2,3,1	4,3,2,1

best arm = 2

best arm = 3

best arm = 4



$$C = (P_{\sigma(1)}, \dots, P_{\sigma(K)})$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^*(C)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

$\nu(\underline{d}, \underline{i}, a) \geq 0$ for all $(\underline{d}, \underline{i}, a)$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1$$

SOME OBSERVATIONS - 1

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^\star(C)}$$

SOME OBSERVATIONS - 1

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^\star(C)}$$

$$\min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu_\epsilon(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a) \geq T^\star(C) - \epsilon$$

QUESTION

SOME OBSERVATIONS - 1

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^\star(C)}$$

$$\min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu_\epsilon(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a) \geq T^\star(C) - \epsilon$$

$$\frac{\# \text{ times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as} \quad n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

SOME OBSERVATIONS - 1

$$\frac{\text{\# times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as} \quad n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

SOME OBSERVATIONS - 1

$$\frac{\text{\# times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as} \quad n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

On the Empirical State-Action Frequencies in Markov Decision Processes Under General Policies

Shie Mannor

Department of Electrical and Computer Engineering, McGill University, 3480 University Street,
Montreal, Québec, Canada H3A 2A7, shie@ece.mcgill.ca, www.ece.mcgill.ca/~shie/

John N. Tsitsiklis

Laboratory for Information and Decision Systems, Massachusetts Institute of Technology,
Cambridge, Massachusetts 02139, jnt@mit.edu, web.mit.edu/~jnt/www/home.html

SOME OBSERVATIONS - 1

$$\frac{\text{\# times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as } n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

$$\left\{ \sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S} \right\}$$

On the Empirical State-Action Frequencies in Markov Decision Processes Under General Policies

Shie Mannor

Department of Electrical and Computer Engineering, McGill University, 3480 University Street,
Montreal, Québec, Canada H3A 2A7, shie@ece.mcgill.ca, www.ece.mcgill.ca/~shie/

John N. Tsitsiklis

Laboratory for Information and Decision Systems, Massachusetts Institute of Technology,
Cambridge, Massachusetts 02139, jnt@mit.edu, web.mit.edu/~jnt/www/home.html

SOME OBSERVATIONS - 1

$$\frac{\text{\# times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as } n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

$$\left\{ \sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S} \right\}$$

On the Empirical State-Action Frequencies in Markov Decision Processes Under General Policies

Shie Mannor

Department of Electrical and Computer Engineering, McGill University, 3480 University Street,
Montreal, Québec, Canada H3A 2A7, shie@ece.mcgill.ca, www.ece.mcgill.ca/~shie/

John N. Tsitsiklis

Laboratory for Information and Decision Systems, Massachusetts Institute of Technology,
Cambridge, Massachusetts 02139, jnt@mit.edu, web.mit.edu/~jnt/www/home.html

($\underline{d}', \underline{i}'$)

SOME OBSERVATIONS - 1

$$\frac{\text{\# times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as } n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

$$\left\{ \sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S} \right\}$$

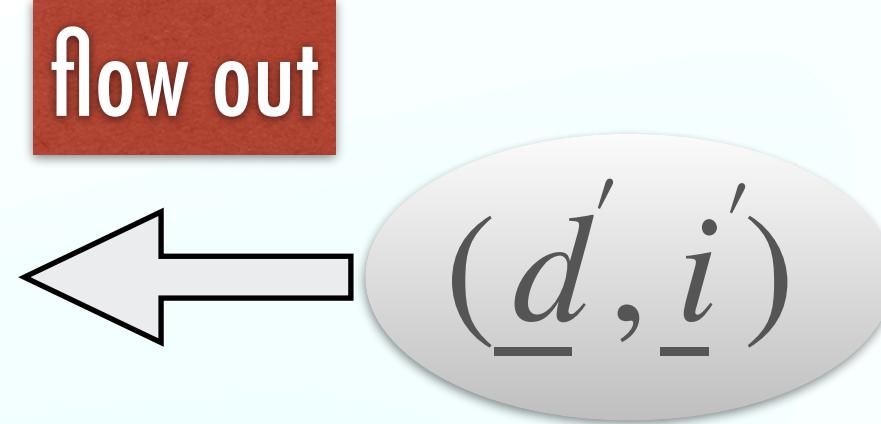
On the Empirical State-Action Frequencies in Markov Decision Processes Under General Policies

Shie Mannor

Department of Electrical and Computer Engineering, McGill University, 3480 University Street,
Montreal, Québec, Canada H3A 2A7, shie@ece.mcgill.ca, www.ece.mcgill.ca/~shie/

John N. Tsitsiklis

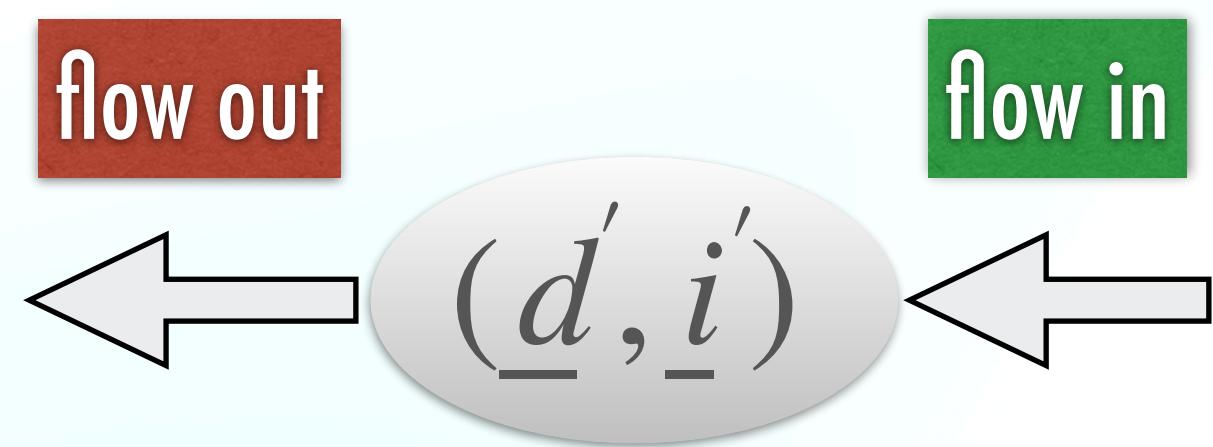
Laboratory for Information and Decision Systems, Massachusetts Institute of Technology,
Cambridge, Massachusetts 02139, jnt@mit.edu, web.mit.edu/~jnt/www/home.html



SOME OBSERVATIONS - 1

$$\frac{\text{\# times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as } n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

$$\left\{ \sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S} \right\}$$



On the Empirical State-Action Frequencies in Markov Decision Processes Under General Policies

Shie Mannor

Department of Electrical and Computer Engineering, McGill University, 3480 University Street,
Montreal, Québec, Canada H3A 2A7, shie@ece.mcgill.ca, www.ece.mcgill.ca/~shie/

John N. Tsitsiklis

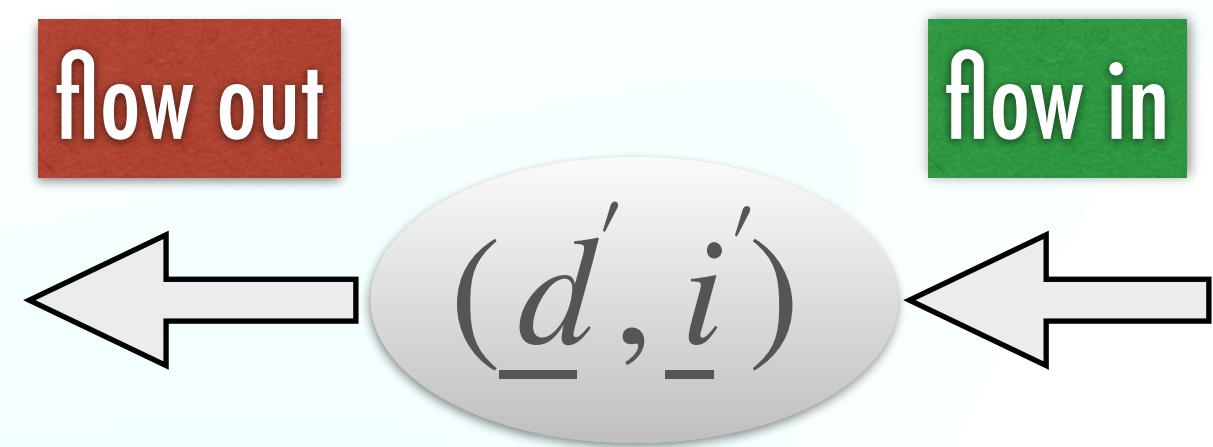
Laboratory for Information and Decision Systems, Massachusetts Institute of Technology,
Cambridge, Massachusetts 02139, jnt@mit.edu, web.mit.edu/~jnt/www/home.html

SOME OBSERVATIONS - 1

$$\frac{\text{\# times } (\underline{d}, \underline{i}, a) \text{ is observed up to time } n}{n} \longrightarrow \nu_\epsilon(\underline{d}, \underline{i}, a) \quad \text{as} \quad n \rightarrow \infty \quad \forall (\underline{d}, \underline{i}, a)$$

flow constraint

$$\left\{ \sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S} \right\}$$



On the Empirical State-Action Frequencies in Markov Decision Processes Under General Policies

Shie Mannor

Department of Electrical and Computer Engineering, McGill University, 3480 University Street, Montreal, Québec, Canada H3A 2A7, shie@ece.mcgill.ca, www.ece.mcgill.ca/~shie/

John N. Tsitsiklis

Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, jnt@mit.edu, web.mit.edu/~jnt/www/home.html

SOME OBSERVATIONS - 2

state space

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

action space

$$\{1, \dots, K\}$$

MDP

state at time t

$$(\underline{d}(t), \underline{i}(t))$$

action at time t

$$A_t$$

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

SOME OBSERVATIONS - 2

state space

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

countably infinite

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

action space

$$\{1, \dots, K\}$$

MDP

state at time t

$$(\underline{d}(t), \underline{i}(t))$$

action at time t

$$A_t$$

SOME OBSERVATIONS - 2

state space

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

countably infinite

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

action space

$$\{1, \dots, K\}$$

MDP

state at time t

$$(\underline{d}(t), \underline{i}(t))$$

action at time t

$$A_t$$

$$d_a(t) \leq R \quad \forall t, a$$

SOME OBSERVATIONS - 2

state space

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

countably infinite

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

action space

$$\{1, \dots, K\}$$

MDP

state at time t

$$(\underline{d}(t), \underline{i}(t))$$

action at time t

$$A_t$$

$$d_a(t) \leq R \quad \forall t, a$$

$$\left\{ \nu(\underline{d}, \underline{i}, a) = \sum_{a'=1}^K \nu(\underline{d}, \underline{i}, a') \quad \text{for all } (\underline{d}, \underline{i}) : d_a = R \right\}$$

SOME OBSERVATIONS - 2

state space

$$\mathbb{S} = \{(\underline{d}, \underline{i})\}$$

countably infinite

$$T^\star(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

action space

$$\{1, \dots, K\}$$

MDP

state at time t

$$(\underline{d}(t), \underline{i}(t))$$

$$A_t$$

R-max delay constraint

$$d_a(t) \leq R \quad \forall t, a$$

$$\left\{ \nu(\underline{d}, \underline{i}, a) = \sum_{a'=1}^K \nu(\underline{d}, \underline{i}, a') \quad \text{for all } (\underline{d}, \underline{i}) : d_a = R \right\}$$

MODIFIED OPTIMISATION

$$T_R^\star(P_{\sigma(1)}, \dots, P_{\sigma(K)}) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a),$$

subject to

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}_R,$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1,$$

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \quad \text{for all } (\underline{d}, \underline{i}, a) \in \mathbb{S}_R \times \mathcal{A},$$

$$\nu(\underline{d}, \underline{i}, a) = \sum_{a'=1}^K \nu(\underline{d}, \underline{i}, a') \quad \text{for all } (\underline{d}, \underline{i}) \in \mathbb{S}_{R,a}, \quad a \in \mathcal{A}.$$

MODIFIED OPTIMISATION

$$T_R^\star(P_{\sigma(1)}, \dots, P_{\sigma(K)}) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a),$$

subject to

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}_R,$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1,$$

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \quad \text{for all } (\underline{d}, \underline{i}, a) \in \mathbb{S}_R \times \mathcal{A},$$

$$\nu(\underline{d}, \underline{i}, a) = \sum_{a'=1}^K \nu(\underline{d}, \underline{i}, a') \quad \text{for all } (\underline{d}, \underline{i}) \in \mathbb{S}_{R,a}, \quad a \in \mathcal{A}.$$

MODIFIED OPTIMISATION

$$T_R^{\star}(P_{\sigma(1)}, \dots, P_{\sigma(K)}) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a),$$

subject to

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}_R,$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1,$$

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \quad \text{for all } (\underline{d}, \underline{i}, a) \in \mathbb{S}_R \times \mathcal{A},$$

$$\nu(\underline{d}, \underline{i}, a) = \sum_{a'=1}^K \nu(\underline{d}, \underline{i}, a') \quad \text{for all } (\underline{d}, \underline{i}) \in \mathbb{S}_{R,a}, \quad a \in \mathcal{A}.$$

MODIFIED OPTIMISATION

$$T_R^{\star}(P_{\sigma(1)}, \dots, P_{\sigma(K)}) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a),$$

subject to

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}_R,$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}_R} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1,$$

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \quad \text{for all } (\underline{d}, \underline{i}, a) \in \mathbb{S}_R \times \mathcal{A},$$

$$\nu(\underline{d}, \underline{i}, a) = \sum_{a'=1}^K \nu(\underline{d}, \underline{i}, a') \quad \text{for all } (\underline{d}, \underline{i}) \in \mathbb{S}_{R,a}, \quad a \in \mathcal{A}.$$

$$\nu_R^{\star} = \{\nu_R^{\star}(\underline{d}, \underline{i}, a)\}$$

ACHIEVABILITY

ALGORITHM FOR BAI - 1

ALGORITHM FOR BAI - 1

$$\lambda_{unif}(a \mid \underline{d}, \underline{i}) = \begin{cases} \frac{1}{K}, & d_a < \textcolor{violet}{R} \quad \forall a, \\ 1, & d_a = \textcolor{violet}{R}. \end{cases}$$

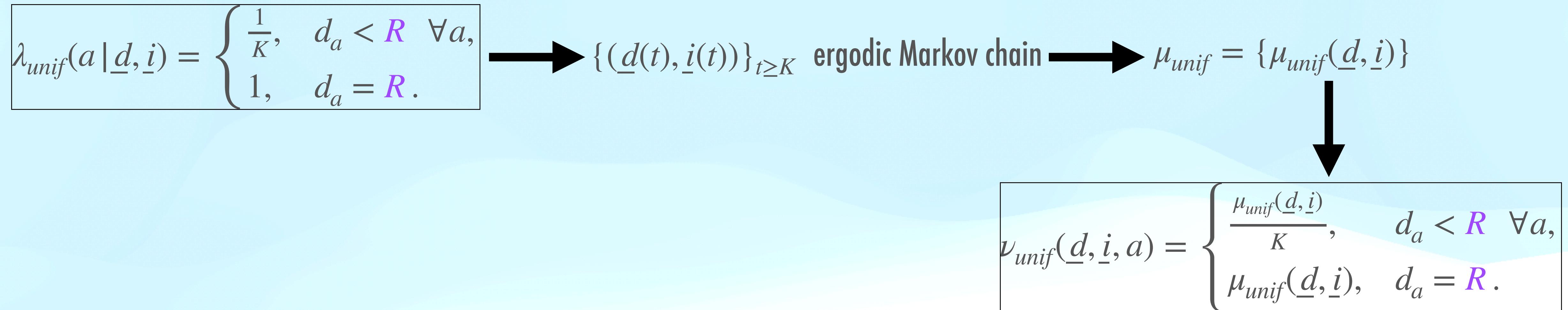
ALGORITHM FOR BAI - 1

$$\lambda_{unif}(a \mid \underline{d}, \underline{i}) = \begin{cases} \frac{1}{K}, & d_a < \textcolor{violet}{R} \quad \forall a, \\ 1, & d_a = \textcolor{violet}{R}. \end{cases} \rightarrow \{(\underline{d}(t), \underline{i}(t))\}_{t \geq K} \text{ ergodic Markov chain}$$

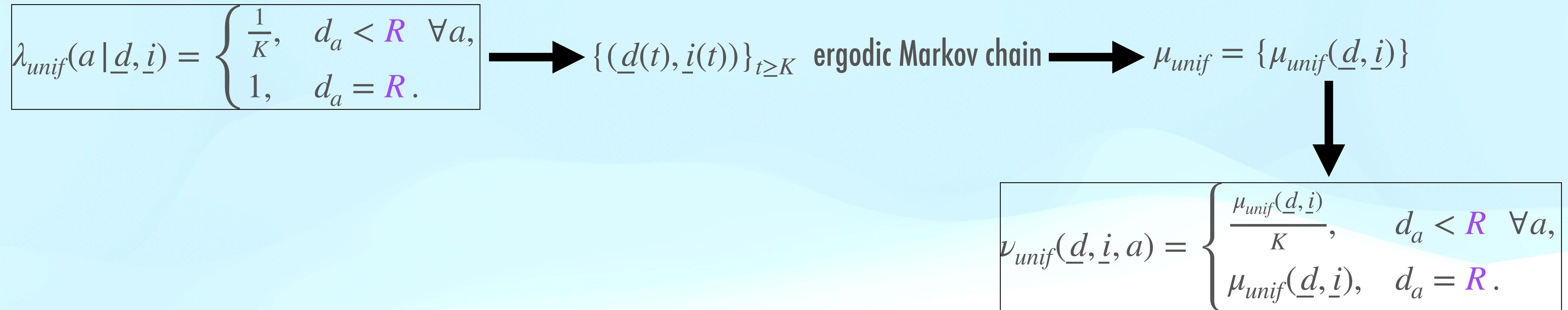
ALGORITHM FOR BAI - 1

$$\lambda_{unif}(a \mid \underline{d}, \underline{i}) = \begin{cases} \frac{1}{K}, & d_a < \textcolor{blue}{R} \quad \forall a, \\ 1, & d_a = \textcolor{blue}{R}. \end{cases} \rightarrow \{(\underline{d}(t), \underline{i}(t))\}_{t \geq K} \text{ ergodic Markov chain} \rightarrow \mu_{unif} = \{\mu_{unif}(\underline{d}, \underline{i})\}$$

ALGORITHM FOR BAI - 1

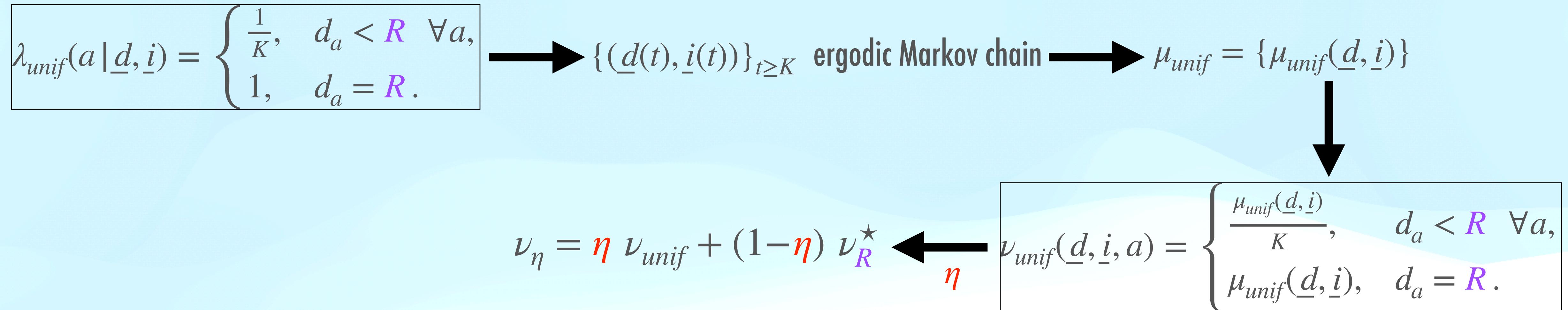


ALGORITHM FOR BAI - 1



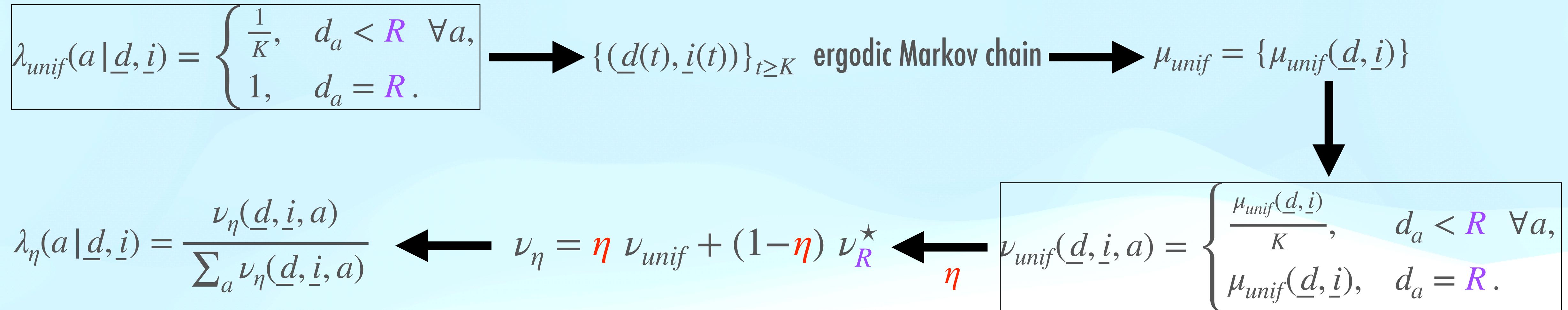
$$\nu_R^\star = \{\nu_R^\star(\underline{d}, \underline{i}, a)\}$$

ALGORITHM FOR BAI - 1



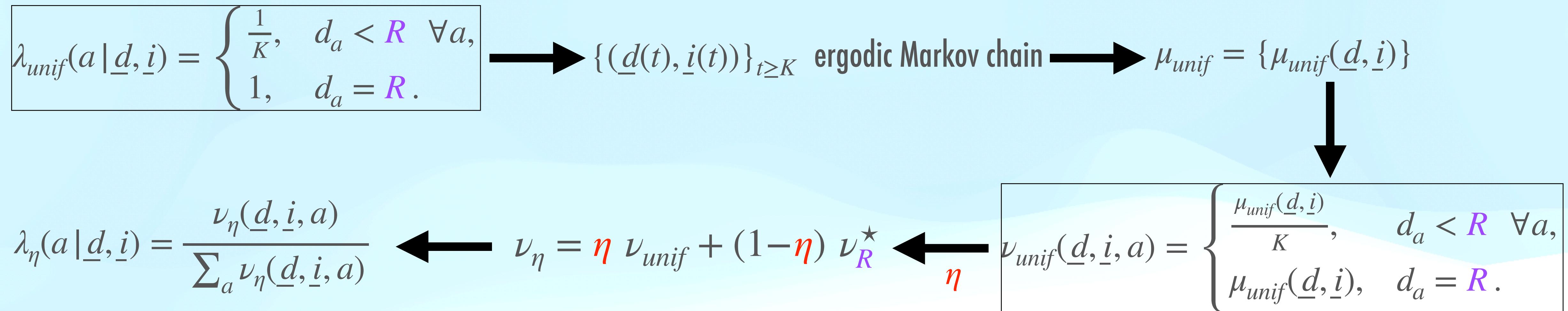
$$\nu_R^\star = \{\nu_R^\star(\underline{d}, \underline{i}, a)\}$$

ALGORITHM FOR BAI - 1



$$\nu_R^\star = \{\nu_R^\star(\underline{d}, \underline{i}, a)\}$$

ALGORITHM FOR BAI - 1



- $\lambda_\eta(a | \underline{d}, \underline{i}) > 0 \implies \lambda_\eta(a | \underline{d}, \underline{i}) \geq \alpha_\eta > 0$
- λ_η rule makes $\{\underline{d}(t), \underline{i}(t)\}$ an ergodic Markov chain

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a \mid \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ , η , R

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a | \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ , η , R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a | \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ , η , R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$

estimate of ground truth

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a | \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ, η, R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$
- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$

estimate of ground truth

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a | \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ, η, R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$
- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$

estimate of ground truth

$\hat{C}(n)$ is the ground truth

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a | \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ , η , R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$

estimate of ground truth

- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$

- STOP

$\hat{C}(n)$ is the ground truth

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a | \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ , η , R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$

estimate of ground truth

- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$

• STOP

$\hat{C}(n)$ is the ground truth

• Declare best arm in $\hat{C}(n)$

ALGORITHM FOR BAI - 2

$$\lambda_\eta(a | \underline{d}, \underline{i}) = \frac{\nu_\eta(\underline{d}, \underline{i}, a)}{\sum_a \nu_\eta(\underline{d}, \underline{i}, a)}$$

inputs

δ, η, R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$

estimate of ground truth

- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$

- STOP

$\hat{C}(n)$ is the ground truth

- Declare best arm in $\hat{C}(n)$

- Else, pull A_{n+1} according to $\lambda_\eta(\cdot | \underline{d}(n), \underline{i}(n))$ assuming $\hat{C}(n)$ to be the ground truth

PERFORMANCE

PERFORMANCE OF THE ALGORITHM (π^*)

inputs

δ, η, R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$ estimate of ground truth
- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$
 - STOP $\hat{C}(n)$ is the ground truth
 - Declare best arm in $\hat{C}(n)$
 - Else, pull A_{n+1} according to $\lambda_\eta(\cdot | \underline{d}(n), \underline{i}(n))$ assuming $\hat{C}(n)$ to be the ground truth

PERFORMANCE OF THE ALGORITHM (π^*)

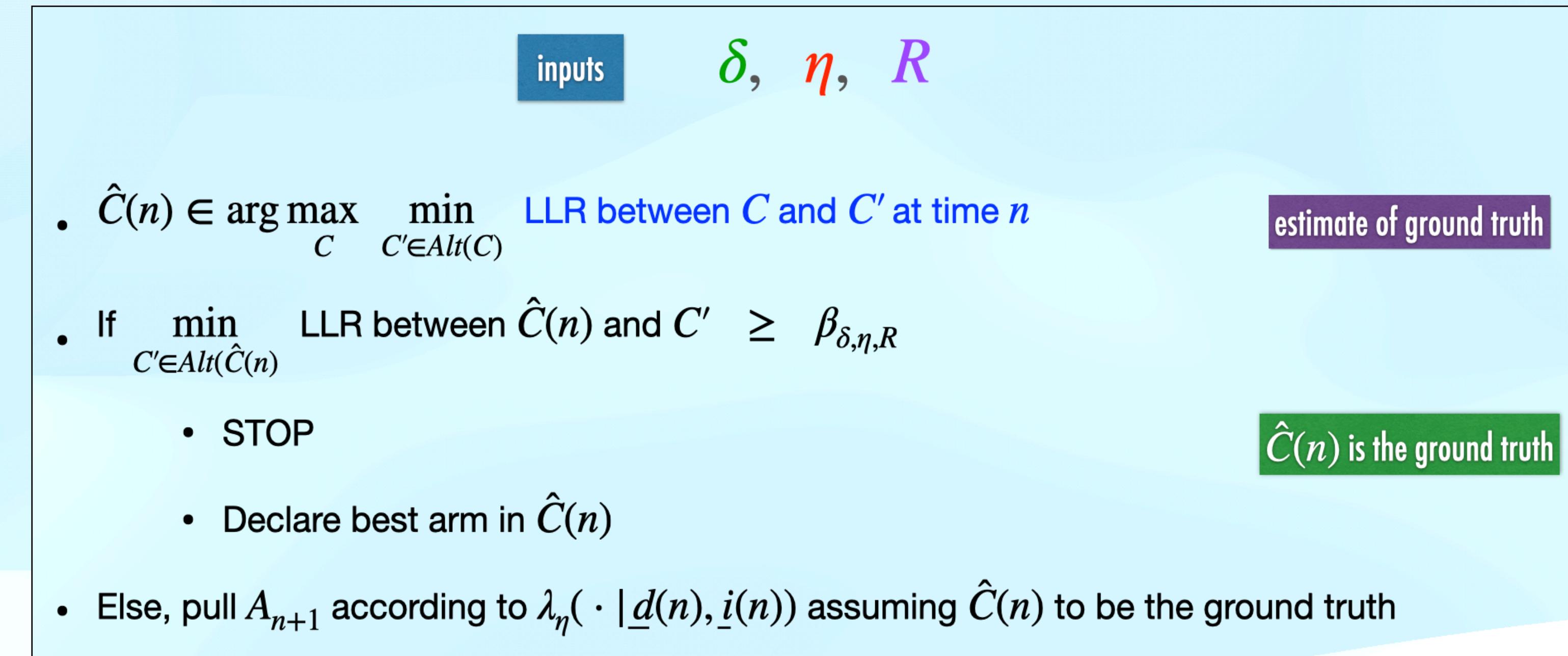
inputs

δ, η, R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$ estimate of ground truth
- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$
 - STOP
 - Declare best arm in $\hat{C}(n)$
 - Else, pull A_{n+1} according to $\lambda_\eta(\cdot | \underline{d}(n), \underline{i}(n))$ assuming $\hat{C}(n)$ to be the ground truth

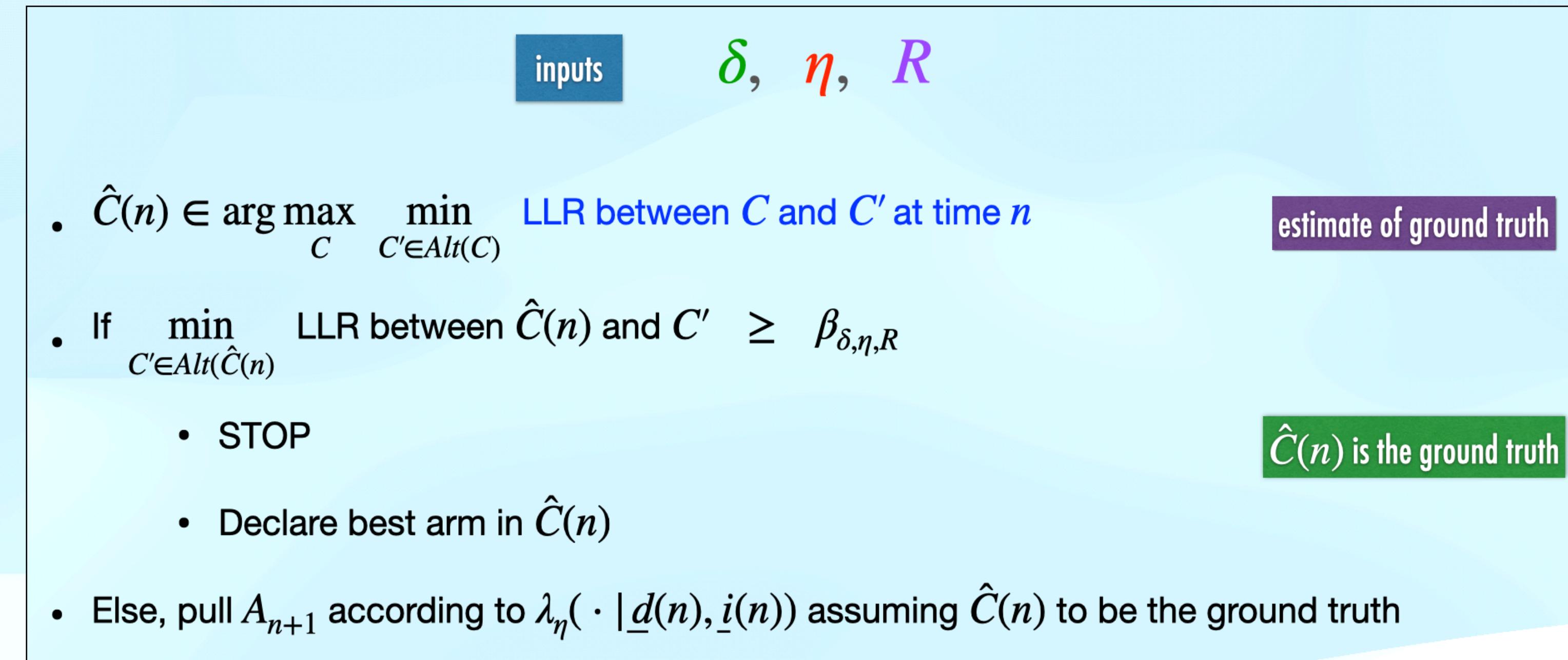
- Stops in finite time w.p.1

PERFORMANCE OF THE ALGORITHM (π^*)



- Stops in finite time w.p.1
- Achieves error probability $\leq \delta$

PERFORMANCE OF THE ALGORITHM (π^*)



- Stops in finite time w.p.1
- Achieves error probability $\leq \delta$

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\delta \downarrow 0} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\eta T_{\text{unif}}(P_1, \dots, P_K) + (1-\eta) T_R^*(P_1, \dots, P_K)}$$

PERFORMANCE OF THE ALGORITHM (π^*)

inputs

δ, η, R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$ estimate of ground truth
- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$
 - STOP
 - Declare best arm in $\hat{C}(n)$
- Else, pull A_{n+1} according to $\lambda_\eta(\cdot | \underline{d}(n), \underline{i}(n))$ assuming $\hat{C}(n)$ to be the ground truth

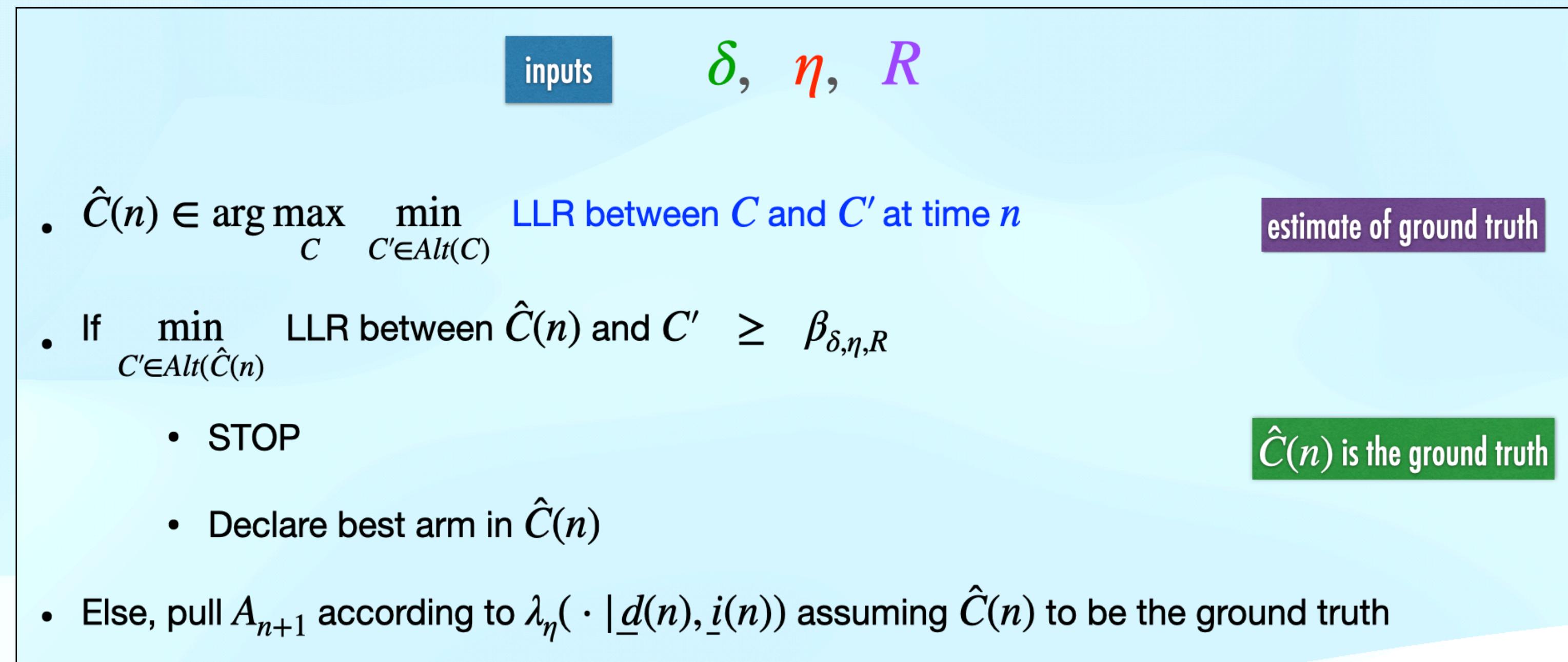
PERFORMANCE OF THE ALGORITHM (π^*)

inputs δ, η, R

- $\hat{C}(n) \in \arg \max_C \min_{C' \in Alt(C)} \text{LLR between } C \text{ and } C' \text{ at time } n$ estimate of ground truth
- If $\min_{C' \in Alt(\hat{C}(n))} \text{LLR between } \hat{C}(n) \text{ and } C' \geq \beta_{\delta, \eta, R}$
 - STOP
 - Declare best arm in $\hat{C}(n)$
- Else, pull A_{n+1} according to $\lambda_\eta(\cdot | \underline{d}(n), \underline{i}(n))$ assuming $\hat{C}(n)$ to be the ground truth

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

PERFORMANCE OF THE ALGORITHM (π^*)



$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{\mathbb{E}^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

RELATION BETWEEN T^* AND T_R^*

RELATION BETWEEN T^* AND T_R^*

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

RELATION BETWEEN T^* AND T_R^*

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

A key monotonicity property

RELATION BETWEEN T^* AND T_R^*

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

A key monotonicity property

$$T_R^*(P_1, \dots, P_K) \leq T_{R+1}^*(P_1, \dots, P_K) \leq T^*(P_1, \dots, P_K) \quad \forall R$$

RELATION BETWEEN T^* AND T_R^*

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

A key monotonicity property

$$T_R^*(P_1, \dots, P_K) \leq T_{R+1}^*(P_1, \dots, P_K) \leq T^*(P_1, \dots, P_K) \quad \forall R$$

$$\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K) \text{ exists}$$

RELATION BETWEEN T^* AND T_R^*

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

A key monotonicity property

$$T_R^*(P_1, \dots, P_K) \leq T_{R+1}^*(P_1, \dots, P_K) \leq T^*(P_1, \dots, P_K) \quad \forall R$$

$$\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K) \text{ exists}$$

IID obs.

RELATION BETWEEN T^* AND T_R^*

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

A key monotonicity property

$$T_R^*(P_1, \dots, P_K) \leq T_{R+1}^*(P_1, \dots, P_K) \leq T^*(P_1, \dots, P_K) \quad \forall R$$

$$\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K) \text{ exists}$$

IID obs.

$$\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K) = T^*(P_1, \dots, P_K)$$

RELATION BETWEEN T^* AND T_R^*

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

A key monotonicity property

$$T_R^*(P_1, \dots, P_K) \leq T_{R+1}^*(P_1, \dots, P_K) \leq T^*(P_1, \dots, P_K) \quad \forall R$$

$$\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K) \text{ exists}$$

IID obs.

$$\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K) = T^*(P_1, \dots, P_K)$$

IID obs.

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0 \\ R \rightarrow \infty}} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

matching bounds

OUR CONTRIBUTIONS

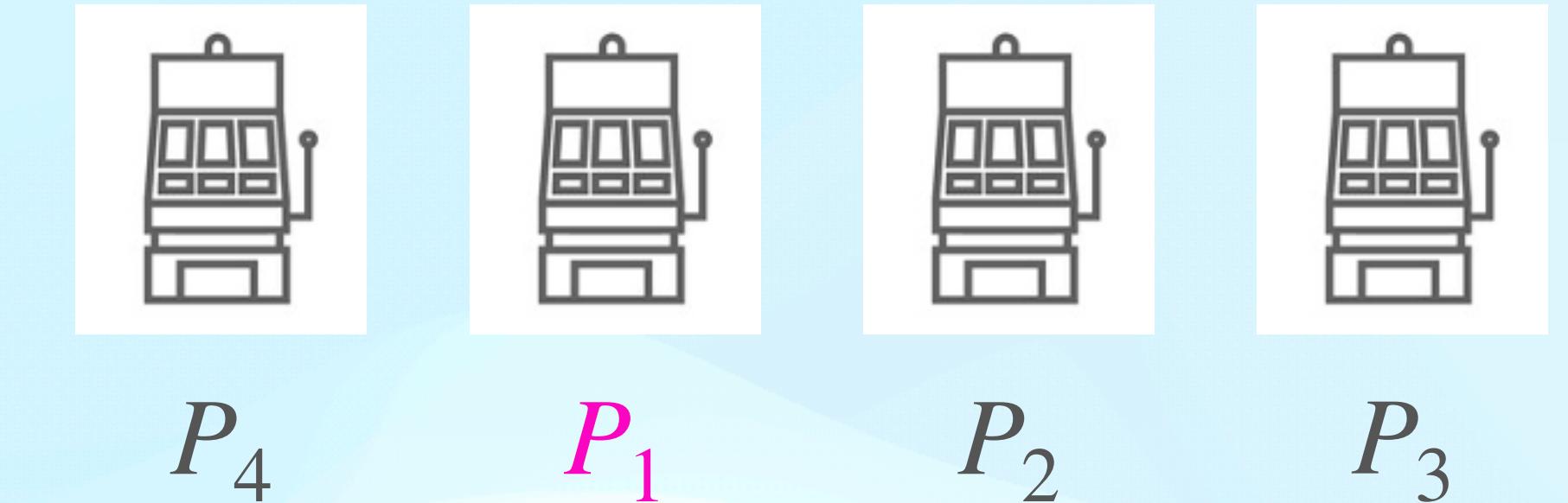
- Converse (lower bound)

$$\liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \geq \frac{1}{T^*(P_1, \dots, P_K)}$$

- Precise characterisation of $T^*(P_1, \dots, P_K)$
- Achievability: an algorithm for finding the best arm (π^*)
- Features of π^* :

- $\pi^* \in \Pi(\delta)$ for any $\delta > 0$
- Practically implementable

- Parameterised upper bound on the expected stopping time of π^*
- $T_R^*(P_1, \dots, P_K)$ is monotone increasing in the parameter R



R-upper bound

$$\limsup_{\substack{\delta \downarrow 0 \\ \eta \downarrow 0}} \frac{\mathbb{E}^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{T_R^*(P_1, \dots, P_K)}$$

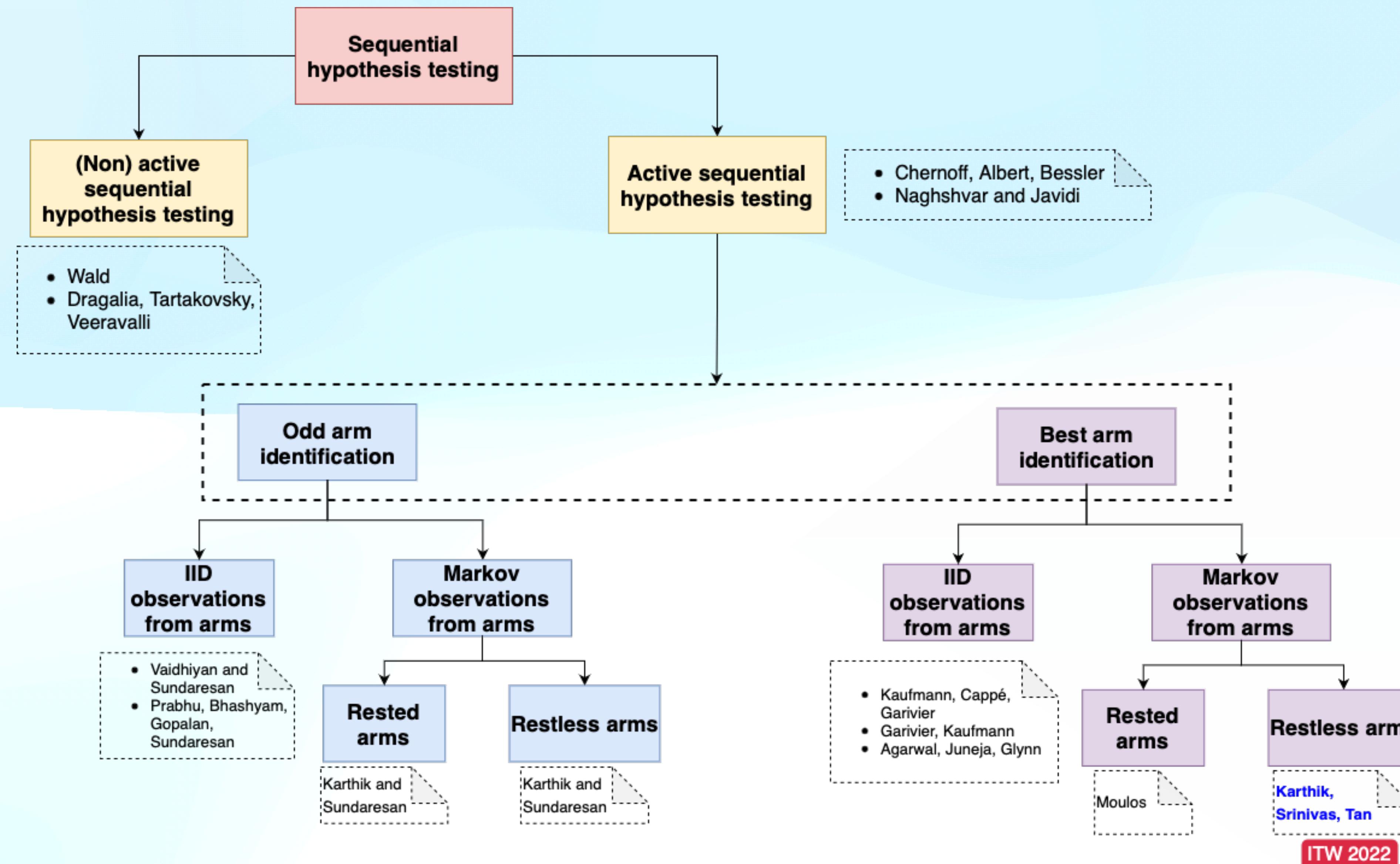
for i.i.d. arms

$$\limsup_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K) = T^*(P_1, \dots, P_K)$$

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \lim_{R \rightarrow \infty} \limsup_{\delta \downarrow 0} \frac{\mathbb{E}^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

WHERE WE STAND

RELATED WORK



FUTURE WORK

FUTURE WORK

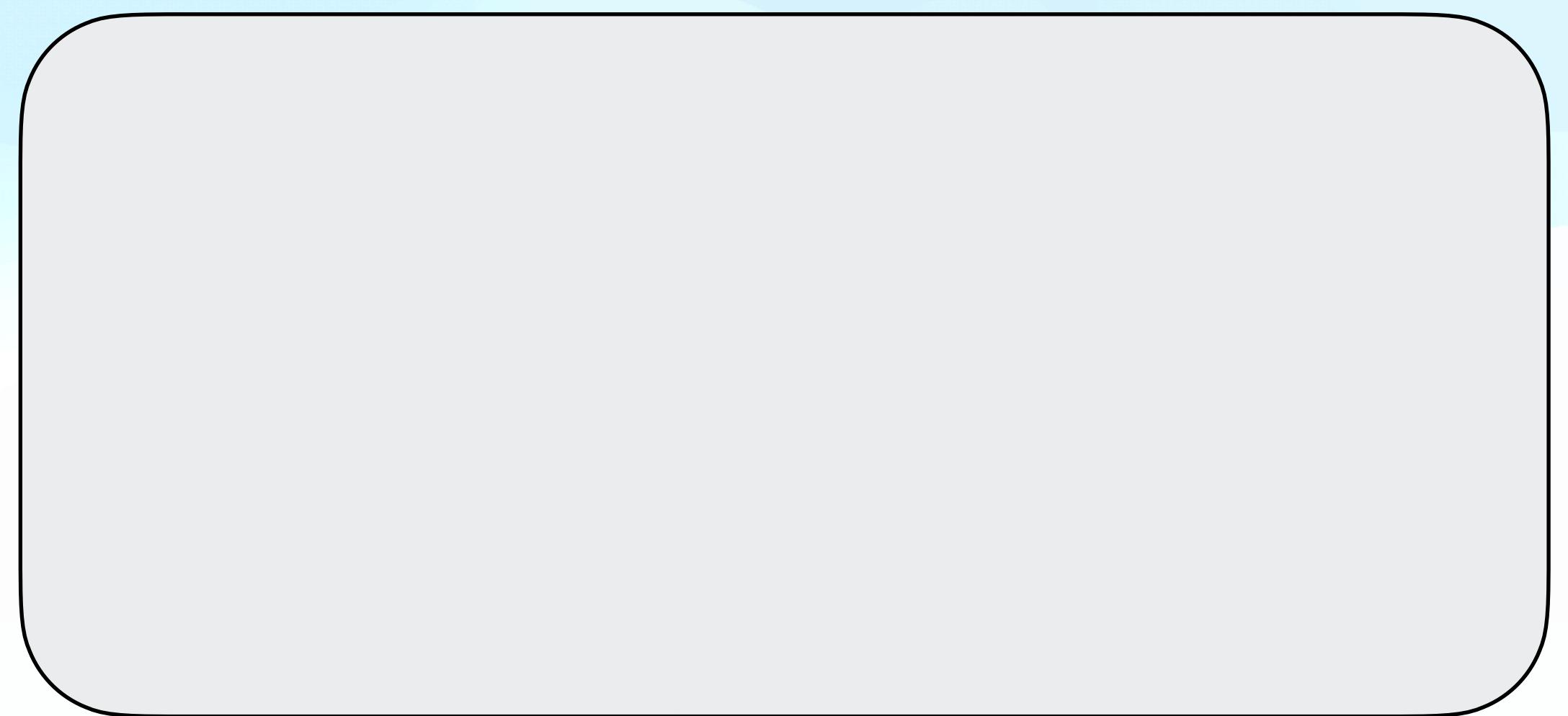
$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \lim_{R \rightarrow \infty} \limsup_{\delta \downarrow 0} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

FUTURE WORK

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \lim_{R \rightarrow \infty} \limsup_{\delta \downarrow 0} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$



$\nu(\underline{d}, \underline{i}, a) \geq 0$ for all $(\underline{d}, \underline{i}, a)$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1$$

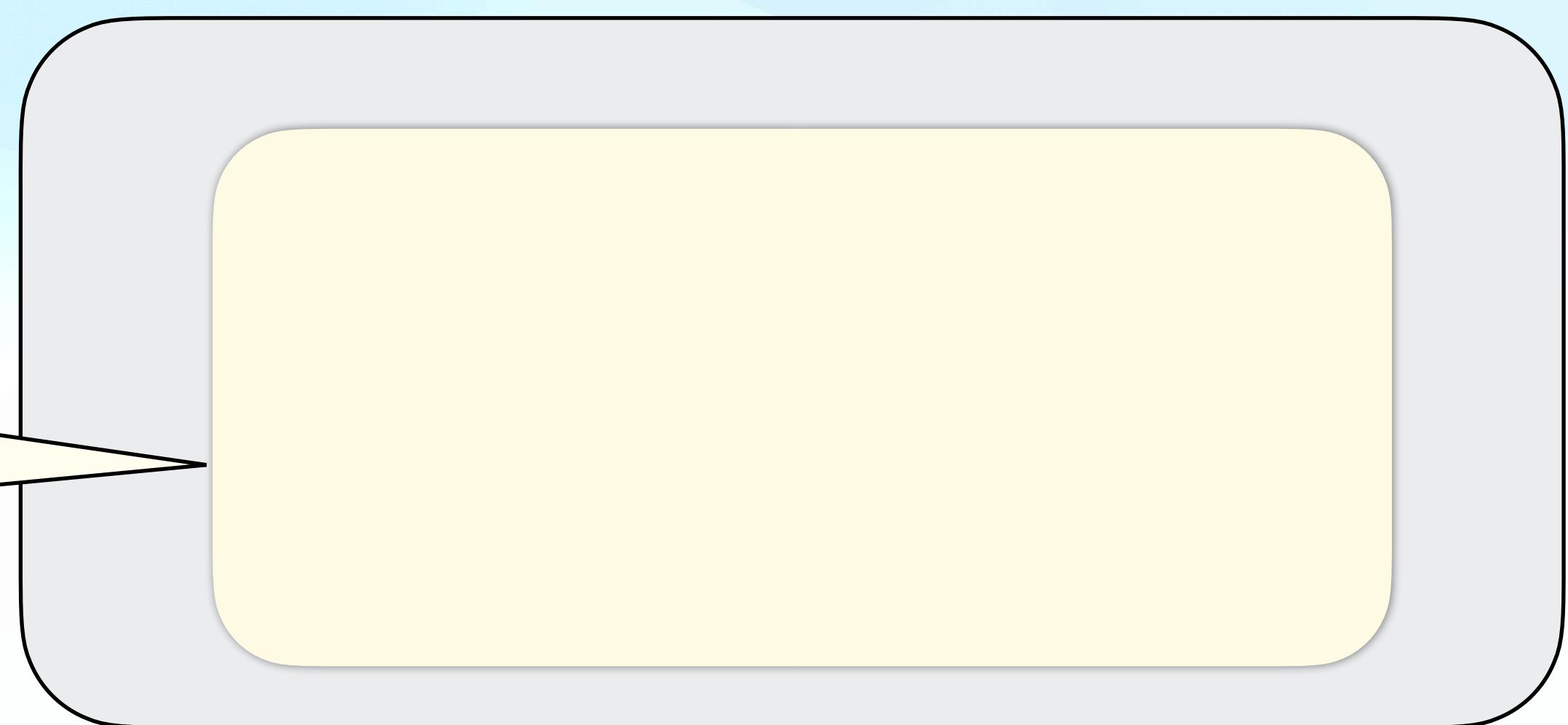
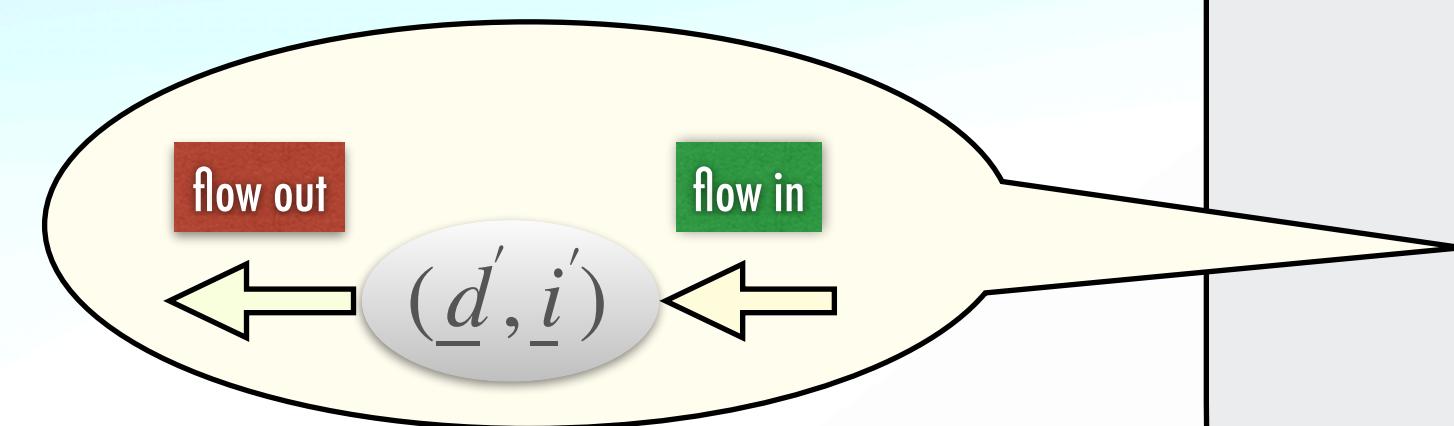
FUTURE WORK

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \lim_{R \rightarrow \infty} \limsup_{\delta \downarrow 0} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

flow constraint

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}$$



$\nu(\underline{d}, \underline{i}, a) \geq 0$ for all $(\underline{d}, \underline{i}, a)$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1$$

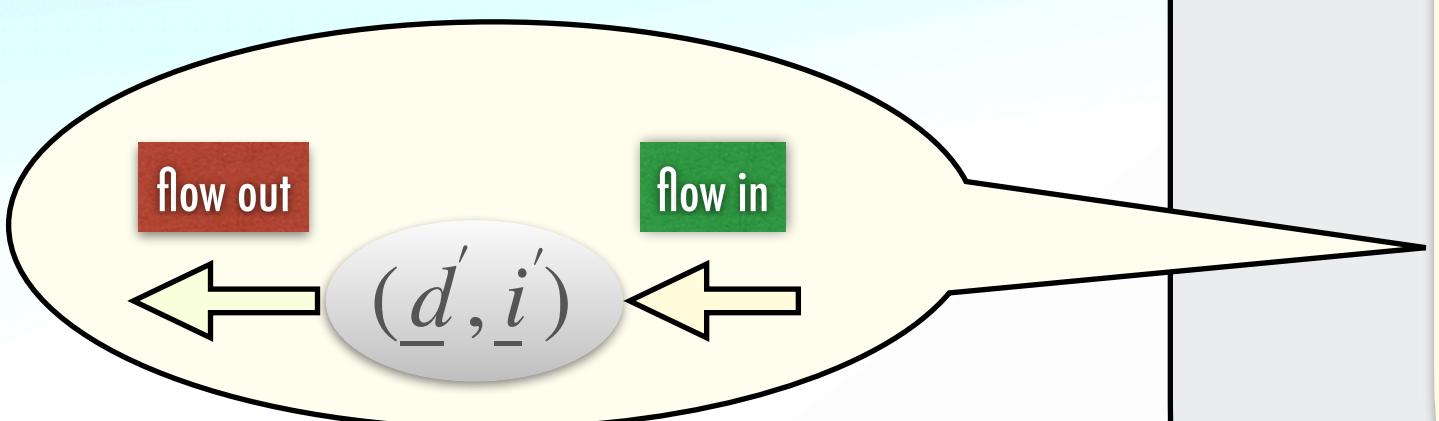
FUTURE WORK

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \lim_{R \rightarrow \infty} \limsup_{\delta \downarrow 0} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

flow constraint

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}$$



$$\nu(\underline{d}, \underline{i}, a) = \mu(\underline{d}, \underline{i}) \cdot \lambda(a | \underline{d}, \underline{i})$$

Ergodic state-action occupancy measures

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \text{ for all } (\underline{d}, \underline{i}, a)$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1$$

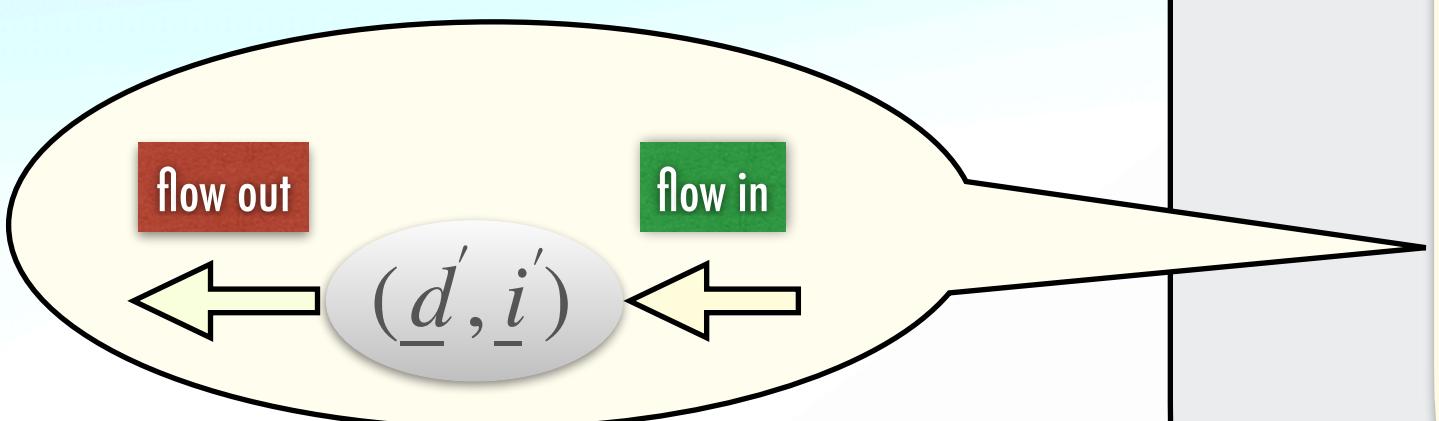
FUTURE WORK

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \lim_{R \rightarrow \infty} \limsup_{\delta \downarrow 0} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

flow constraint

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}$$



??

$$\nu(\underline{d}, \underline{i}, a) = \mu(\underline{d}, \underline{i}) \cdot \lambda(a | \underline{d}, \underline{i})$$

Ergodic state-action occupancy measures

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \text{ for all } (\underline{d}, \underline{i}, a)$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1$$

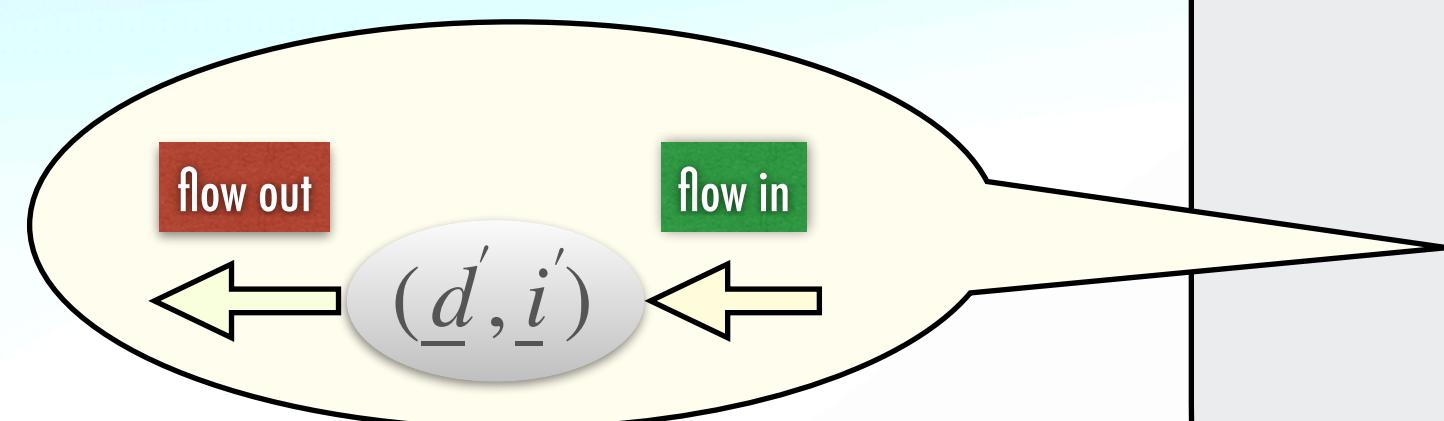
FUTURE WORK

$$\frac{1}{T^*(P_1, \dots, P_K)} \leq \liminf_{\delta \downarrow 0} \inf_{\pi \in \Pi(\delta)} \frac{\mathbb{E}^\pi[\text{stopping time under } \pi]}{\log(1/\delta)} \leq \lim_{R \rightarrow \infty} \limsup_{\delta \downarrow 0} \frac{E^{\pi^*}[\text{stopping time under } \pi^*]}{\log(1/\delta)} \leq \frac{1}{\lim_{R \rightarrow \infty} T_R^*(P_1, \dots, P_K)}$$

$$T^*(C) = \sup_{\nu} \min_{C' \in Alt(C)} \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) k_{CC'}(\underline{d}, \underline{i}, a)$$

flow constraint

$$\sum_{a=1}^K \nu(\underline{d}', \underline{i}', a) = \sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) \mathbb{P}((\underline{d}, \underline{i}, a) \rightarrow (\underline{d}', \underline{i}')) \quad \text{for all } (\underline{d}', \underline{i}') \in \mathbb{S}$$



??

$$\nu(\underline{d}, \underline{i}, a) = \mu(\underline{d}, \underline{i}) \cdot \lambda(a | \underline{d}, \underline{i})$$

Ergodic state-action occupancy measures

- Unknown TPMs
- Hidden Markov observations
- Second-order asymptotics

$$\nu(\underline{d}, \underline{i}, a) \geq 0 \text{ for all } (\underline{d}, \underline{i}, a)$$

$$\sum_{(\underline{d}, \underline{i}) \in \mathbb{S}} \sum_{a=1}^K \nu(\underline{d}, \underline{i}, a) = 1$$

ACKNOWLEDGEMENTS

- Srinivas and Vincent for their enthusiasm
- National Research Foundation, Singapore, and DSO National Laboratories for the generous sponsorship (award no: AISG2-RP-2020-018)



Srinivas Reddy Kota
ksreddy@nus.edu.sg

<https://sites.google.com/view/srinivas-reddy-kota/home>



Vincent Y. F. Tan
vtan@nus.edu.sg
<https://vyftan.github.io/>

ACKNOWLEDGEMENTS

- Srinivas and Vincent for their enthusiasm
- National Research Foundation, Singapore, and DSO National Laboratories for the generous sponsorship (award no: AISG2-RP-2020-018)



Srinivas Reddy Kota
ksreddy@nus.edu.sg

<https://sites.google.com/view/srinivas-reddy-kota/home>



Vincent Y. F. Tan
vtan@nus.edu.sg
<https://vyftan.github.io/>

THANK YOU!

ACKNOWLEDGEMENTS

- Srinivas and Vincent for their enthusiasm
- National Research Foundation, Singapore, and DSO National Laboratories for the generous sponsorship (award no: AISG2-RP-2020-018)



Srinivas Reddy Kota
ksreddy@nus.edu.sg

<https://sites.google.com/view/srinivas-reddy-kota/home>



Vincent Y. F. Tan
vtan@nus.edu.sg
<https://vyftan.github.io/>

THANK YOU!

Questions? Let's talk!

Contact: karthik@nus.edu.sg
<https://karthikpn.com>