# Dynamic resource allocation problems in communication networks:

## Introduction and the Finite Horizon Restless Bandit problem

Alexandre Reiffers-Masson

Equipe Maths&Net, IMT Atlantique, CS department
LabSTICC (UMR CNRS 6285)

June 25, 2024

# Acknowledgement

This course has been also elaborate during the project Ramonaas[1] (Regional Program STIC-AmSud):

- a STIC/AMSUD project between CAPES/BR (88881.694462/2022-01);
- Ministry for Europe and Foreign Affairs/FR;
- Campus France/FR and the National Agency for Research;
- Innovation/UY (MOV_CO_2022_1_1011515)

IMT Atlantique : -Engineering institution under the tutelage of the French Ministry of Industry
- 3 campuses North-West of France
- Member of the Institut Mines-Télécom

**Track MLA:**
*where?* on Brest campus

*when?*
**from 30 march to 2nd june 2023**

**IMT Atlantique**
Bretagne-Pays de la Loire
École Mines-Télécom

Brest
Rennes
Nantes

# What do in Brest?

# Agenda of the course (Room B03-036)

- **Day I:** Resource allocation problem and Restless Bandit.
    1. Course, 9h - 12h, teacher: Alexandre Reiffers-Masson
    2. Lab, 13h30 - 16h30, teacher: Lucas Lopes, Alexandre Reiffers-Masson

# Agenda of the course (Room B03-036)

- **Day I:** Resource allocation problem and Restless Bandit.
  1. Course, 9h - 12h, teacher: Alexandre Reiffers-Masson
  2. Lab, 13h30 - 16h30, teacher: Lucas Lopes, Alexandre Reiffers-Masson
- **Day II**
  1. Course on Resource allocation in optical networks, 9h - 12h, teacher: Luiz Anet Neto
  2. Course + Lab on aggregation techniques in MDP, 13h30 - 16h30, teacher: Olivier Tsemogne

# Agenda of the course (Room B03-036)

- **Day I:** Resource allocation problem and Restless Bandit.
  1. Course, 9h - 12h, teacher: Alexandre Reiffers-Masson
  2. Lab, 13h30 - 16h30, teacher: Lucas Lopes, Alexandre Reiffers-Masson
- **Day II**
  1. Course on Resource allocation in optical networks, 9h - 12h, teacher: Luiz Anet Neto
  2. Course + Lab on aggregation techniques in MDP, 13h30 - 16h30, teacher: Olivier Tsemogne
- **Day III:** Deep-learning and deep reinforcement learning applied to Resource Allocation Problems.
  1. Course, 9h - 12h, teacher: Alexandre Reiffers-Masson
  2. Lab, 13h30 - 16h30, teacher: Claudina Rattaro, Lucas Inglés, Alexandre Reiffers-Masson

# Agenda of Day I

**Provably efficient heuristics for solving large-scale resource allocation problems.**

- Introduction to resource allocation problem, Markov Decision Process, Restless Bandit in Finite Horizon and Infinite Horizon.

- Weakly coupled MDP and the resolving heuristic.

- Constrained Finite Horizon Stochastic Optimization Problems.

# Objectives of the course

Provably efficient heuristics for solving large-scale resource allocation problems

1. Design heuristics and prove asymptotically optimal properties.
2. Code the heuristic in Python using *cvxpy*.

# Machine Maintenance[2]

- **Scenario:** A collection of $N$ machines which deteriorate under usage is maintained by a set of $\alpha$ repairmen. Maintenance interventions will improve a machine's condition and may preempt costly breakdowns.

[2]Glazebrook, K. D., Mitchell, H. M., & Ansell, P. S. (2005). Index policies for the maintenance of a collection of machines by a set of repairmen. European Journal of Operational Research, 165(1), 267-284.

# Machine Maintenance[2]

- **Scenario:** A collection of $N$ machines which deteriorate under usage is maintained by a set of $\alpha$ repairmen. Maintenance interventions will improve a machine's condition and may preempt costly breakdowns.

- **Challenge:** How to *allocate* repairmen at each instant $t$, *knowing* the state of each machine, *to minimize the total expected deterioration of machines*.

[2]Glazebrook, K. D., Mitchell, H. M., & Ansell, P. S. (2005). Index policies for the maintenance of a collection of machines by a set of repairmen. European Journal of Operational Research, 165(1), 267-284.

# Machine Maintenance[2]

- **Scenario:** A collection of $N$ machines which deteriorate under usage is maintained by a set of $\alpha$ repairmen. Maintenance interventions will improve a machine's condition and may preempt costly breakdowns.

- **Challenge:** How to *allocate* repairmen at each instant $t$, *knowing* the state of each machine, *to minimize the total expected deterioration of machines*.

---

[2]Glazebrook, K. D., Mitchell, H. M., & Ansell, P. S. (2005). Index policies for the maintenance of a collection of machines by a set of repairmen. European Journal of Operational Research, 165(1), 267-284.

# Model: Evolution of the state of a machine

- Let $S_k(t) \in \{1, \dots, S\}$ be the state of the machine of the machine $k$.

# Model: Evolution of the state of a machine

- Let $S_k(t) \in \{1, \ldots, S\}$ be the state of the machine of the machine $k$.
- The greater the value of $S_k(t)$, the worse the machine $k$ performs.

# Model: Evolution of the state of a machine

- Let $S_k(t) \in \{1, \ldots, S\}$ be the state of the machine of the machine $k$.
- The greater the value of $S_k(t)$, the worse the machine $k$ performs.
- **Remark:** Usually the evolution of the state of machine $k$ is modelled by a *controlled Markov chain*.

# Model: Evolution of the state of a machine

- Let $S_k(t) \in \{1, \ldots, S\}$ be the state of the machine of the machine $k$.

- The greater the value of $S_k(t)$, the worse the machine $k$ performs.

- **Remark:** Usually the evolution of the state of machine $k$ is modelled by a *controlled Markov chain*.

- At each instant $t = 0, \ldots, N-1$ a repairman is checking the machine $k$ (action $a = 1$) or no one is checking (action $a = 0$). We can assume that

# Model: Evolution of the state of a machine

- Let $S_k(t) \in \{1, \ldots, S\}$ be the state of the machine of the machine $k$.

- The greater the value of $S_k(t)$, the worse the machine $k$ performs.

- **Remark:** Usually the evolution of the state of machine $k$ is modelled by a *controlled Markov chain*.

- At each instant $t = 0, \ldots, N - 1$ a repairman is checking the machine $k$ (action $a = 1$) or no one is checking (action $a = 0$). We can assume that

$$\mathbb{P}(S(t+1) = 0 | S(t) = s, A(t) = 1) \quad = \quad 1$$

# Model: Evolution of the state of a machine

- Let $S_k(t) \in \{1, \ldots, S\}$ be the state of the machine of the machine $k$.

- The greater the value of $S_k(t)$, the worse the machine $k$ performs.

- **Remark:** Usually the evolution of the state of machine $k$ is modelled by a *controlled Markov chain*.

- At each instant $t = 0, \ldots, N-1$ a repairman is checking the machine $k$ (action $a = 1$) or no one is checking (action $a = 0$). We can assume that

$$
\begin{array}{rcl}
\mathbb{P}(S(t+1) = 0 | S(t) = s, A(t) = 1) & = & 1 \\
\mathbb{P}(S(t+1) = s' | S(t) = s, A(t) = 0) & = & I\{s' \geq s\} p_{ss'}
\end{array}
\tag{1}
$$

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.
- **System dynamic:**

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.
- **System dynamic:**
    1. When a charging spot is available, a new vehicle may join the system and occupy the spot.

---

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.

- **System dynamic:**
    1. When a charging spot is available, a new vehicle may join the system and occupy the spot.
    2. Upon occupying the spot, the vehicle announces the time that it will leave the station and the amount of electricity that it needs.

---

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.
- **System dynamic:**
  1. When a charging spot is available, a new vehicle may join the system and occupy the spot.
  2. Upon occupying the spot, the vehicle announces the time that it will leave the station and the amount of electricity that it needs.

---

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.
- **System dynamic:**
  1. When a charging spot is available, a new vehicle may join the system and occupy the spot.
  2. Upon occupying the spot, the vehicle announces the time that it will leave the station and the amount of electricity that it needs.
- **Reward and cost:**

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.

- **System dynamic:**
  1. When a charging spot is available, a new vehicle may join the system and occupy the spot.
  2. Upon occupying the spot, the vehicle announces the time that it will leave the station and the amount of electricity that it needs.

- **Reward and cost:**
  1. The charging station obtains a reward of $1 - c$ for each unit of electricity provided.

---

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

# Deadline Scheduling (for charging electric vehicles)[3]

- **Set-up:** Charging station has $N$ charging spots and enough power to charge $M$ vehicles at each round.
- **System dynamic:**
    1. When a charging spot is available, a new vehicle may join the system and occupy the spot.
    2. Upon occupying the spot, the vehicle announces the time that it will leave the station and the amount of electricity that it needs.
- **Reward and cost:**
    1. The charging station obtains a reward of $1 - c$ for each unit of electricity provided.
    2. If the station cannot fully charge the vehicle by the time it leaves, the station needs to pay a penalty proportional to the amount of the unfulfilled charge.

---

[3]Yu, Zhe, Yunjian Xu, and Lang Tong. "Deadline scheduling as restless bandits." IEEE Transactions on Automatic Control 63, no. 8 (2018): 2343-2358.

## Model

**States:** Let $T_k(t) = d_k - t$ be the lead time to deadline $d_k$ and $B_k(t)$ be the amount of electricity. The state of the $k$-th spot in defined as

# Model

**States:** Let $T_k(t) = d_k - t$ be the lead time to deadline $d_k$ and $B_k(t)$ be the amount of electricity. The state of the $k$-th spot in defined as

$$S_k(t) = \begin{cases} (0,0) & \text{if no job is at the } k\text{-th position,} \\ (T_k(t), B_k(t)) & \text{otherwise.} \end{cases}$$

## Model

**States:** Let $T_k(t) = d_k - t$ be the lead time to deadline $d_k$ and $B_k(t)$ be the amount of electricity. The state of the $k$-th spot in defined as

$$S_k(t) = \begin{cases} (0,0) & \text{if no job is at the } k\text{-th position,} \\ (T_k(t), B_k(t)) & \text{otherwise.} \end{cases}$$

**Evolution:** The (Markovian) evolution of the state is given by:

# Model

**States:** Let $T_k(t) = d_k - t$ be the lead time to deadline $d_k$ and $B_k(t)$ be the amount of electricity. The state of the $k$-th spot in defined as

$$S_k(t) = \begin{cases} (0,0) & \text{if no job is at the } k\text{-th position,} \\ (T_k(t), B_k(t)) & \text{otherwise.} \end{cases}$$

**Evolution:** The (Markovian) evolution of the state is given by:

$$S_k(t+1) = \begin{cases} (T_k(t)-1, [B_k(t)-a_k(t)]_+) & \text{if } T_k(t) > 1, \\ (T, B) \text{ with prob. } Q(T, B) & \text{otherwise,} \end{cases}$$

where $a_k(t)$ is the amount of electricity given to spot $k$ at instant $t$.

## Other applications

- Wireless Communication;
- Web Crawling;
- Congestion Control;
- Queuing Systems;
- Cluster and Cloud computing;
- Target Tracking;
- Clinical Trials.

# A quick recall on Markov chain

**Definition**
A *Markov Chain* consists of a finite set $\mathcal{S}$ (called the state space) together with a countable family of random variables $S(0), S(1), S(2), \ldots$ with values in $\mathcal{S}$ such that for all $t \geq 0$:

# A quick recall on Markov chain

### Definition
A *Markov Chain* consists of a finite set $\mathcal{S}$ (called the state space) together with a countable family of random variables $S(0), S(1), S(2), \ldots$ with values in $\mathcal{S}$ such that for all $t \geq 0$:

$$\mathbb{P}(S(t+1) = s'|S(0) = s_0, \ldots, S(t) = s_t) = \mathbb{P}(S(t+1) = s'|S(t) = s_t), \tag{2}$$

# A quick recall on Markov chain

### Definition
A *Markov Chain* consists of a finite set $\mathcal{S}$ (called the state space) together with a countable family of random variables $S(0), S(1), S(2), \ldots$ with values in $\mathcal{S}$ such that for all $t \geq 0$:

$$\mathbb{P}(S(t+1) = s'|S(0) = s_0, \ldots, S(t) = s_t) = \mathbb{P}(S(t+1) = s'|S(t) = s_t), \tag{2}$$

- We refer to this fundamental equation as the *Markov property*.

# A quick recall on Markov chain

### Definition
A *Markov Chain* consists of a finite set $\mathcal{S}$ (called the state space) together with a countable family of random variables $S(0), S(1), S(2), \ldots$ with values in $\mathcal{S}$ such that for all $t \geq 0$:

$$\mathbb{P}(S(t+1) = s' | S(0) = s_0, \ldots, S(t) = s_t) = \mathbb{P}(S(t+1) = s' | S(t) = s_t), \tag{2}$$

- We refer to this fundamental equation as the *Markov property*.
- We set for every $s, s' \in \mathcal{S}$:

$$p_{ss'} := \mathbb{P}(S(t+1) = s' | S(t) = s).$$

# A quick recall on Markov chain

> ### Definition
> A *Markov Chain* consists of a finite set $\mathcal{S}$ (called the state space) together with a countable family of random variables $S(0), S(1), S(2), \ldots$ with values in $\mathcal{S}$ such that for all $t \geq 0$:
>
> $$\mathbb{P}(S(t+1) = s'|S(0) = s_0, \ldots, S(t) = s_t) = \mathbb{P}(S(t+1) = s'|S(t) = s_t), \tag{2}$$

- We refer to this fundamental equation as the *Markov property*.
- We set for every $s, s' \in \mathcal{S}$:

$$p_{ss'} := \mathbb{P}(S(t+1) = s'|S(t) = s).$$

The matrix $P := [[p_{ss'}]]_{s,s'}$ is called the *transition matrix*.

# One arm Restless Bandit

A **One Arm Restless Bandit** is a MDP defined on:

# One arm Restless Bandit

A **One Arm Restless Bandit** is a MDP defined on:

- A finite state space $\mathcal{S} := \{1, \ldots, d\}$,

# One arm Restless Bandit

A **One Arm Restless Bandit** is a MDP defined on:

- A finite state space $\mathcal{S} := \{1, \ldots, d\}$,
- A finite action space $\mathcal{A} := \{0, 1\}$,

# One arm Restless Bandit

A **One Arm Restless Bandit** is a MDP defined on:

- A finite state space $\mathcal{S} := \{1, \ldots, d\}$,
- A finite action space $\mathcal{A} := \{0, 1\}$,

and where:

# One arm Restless Bandit

A **One Arm Restless Bandit** is a MDP defined on:

- A finite state space $\mathcal{S} := \{1, \ldots, d\}$,
- A finite action space $\mathcal{A} := \{0, 1\}$,

and where:

- $S(t) \in \mathcal{S}$ is the state of the bandit at the discrete decision time $t \in \{0, \cdots, T\}$,

# One arm Restless Bandit

A **One Arm Restless Bandit** is a MDP defined on:

- A finite state space $\mathcal{S} := \{1, \ldots, d\}$,
- A finite action space $\mathcal{A} := \{0, 1\}$,

and where:

- $S(t) \in \mathcal{S}$ is the state of the bandit at the discrete decision time $t \in \{0, \cdots, T\}$,
- $A(t) \in \mathcal{A}$ is the action taken by the decision maker at the discrete decision time $t \in \{0, \cdots, T\}$.

# One arm Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

# One arm Restless Bandit

*For each time-step $t = 0, \ldots, T - 1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) \in \mathcal{S}$;

# One arm Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) \in \mathcal{S}$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) \in \mathcal{A}$;

# One arm Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) \in \mathcal{S}$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) \in \mathcal{A}$;

3. The decision-maker collects the reward $r_{S(t)}^{A(t)}$;

# One arm Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) \in \mathcal{S}$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) \in \mathcal{A}$;

3. The decision-maker collects the reward $r_{S(t)}^{A(t)}$;

4. The Markov process evolves to $S(t+1) = s'$ with probability $p_{S(t),s'}^{A(t)}$.

# One arm Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) \in \mathcal{S}$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) \in \mathcal{A}$;

3. The decision-maker collects the reward $r_{S(t)}^{A(t)}$;

4. The Markov process evolves to $S(t+1) = s'$ with probability $p_{S(t),s'}^{A(t)}$.

**Objective:** Maximize the expected total sum of rewards over the $T$ time-steps.

# One arm Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) \in \mathcal{S}$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) \in \mathcal{A}$;

3. The decision-maker collects the reward $r_{S(t)}^{A(t)}$;

4. The Markov process evolves to $S(t+1) = s'$ with probability $p_{S(t),s'}^{A(t)}$.

**Objective:** Maximize the expected total sum of rewards over the $T$ time-steps.

´ **Knowns parameters:** $\mathcal{S}$, $\mathcal{A}$, reward $R := [[r_s^a]]_{s,a}$, Horizon $T$, transition matrix $P^a := [[p_{s,s'}^a]]_{s,s'}$.

# Control Policies

### Definition

- A **general control policy** $\pi$ is a mapping from each *possible history* $h_t = \{s_\tau, a_\tau\}_{0 \leq \tau \leq t}$ to $a_t = \pi_t(h_t)$.

# Control Policies

### Definition

- A **general control policy** $\pi$ is a mapping from each *possible history* $h_t = \{s_\tau, a_\tau\}_{0 \le \tau \le t}$ to $a_t = \pi_t(h_t)$.

- A **Markov control policy** $\pi$ depends on the current state and time only: $a_t = \pi_t(s_t)$ .

# Control Policies

### Definition

- A **general control policy** $\pi$ is a mapping from each *possible history* $h_t = \{s_\tau, a_\tau\}_{0 \le \tau \le t}$ to $a_t = \pi_t(h_t)$.

- A **Markov control policy** $\pi$ depends on the current state and time only: $a_t = \pi_t(s_t)$ .

- A **randomized control policy**, the action $a_t$ is chosen according to a probability distribution $\pi_t(a|h_t)$ over $A_t$ .

# Control Policies

## Definition

- A **general control policy** $\pi$ is a mapping from each *possible history* $h_t = \{s_\tau, a_\tau\}_{0 \leq \tau \leq t}$ to $a_t = \pi_t(h_t)$.
- A **Markov control policy** $\pi$ depends on the current state and time only: $a_t = \pi_t(s_t)$ .
- A **randomized control policy**, the action $a_t$ is chosen according to a probability distribution $\pi_t(a|h_t)$ over $A_t$ .

**Remarks:**

1. Randomized Markov policies suffice to achieve the optimum in classical constrained MDP. *In this course, we will only use such policies.*

# Control Policies

## Definition

- A **general control policy** $\pi$ is a mapping from each *possible history* $h_t = \{s_\tau, a_\tau\}_{0 \leq \tau \leq t}$ to $a_t = \pi_t(h_t)$.
- A **Markov control policy** $\pi$ depends on the current state and time only: $a_t = \pi_t(s_t)$ .
- A **randomized control policy**, the action $a_t$ is chosen according to a probability distribution $\pi_t(a|h_t)$ over $A_t$ .

**Remarks:**

1. Randomized Markov policies suffice to achieve the optimum in classical constrained MDP. *In this course, we will only use such policies.*

2. When a randomized Markov policy $\pi_t$ is used, the probability that the Markov process evolves to $S(t+1) = s'$ and action $A(t) = a$, knowing $S(t) = s$ is given by $p^a_{s,s'} \pi^a(s)$.

# Mathematical formulation of the problem

For a given randomized Markov policy $\pi := [\pi_t]_{0 \leq t \leq T-1}$, we define the cumulative reward:

# Mathematical formulation of the problem

For a given randomized Markov policy $\pi := [\pi_t]_{0 \leq t \leq T-1}$, we define the cumulative reward:

$$V_1^\pi(x^0, T) := \mathbb{E} \left[ \sum_{t=0}^{T-1} R_{S_t}^{A_t} \mid S_0 = x^0 \right].$$

# Mathematical formulation of the problem

For a given randomized Markov policy $\pi := [\pi_t]_{0 \leq t \leq T-1}$, we define the cumulative reward:

$$V_1^\pi(x^0, T) := \mathbb{E}\left[\sum_{t=0}^{T-1} R_{S_t}^{A_t} \mid S_0 = x^0\right].$$

The **Value function** is given by:

$$V_1^*(x^0, T) = \min_\pi V_1^\pi(x^0, T)$$

## LP formulation

Let us define the following LP problem:

$$
\begin{aligned}
\min_{y \geq 0} \quad & \sum_{t=0}^{T-1} \sum_{s,a} R_s^a y_{a,s}(t) \\
\text{s.t.} \quad & y_{s,0}(t) + y_{s,1}(t) = x_s(t), \ \forall t \in [[0, T-1]], \ \forall s \in \mathcal{S}, \\
& x_s(t) = \sum_{s'} \sum_a y_{s',a}(t-1) p_{s',s}^a \ \forall t \in [[1, T-1]], \ \forall s \in \mathcal{S}, \\
& x_s(0) = x^0, \ \forall s \in \mathcal{S}
\end{aligned}
\tag{3}
$$

# Equivalence

**Lemma:** Let $y^*$ be a solution of (3). If for all $0 \leq t \leq T-1$, for all $s \in \mathcal{S}$ and for all $a \in \mathcal{A}$, we define

$$\pi_t(a|s) = \left\{ \begin{array}{ll} y_s^a(t)/x_s^a(t), & \text{if } x_s^a(t) > 0, \\ 0 & \text{otherwise,} \end{array} \right.$$

then

$$V_1^*(x^0, T) = V_1^\pi(x^0, T).$$

# $N$-Arms Restless Bandit

A $N$-**arms Restless Bandit** is composed of $N$ statistically MDPs where:

# $N$-Arms Restless Bandit

A $N$-**arms Restless Bandit** is composed of $N$ statistically MDPs where:

- $S_k(t) \in \mathcal{S}$ is the state of the arm $k$ at the discrete decision time $t \in \{0, \cdots, T\}$,

# $N$-Arms Restless Bandit

A $N$-**arms Restless Bandit** is composed of $N$ statistically MDPs where:

- $S_k(t) \in \mathcal{S}$ is the state of the arm $k$ at the discrete decision time $t \in \{0, \cdots, T\}$,
- $A_k(t) \in \{0, 1\}$ is the action taken by the decision maker at the discrete decision time $t \in \{0, \cdots, T\}$.

# $N$-Arms Restless Bandit

A $N$-**arms Restless Bandit** is composed of $N$ statistically MDPs where:

- $S_k(t) \in \mathcal{S}$ is the state of the arm $k$ at the discrete decision time $t \in \{0, \cdots, T\}$,
- $A_k(t) \in \{0, 1\}$ is the action taken by the decision maker at the discrete decision time $t \in \{0, \cdots, T\}$.
- We assume that the decision maker chooses a fraction $0 < \alpha < 1$ of the $N$ arms to be activated.

# $N$-Arms Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

# $N$-Arms Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) := [S_1(t), \ldots, S_N(t)] \in \mathcal{S}^N$;

# $N$-Arms Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) := [S_1(t), \ldots, S_N(t)] \in \mathcal{S}^N$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) := [A_1(t), \ldots, A_N(t)] \in \{0, 1\}^N$, such that $\sum_k A_k(t) \leq N\alpha$;

# $N$-Arms Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) := [S_1(t), \ldots, S_N(t)] \in \mathcal{S}^N$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) := [A_1(t), \ldots, A_N(t)] \in \{0, 1\}^N$, such that $\sum_k A_k(t) \leq N\alpha$;

3. The decision-maker collects the reward $\sum_k r_{S_k(t)}^{A_k(t)}$;

# $N$-Arms Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) := [S_1(t), \ldots, S_N(t)] \in \mathcal{S}^N$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) := [A_1(t), \ldots, A_N(t)] \in \{0, 1\}^N$, such that $\sum_k A_k(t) \leq N\alpha$;

3. The decision-maker collects the reward $\sum_k r_{S_k(t)}^{A_k(t)}$;

4. For every $k$, the arm $k$ evolves to $S_k(t+1) = s'$ with probability $p_{S_k(t), s'}^{A_k(t)}$.

# $N$-Arms Restless Bandit

*For each time-step $t = 0, \ldots, T - 1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) := [S_1(t), \ldots, S_N(t)] \in \mathcal{S}^N$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) := [A_1(t), \ldots, A_N(t)] \in \{0, 1\}^N$, such that $\sum_k A_k(t) \le N\alpha$;

3. The decision-maker collects the reward $\sum_k r_{S_k(t)}^{A_k(t)}$;

4. For every $k$, the arm $k$ evolves to $S_k(t + 1) = s'$ with probability $p_{S_k(t), s'}^{A_k(t)}$.

**Objective:** Maximize the expected total sum of rewards over the $T$ time-steps.

# $N$-Arms Restless Bandit

*For each time-step $t = 0, \ldots, T-1$:*

1. The decision-maker gets full knowledge of the current system state $S(t) := [S_1(t), \ldots, S_N(t)] \in \mathcal{S}^N$;

2. Once $S(t)$ has been observed, the decision-maker chooses a control $A(t) := [A_1(t), \ldots, A_N(t)] \in \{0, 1\}^N$, such that $\sum_k A_k(t) \le N\alpha$;

3. The decision-maker collects the reward $\sum_k r_{S_k(t)}^{A_k(t)}$;

4. For every $k$, the arm $k$ evolves to $S_k(t+1) = s'$ with probability $p_{S_k(t), s'}^{A_k(t)}$.

**Objective:** Maximize the expected total sum of rewards over the $T$ time-steps.

´ **Knowns parameters:** $\mathcal{S}$, reward $R := [[r_s^a]]_{s,a}$, Horizon $T$, transition matrix $P^a := [[p_{s,s'}^a]]_{s,s'}$.

# A new state representation

To simplify the mathematical formulation of the problem, we make the following observation:

# A new state representation

To simplify the mathematical formulation of the problem, we make the following observation:

- **Arms are exchangeable.** Two arms in the same state and for which the same action is chosen provide the *same reward* and have the *same transition probabilities*.

# A new state representation

To simplify the mathematical formulation of the problem, we make the following observation:

- **Arms are exchangeable.** Two arms in the same state and for which the same action is chosen provide the *same reward* and have the *same transition probabilities*.
- **Implication:** It is equivalent to use a control on $S(t)$ or on the *empirical distribution of states* at every instant $t$.

# A new state representation

To simplify the mathematical formulation of the problem, we make the following observation:

- **Arms are exchangeable.** Two arms in the same state and for which the same action is chosen provide the *same reward* and have the *same transition probabilities*.
- **Implication:** It is equivalent to use a control on $S(t)$ or on the *empirical distribution of states* at every instant $t$.

**New notations:**

# A new state representation

To simplify the mathematical formulation of the problem, we make the following observation:

- **Arms are exchangeable.** Two arms in the same state and for which the same action is chosen provide the *same reward* and have the *same transition probabilities*.

- **Implication:** It is equivalent to use a control on $S(t)$ or on the *empirical distribution of states* at every instant $t$.

**New notations:**

- $M_s^{(N)}(t) :=$ the fraction of arms in state $s$ at time $t$.
  $M^{(N)}(t) := [M_s^{(N)}(t)]_{s \in \mathcal{S}}$ is the associated vector.

# A new state representation

To simplify the mathematical formulation of the problem, we make the following observation:

- **Arms are exchangeable.** Two arms in the same state and for which the same action is chosen provide the *same reward* and have the *same transition probabilities*.

- **Implication:** It is equivalent to use a control on $S(t)$ or on the *empirical distribution of states* at every instant $t$.

**New notations:**

- $M_s^{(N)}(t) :=$ the fraction of arms in state $s$ at time $t$. $M^{(N)}(t) := [M_s^{(N)}(t)]_{s \in \mathcal{S}}$ is the associated vector.

- $Y_{s,a}^{(N)}(t) :=$ the fraction of arms in state $s$ at time $t$ for which decision $a$ is taken. $Y^{(N)}(t) := [Y_{s,a}^{(N)}(t)]_{s \in \mathcal{S}, a \in \{0,1\}}$ is the associated vector.

## Mathematical Formulation

$$\min_{\pi} \quad \mathbb{E} \sum_{t=0}^{T-1} \sum_{s,a} r_s^a Y_{a,s}^{(N)}(t) := V_{opt}^{(N)}(m(0), T) \tag{4a}$$

s.t.   Arms follow the Markovian evolution generated by $\Pi_n p_{s_n, s_n'}^{a_n}$,
$$\tag{4b}$$

$$Y_{0,s}^{(N)}(t) + Y_{1,s}^{(N)}(t) = M_s^{(N)}(t), \ \forall t \in [[0, T-1]], \ \forall s \in \mathcal{S}, \tag{4c}$$

$$\sum_s Y_{s,1}^{(N)}(t) \leq \alpha \ \forall t \in [[0, T-1]],, \tag{4d}$$

$$M_s^{(N)}(0) = m_s(0), \ \forall s \in \mathcal{S}, \tag{4e}$$

where $m_s(0) = \frac{1}{N} \sum_{k=1}^{N} I\{S_k(0) = s\}$ , for all $s \in \mathcal{S}$.

# Difficulty

The key difficulty of the $N$-Arms Restless Bandit problem is coming from:

$$\sum_s Y_{s,1}^{(N)}(t) \le \alpha \;\forall t \in [[0, T-1]],$$

which couples all the arms together.

**Challenge:**
How to design an efficient heuristic to solve such problem?

# Outline of the approach

1. **Relaxation:** Classical approach is to relax this constraint and consider a problem where this constraint has to be satisfied only in expectation:

# Outline of the approach

1. **Relaxation:** Classical approach is to relax this constraint and
   consider a problem where this constraint has to be satisfied
   only in expectation:

$$\sum_s \mathbb{E}[Y_{s,1}^{(N)}(t)] \le \alpha, \ \forall t \in [[0, T-1]].$$

# Outline of the approach

1. **Relaxation:** Classical approach is to relax this constraint and consider a problem where this constraint has to be satisfied only in expectation:

$$\sum_s \mathbb{E}[Y_{s,1}^{(N)}(t)] \leq \alpha, \ \forall t \in [[0, T-1]].$$

2. **Interpolation:** Construct a sequence of decision rules $\pi_t : \Delta^d \to \Delta^{2d}$ which is optimal for the relaxed problem.

# Relaxed problem

$$\min_{\pi} \quad \sum_{t=0}^{T-1} \mathbb{E} \sum_{s,a} r_s^a Y_{a,s}^{(N)}(t) =: V_{rel}^{(N)}(m(0), T) \tag{5a}$$

s.t.    Arms follow the Markovian evolution, (5b)

$$Y_{0,s}^{(N)}(t) + Y_{1,s}^{(N)}(t) = M_s^{(N)}(t), \ \forall t \in [[0, T-1]], \ \forall s \in \mathcal{S}, \tag{5c}$$

$$\sum_s \mathbb{E}[Y_{s,1}^{(N)}(t)] \le \alpha \ \forall t \in [[0, T-1]],, \tag{5d}$$

$$M_s^{(N)}(0) = m_s(0), \ \forall s \in \mathcal{S}, \tag{5e}$$

## LP formulation

Let us define the following LP problem:

$$
\begin{aligned}
\min_{y \geq 0} \quad & \sum_{t=0}^{T-1} \sum_{s,a} r_s^a y_{s,a}(t) =: V_{LP}(m(0), T) \\
\text{s.t.} \quad & y_{s,0}(t) + y_{s,1}(t) = m_s(t), \ \forall t \in [[0, T-1]], \ \forall s \in \mathcal{S}, \\
& m_s(t) = \sum_{s'} \sum_{a} y_{s',a}(t-1) p_{s',s}^a \ \forall t \in [[1, T-1]], \ \forall s \in \mathcal{S}, \\
& \sum_{s} y_{s,1}(t) \leq \alpha, \ \forall t \in [[0, T-1]],, \\
& m_s(0) = m^0, \ \forall s \in \mathcal{S}
\end{aligned}
$$

(6)

## LP formulation

Let us define the following LP problem:

$$
\min_{y \geq 0} \quad \sum_{t=0}^{T-1} \sum_{s,a} r_s^a y_{s,a}(t) =: V_{LP}(m(0), T)
$$

$$
\text{s.t.} \quad y_{s,0}(t) + y_{s,1}(t) = m_s(t), \ \forall t \in [[0, T-1]], \ \forall s \in \mathcal{S},
$$

$$
m_s(t) = \sum_{s'} \sum_{a} y_{s',a}(t-1) p_{s',s}^a \ \forall t \in [[1, T-1]], \ \forall s \in \mathcal{S},
$$

$$
\sum_s y_{s,1}(t) \leq \alpha, \ \forall t \in [[0, T-1]],,
$$

$$
m_s(0) = m^0, \ \forall s \in \mathcal{S}
$$

$$
\tag{6}
$$

We denote by $y^* := [[[y_{s,a}^*(t)]]]_{s,a,t}$ the optimal solution of (6) and we also define $m^* := [[m_s(t) := \sum_a y_{s,a}^*(t)]]_{s,t}$.

# Equivalence

**Lemma:**

$$
\begin{aligned}
V_{rel}(m^0, T) &= V_{LP}(m^0, T), \\
V_{opt}^{(N)}(m(0), T) &\geq V_{LP}(m^0, T).
\end{aligned}
$$

## Projection

We define the set of feasible control at time $t$ by:

$$\mathcal{Y}(M^{(N)}(t)) := \left\{ y \in \mathbb{R}_+^{2d} \mid \sum_a y_{s,a} = M_s^{(N)}(t) \ \forall s \in \mathcal{S}; \ \sum_s y_{s,1} \leq \alpha \right\}$$

## Projection

We define the set of feasible control at time $t$ by:

$$\mathcal{Y}(M^{(N)}(t)) := \left\{ y \in \mathbb{R}_+^{2d} \mid \sum_a y_{s,a} = M_s^{(N)}(t) \ \forall s \in \mathcal{S}; \ \sum_s y_{s,1} \leq \alpha \right\}$$

Some observations:

## Projection

We define the set of feasible control at time $t$ by:

$$\mathcal{Y}(M^{(N)}(t)) := \left\{ y \in \mathbb{R}_+^{2d} \mid \sum_a y_{s,a} = M_s^{(N)}(t) \; \forall s \in \mathcal{S}; \; \sum_s y_{s,1} \leq \alpha \right\}$$

Some observations:

1. In general $y^*(t) \notin \mathcal{Y}(M^{(N)}(t))$;

## Projection

We define the set of feasible control at time $t$ by:

$$\mathcal{Y}(M^{(N)}(t)) := \left\{ y \in \mathbb{R}_+^{2d} \mid \sum_a y_{s,a} = M_s^{(N)}(t) \; \forall s \in \mathcal{S}; \; \sum_s y_{s,1} \leq \alpha \right\}$$

Some observations:

1. In general $y^*(t) \notin \mathcal{Y}(M^{(N)}(t))$;
2. In general $y^*(t) \in \mathcal{Y}(m^*(t))$.

## Projection

We define the set of feasible control at time $t$ by:

$$\mathcal{Y}(M^{(N)}(t)) := \left\{ y \in \mathbb{R}_+^{2d} | \sum_a y_{s,a} = M_s^{(N)}(t) \; \forall s \in \mathcal{S}; \; \sum_s y_{s,1} \leq \alpha \right\}$$

Some observations:

1. In general $y^*(t) \notin \mathcal{Y}(M^{(N)}(t))$;
2. In general $y^*(t) \in \mathcal{Y}(m^*(t))$.

---

We define the following **projection operator**:

$$\pi_t^{Proj}(M^{(N)}) := \mathsf{Proj}_t(M^{(N)}) := \mathsf{argmin}_{y \in \mathcal{Y}(M^{(N)}(t))} \|y - y^*(t)\|_2^2. \tag{7}$$

---

# Algorithm

**The Projection Policy**

- **Input:** Initial system configuration vector $m(0)$ and time horizon $T$.
- **Solve** The LP to obtain $y^*$;
- **Set** $\hat{M} := \mathsf{m}(0)$;
- **For** $t = 0, 2, \ldots, T-1$ **do:**
    1. *Projection step:* Compute $\hat{y}(t) := \mathsf{Proj}_t(\hat{M})$;
    2. *Rounding step:* For all $s \in \mathcal{S}$, set:

    $$\hat{Y}_{s,a}^{(N)}(t) = \left\{ \begin{array}{ll} N^{-1} \lfloor N \hat{y}_{s,1}(t) \rfloor & \text{if } a = 1, \\ \hat{M}_s - N^{-1} \lfloor N \hat{y}_{s,1}(t) \rfloor & \text{otherwise.} \end{array} \right.$$

    3. Use control $\hat{Y}^{(N)}$ to advance to the next time-step ;
    4. Set $\hat{M} :=$ current empirical distribution;

# Policy construction scheme

A simple scheme can explain easily our algorithm:

# Policy construction scheme

A simple scheme can explain easily our algorithm:

$$\hat{M}(t) \xrightarrow[\text{Proj. step}]{\mathsf{Proj}_t(\hat{M}(t))} \hat{y}(t) \xrightarrow[\text{Roun. step}]{} \hat{Y}_{s,a}^{(N)}(t) \xrightarrow[\text{Trans. step}]{} \hat{M}(t+1)$$

**Remark:** We will see in the next theorem that the projection step can be replaced by map $\pi(\cdot)$ such that:

# Policy construction scheme

A simple scheme can explain easily our algorithm:

$$\hat{M}(t) \xrightarrow[\text{Proj. step}]{\text{Proj}_t(\hat{M}(t))} \hat{y}(t) \xrightarrow[\text{Roun. step}]{} \hat{Y}^{(N)}_{s,a}(t) \xrightarrow[\text{Trans. step}]{} \hat{M}(t+1)$$

**Remark:** We will see in the next theorem that the projection step can be replaced by map $\pi(\cdot)$ such that:

1. *Admissible policy:* $\pi_t(M^{(N)}(t)) \in \mathcal{Y}(M^{(N)}(t))$,

# Policy construction scheme

A simple scheme can explain easily our algorithm:

$$\hat{M}(t) \xrightarrow[\text{Proj. step}]{\text{Proj}_t(\hat{M}(t))} \hat{y}(t) \xrightarrow[\text{Roun. step}]{} \hat{Y}_{s,a}^{(N)}(t) \xrightarrow[\text{Trans. step}]{} \hat{M}(t+1)$$

**Remark:** We will see in the next theorem that the projection step can be replaced by map $\pi(\cdot)$ such that:

1. *Admissible policy:* $\pi_t(M^{(N)}(t)) \in \mathcal{Y}(M^{(N)}(t))$,
2. *LP-compatible policy:* $\pi_t(m^*(t)) = y^*(t)$.

# Main result

## Theorem
*Let $\pi := \{\pi_t\}_{0 \leq t \leq T-1}$ be a continuous an admissible and continuous policy then*

# Main result

## Theorem
Let $\pi := \{\pi_t\}_{0 \le t \le T-1}$ be a continuous an admissible and continuous policy then

$$\lim_{N \to +\infty} V_{opt,\pi}^{(N)}(m(0), T) = V_{rel,\pi}(m(0), T).$$

# Main result

### Theorem
*Let $\pi := \{\pi_t\}_{0 \le t \le T-1}$ be a continuous an admissible and continuous policy then*

$$\lim_{N \to +\infty} V_{opt,\pi}^{(N)}(m(0), T) = V_{rel,\pi}(m(0), T).$$

*Moreover if $\pi$ is LP-compatible then*

# Main result

## Theorem
*Let $\pi := \{\pi_t\}_{0 \le t \le T-1}$ be a continuous an admissible and continuous policy then*

$$\lim_{N \to +\infty} V_{opt,\pi}^{(N)}(m(0), T) = V_{rel,\pi}(m(0), T).$$

*Moreover if $\pi$ is LP-compatible then*

$$\lim_{N \to +\infty} V_{opt,\pi}^{(N)}(m(0), T) = V_{LP}(m(0), T).$$

# Infinite Restless Bandit

The initial Restless Bandit was defined as follows:

$$\min_{\pi \in \Pi} \quad \lim_{T \to +\infty} \frac{1}{T} \mathbb{E} \sum_{t=0}^{T-1} \sum_{s,a} r_s^a Y_{a,s}^{(N)}(t) =: V_{opt}^{(N)}(\infty) \tag{8a}$$

s.t. Arms follow the Markovian evolution generated by $\Pi_n p_{s_n, s_n'}^{a_n}$, (8b)

$$Y_{0,s}^{(N)}(t) + Y_{1,s}^{(N)}(t) = M_s^{(N)}(t), \ \forall t \in [[0, T-1]], \ \forall s \in \mathcal{S}, \tag{8c}$$

$$\sum_s Y_{s,1}^{(N)}(t) \leq \alpha \ \forall t \in [[0, T-1]],, \tag{8d}$$

$$M_s^{(N)}(0) = m_s(0), \ \forall s \in \mathcal{S}, \tag{8e}$$

where $m_s(0) = \frac{1}{N} \sum_{k=1}^N I\{S_k(0) = s\}$ , for all $s \in \mathcal{S}$ and $\Pi$ is the set of Markovian policy.

# The associated LP

We next relax the constraints $\sum_s Y_{s,1}^{(N)}(t) \le \alpha, \ \forall t \in [[0, T-1]]$ into:

## The associated LP

We next relax the constraints $\sum_s Y_{s,1}^{(N)}(t) \leq \alpha, \ \forall t \in [[0, T-1]]$ into:

$$\lim_{T \to +\infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_s \mathbb{E}_\pi[Y_{s,1}^{(N)}(t)] \leq \alpha. \tag{9}$$

## The associated LP

We next relax the constraints $\sum_s Y_{s,1}^{(N)}(t) \leq \alpha, \ \forall t \in [[0, T-1]]$ into:

$$\lim_{T \to +\infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_s \mathbb{E}_\pi[Y_{s,1}^{(N)}(t)] \leq \alpha. \qquad (9)$$

By defining $y_{s,a} = \lim_{T \to +\infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_s \mathbb{E}_\pi[Y_{s,a}^{(N)}(t)]$, for all $a$ and $s$, we then obtain the following linear program:

## The associated LP

We next relax the constraints $\sum_s Y_{s,1}^{(N)}(t) \le \alpha, \ \forall t \in [[0, T-1]]$ into:

$$\lim_{T \to +\infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_s \mathbb{E}_\pi[Y_{s,1}^{(N)}(t)] \le \alpha. \tag{9}$$

By defining $y_{s,a} = \lim_{T \to +\infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_s \mathbb{E}_\pi[Y_{s,a}^{(N)}(t)]$, for all $a$ and $s$, we then obtain the following linear program:

$$\begin{aligned} \min_{y \ge 0} \quad & \sum_{s,a} r_s^a y_{s,a} =: V_{LP}(\infty) \\ \text{s.t.} \quad & y_{s,0} + y_{s,1} = \sum_{s'} \sum_a y_{s',a} p_{s',s}^a, \ \forall s \in \mathcal{S}, \\ & \sum_s y_{s,1} \le \alpha, \sum_{s,a} y_{s,a} = 1. \end{aligned} \tag{10}$$

# Policies

As before we need to define a policy such that we can transfer the solution $y^*$ of the LP to $N$-arms problems. Three solutions:

- LP-priority policy[4],
- LP-index policy[5],
- Whittle indices[6].

---

[4]Verloop M (2016) Asymptotically optimal priority policies for indexable and nonindexable restless bandits. Annals of Applied Probability 26(4):1947-1995.

[5]Gast, Nicolas, Bruno Gaujal, and Chen Yan. "Linear Program-Based Policies for Restless Bandits: Necessary and Sufficient Conditions for (Exponentially Fast) Asymptotic Optimality." Mathematics of Operations Research (2023).

[6]Weber RR, Weiss G (1990) On an index policy for restless bandits. Journal of Applied Probability 27(3):637-648, ISSN 00219002

# LP-priority

We define the following four sets, which form a partition of $S$

$$\mathcal{S}^+ = \{s \in \mathcal{S} | y_{s,1}^* > 0, \ y_{s,0}^* = 0\}, \tag{11}$$

$$\mathcal{S}^0 = \{s \in \mathcal{S} | y_{s,1}^* > 0, \ y_{s,0}^* > 0\}, \tag{12}$$

$$\mathcal{S}^- = \{s \in \mathcal{S} | y_{s,1}^* = 0, \ y_{s,0}^* > 0\}, \tag{13}$$

$$\mathcal{S}^\emptyset = \{s \in \mathcal{S} | y_{s,1}^* = 0, \ y_{s,0}^* = 0\}. \tag{14}$$

**Definition:** The set of **LP-priorities** are defined as
$\Sigma := \cup_{y^*} \Sigma(y^*)$, where $\Sigma(y^*)$ is the set of permutations
$\sigma = \sigma_1 \ldots \sigma_d$ of the $d$ states such that any state in $\mathcal{S}^+$ appears
before any state in $\mathcal{S}^0$, and any state in $\mathcal{S}^0$ appears before any
state in $\mathcal{S}^-$.

# LP-indices

By strong duality, there exists Lagrange multiplier $\gamma^* \in \mathbb{R}$ such that $y^*$ is also an optimal solution to the following linear program:

## LP-indices

By strong duality, there exists Lagrange multiplier $\gamma^* \in \mathbb{R}$ such that $y^*$ is also an optimal solution to the following linear program:

$$
\begin{aligned}
g(\gamma^*) = \min_{y \geq 0} \quad & \sum_{s,a} r_s^a y_{s,a} + \gamma^* \sum_s y_{s,1} \\
\text{s.t.} \quad & y_{s,0} + y_{s,1} = \sum_{s'} \sum_a y_{s',a} p_{s',s}^a, \ \forall s \in \mathcal{S}, \\
& \sum_{s,a} y_{s,a} = 1.
\end{aligned}
$$

# LP-index policy

We can transform this LP into an MDP, with the value function $V^*(s)$ satisfies the Bellman equation:

$$g(\gamma^*)+V^*(s) = \min\{\underbrace{r_s^1 + \gamma^* r_s^1 + \sum_{s'} p_{s,s'}^1 V^*(s')}_{=:Q_s^1}, \underbrace{r_s^0 + \sum_{s'} p_{s,s'}^0 V^*(s')}_{=:Q_s^0}\}.$$

- The LP indices for the infinite horizon are defined as $I_s := Q_s^1 - Q_s^0$ for state s.
- The **LP-index policy** is the strict priority policy by using the values $I_s$ as a priority order to rank states within $\mathcal{S}^+$, $\mathcal{S}^-$ and $\mathcal{S}^0$ at each decision epoch.

# Whittle indices

- Let us define for each value $\gamma \in \mathbb{R}$, the value function $V_s(\gamma)$ for state $s$ satisfies the Bellman equation:

$$g(\gamma) + V^*(s, \gamma) = \min\{\underbrace{r_s^1 + \gamma^* r_s^1 + \sum_{s'} p_{s,s'}^1 V^*(s', \gamma)}_{Q_s^1(\gamma)},$$
$$\underbrace{r_s^0 + \sum_{s'} p_{s,s'}^0 V^*(s', \gamma)}_{Q_s^0(\gamma)}\}.$$

- Let us also define the set for which the arg min of the parametrized Bellman equation:

$$\mathcal{S}(\gamma) := \{s \in \mathcal{S} | Q_s^1(\gamma) > Q_s^0(\gamma)\}.$$

# Whittle indices (cont'd)

- We say that the Restless Bandit is **indexable** if $\mathcal{S}(\gamma)$ expands monotonically from to the full set $\mathcal{S}$ when $\gamma$ is decreased from $+\infty$ to $-\infty$.

- The **Whittle index** $\gamma_s$ for state $s$ is defined to be the supremum value of $\gamma$ for which $s$ belongs to $\mathcal{S}(\gamma)$.

- **Whittle index policy** is the strict priority policy by using the values $\gamma_s$ as a priority score to rank states within $\mathcal{S}^+$, $\mathcal{S}^-$ and $\mathcal{S}^0$ at each decision epoch.

# Link between the policies

### Theorem
*Assume that the infinite horizon RB is unichain, so that $\mathcal{S}^{\emptyset} = \emptyset$ .*
*Then:*

- $s \in \mathcal{S}^+ \Rightarrow I_s > 0$, $s \in \mathcal{S}^- \Rightarrow I_s < 0$, $s \in \mathcal{S}^0 \Rightarrow I_s = 0$.
- *If we assume furthermore that the infinite horizon RB is indexable in Whittle's sense, then their Whittle indices $\gamma(s)$ satisfy:* $s \in \mathcal{S}^+ \Rightarrow \gamma(s) > \gamma^*$, $s \in \mathcal{S}^- \Rightarrow \gamma(s) < \gamma^*$, $s \in \mathcal{S}^0 \Rightarrow \gamma(s) = \gamma^*$.

# Bibliography

- The proof of the main theorem and more advance theorem can be found here: Gast, Nicolas, Bruno Gaujal, and Chen Yan. "Linear Program-Based Policies for Restless Bandits: Necessary and Sufficient Conditions for (Exponentially Fast) Asymptotic Optimality." Mathematics of Operations Research (2023).

- If you want to find a lot of different applications, you can have a look at: Avrachenkov, Konstantin E., and Vivek S. Borkar. "Whittle index based Q-learning for restless bandits with average reward." Automatica 139 (2022): 110186.