



HAL
open science

Testing Indexability and Computing Whittle and Gittins Index in Subcubic Time

Nicolas Gast, Bruno Gaujal, Kimang Khun

► **To cite this version:**

Nicolas Gast, Bruno Gaujal, Kimang Khun. Testing Indexability and Computing Whittle and Gittins Index in Subcubic Time. 2022. hal-03602458v3

HAL Id: hal-03602458

<https://inria.hal.science/hal-03602458v3>

Preprint submitted on 8 Dec 2022 (v3), last revised 20 Jun 2023 (v5)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Testing Indexability and Computing Whittle and Gittins Index in Subcubic Time

Nicolas Gast, Bruno Gaujal and Kimang Khun

Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP*, LIG, 38000, Grenoble, France.

*Institute of Engineering Univ. Grenoble Alpes.

Contributing authors: nicolas.gast@inria.fr; bruno.gaujal@inria.fr; kimang.khun@inria.fr;

Abstract

Whittle index is a generalization of Gittins index that provides very efficient allocation rules for restless multi-armed bandits. In this work, we develop an algorithm to test the indexability and compute the Whittle indices of any finite-state restless bandit arm. This algorithm works in the discounted and non-discounted cases, and can compute Gittins index. Our algorithm builds on three tools: (1) a careful characterization of Whittle index that allows one to compute recursively the k th smallest index from the $(k - 1)$ th smallest, and to test indexability, (2) the use of the Sherman-Morrison formula to make this recursive computation efficient, and (3) a sporadic use of the fastest matrix inversion and multiplication methods to obtain a subcubic complexity. We show that an efficient use of the Sherman-Morrison formula leads to an algorithm that computes Whittle index in $(2/3)n^3 + o(n^3)$ arithmetic operations, where n is the number of states of the arm. The careful use of fast matrix multiplication leads to the first subcubic algorithm to compute Whittle or Gittins index: By using the current fastest matrix multiplication, the theoretical complexity of our algorithm is $O(n^{2.5286})$. We also develop an efficient implementation of our algorithm that can compute indices of Markov chains with several thousands of states in less than a few seconds.

Keywords: Whittle Index, Gittins Index, Restless Bandit, Multi-armed Bandit, Sherman-Morrison, Markov Decision Process, Fast Matrix Multiplication

1 Introduction

Markovian bandits form a subclass of multi-armed bandit problems in which each arm has an internal state that evolves over time in a Markovian manner, as a function of the decision maker's actions. In such a problem, at each time step, the decision maker observes the state of all arms and chooses which one to activate. When the state of an arm evolves only when this arm is chosen, one falls into the category of *rested* Markovian bandits for which an optimal policy (in the discounted case) was found by Gittins [1]. When the state of an arm can also evolve when the arm is not chosen, the problem is called a *restless bandit* problem, and computing an optimal policy is computationally difficult [2].

In his seminal paper [3], Whittle proposed a very efficient heuristic: For each arm, an index function maps each state of the arm to a real number. The Whittle index policy then consists in activating the arms having the highest index first. This heuristic generalizes Gittins index to restless bandits. Contrary to the rested case, the Whittle index policy is in general not optimal. Yet, this policy has been proven to be very efficient over the years: up to a condition called *indexability*, Whittle index has been shown to be (in the undiscounted case) asymptotically optimal as the number of arms grows to infinity under certain technical assumptions [4–6]. Moreover, the heuristic performs extremely well in practice [7–9]. Restless bandits and Whittle index have been applied to many scheduling and resource allocation problems such as wireless communication [10, 11], web crawling [12, 13], congestion control [14, 15], queueing systems [16–22], and clinical trials [23].

The above examples show that, when a problem is indexable, computing Whittle index is a very efficient way to construct a nearly-optimal heuristic. This raises a few important questions, that we study in this paper:

- Is testing indexability computationally hard?
- Is there an efficient algorithm to compute Whittle index?
- Is Whittle index harder to compute than Gittins index?

Related work. The computation of Gittins index has received a lot of attention in the past, see for instance [24–27] and the recent survey [28]. For a n -state arm, the algorithms having the smallest complexity perform $(2/3)n^3 + O(n^2)$ arithmetic operations [28]. Note that in page 4 of [29] the author claims that it is unlikely that this complexity can be improved. As we see later, we do improve upon this complexity.

Concerning Whittle index, to the best of our knowledge, there is no efficient general purpose algorithm to test indexability, and most papers studying Whittle index either assume that the studied model is indexable or focus on specific classes of restless bandits for which the structure of arms can be used to show indexability, see *e.g.* [17, 19, 30, 31]. Assuming indexability, the computation of Whittle index has been considered by a few papers.

For restless bandits with discount $\beta \in (0, 1)$, the most efficient numerical algorithm to compute Whittle index was recently presented in [29]. This algorithm, called *fast-pivoting* algorithm, performs $(2/3)n^3 + O(n^2)$ arithmetic

operations¹ if the initialization phase is excluded from the count. This is done by using the parametric simplex method and exploiting the special structure of this linear system to reduce the complexity of simplex pivoting steps. This fast-pivoting algorithm is an efficient implementation of adaptive-greedy algorithm [32] which outputs Whittle index if and only if the arm satisfies a technical condition called partial conservation law (PCL), which is more restrictive than just being indexable. So, it is not applicable for all indexable restless bandits. Based on a geometric interpretation of Whittle index, the authors in [33] propose a refinement of the adaptive-greedy algorithm of [32] to compute Whittle indices of all indexable restless bandits. For a n -state arm, the refined algorithm achieves a $O(n^3)$ complexity by using the Sherman-Morrison formula. The authors also propose a few checkable conditions to test indexability. However, those conditions are not necessary for indexability, which means that if an arm does not verify the conditions, we cannot conclude that the arm is non-indexable and an algorithm to check indexability is still needed. Also, no detailed description is given for adapting those conditions and their algorithm to restless bandit without discount. A thorough comparison between our algorithm and [29, 33] is given in Appendix E. While computing the Whittle indices of a known arm's model is still a challenge, there is interesting work in trying to learn Whittle index when only the arm's simulator is given and the arm's model is unknown. For instance, [12, 34, 35] use Q-learning algorithm to estimate Whittle index as time evolves in finite-state restless bandits. Moreover, the work of [36] uses deep reinforcement learning framework to estimate the Whittle indices of the arms with large state space or convoluted transition kernel, assuming a notion of strong indexability.

For non-discounted restless bandit $\beta = 1$, the author of [29] only provides a brief description of how the fast-pivoting algorithm given in the discounted case can be adapted although no explicit algorithm is given in that paper. For continuous-time n -state restless bandits, the work of [37] proposes an algorithm to check indexability and compute Whittle index with a complexity exponential in the number of states n of each arm. According to Remark 4.1 of that paper, this complexity can be reduced to $O(n^5)$ if the restless bandit is known to be indexable and threshold-based policies are optimal. It is stated that their approach is not applicable for discounted restless bandits. For learning aspect, the work of [34] shows as to learn Whittle index in non-discounted case by maintaining two Q-functions, updating them using Q-learning algorithm, and deducing Whittle index from them when needed. The way that Whittle indices are computed is very close to our work but less efficient than our algorithm since the authors are more interested in learning the index.

Contributions. In this paper, we investigate Whittle index computation in restless multi-armed bandit problems and present four main contributions.

Our first contribution is to propose a univocal definition of indexability. This definition clarifies ambiguities that exist in the classical definitions: Classical definitions assume that an arm is indexable if the optimal policy is a

¹multiplications and additions of real numbers, regardless of their values

non-decreasing function of some penalty term λ . While this definition works for most practical cases, it is not always precise enough because the optimal policy is in general not unique. In our definition, we specify the notion of increasingness that should be used. Our definition guarantees the uniqueness of Whittle indices.

Our second contribution is to propose a unified algorithm that computes the Whittle indices for both discounted and non-discounted restless bandits. Our algorithm, which can be viewed as a refinement of the algorithm in [33], tests whether the input arm is indexable or not, and computes Whittle index if the arm is. As a byproduct, our algorithm can compute Gittins index in rested bandits which are a subclass of restless bandits. This algorithm computes the indices in increasing order, and relies on an efficient use of the Sherman-Morrison formula to compute Whittle index in $(2/3)n^3 + O(n^2)$ plus subcubic time [38] to solve a linear system of order n . This algorithm can detect on the fly if a computed index violates the indexability condition, which adds an extra $(1/3)n^3 + O(n^2)$ arithmetic operations. This later test is optional: the complexity of our algorithm is $n^3 + o(n^3)$ when testing indexability and $(2/3)n^3 + o(n^3)$ without the test. For discounted problems, our algorithm works for any finite-state arm. For non-discounted problems, our algorithm takes as an input any weakly communicating arm and can output three results: the arm is indexable, non-indexable, or multichain. We show the correctness of the algorithm which proves that for unichain arms, our algorithm can always test if the arm is indexable or not. The possible outputs of our algorithm are summarized in Figure 1.

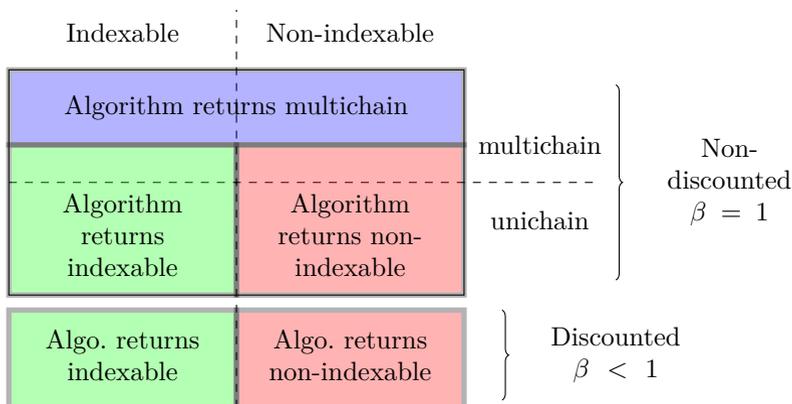


Fig. 1: Possible outputs of our algorithm: For unichain or discounted problems, our algorithm tests indexability and returns the index if and only if the problem is indexable. For some multichain problems, the algorithm can test indexability. For the others, it only returns that the problem is multichain.

Our third contribution is to show how to reduce the complexity of the above algorithm to obtain the first subcubic algorithm to compute Whittle index. This improvement is made possible by the fact that a linear system can be solved in subcubic time. By carefully reordering the computations, we show that it is possible to reduce the use of the Sherman-Morrison formula at the price of solving more linear systems. The subcubic complexity comes by striking a good balance between having too many or too few linear systems to solve. By using the current fastest matrix multiplication method, our algorithm can test indexability and compute Whittle index in $O(n^{2.5286})$. Our algorithm is also the first subcubic algorithm to compute Gittins index.

Our fourth and last contribution is to provide an open-source implementation of our algorithm and present an empirical evaluation of the performance of our implementation. Our results show that our algorithm is very efficient in computing Whittle index and testing indexability. Moreover, our simulations indicate that the subcubic version of our algorithm not only has an asymptotically small complexity but is also faster in practice than our original $(2/3)n^3$ algorithm. Testing the indexability and computing indices takes less than one second for $n = 1000$ states and less than 10 minutes for $n = 15000$ states. This is 15 to 20 times faster than the computation times reported in [29].

Road map. The paper is organized as follows. We introduce the problem and the definition of Whittle index in Section 2. In Section 3, we characterize Whittle index and provide a general idea of how to compute Whittle index. Then, we show, in Section 4, how to use the Sherman-Morrison formula to compute the indices efficiently. We then show how to reduce the complexity of the algorithm by using fast matrix multiplication method in Section 5. We compare the numerical result of different variants of our algorithm in Section 6. We show how to adapt this approach to the discounted case in Section 7. Finally, we conclude in Section 8.

2 Restless bandits and indexability

2.1 Restless bandit arms and multi-armed bandit

In this paper, a restless bandit arm (that we denote later by RB) is a Markov decision process (MDP) with discrete state space $[n] := \{1, \dots, n\}$ and binary action space $\{0, 1\}$, where 0 denotes the action “rest” and 1 denotes the action “activate”. The time is discrete and the evolution is Markovian: If the MDP is in state i and action a is chosen, the decision maker earns an instantaneous reward r_i^a and the arm transitions to a new state j with probability P_{ij}^a . We denote this MDP by the pair (r, P) . The name “restless” comes from the fact that an arm put at rest may still transition to a new state.

A restless multi-armed bandit (RMAB) problem is a finite collection of $M \in \mathbb{N}^*$ independent RB arms. At time t , the decision maker observes the state of all arms, and can choose up to m arms to activate, where $m \leq M$ is a fixed constant. The decision maker then earns a reward that is the sum of the

rewards of all arms. The objective of the decision maker is to identify an allocation rule that maximizes the average reward earned over an infinite number of time steps. Such a problem is notoriously difficult to solve, as its complexity grows exponentially with the number of arms [2]. In his seminar paper [3], Whittle proposes the following approach: *if all arms verify a technical condition known as indexability, then each state i of each arm is associated with a real number λ_i , that is now known as the Whittle index of state i . At each time, the decision maker activates the m arms whose Whittle index of their current states are the m greatest indices.* As mentioned earlier, this heuristic performs extremely well in practice, see *e.g.*, [7–9]. This shows that if the M arms are all indexable, then one can derive a very efficient allocation rule for RMAB problems by computing the Whittle indices of all arms. The Whittle indices of an arm do not depend on the other arms. This shows that the computational cost of Whittle index policy is linear in the number of arms multiplied by the time to compute the indices for a single arm. Hence, in the remaining of the paper, we focus on a single arm and present a new algorithm to test indexability and compute the index of a given arm.

2.2 Indexability and Whittle index

For the remaining of the paper, we consider a single arm (r, P) that has n states and that is assumed to be *weakly communicating* (see definition below). In this section, we introduce the notion of indexability and discuss some ambiguities that we have found when using the definition of indexability defined in previous works.

2.2.1 Policy and arm structure

An arm is a two-action MDP. Hence, a policy π is a subset of the state space, $\pi \subseteq [n]$, such that the policy chooses to activate the arm in state i if $i \in \pi$. We say that π is the set of *active* states and we say that state i is *passive* if $i \notin \pi$. By abuse of notation, we will write $\pi_i = 1$ if $i \in \pi$ and $\pi_i = 0$ if $i \notin \pi$, and we denote by \mathbf{P}^π the transition matrix corresponding to the policy π , *i.e.*, $P_{ij}^\pi = P_{ij}^{\pi_i}$.

Following classical definitions in the literature [39], we say that:

- A policy π is *unichain* if the transition matrix \mathbf{P}^π induced by π has a unique recurrent class. A policy that is not unichain is called *multichain*.
- An arm is unichain if all policies $\pi \subseteq [n]$ are unichain. An arm is multichain if it is not unichain, *i.e.*, if there exists a policy $\pi \subseteq [n]$ that is multichain.
- An arm is *weakly communicating* if there exists a set of states $\mathcal{S} \subseteq [n]$ with each state in \mathcal{S} accessible from every other state of \mathcal{S} and each state in $[n] \setminus \mathcal{S}$ is transient under every policy.

These definitions are important when working with average reward MDPs. In particular, the notion of weakly communicating guarantees that the average reward of the optimal policy, that we will define in the following section, does

not dependent on the initial state [39]. This notion is strictly weaker than the notion of unichain. Indeed, weakly communicating characterizes patterns of states that are accessible from each other under some policy and not the chain structure induced by the policy. So, a weakly communicating MDP can be unichain or multichain.

Testing if a policy is unichain is computationally feasible and can be solved in $O(n^2)$ by using Tarjan's strongly connected component algorithm. Testing if an arm is unichain is, however, NP-hard [40]. Testing if an arm is weakly communicating can also be done in $O(n^2)$: Indeed, an arm is weakly communicating if and only if² the transition matrix $(\mathbf{P}^0 + \mathbf{P}^1)/2$ is unichain.

2.2.2 Gain optimality and Bellman optimality

Following a policy π , we denote by g_i^π the long-term average reward that a decision maker would obtain when starting in state i . In the remainder of the paper, we use the term “gain” to denote the long-term average reward. Let $g_i^* = \max_\pi g_i^\pi$ be the maximal gain starting from state i . As we assume that the arm is weakly communicating, it is shown in [39, Chapter 8] that the quantity g_i^* does not depend on i and we simply denote it by g^* . We say that a policy is gain optimal if $g_i^\pi = g^*$ for all state i .

It is shown in [39, Chapter 8] that g^* is the optimal gain if and only if there exists a vector \mathbf{h}^* , called the bias vector that satisfies the Bellman optimality equation: for all $i \in [n]$,

$$g^* + h_i^* = \max_{a \in \{0,1\}} \left(r_i^a + \sum_{j=1}^n P_{ij}^a h_j^* \right). \quad (1)$$

We say that a policy π is *Bellman optimal* if there exists³ a bias vector \mathbf{h}^* that is a solution of (1) and such that π attains the maximum in (1), *i.e.*: for all i , $\pi_i = \arg \max_{a \in \{0,1\}} (r_i^a + \sum_{j=1}^n P_{ij}^a h_j^*)$.

The notion of Bellman optimality is stronger than the notion of gain optimality: A Bellman optimal policy is gain optimal but the converse is not true in general. Note that the distinction between gain optimal and Bellman optimal policies is only important for the average reward criterion. This distinction disappears for the discounted case that we discuss in Section 7.

2.2.3 λ -penalized MDP and definition of indexability

For each $\lambda \in \mathbb{R}$, we define a λ -penalized MDP⁴ whose transition matrices are the same as in the original MDP and whose reward at time $t \geq 0$ when taking

²This comes from two facts. First, a state is accessible from another state if it is accessible in a Markov chain whose transition matrix is $(\mathbf{P}^0 + \mathbf{P}^1)/2$. Second, a state is transient for any policy if it is transient in a Markov chain whose transition matrix is $(\mathbf{P}^0 + \mathbf{P}^1)/2$.

³If the MDP is unichain, then the bias vector \mathbf{h}^* is unique up to an additive constant. This is in general not the case for multichain MDPs (even if the MDP is weakly communicating).

⁴not to be confused with β -discounted MDPs, where the discount is on rewards and not on actions.

action a_t in state s_t is $r_{s_t}^{a_t} - \lambda a_t$. The quantity λ is a penalty for taking action “activate”. For λ -penalized MDPs, we define the gain and bias functions as in Section 2.2.2, but these quantities now depend on λ . Hence, we will write them as functions of λ : For instance, the optimal gain is $g^*(\lambda)$, and we will use the notation $\mathbf{h}^*(\lambda)$ to denote an optimal bias, and $\pi^*(\lambda)$ to denote an optimal policy.

The classical definition of indexability use in the literature [29, 33, 34, 36] says that an arm is indexable if and only if the optimal policy $\pi^*(\lambda)$ is non-increasing in λ (for the inclusion order). If an arm is indexable, these papers define the Whittle index of a state i as a real number λ_i such that $\pi^*(\lambda) = \{i \in [n] : \lambda_i > \lambda\}$. This definition is ambiguous for two reasons: First, optimal policies are in general not unique. Hence, the notion of $\pi^*(\lambda)$ being non-increasing is unclear: should all optimal policies be non-increasing or at least one? Second, the notion of optimality for a policy is also unclear: should it mean “gain optimal”, “bias optimal” or another notion of optimality?

To solve these ambiguities, in this paper, we use the following definition of indexability.

Definition 1. *Given a finite-state weakly communicating arm, let $\Pi^*(\lambda)$ be the set of Bellman optimal policies for a penalty λ . We say that the arm is indexable if for all $\lambda < \lambda'$, and all policies $\pi \in \Pi^*(\lambda)$ and $\pi' \in \Pi^*(\lambda')$, then $\pi \supseteq \pi'$.*

This definition says that the function $\pi^*(\lambda) \supseteq \pi^*(\lambda')$ regardless of the choice of Bellman optimal policies. As we show next, it guarantees that the Whittle indices are uniquely defined when they exist. As we detail in Appendix A.1, this is not necessarily the case when we consider other interpretations of the classical definition.

2.3 Definition of Whittle index and characterization of indexability

The proposition below shows that Definition 1 implies that Whittle index is well defined and proposes a characterization of any indexable arm, that we will later use to derive our algorithm.

Lemma 1. *In a n -state weakly communicating arm, the following three properties are equivalent:*

- (i) *The arm is indexable.*
- (ii) *For all state $i \in [n]$, there exists a unique penalty λ_i – called the Whittle index of state i – such that if $\pi \in \Pi^*(\lambda)$ is any Bellman optimal policy for the penalty λ , then $\pi_i = 1$ if $\lambda < \lambda_i$ and $\pi_i = 0$ if $\lambda > \lambda_i$.*
- (iii) *There is a non decreasing sequence of penalties $\mu_{\min}^0 := -\infty \leq \mu_{\min}^1 \leq \mu_{\min}^2 \leq \dots \leq \mu_{\min}^n \leq \mu_{\min}^{n+1} := +\infty$ and a sequence of policies $\pi^1 := [n] \supseteq \pi^2 \supseteq \dots \supseteq \pi^{n+1} := \emptyset$ such that:*
 - *If $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$, there exists a unique Bellman optimal policy π^k .*

- If k is such that $\mu_{\min}^{k-1} < \mu_{\min}^k$, then all Bellman optimal policies for the penalty μ_{\min}^{k-1} contain π^k , and π^k contains all Bellman optimal policies for the penalty μ_{\min}^k .

In the above lemma, we use a subscript “min” in the penalties μ_{\min}^k in order to be consistent with the same quantities used in Algorithm 1 and 2. The signification of this “min” is because it will be a minimum of values of the form μ_i^k . We should stress that these quantities (as well as the Whittle index λ_i) can either be finite or infinite. When we say that “a policy π is optimal for the penalty $+\infty$ ”, this means “there exists a penalty $\bar{\lambda}$ such that π is optimal for all $\lambda \geq \bar{\lambda}$ ”. Also, the last part of the lemma implies that π^k is the unique Bellman optimal policy for all penalty $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$.

Proof The lemma is a direct consequence of the definition of indexability.

(i) \Rightarrow (ii) – Assume first that the arm is indexable. Let $i \in [n]$ be a state and let $\lambda_i = \sup\{\lambda : \exists \pi \in \Pi^*(\lambda) \text{ such that } \pi_i = 0\}$. By Definition 1, if π' is a Bellman optimal policy for a penalty $\lambda > \lambda_i$, then $\pi' \subseteq \pi$, which in turn implies that $\pi'_i = 0$. Similarly, if $\lambda < \lambda_i$, then $\pi'_i = 1$. This implies (ii).

(ii) \Rightarrow (iii) – Assume (ii) and let σ^k be the state with the k th smallest index (where ties are broken arbitrarily). Let $\mu_{\min}^k := \lambda_{\sigma^k}$ be the index of the state σ^k and let $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$. By (ii), any Bellman optimal policy for the penalty $\lambda < \lambda_{\sigma^{k-1}}$ contains $\pi^k := [n] \setminus \{\sigma^1, \dots, \sigma^{k-1}\}$. Similarly, π^k contains any Bellman optimal policy for the penalty $\lambda > \lambda_{\sigma^{k-1}}$. This implies that the policy π^k is the unique optimal policy for all $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$.

(iii) \Rightarrow (i) – The property (iii) implies that implies that π^k is the unique Bellman optimal policy for all $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$. \square

3 Condition for indexability and basic algorithm

This section aims at providing a basic algorithm to detect whether an arm is indexable or not and if it is the case, to compute the Whittle index of all states. This algorithm tries to construct a sequence of *unichain* policies $\pi^1 \supseteq \pi^2 \supseteq \dots$ that satisfy the conditions of Lemma 1. We prove the correctness of our algorithm: if it can construct such a sequence, then the problem is indexable and the computed indices are correct. If the algorithm cannot compute such a sequence of policies, this is either because the problem is not indexable, or because the arm is multichain.

3.1 Condition for optimality

In this section, we provide two technical lemmas that we will use in our algorithm. They provide conditions to verify when a given unichain policy is Bellman optimal and if yes, when it is the unique Bellman optimal policy.

Let $\pi \subseteq [n]$ be a unichain policy and fix a penalty λ . By [39, Chapter 8], there exists a gain g^π and a bias vector \mathbf{h}^π such that

$$g^\pi(\lambda) + h_i^\pi(\lambda) = r_i^{\pi_i} - \lambda\pi_i + \sum_{j=1}^n P_{ij}^{\pi_i} h_j^\pi(\lambda). \quad (2)$$

We denote by α_i^π the *active advantage* in state i under policy π , which is the difference between the value in state i of action activate and the one of action rest. It is defined by:

$$\alpha_i^\pi(\lambda) := r_i^1 - r_i^0 - \lambda + \sum_{j=1}^n (P_{ij}^1 - P_{ij}^0) h_j^\pi(\lambda). \quad (3)$$

For a unichain policy, Equation (2) uniquely determines the vector $\mathbf{h}^\pi(\lambda)$ up to an additive constant $c\mathbf{1}$ (see [39, Chapter 8]). Hence the active advantage vector $\boldsymbol{\alpha}^\pi(\lambda)$ is uniquely determined for a unichain policy π . As we will see later, the function $\boldsymbol{\alpha}^\pi(\lambda)$ is affine in λ .

Our algorithm computes the Whittle index by increasing order, by trying to eliminate states one by one. The following lemma shows that to compute the next Bellman optimal policy, one should look at when the active advantage of a state is equal to 0. In this lemma, $\pi \ominus \{i\}$ denotes the symmetric difference between π and $\{i\}$, *i.e.*, $\pi \ominus \{i\} = \pi \setminus \{i\}$ if $i \in \pi$ and $\pi \ominus \{i\} = \pi \cup \{i\}$ if $i \notin \pi$. Also, the active advantage provides necessary and sufficient condition for a unichain policy to be Bellman optimal, and/or to be the unique Bellman optimal policy as shown in the following lemma.

Lemma 2. *In a weakly communicating arm (r, P) , let π be a unichain policy. Then, for any penalty λ :*

- (i) π is Bellman optimal if and only if $\alpha_i^\pi(\lambda) \geq 0$ for all $i \in \pi$ and $\alpha_i^\pi(\lambda) \leq 0$ for all $i \notin \pi$.
- (ii) if π is Bellman optimal and $\alpha_i^\pi(\lambda) = 0$. Then, $\pi \ominus \{i\}$ is also Bellman optimal. Moreover, if $\pi \ominus \{i\}$ is unichain, then $\boldsymbol{\alpha}^\pi(\lambda) = \boldsymbol{\alpha}^{\pi \ominus \{i\}}(\lambda)$.
- (iii) π is the unique Bellman optimal policy if and only if $\alpha_i^\pi(\lambda) > 0$ for all $i \in \pi$ and $\alpha_i^\pi(\lambda) < 0$ for all $i \notin \pi$.

Proof The first two points are direct consequences of the definition of Bellman optimality: a policy is Bellman optimal if and only if its bias \mathbf{h}^π satisfies (1). The definition of the active advantage implies that the active advantage of a Bellman optimal policy should be non-negative for active states $i \in \pi$, and non-positive for passive states $i \notin \pi$. Moreover, if there exists a state i whose active advantage is null, $\alpha_i^\pi(\lambda) = 0$, then the policy $\pi \ominus \{i\}$ is also Bellman optimal.

Points (i) and (ii) also show one direction of the equivalence of (iii): If policy π is the unique Bellman optimal policy, then for all state i , $\alpha_i^\pi(\lambda) \neq 0$. The non-trivial property is the other direction of the equivalence. This is a consequence of Lemma 6 that we prove in Appendix B.1. \square

Note that the main difficulty in proving Lemma 6 is that we do not assume the arm to be unichain: for a unichain arm, the bias of the optimal policy is unique up to a constant vector (see [41] and [39, Section 8.4]). This implies that if π is an optimal policy and $\alpha_i^\pi(\lambda) \neq 0$ for all i , then π is the unique optimal policy. The proof of Lemma 6 that we do in Appendix B.1 does not require the MDP to be unichain, but only the policy π to be Bellman optimal and unichain. In Appendix B.1, we prove that this holds not only for two-action MDPs but also for any MDP with finite state and action spaces.

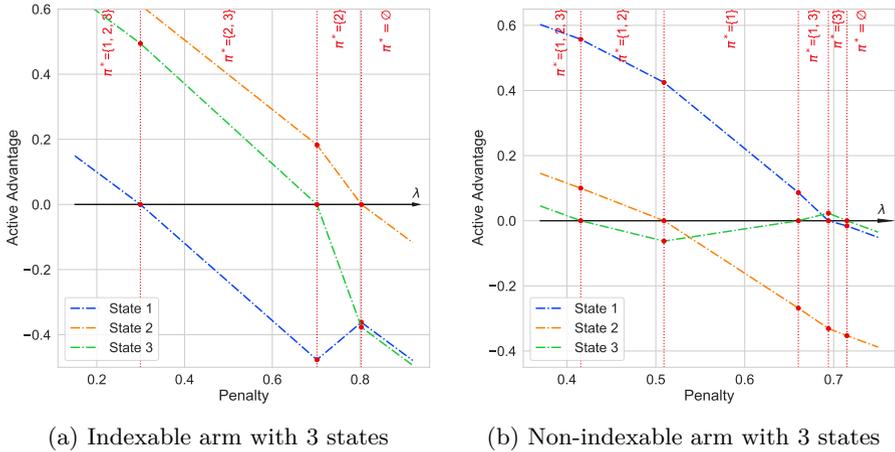


Fig. 2: The active advantage as a function of penalty for two unichain examples, one is indexable (Figure 2a) and the other is not (Figure 2b). The red dots mark where the lines change their slope. The parameters of both examples are provided in Appendix A.3.

To illustrate Lemma 1 and Lemma 2, we consider two three-state arms. For each model, we plot in Figure 2 the active advantage $\alpha_i^{\pi^*(\lambda)}(\lambda)$ as a function of the penalty λ for all states $i \in \{1, 2, 3\}$. By Lemma 2, we know that an optimal policy should activate all states having a positive advantage and rest all states having a negative advantage. Combined with the characterization of Lemma 1, this shows that:

- The model presented in Figure 2a is indexable: the optimal policy is a non-increasing function of λ and the indices are $\lambda_1 \approx 0.3$, $\lambda_2 \approx 0.8$ and $\lambda_3 \approx 0.7$.
- The model presented in Figure 2b is not indexable: the optimal policy $\pi^*(0.7) = \{3\}$ is not included in $\pi^*(0.6) = \{1\}$.

3.2 Overview of the algorithm

Our algorithm computes Whittle index in increasing order by navigating through unichain Bellman optimal policies and using the characterization provided by Lemma 1. It follows the graphical construction given in Figure 2a. It uses the following facts:

- If policy $\pi^1 := [n]$ is unichain, then α^{π^1} is decreasing in λ
- Similarly, if policy $\pi^{n+1} := \emptyset$ is unichain, then $\alpha^{\pi^{n+1}}$ is decreasing in λ
- For an indexable arm, computing the index can be done by a greedy algorithm that constructs a sequence of penalties $\mu_{\min}^1 \leq \mu_{\min}^2 \leq \dots \leq \mu_{\min}^n$ and a sequence of unichain policies $\pi^1 \supseteq \dots \supseteq \pi^n$ by looking at where $\alpha_i^{\pi^k}(\lambda)$ intersects horizontal axis for all $i \in \pi^k$.
- The arm is indexable if and only if for all k such that $\mu_{\min}^{k-1} < \mu_{\min}^k$, the constructed π^k is the largest Bellman optimal policy for the penalty μ_{\min}^k .

In order to compute the Whittle index and test the indexability, our algorithm needs that all policies π^k constructed by the algorithm to be unichain. It does not require the arm to be unichain.

This leads to Algorithm 1, that we write in pseudo-code. This algorithm relies on two subroutines: on Line 7, to compute the next index and on Line 8 to test if a policy is Bellman optimal. We will describe later in the paper how to implement these functions in an efficient manner. Note that in all the paper, we use the superscript k (e.g., π^k, μ^k, σ^k) to refer to the quantities computed at iteration k . We use the subscripts i or j (e.g., $\pi_i, \lambda_i, \mu_i, \pi_j$) to refer to the quantities related to states i or j .

Note that when $\mu_{\min}^k = +\infty$, the quantity $\alpha_i^{\pi^k}(\mu_{\min}^k)$ defined in Line 7 of Algorithm 1 should be understood as $\lim_{\lambda \rightarrow \infty} \alpha_i^{\pi^k}(\lambda) \in \mathbb{R} \cup \{-\infty, +\infty\}$. These limits are well defined because the functions α s are affine in λ .

To illustrate how the algorithm works, we plot in Figure 3 the values computed by the algorithm for the two arms represented in Figure 2. For both models (indexable and non-indexable), the algorithm starts with the policy $[n]$ for which the derivative of the active advantage with respect to λ is -1 for all states. It then computes μ_{\min}^1 which is the potential index of State $\sigma^1 = 1$ for 3a and of State $\sigma^1 = 3$ for 3b. The algorithm then moves to iteration 2 and computes $\mu_{\min}^2 > \mu_{\min}^1$ for both models and observes that for both models $\alpha_{\sigma^1}^{\pi^2}(\mu_{\min}^2) < 0$ for both models. Then the algorithm moves to iteration 3 and computes $\mu_{\min}^3 > \mu_{\min}^2$. There are now two cases:

- For 3a, the algorithm verifies that $\pi^4 := \emptyset$ is unichain (α_i^{\emptyset} is decreasing in λ for all i), terminates, and returns that the arm is indexable.
- For 3b, the algorithm realizes that $\alpha_{\sigma^1}^{\pi^3}(\mu_{\min}^3) > 0$ which shows that this model is not indexable.

Note that for the indexable example of Figure 3a, the active advantage function is not a decreasing function of λ . Hence, this example is neither PCL-indexable (defined in [29, Definition 3]) nor strongly-indexable (defined in [36]). However, this does not prevent our algorithm from working.

Algorithm 1 Given a n -state weakly communicating arm, test indexability and compute Whittle index (if indexable).

```

1: Set  $\pi^1 := [n]$ ,  $\mu_{\min}^0 := -\infty$ 
2: if  $\pi^1$  is multichain then
3:   return the arm is multichain
4: end if
5: for  $k = 1$  to  $n$  do
6:   Compute  $\alpha^{\pi^k}(\lambda)$ 
7:   Let  $\mu_{\min}^k := \inf\{\lambda \geq \mu_{\min}^{k-1} : \exists i \in \pi^k, \alpha_i^{\pi^k}(\lambda) = 0\}$ 
8:   if  $\mu_{\min}^{k-1} < \mu_{\min}^k$  and for some  $i \notin \pi^k$ ,  $\alpha_i^{\pi^k}(\mu_{\min}^k) \geq 0$  then
9:     return the arm is not indexable
10:  end if
11:  if  $\mu_{\min}^k = +\infty$  then
12:    Set  $\lambda_i := +\infty$  for all  $i \in \pi^k$ 
13:    return the arm is indexable and the indices are  $\{\lambda_i\}_{i \in [n]}$ .
14:  end if
15:  Let  $\sigma^k \in \pi^k$  be such that  $\alpha_{\sigma^k}^{\pi^k}(\mu_{\min}^k) = 0$  and  $\lambda_{\sigma^k} = \mu_{\min}^k$ 
16:  Set  $\pi^{k+1} := \pi^k \setminus \{\sigma^k\}$ 
17:  if  $\pi^{k+1}$  is multichain then
18:    return the arm is multichain
19:  end if
20: end for
21: return the arm is indexable and the indices are  $\{\lambda_i\}_{i \in [n]}$ .

```

3.3 Correctness of Algorithm 1

The following result shows that Algorithm 1 is correct.

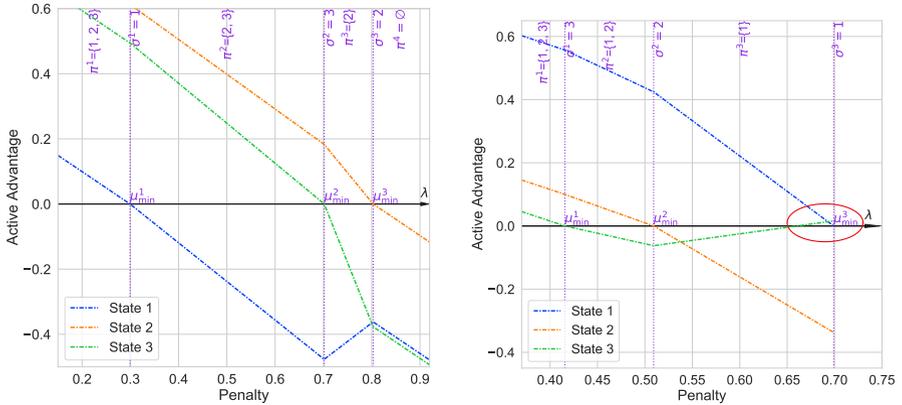
Theorem 1. *Given a n -state weakly communicating arm:*

- (i) *if Algorithm 1 outputs “the arm is indexable and the indices are $\{\lambda_i\}_{i \in [n]}$ ”, then the arm is indexable and each λ_i is the Whittle index of state i ;*
- (ii) *if Algorithm 1 outputs “non-indexable”, then the arm is non-indexable;*
- (iii) *if Algorithm 1 outputs “multichain”, then the arm is multichain.*

A direct consequence of Theorem 1 is that for unichain⁵ arms, Algorithm 1 provides a full characterization of indexability.

Corollary 2. *Given a unichain arm with finite states, Algorithm 1 outputs “the arm is indexable” if and only if it is indexable.*

⁵If a MDP is unichain, then it is weakly communicating.



(a) Indexable arm with 3 states. Policy $\pi^4 = \emptyset$ is unichain which implies that $\alpha_i^{\pi^4}$ is decreasing in λ for each i .

(b) Non-indexable arm with 3 states. The algorithm stops at iteration 3 because $\alpha_3^{\pi^3}(\mu_{\min}^3) > 0$ (the green line at the zone circled with red ellipse).

Fig. 3: The active advantage $\alpha_i^{\pi^k}(\lambda)$ computed by the algorithm, for the two examples of Figure 2.

Note that Algorithm 1 does not require the arm to be unichain to work. The required condition is that the Bellman optimal policies $\{\pi^k\}_{k \geq 1}$ that Algorithm 1 uses are all unichain. In particular, there exist examples of arms that are multichain and indexable and for which the algorithm returns “indexable”. Similarly, there exist examples of arms that are multichain and non-indexable and for which the algorithm returns “non-indexable”. In fact, our algorithm works by exploring solely unichain policies. It does so because (in general) the bias of multichain policy is not unique. The characterization of Bellman optimal policies is much more difficult for multichain models and the notion of indexability becomes more elusive (see Example A1 in Appendix A.1). When Algorithm 1 returns “multichain”, it means that the algorithm is unable to decide whether the arm is indexable or not (but the algorithm knows that the arm is multichain because it has just found a multichain policy).

Proof of Theorem 1 Proof of (i) – We first prove by induction on k that:

- If Algorithm 1 completes iteration $k \geq 1$, then π^k and π^{k+1} are unichain and
- (A) π^k is the unique Bellman optimal policy for all $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$;
 - (B) π^{k+1} is Bellman optimal for μ_{\min}^k and $\alpha^{\pi^{k+1}}(\mu_{\min}^k) = \alpha^{\pi^k}(\mu_{\min}^k)$.

Base case $k = 1$: As we prove later in (17), $\frac{\partial \alpha_i^{\pi^1}}{\partial \lambda} = -1$. So, for each $i \in \pi^1$, $\alpha_i^{\pi^1}(\lambda)$ is decreasing in λ . By definition, μ_{\min}^1 is the smallest λ such that one of the $\alpha_i^{\pi^1}(\lambda) = 0$. Hence, for all $\lambda < \mu_{\min}^1$: $\alpha_i^{\pi^1}(\lambda) > 0$. By Lemma 2, this shows that π^1 is the unique Bellman optimal policy for all $\lambda \in (\mu_{\min}^0, \mu_{\min}^1)$, so (A) is true. Moreover, since π^1 is

Bellman optimal for the penalty μ_{\min}^1 and π^2 is unichain, Lemma 2 implies that π^2 is Bellman optimal for the penalty μ_{\min}^1 and $\alpha^{\pi^2}(\mu_{\min}^1) = \alpha^{\pi^1}(\mu_{\min}^1)$. This shows (B).

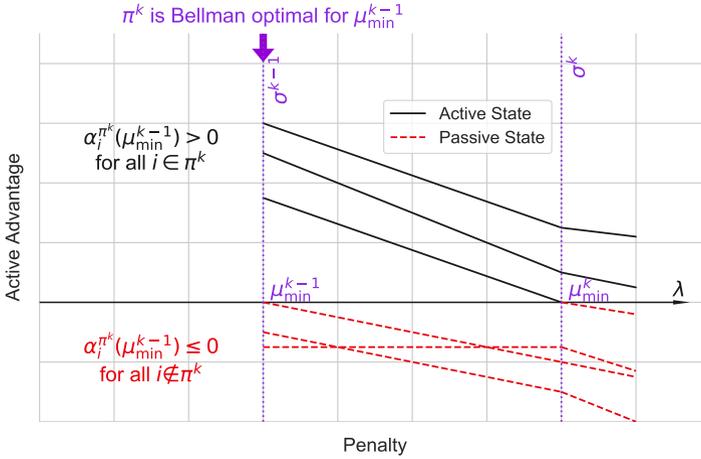


Fig. 4: Illustration of what happens when $\mu_{\min}^{k-1} < \mu_{\min}^k$ and when the test of Line 8 is successful (recall that the function $\alpha_i^{\pi^k}(\lambda)$ is affine in λ). The black lines are the advantage functions $\alpha_i^{\pi^k}(\lambda)$ of active state $i \in \pi^k$. The dashed red lines are the advantage functions $\alpha_i^{\pi^k}(\lambda)$ of passive state $i \notin \pi^k$.

Suppose that the induction is true until iteration $k - 1$ and that the algorithm completes iteration k . If $\mu_{\min}^{k-1} = \mu_{\min}^k$, (A) is trivial and (B) is a direct consequence of the definition of μ_{\min}^k and of Lemma 2. Consider now that $\mu_{\min}^{k-1} < \mu_{\min}^k$ and observe what happens in Figure 4. By the induction hypothesis, π^k is Bellman optimal for the penalty μ_{\min}^{k-1} . Moreover, by definition of μ_{\min}^k , together with $\mu_{\min}^{k-1} < \mu_{\min}^k$, $\alpha_i^{\pi^k}(\mu_{\min}^{k-1}) \neq 0$ for all $i \in \pi^k$. Hence:

$$\alpha_i^{\pi^k}(\mu_{\min}^{k-1}) > 0 \text{ for } i \in \pi^k \quad \text{and} \quad \alpha_i^{\pi^k}(\mu_{\min}^{k-1}) \leq 0 \text{ for } i \notin \pi^k.$$

Finally, thanks to the test on Line 8 of the algorithm, we have:

$$\alpha_i^{\pi^k}(\mu_{\min}^k) \geq 0 \text{ for } i \in \pi^k \quad \text{and} \quad \alpha_i^{\pi^k}(\mu_{\min}^k) < 0 \text{ for } i \notin \pi^k. \quad (4)$$

In consequence, for each $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$, $\alpha_i^{\pi^k}(\lambda) > 0$ for $i \in \pi^k$ and $\alpha_i^{\pi^k}(\lambda) < 0$ for $i \notin \pi^k$. Lemma 2 implies that π^k is the unique Bellman optimal policy for each $\lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k)$. This shows (A). Also, (4) implies that π^k is Bellman optimal for the penalty μ_{\min}^k . Combine this with the fact that π^{k+1} is unichain, Lemma 2 implies that π^{k+1} is Bellman optimal for μ_{\min}^k and $\alpha^{\pi^{k+1}}(\mu_{\min}^k) = \alpha^{\pi^k}(\mu_{\min}^k)$. This shows (B). So, the induction is also true for iteration k .

This shows that the induction property is true for all $k \in [n]$. In particular, when $\pi^{n+1} := \emptyset$ is unichain, $\alpha_i^{\pi^{n+1}}$ is decreasing in λ as we prove later in (18) that $\frac{\partial \alpha_i^{\pi^{n+1}}}{\partial \lambda} = -1$. Combine this with the fact that $\alpha_i^{\pi^{n+1}}(\mu_{\min}^n) = \alpha_i^{\pi^n}(\mu_{\min}^n) \leq 0$

for all i , Lemma 2 implies that π^{n+1} is the unique Bellman optimal policy for $\lambda \in (\mu_{\min}^n, \mu_{\min}^{n+1})$. To sum up, there are two cases for which the algorithm outputs that the arm is indexable:

1. if the algorithm goes until the end of iteration n , then the sequence of values $\{\mu_{\min}^k\}_{k \in [n+1]}$ and of policies $\{\pi^k\}_{k \in [n+1]}$ satisfies the conditions of Lemma 1(iii) and the arm is indexable.
2. if the algorithm stops at iteration k because $\mu_{\min}^k = +\infty$, then one can set $\mu_{\min}^{k+1} := \dots := \mu_{\min}^{n+1} := +\infty$ and define a sequence of policies $\pi^{k+1} \supseteq \dots \supseteq \pi^{n+1} := \emptyset$ by eliminating all states of π^k in an arbitrary order. These sequences satisfies the conditions of Lemma 1(iii) and the arm is indexable.

Proof of (ii) – Our algorithm outputs non-indexable if there exists an iteration k and a state $j \notin \pi^k$, such that $\alpha_j^{\pi^k}(\mu_{\min}^k) \geq 0$ when $\mu_{\min}^{k-1} < \mu_{\min}^k$. We know that π^k is Bellman optimal for μ_{\min}^{k-1} , otherwise the algorithm would have stopped before. Assume that:

$$\text{All Bellman optimal policies for any } \lambda \in (\mu_{\min}^{k-1}, \mu_{\min}^k] \text{ are included in } \pi^k. \quad (5)$$

We will see that this leads to a contradiction. We distinguish two possibilities:

1. π^k is Bellman optimal for the penalty μ_{\min}^k – This implies that $\alpha_i^{\pi^k}(\mu_{\min}^k) \leq 0$ for all $i \notin \pi^k$ which, together with $\alpha_j^{\pi^k}(\mu_{\min}^k) \geq 0$, implies that $\alpha_j^{\pi^k}(\mu_{\min}^k) = 0$. By Lemma 2, this would imply that $\pi^k \cup \{j\}$ is Bellman optimal for μ_{\min}^k . This is in contradiction with (5).
2. π^k is not Bellman optimal for μ_{\min}^k – In this case, we denote by $\tilde{\lambda}$ as the smallest penalty $\lambda \in [\mu_{\min}^{k-1}, \mu_{\min}^k]$ such that there exists $\pi \subsetneq \pi^k$ that is Bellman optimal for $\tilde{\lambda}$ (it exists because π^k is not Bellman optimal for μ_{\min}^k and we assumed (5)). By definition of $\tilde{\lambda}$, π and π^k are both Bellman optimal for the penalty $\tilde{\lambda}$. Let $i \in \pi^k \setminus \pi$. By Lemma 2, this implies that $\alpha_i^{\pi^k}(\tilde{\lambda}) = 0$. The problem is that by definition, μ_{\min}^k is the smallest penalty λ for which there exists $i \in \pi^k$ such that $\alpha_i^{\pi^k}(\lambda) = 0$. This implies that $\tilde{\lambda} = \mu_{\min}^k$ which in turn implies that π^k is optimal for μ_{\min}^k . This leads to a contradiction.

This shows that neither case 1 nor 2 are possible. So, (5) cannot be true. In consequence, the negation of (5) is true: there exists $\lambda > \mu_{\min}^{k-1}$ and $\pi \not\subseteq \pi^k$ such that π is Bellman optimal for λ . This contradicts Definition 1 and therefore implies that the arm is not indexable.

Proof of (iii) – if our algorithm outputs multichain, then the arm is multichain. This is straightforward based on the definition of multichain MDP. \square

We should note that by Line 15, it is possible to have $\mu_{\min}^k = \mu_{\min}^{k-1}$. This happens when several states have the same value of Whittle index. This is not problematic because we are sure that $\sigma^k \neq \sigma^{k-1}$ by Line 16.

In the proof of Theorem 1(i), we showed that when policy $[n]$ is unichain, the function $\alpha_i^{\pi^1}(\lambda)$ is decreasing in λ which implies that it crosses the line 0 at some finite value μ_i^1 . This implies that for an indexable arm, if $[n]$ is unichain then the Whittle index are all strictly larger than $-\infty$. A symmetric argument shows that if policy \emptyset is unichain, then all Whittle index are strictly smaller than $+\infty$. This implies the following result.

Corollary 3. *Given a unichain arm with n states, if the arm is indexable, then the indices of the n states are finite: $\lambda_i \notin \{-\infty, +\infty\}$ for all $i \in [n]$.*

This is not necessarily true for multichain arms (see the discussion in Appendix A.2.)

3.4 Naive implementation of Algorithm 1 (in $O(n^4)$)

For a given penalty λ , we consider a policy π that is Bellman optimal and unichain. Recall that $g^*(\lambda)$ is the maximal gain, and $\mathbf{h}^\pi(\lambda) \in \mathbb{R}^n$ is a solution (2). We consider $\mathbf{h}^\pi(\lambda)$ such that $h_1^\pi(\lambda) = 0$. Recall from (2) that for all $i \in [n]$:

$$g^*(\lambda) + h_i^\pi(\lambda) = r_i^{\pi_i} - \lambda\pi_i + \sum_{j=1}^n P_{ij}^{\pi_i} h_j^\pi(\lambda). \quad (6)$$

The above system is a system of $n+1$ linear equations with $n+1$ variables (the additional equation begins with $h_1^\pi(\lambda) = 0$). As π is unichain, the maximal gain $g^*(\lambda)$ and $\mathbf{h}^\pi(\lambda)$ are uniquely determined by the system of linear equations (6), together with the condition that $h_1^\pi(\lambda) = 0$. Note that in (6) the sum is for $j = 1$ to n . Since $h_1^\pi(\lambda) = 0$, it can be transformed into a sum from $j = 2$ to n .

Let us define the vector $\mathbf{v}^\pi(\lambda) := [g^*(\lambda) \ h_2^\pi(\lambda) \ \dots \ h_n^\pi(\lambda)]^\top$ which is similar to the vector $\mathbf{h}^\pi(\lambda)$ in which we replaced $h_1^\pi(\lambda)$ by $g^*(\lambda)$. We can write Equation (6) under a matrix form as:

$$\mathbf{A}^\pi \mathbf{v}^\pi(\lambda) = \mathbf{r}^\pi - \lambda\boldsymbol{\pi}, \quad (7)$$

where \mathbf{r}^π is the reward vector under π : $\mathbf{r}^\pi := [r_1^{\pi_1} \ \dots \ r_n^{\pi_n}]^\top$, $\boldsymbol{\pi} := [\pi_1 \ \dots \ \pi_n]^\top$, and \mathbf{A}^π is the following square matrix:

$$\mathbf{A}^\pi := \begin{bmatrix} 1 & & & \\ 1 & 1 & & \\ & & \ddots & \\ 1 & & & 1 \end{bmatrix} - \begin{bmatrix} 0 & P_{12}^{\pi_1} & \dots & P_{1n}^{\pi_1} \\ 0 & P_{22}^{\pi_2} & \dots & P_{2n}^{\pi_2} \\ \vdots & & \ddots & \vdots \\ 0 & P_{n2}^{\pi_n} & \dots & P_{nn}^{\pi_n} \end{bmatrix} = \begin{bmatrix} 1 & -P_{12}^{\pi_1} & \dots & -P_{1n}^{\pi_1} \\ 1 & 1 - P_{22}^{\pi_2} & \dots & -P_{2n}^{\pi_2} \\ \vdots & & \ddots & \vdots \\ 1 & -P_{n2}^{\pi_n} & \dots & 1 - P_{nn}^{\pi_n} \end{bmatrix} \quad (8)$$

As we show in Lemma 7, the matrix \mathbf{A}^π is invertible if and only if policy π is unichain. In consequence, \mathbf{v}^π is an affine function of λ :

$$\mathbf{v}^\pi(\lambda) = (\mathbf{A}^\pi)^{-1}(\mathbf{r}^\pi - \lambda\boldsymbol{\pi}) = (\mathbf{A}^\pi)^{-1}\mathbf{r}^\pi - \lambda(\mathbf{A}^\pi)^{-1}\boldsymbol{\pi} \quad (9)$$

For a state i , let $\delta_i := r_i^1 - r_i^0$, $\Delta_{i1} := 0$ and $\Delta_{ij} := P_{ij}^1 - P_{ij}^0$ for $j \in \{2, \dots, n\}$. By definition of the advantage function of (3), we have:

$$\boldsymbol{\alpha}^\pi(\lambda) = \boldsymbol{\delta} - \lambda\mathbf{1} + \boldsymbol{\Delta}\mathbf{v}^\pi(\lambda).$$

For each active state i , we want to find the smallest penalty $\mu_i^k \geq \mu_{\min}^{k-1}$ such that $\alpha_i^{\pi^k}(\mu_i^k) = 0$. Suppose that π^{k-1} and π^k are unichain and π^{k-1} is Bellman optimal for μ_{\min}^{k-1} . By Lemma 2, $\alpha^{\pi^k}(\mu_{\min}^{k-1}) = \alpha^{\pi^{k-1}}(\mu_{\min}^{k-1})$. Let $\mathbf{d}^{\pi^k} := -(\mathbf{A}^{\pi^k})^{-1}\boldsymbol{\pi}^k$. By (9), $\alpha^{\pi^k}(\lambda)$ is a linear function of λ whose derivative is $(-\mathbf{1} + \Delta \mathbf{d}^{\pi^k})$. In particular, $\alpha^{\pi^k}(\lambda) = \alpha^{\pi^{k-1}}(\mu_{\min}^{k-1}) - (\lambda - \mu_{\min}^{k-1})(\mathbf{1} - \Delta \mathbf{d}^{\pi^k})$. Thus, $\alpha_i^{\pi^k}(\lambda) = 0$ if and only if

$$\alpha_i^{\pi^{k-1}}(\mu_{\min}^{k-1}) = (\lambda - \mu_{\min}^{k-1})\left(1 - \sum_{j=2}^n \Delta_{ij} d_j^{\pi^k}\right). \quad (10)$$

Recall that $\alpha_i^{\pi^{k-1}}(\mu_{\min}^{k-1})$ is non-negative for active state $i \in \pi^k$. The value μ_i^k is the smallest $\lambda \geq \mu_{\min}^{k-1}$ that satisfies Equation (10). There are two cases:

1. if $\alpha_i^{\pi^{k-1}}(\mu_{\min}^{k-1}) = 0$, then $\mu_i^k := \mu_{\min}^{k-1}$;
2. if $\alpha_i^{\pi^{k-1}}(\mu_{\min}^{k-1}) > 0$ and
 - (a) if $1 - \sum_{j=2}^n \Delta_{ij} d_j^{\pi^k} > 0$, then

$$\mu_i^k := \mu_{\min}^{k-1} + \frac{\alpha_i^{\pi^{k-1}}(\mu_{\min}^{k-1})}{1 - \sum_{j=2}^n \Delta_{ij} d_j^{\pi^k}}; \quad (11)$$

- (b) if $1 - \sum_{j=2}^n \Delta_{ij} d_j^{\pi^k} \leq 0$, then $\mu_i^k := +\infty$.

This shows that, for a given k , computing μ_{\min}^k of Line 7 can be done in $O(n^3)$: A first part in $O(n^3)$ to compute the inverse of matrix \mathbf{A}^{π} and to compute \mathbf{d}^{π} , plus some smaller order terms to compute the solutions of (10). Similarly, the test in Line 8 of Algorithm 1 can also be implemented in $O(n^3)$ by using $\alpha^{\pi^k}(\mu_{\min}^k) = \alpha^{\pi^{k-1}}(\mu_{\min}^{k-1}) - (\mu_{\min}^k - \mu_{\min}^{k-1})(\mathbf{1} - \Delta \mathbf{d}^{\pi^k})$ with the convention when $\mu_{\min}^k = +\infty$ that $+\infty \times 0 = 0$ and $+\infty \times x = \text{sign}(x)\infty$ for any $x \neq 0$. This leads to an overall complexity of $O(n^4)$ for Algorithm 1 that contains n loops each having a $O(n^3)$ complexity. If at some iteration k the matrix \mathbf{A}^{π^k} is not invertible, then Lemma 7 implies that π^k is multichain. In consequence, the algorithm outputs multichain and stops. We integrate this in the newer version of our algorithm below.

4 The $(2/3)n^3 + o(n^3)$ algorithm

This section describes a way to implement Algorithm 1 efficiently using $O(n^3)$ operations. The main idea is to use the Sherman-Morrison formula to compute in $O(n^2)$ the active advantage vector $\alpha^{\pi^k}(\lambda)$ associated to π^k from the one associated to π^{k-1} . This leads to a $O(n^3)$ algorithm. Once this main idea is in place, we show how to avoid unnecessary computations to obtain an algorithm that performs $(2/3)n^3 + o(n^3)$ arithmetic operations.

4.1 Additional notations

In order to obtain a more efficient and compact algorithm, for an iteration k and a state i , we define $y_i^k := \sum_{j=2}^n \Delta_{ij} d_j^{\pi^k}$ and $z_i^k := \alpha_i^{\pi^k}(\mu_{\min}^k)$, where $d_j^{\pi^k}$, Δ_{ij} and $\alpha_i^{\pi^k}$ are as in (11). Equation (11) can be rewritten as

$$\mu_i^k = \mu_{\min}^{k-1} + \frac{z_i^{k-1}}{1 - y_i^k}. \quad (12)$$

The above equation can be used to compute μ_i^k and μ_{\min}^k easily from y_i^k and z_i^{k-1} . Indeed, from the previous section, we have $\alpha_i^{\pi^k}(\mu_{\min}^k) = \alpha_i^{\pi^{k-1}}(\mu_{\min}^{k-1}) - (\mu_{\min}^k - \mu_{\min}^{k-1})(1 - \sum_{j=2}^n \Delta_{ij} d_j^{\pi^k})$ which translates into

$$z_i^k = z_i^{k-1} - (\mu_{\min}^k - \mu_{\min}^{k-1})(1 - y_i^k). \quad (13)$$

This shows that the critical values to compute are the variables y_i^k . In the remainder of this section, we show that the quantity y_i^k can be computed efficiently by a recursive formula.

4.2 Application of the Sherman-Morrison formula

To compute $y_i^{k+1} := \sum_{j=2}^n \Delta_{ij} d_j^{\pi^{k+1}}$, we need to compute the quantities $d^{\pi^{k+1}} := -(\mathbf{A}^{\pi^{k+1}})^{-1} \boldsymbol{\pi}^{k+1}$. This requires the inverse of $\mathbf{A}^{\pi^{k+1}}$. By definition of π^k , two policies π^k and π^{k+1} differ by exactly one state: $\boldsymbol{\pi}^{k+1} = \boldsymbol{\pi}^k - \mathbf{e}_{\sigma^k}$ where \mathbf{e}_j denotes the column vector with a 1 in j th coordinate and 0's elsewhere. Also by definition of \mathbf{A}^π in (8), the two matrices \mathbf{A}^{π^k} and $\mathbf{A}^{\pi^{k+1}}$ differ only at the row σ^k :

$$\mathbf{A}^{\pi^{k+1}} = \mathbf{A}^{\pi^k} + \mathbf{e}_{\sigma^k} \boldsymbol{\Delta}_{\sigma^k} \quad (14)$$

where $\boldsymbol{\Delta}_{\sigma^k}$ is a row vector defined as in the previous section.

One can efficiently compute the inverse of matrix $\mathbf{A}^{\pi^{k+1}}$ from the one of \mathbf{A}^{π^k} by using the Sherman-Morrison formula, which says that if $\mathbf{A} \in \mathbb{R}^{n \times n}$ is an invertible square matrix and $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$ are two column vectors, then the matrix $\mathbf{A} + \mathbf{p}\mathbf{q}^\top$ is invertible if and only if $1 + \mathbf{q}^\top \mathbf{A}^{-1} \mathbf{p} \neq 0$ and if $\mathbf{A} + \mathbf{p}\mathbf{q}^\top$ is invertible, then:

$$(\mathbf{A} + \mathbf{p}\mathbf{q}^\top)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{p}\mathbf{q}^\top \mathbf{A}^{-1}}{1 + \mathbf{q}^\top \mathbf{A}^{-1} \mathbf{p}}.$$

Let $X_{ij}^k := \boldsymbol{\Delta}_i (\mathbf{A}^{\pi^k})^{-1} \mathbf{e}_j$. Following (14), we can apply the Sherman-Morrison formula with matrix \mathbf{A}^{π^k} , and vectors $\mathbf{p} = \mathbf{e}_{\sigma^k}$ and $\mathbf{q}^\top = \boldsymbol{\Delta}_{\sigma^k}$. After some

simplification, we get:

$$\begin{aligned} X_{ij}^{k+1} &:= \Delta_i(\mathbf{A}^{\pi^{k+1}})^{-1} \mathbf{e}_j = \Delta_i(\mathbf{A}^{\pi^k} + \mathbf{e}_{\sigma^k} \Delta_{\sigma^k})^{-1} \mathbf{e}_j \\ &= X_{ij}^k - \frac{X_{i\sigma^k}^k}{1 + X_{\sigma^k\sigma^k}^k} X_{\sigma^k j}^k \end{aligned} \quad (15)$$

$$\text{In particular, } X_{i\sigma^k}^{k+1} = \frac{X_{i\sigma^k}^k}{1 + X_{\sigma^k\sigma^k}^k}.$$

Before computing X_{ij}^{k+1} , we need to verify that π^{k+1} is unichain. With the help of Lemma 7 and the Sherman-Morrison formula, this can be done easily: π^{k+1} is unichain if and only if $1 + X_{\sigma^k\sigma^k}^k \neq 0$.

For $y_i^{k+1} := \Delta_i \mathbf{d}^{\pi^{k+1}} = -\Delta_i(\mathbf{A}^{\pi^{k+1}})^{-1} \boldsymbol{\pi}^{k+1}$, we use $\boldsymbol{\pi}^{k+1} = \boldsymbol{\pi}^k - \mathbf{e}_{\sigma^k}$ and apply the Sherman-Morrison formula to get:

$$\begin{aligned} y_i^{k+1} &= -\Delta_i(\mathbf{A}^{\pi^k} + \mathbf{e}_{\sigma^k} \Delta_{\sigma^k})^{-1} (\boldsymbol{\pi}^k - \mathbf{e}_{\sigma^k}) \\ &= -\Delta_i(\mathbf{A}^{\pi^k})^{-1} (\boldsymbol{\pi}^k - \mathbf{e}_{\sigma^k}) + \frac{\Delta_i(\mathbf{A}^{\pi^k})^{-1} \mathbf{e}_{\sigma^k} \Delta_{\sigma^k} (\mathbf{A}^{\pi^k})^{-1}}{1 + \Delta_{\sigma^k} (\mathbf{A}^{\pi^k})^{-1} \mathbf{e}_{\sigma^k}} (\boldsymbol{\pi}^k - \mathbf{e}_{\sigma^k}) \\ &= y_i^k + X_{i\sigma^k}^k + \frac{X_{i\sigma^k}^k (-y_{\sigma^k}^k - X_{\sigma^k\sigma^k}^k)}{1 + X_{\sigma^k\sigma^k}^k} \\ &= y_i^k + \frac{X_{i\sigma^k}^k (1 - y_{\sigma^k}^k)}{1 + X_{\sigma^k\sigma^k}^k} = y_i^k + (1 - y_{\sigma^k}^k) X_{i\sigma^k}^{k+1} \end{aligned} \quad (16)$$

The above formula indicate how to compute \mathbf{y}^{k+1} from \mathbf{y}^k . To complete this analysis, let us show that $\mathbf{y}^1 = \mathbf{0}$. For a given policy π , the vector \mathbf{d}^π satisfies the same equation as Equation (7) but replacing $r_i^{\pi_i} - \lambda \pi_i$ by $-\pi_i$. This implies that for $\boldsymbol{\pi} = \boldsymbol{\pi}^1 = [1 \dots 1]^\top$, one has $\mathbf{d}^\pi = [-1 \ 0 \dots 0]^\top$ as d_1^π is the long-term average reward of a Markov reward process whose reward is negative one in all states and d_2^π, \dots, d_n^π is the bias of this process. This shows that for all i , one has $y_i^1 := \sum_{j=2}^n \Delta_{ij} d_j^\pi = 0$. Moreover, by (9), one has

$$\begin{aligned} \boldsymbol{\alpha}^{\pi^1}(\lambda) &= \boldsymbol{\delta} - \lambda \mathbf{1} + \Delta(\mathbf{A}^{\pi^1})^{-1} \mathbf{r}^{\pi^1} - \underbrace{\lambda \Delta(\mathbf{A}^{\pi^1})^{-1} \boldsymbol{\pi}^1}_{=:\mathbf{y}^1} \\ &= \boldsymbol{\delta} - \lambda \mathbf{1} + \mathbf{X}^1 \mathbf{r}^{\pi^1}. \end{aligned} \quad (17)$$

Finally, for $\boldsymbol{\pi}^{n+1} = [0 \dots 0]^\top$, one has

$$\boldsymbol{\alpha}^{\pi^{n+1}}(\lambda) = \boldsymbol{\delta} - \lambda \mathbf{1} + \Delta(\mathbf{A}^{\pi^{n+1}})^{-1} \mathbf{r}^{\pi^{n+1}}. \quad (18)$$

4.3 Detailed algorithm

Equation (12) shows how to compute μ_i^k from the values of y_i^k and z_i^{k-1} while (16), (15) and (13) show how to compute the values of \mathbf{y} , \mathbf{z} and \mathbf{X} recursively in k . In order to compute μ_{\min}^k and σ^k , one only needs to compute the values μ_i^k for $i \in \pi^k$. Once $\mu_{\min}^k = \min_{i \in \pi^k} \mu_i^k$ is determined, if $\mu_{\min}^{k-1} < \mu_{\min}^k$, then Line 8 of Algorithm 1 can be performed, based on (13), by checking if $z_i^k \geq 0$ for some $i \in [n] \setminus \pi^k$.

Algorithm 2 Given a n -state weakly communicating arm, test indexability and compute Whittle indices (if indexable).

```

1: Set  $\pi^1 := [n]$ ,  $k_0 = 1$ ,  $\Delta_{i1} := 0$ ,  $\Delta_{ij} := P_{ij}^1 - P_{ij}^0, \forall i \in [n], j \in \{2, \dots, n\}$ ,
    $\mathbf{y}^1 = \mathbf{0}$  and  $\mathbf{A}^{\pi^1}$  is defined by (8).
2: if  $\pi^1$  is multichain then
3:   return the arm is multichain
4: end if
5: Set  $\mathbf{X}^1 := \Delta(\mathbf{A}^{\pi^1})^{-1}$ 
6: Set  $\boldsymbol{\mu}^1 := \mathbf{r}^1 - \mathbf{r}^0 + \mathbf{X}^1 \mathbf{r}^1$ 
7: Let  $\sigma^1 := \arg \min_{i \in \pi^1} \mu_i^1$ ,  $\lambda_{\sigma^1} = \mu_{\sigma^1}^1 := \mu_{\sigma^1}^1$ , and  $\pi^2 := \pi^1 \setminus \{\sigma^1\}$ 
8: Set  $\mathbf{z}^1 = \boldsymbol{\mu}^1 - \mu_{\min}^1 \mathbf{1}$ 
9: for  $k = 2$  to  $n$  do
10:   Update  $\mathbf{X}(k-1)$  ▷ Here we call Subroutine 3 or Subroutine 4
11:   Set  $\mathbf{y}^k = \mathbf{y}^{k-1} + (1 - y_{\sigma^{k-1}}^{k-1}) \mathbf{X}_{:\sigma^{k-1}}^k$ 
12:   for  $i \in \pi^k$  do
13:     Set  $\mu_i^k := \begin{cases} \mu_{\min}^{k-1}, & \text{if } z_i^{k-1} = 0 \\ \mu_{\min}^{k-1} + \frac{z_i^{k-1}}{1 - y_i^k}, & \text{if } z_i^{k-1} > 0 \text{ and } 1 - y_i^k > 0 \\ +\infty, & \text{otherwise} \end{cases}$ 
14:   end for
15:   Let  $\sigma^k := \arg \min_{i \in \pi^k} \mu_i^k$  and  $\mu_{\min}^k := \mu_{\sigma^k}^k$ 
16:   Set  $\mathbf{z}^k = \mathbf{z}^{k-1} - (\mu_{\min}^k - \mu_{\min}^{k-1})(\mathbf{1} - \mathbf{y}^k)$ 
17:   if  $\mu_{\min}^{k-1} < \mu_{\min}^k$  and  $z_i^k \geq 0$  for some  $i \in [n] \setminus \pi^k$  then
18:     return the arm is not indexable
19:   end if
20:   if  $\mu_{\min}^k = +\infty$  then
21:     Set  $\lambda_i := +\infty$  for all  $i \in \pi^k$ 
22:     return the arm is indexable and the indices are  $\{\lambda_i\}_{i \in [n]}$ .
23:   end if
24:   Set  $\lambda_{\sigma^k} = \mu_{\min}^k$  and  $\pi^{k+1} := \pi^k \setminus \{\sigma^k\}$ 
25: end for
26: if  $\pi^{n+1}$  is multichain then
27:   return the arm is multichain
28: end if
29: return the arm is indexable and the indices are  $\{\lambda_i\}_{i \in [n]}$ .

```

Subroutine 3 Update_X(k)

```

1: for  $\ell = 1$  to  $k - 1$  do
2:   for  $i \in [n]$  do ▷ or  $i \in \pi^{\ell+1}$  if we do not test indexability.
3:      $X_{i\sigma^k}^{\ell+1} = X_{i\sigma^k}^\ell - X_{i\sigma^\ell}^{\ell+1} X_{\sigma^\ell \sigma^k}^\ell$ 
4:   end for
5: end for
6: if  $1 + X_{\sigma^k \sigma^k}^k = 0$  then
7:   return the arm is multichain
8: end if
9: for  $i \in [n]$  do ▷ or  $i \in \pi^{k+1}$  if we do not test indexability.
10:   $X_{i\sigma^k}^{k+1} = \frac{X_{i\sigma^k}^k}{1 + X_{\sigma^k \sigma^k}^k}$ 
11: end for

```

This leads to Algorithm 2 that can be decomposed as follows:

1. In Lines 1 to 8, we initialize the various variables. The main complexity of this part is to compute the matrix \mathbf{X}^1 , which is equivalent to solving the linear system $\mathbf{X}\mathbf{A}^{\pi^1} = \mathbf{\Delta}$. It can be done by inverting the matrix \mathbf{A}^{π^1} and multiplying this by the matrix $\mathbf{\Delta}$. This can be done in a subcubic complexity by using for instance Strassen's algorithm [38].
2. We then enter the main loop:
 - We update the vectors $\boldsymbol{\mu}$, \mathbf{z} by using Equations (12) and (13), and test indexability. This costs $O(n)$ operations per iteration, thus $O(n^2)$ in total.
 - We update the vector \mathbf{X} according to (15). The “naive” way to do so is to use Subroutine 3. At iteration k this costs $2kn$ arithmetic operations if we test indexability, and $2 \sum_{l=1}^k (n-l)$ if we do not test indexability. The total complexity of computing \mathbf{X} is $n^3 + O(n^2)$ arithmetic operations if we test indexability and $(2/3)n^3 + O(n^2)$ if we do not. A detailed study of the arithmetic complexity is provided in Appendix C, where we also provide details on how to efficiently implement the algorithm, including how to optimize the cost of memory access.
 - In Line 11, we update the vector \mathbf{y} by using Equation (16), which costs $O(n)$ per iteration.
3. Testing if π^{n+1} is unichain can be done in $O(n^2)$ by Tarjan's strongly connected component algorithm.

Hence, the total complexity of this algorithm is $n^3 + o(n^3)$ if we test indexability and $(2/3)n^3 + o(n^3)$ if we do not test indexability. Without testing the indexability, our algorithm has the same main complexity as [29]. However, the algorithm of [29] computes Whittle index only for an arm that is PCL-indexable. Hence we can claim that our algorithm is the first algorithm that computes Whittle index with cubic complexity for all indexable restless bandits. It is also the first algorithm with cubic complexity that detects

non-indexability of restless bandits. Note that the ties are broken arbitrarily for Lines 7 and 15. Moreover, as mentioned above, the weakly communicating assumption can be verified in $O(n^2)$ by using Tarjan's strongly connected component algorithm for transition matrix $(\mathbf{P}^0 + \mathbf{P}^1)/2$.

Remark: Equivalently, instead of calling Subroutine 3 at Line 10, one could do the following update for all $j \in \pi^{k+1}$ and $i \in [n]$ (or $i \in \pi^{k+1}$ if we do not test indexability):

$$X_{ij}^{k+1} = X_{ij}^k - \frac{X_{i\sigma^k}^k}{1 + X_{\sigma^k\sigma^k}^k} X_{\sigma^k j}^k. \quad (19)$$

This iterative update is very close to the one used in [29, 33] for discounted restless bandit. This results in an algorithm that has the same total complexity as Subroutine 3 (of $n^3 + O(n^2)$ or $(2/3)n^3 + O(n^2)$ with or without the indexability test) because both algorithms will have computed the same values of X_{ij}^k . The reason to use Subroutine 3 is that, as we will see in the next section, not all values of X_{ij}^k are needed at iteration k : in particular, for $\ell < k$, the computation of the value $X_{i\sigma^{\ell-1}}^\ell$ has no interest per say and is only useful because it allows to recursively compute $X_{i\sigma^{\ell-1}}^k$. In the section below, we show how to reduce the cost by avoiding the computation of $X_{i\sigma^{\ell-1}}^\ell$ when ℓ is much smaller than k . We comment more on the differences with [29, 33] in Appendix E and in particular we explain why our approach can be tuned into a subcubic algorithm while (19) cannot.

5 The subcubic algorithm

5.1 Main idea: recomputing \mathbf{X}^k from \mathbf{A}^{π^k} periodically

The main computational burden of Algorithm 2 is concentrated on two lines: on Line 5 where we compute \mathbf{X}^1 by solving a linear system, and on Line 10 where we compute the column vector $\mathbf{X}_{\cdot,\sigma^{k-1}}^k$ from $\mathbf{X}_{\cdot,\sigma^{k-1}}^1$. The remainder of the code runs in $O(n^2)$ operations and is therefore negligible for large matrices. In fact, the computation of \mathbf{X}^1 is done by solving a linear system, which can be computed by using a subcubic algorithm (like Strassen [38]). In this section, we show how to optimize our algorithm by reducing the complexity of the update_X() function, at the price of recomputing the full matrix \mathbf{X}^{k_0} from $\mathbf{A}^{\pi^{k_0}}$ periodically.

At iteration k of Algorithm 2, the quantities $X_{i\sigma^{k-1}}^k$ are used at Line 11 to obtain the values y_i^k . The matrix \mathbf{X}^k is defined as $\mathbf{X}^k := \mathbf{\Delta}(\mathbf{A}^{\pi^k})^{-1}$. It also satisfies Equation (15), that is, for an iteration ℓ and states i and j , we have:

$$X_{ij}^{\ell+1} = X_{ij}^\ell - X_{i\sigma^\ell}^{\ell+1} X_{\sigma^\ell j}^\ell. \quad (20)$$

The way Subroutine 3 is implemented is to initialize $\mathbf{X}^1 := \Delta(\mathbf{A}^{\pi^1})^{-1}$ and then use (20) recursively to compute the column vector $\mathbf{X}_{:\sigma^{k-1}}^k$ from $\mathbf{X}_{:\sigma^{k-1}}^1$ at each iteration.

Here, we propose an alternative formulation which consists in recomputing the whole matrix $\mathbf{X}^k := \Delta(\mathbf{A}^{\pi^k})^{-1}$ every K iterations. In the meantime, we use (20) to compute the values $X_{i\sigma^{k-1}}^\ell$ for $\ell \in \{k_0 + 1, \dots, k\}$ where k_0 is the iteration at which we recomputed the whole matrix $\mathbf{X}^{k_0} := \Delta(\mathbf{A}^{\pi^{k_0}})^{-1}$. This can be implemented by replacing the call to Subroutine 3 at Line 10 with a call to Subroutine 4.

Subroutine 4 Update_X_FMM(k)

```

1: if  $k$  is an iteration at which we recompute the whole  $\mathbf{X}^{k+1}$  then
2:   Set  $k_0 := k + 1$ 
3:   if  $\mathbf{A}^{\pi^{k_0}}$  is not invertible then
4:     return the arm is multichain
5:   end if
6:   Set  $\mathbf{X}^{k_0} := \Delta(\mathbf{A}^{\pi^{k_0}})^{-1}$ 
7: else
8:   for  $\ell = k_0$  to  $k - 1$  do
9:     for  $i \in [n]$  do ▷ or  $i \in \pi^{\ell+1}$  if we do not test indexability.
10:       $X_{i\sigma^\ell}^{\ell+1} = X_{i\sigma^\ell}^\ell - X_{i\sigma^\ell}^{\ell+1} X_{\sigma^\ell \sigma^\ell}^\ell$ 
11:    end for
12:  end for
13:  if  $1 + X_{\sigma^k \sigma^k}^k = 0$  then
14:    return the arm is multichain
15:  end if
16:  for  $i \in [n]$  do ▷ or  $i \in \pi^{k+1}$  if we do not test indexability.
17:     $X_{i\sigma^k}^{k+1} = \frac{X_{i\sigma^k}^k}{1 + X_{\sigma^k \sigma^k}^k}$ 
18:  end for
19: end if

```

To see why this can be more efficient, we illustrate in Figure 5 the pairs (ℓ, σ^{k-1}) for which we compute the vector $\mathbf{X}_{:\sigma^{k-1}}^\ell$, either for Subroutine 3 or Subroutine 4. In each case, a vertical blue line indicates that we recompute the whole matrix \mathbf{X}^k by solving a linear system. The gray zone corresponds to the values (ℓ, σ^{k-1}) for which we compute $\mathbf{X}_{:\sigma^{k-1}}^\ell$ using Equation (20) and the red squares represent the vector $\mathbf{X}_{:\sigma^{k-1}}^k$ used at Line 11 of Algorithm 2. For Subroutine 3, we do one matrix inversion at the beginning and then compute for all (ℓ, σ^{k-1}) with $\ell \leq k$ because the red square at value (k, σ^{k-1}) is computed starting from the vertical blue line at value $(1, \sigma^{k-1})$. For Subroutine 4, we do $\lfloor n/K \rfloor$ (here $\lfloor n/K \rfloor = 4$) full recomputation of \mathbf{X}^k , which

correspond to the vertical blue lines. We gain in terms of operations because the surface of the gray zone to compute is divided by $\lfloor n/K \rfloor$.

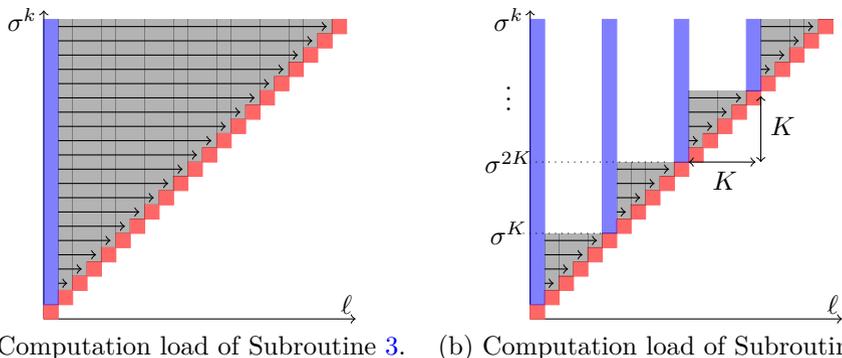


Fig. 5: Illustration of the improvement proposed by replacing Subroutine 3 with Subroutine 4 when we do not test the indexability. For Subroutine 4, we solve more linear systems (each vertical blue line corresponds to solving a linear system) but we reduce the gray zone to compute. A linear arrow corresponds to the internal loop of Line 8 of Subroutine 4.

Note that the y -axis of Figure 5 is ordered by increasing value of σ^k (and not by increasing value of k). The value of σ^k is computed at iteration k but unknown before iteration k . This explains why in Subroutine 4, when we recompute the matrix (X_{ij}^k) at an iteration k (vertical blue lines in Figure 5(b)), we recompute it for all i, j and not just $i, j \in \{\sigma^k, \dots, \sigma^K\}$ (which are the only values that we will use): Indeed, $\sigma^k, \dots, \sigma^K$ are unknown at iteration k .

5.2 A subcubic algorithm for Whittle index

We now assume to have access to a subcubic matrix multiplication algorithm that satisfies the following property:

(FMM) There exists an algorithm to multiply a matrix of size $n \times n$ by a matrix of size⁶ $n \times n^\gamma$ that runs in $O(n^{\omega(\gamma)})$, where $\omega : [0, 1] \rightarrow [2, 3]$ is a non-decreasing function.

Going back to Algorithm 2 where Line 10 is Subroutine 4, we now assume that we recompute the whole matrix \mathbf{X}^k every $O(n^\gamma)$ iterations. The new algorithm has a subcubic complexity:

Theorem 4. *Given a n -state weakly communicating arm, Algorithm 2 with Subroutine 4 checks indexability and computes Whittle (and Gittins) index in time at least $\Omega(n^{2.5})$ and at most $O(n^{2.5286})$ when choosing $\gamma = 0.5286$.*

⁶We write n^γ which is possibly non-integer. For the sake of simplicity, we write n^γ and $n^{1-\gamma}$ but they should be understood as $\lfloor n^\gamma \rfloor$ and $\lfloor n^{1-\gamma} \rfloor$ respectively.

We believe that Theorem 4 is the first theoretical result that shows that Whittle index can be computed in subcubic time. As we show in Section 7, this algorithm can be directly extended to discounted index. As a byproduct, we also obtain the first subcubic algorithm to compute Gittins index.

Proof The algorithm starts by computing \mathbf{X}^1 which can be done in $O(n^{\omega(1)})$. Then, there are $n^{1-\gamma}$ times that we do:

1. We fill the “gray” mini matrices by using (20). This amounts to three for loops of size n (for i), n^γ (for k) and n^γ (for ℓ). Hence, each small gray matrix costs $O(n^{1+2\gamma})$.
2. At the end of a cycle, we recompute the full inverse by updating $(\mathbf{A}^{\pi^k})^{-1}$ from $(\mathbf{A}^{\pi^{k-n^\gamma}})^{-1}$. As we show in Lemma 3 (stated below), this can be done in $O(n^{\omega(\gamma)})$.

This implies that the algorithm has a complexity:

$$O(n^{\omega(\gamma)}) + n^{1-\gamma} \left(O(n^{1+2\gamma}) + O(n^{\omega(\gamma)}) \right) = O(n^{\max\{2+\gamma, 1-\gamma+\omega(\gamma)\}}).$$

To compute the optimal value of γ minimizing this expression requires the knowledge of the function $\omega(\gamma)$ which is not known. The current state of the art only gives a lower bound ($\omega(\gamma) \geq 2$) and an upper bound described in [42].

It is shown in [42] that $\gamma = 0.5286$ is the smallest currently known value of γ for which $\omega(\gamma) < 1 + 2\gamma$. This implies that the complexity is at most $O(n^{2.5286})$.

As for the lower bound, $\omega(\gamma) \geq 2$ implies that the complexity of the algorithm is at least $\Omega(n^{2.5})$. \square

In the next lemma, \mathbf{B} plays the role of \mathbf{A}^{π^k} and \mathbf{A} the role of $\mathbf{A}^{\pi^{k-n^\gamma}}$. Note that as required in the lemma, exactly n^γ rows and columns are changed between the two.

Lemma 3. *Assume (FMM). Let \mathbf{A} be a square matrix whose inverse \mathbf{A}^{-1} has already been computed, and let \mathbf{B} be an invertible square matrix such that $\mathbf{A} - \mathbf{B}$ is of rank smaller than n^γ . Then, it is possible to compute the inverse of \mathbf{B} in $O(n^{\omega(\gamma)})$.*

Proof The matrix \mathbf{B} can be written as $\mathbf{B} = \mathbf{A} + \mathbf{UCV}$ where \mathbf{U} is a $n \times n^\gamma$ matrix, \mathbf{C} is $n^\gamma \times n^\gamma$ and \mathbf{V} is $n^\gamma \times n$. The Sherman–Morrison–Woodbury formula [43] states that

$$\mathbf{B}^{-1} = (\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}.$$

This shows that \mathbf{B}^{-1} can be computed by:

- Computing $\mathbf{D} := \mathbf{A}^{-1}\mathbf{U}$ and $\mathbf{E} := \mathbf{VA}^{-1}$: this takes $O(n^{\omega(\gamma)})$.
- Computing $\mathbf{F} := (\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}$: as this is the inversion of a $n^\gamma \times n^\gamma$ matrix, it can be done in $O(n^{\gamma\omega(1)})$ where $\gamma\omega(1) \leq \omega(\gamma)$.
- Computing $\mathbf{G} := \mathbf{DF}$ and then \mathbf{GE} : this again takes $O(n^{\omega(\gamma)})$.

Hence, computing \mathbf{B}^{-1} can be done in $O(n^{\omega(\gamma)})$ operations for the inversion and all multiplications plus an additional $O(n^2)$ term for the subtraction and the addition. As $\omega(\gamma) \geq 2$, this concludes the proof of the lemma. \square

5.3 The subcubic algorithm in practice

The complexity of $O(n^{2.5286})$ given in Theorem 4 is mainly of theoretical interest. The value $\gamma = 0.5286$ is obtained by using the best upper bound on $\omega(\gamma)$ known today which is based on the Coppersmith-Winograd algorithm and its variants. The Coppersmith-Winograd algorithm (or its variants) are, however, known as a *galactic algorithm*: the hidden constant in the $O()$ is so large that their runtime is prohibitive for any reasonable value of n . Hence, the existence of these algorithms is of theoretical interest but has limited applicability.

This does not discard the practical improvement provided by Subroutine 4 which is based on the mere fact that multiplying two matrices (or inverting a matrix) is faster than three nested loops even for matrices of moderate size. To verify this, we launched a detailed profiling of the code of Algorithm 2 with the non-optimized Subroutine 3. It shows that for a problem of dimensions 5000, the update of Line 10 takes more than 90% of the computation time, the initialization of \mathbf{X}^1 on Line 5 takes about 5% of the time and the rest of the code takes less than 1% of the running time.

Now, if inverting the full matrix takes about 5% of the execution time, and updating the gray zone takes 95%, then by doing 5 updates, one can hope to obtain an algorithm whose running time is roughly $5 \times 5 + 95/5 \approx 43\%$ the one of the original implementation. As we observe in Section 6, this is close to the gain that we obtain in practice. A general way to choose the best number of updates is used in the numerical section. It is based on the following reasoning. For large matrices (say $n \geq 10^3$), the fastest implementations of matrix multiplication and inversion are based on Strassen's algorithm [44, 45]. As we report in Appendix D, the time to solve a linear system of size n by using the default installation of `scipy` seems to run in $O(n^{2.8})$. By replacing the function $\omega(\gamma)$ used in Theorem 4 by a more practical bound ($\omega(\gamma) = 2.8$), the best value for γ becomes $\gamma = 0.9$. This indicates that our algorithm can be implemented in $O(n^{2.9})$ by doing $O(n^{0.1})$ recomputation of \mathbf{X}^k from \mathbf{A}^{π^k} . Note that even for very large values of n (like $n = 15000$), $n^{0.1}$ remains quite small, e.g., $15000^{0.1} \approx 2.6$. In practice, we observe that updating $\text{int}(2n^{0.1})$ times (the notation $\text{int}(x)$ indicates that it is rounded to the closest integer) gives the best performance among all algorithms, as reported in the next section.

6 Numerical experiments

In complement to our theoretical analysis, we developed a python package that implements Algorithm 2 and gives the choice of using the variant of Subroutine 3 or of Subroutine 4 to do the “update_X()” function. This package relies on three python libraries: `scipy` and `numpy` for matrix operations, and `numba` to compile the python code. To facilitate its usage, this package can be installed by using `pip install markovianbandit-pkg`.

All experiments were conducted on a laptop (Macbook Pro 2020) with an Intel Core i9 CPU at 2.3 GHz with 16GB of Memory using Python 3.6.9 ::

Anaconda custom (64-bit) under macOS Big Sur version 11.6.2. The version of the packages are `scipy` version 1.5.4, `numpy` version 1.19.5 and `numba` version 0.53.1. The code of all experiments is available at <https://gitlab.inria.fr/markovianbandit/efficient-whittle-index-computation>.

Finally, in all of our experiments, the generated arms are all unichain.

6.1 Time to compute Whittle indices

To test the implementation of our algorithm, we randomly generate restless bandit arms with n states where $n \in \{100, 1000, \dots, 15000\}$. In each case, both transition matrices are uniform probabilistic matrices: for each row of each matrix, we generate n i.i.d. entries following the exponential distribution and divide the row by its sum. This means that all matrices are dense. We use dense matrix since it is the worst case for computational interest. By running our algorithms on sparse matrix, we would expect to have faster running time. Note that all tested matrices are indexable. This is coherent with [32, 46] that report that for uniform matrices, the probability of finding a non-indexable example decreases very rapidly with the dimension n . Finally, reward vectors were generated from random Uniform[0,1) entries.

Table 1: Running time (in seconds) of the variants of Algorithm 2

n	$O(n^3)$ algorithm (Subroutine 3)		Subcubic algorithm (Subroutine 4)	
	(a) With index. test	(b) Without test	(c) With index. test	(d) Without test
100	0.006	0.004	0.007	0.005
1000	0.2	0.2	0.2	0.2
2000	1.9	1.2	1.1	1.0
3000	7.2	5.0	3.2	2.6
4000	17	12	8	6
5000	34	25	16	12
6000	60	43	27	20
7000	95	68	42	33
8000	142	99	63	49
9000	199	141	92	70
10000	275	191	122	95
11000	361	257	164	133
12000	471	339	225	190
13000	620	428	286	243
14000	790	553	408	317
15000	965	685	501	403

We record the runtime of the different variants of our algorithm and report the results in Table 1. Note that these results present the whole execution time of the algorithm, including the initialization phase in which \mathbf{X}^1 is computed. For each value of n , we run Algorithm 2 with four variants:

- The first two columns correspond to the $O(n^3)$ algorithm (that uses Subroutine 3 for Line 10), either (a) with the indexability test, or (b) without the indexability test.

- The last two columns correspond to the subcubic algorithm (that uses Subroutine 4 for Line 10 with $\text{int}(2n^{0.1})$ updates), either (c) with the indexability test, or (d) without the indexability test.

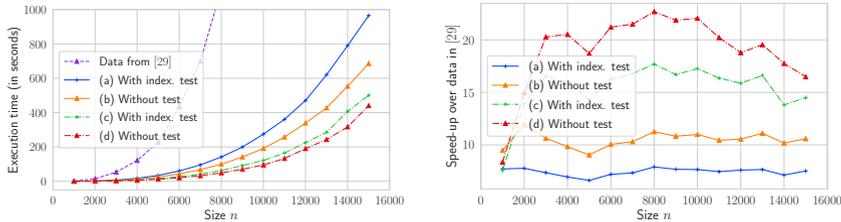
Our numbers show that our algorithm can compute the Whittle index in less than one second for $n = 1000$ states and slightly less than 7 minutes for $n = 15000$ states with variant (d). As expected, not doing the indexability test does improve the performance compared to doing the indexability test (here by a factor approximately 1/3 for Subroutine 3 and 1/5 for Subroutine 4). More importantly, this table shows that the time when using the subcubic variant, Subroutine 4, diminishes the computation time by about 40% to 50% compared to when using Subroutine 3. Note that for $n = 100$, using the Subroutine 3 is slightly faster than using the Subroutine 4 (while both takes only a few milliseconds). This indicates that the subcubic algorithm becomes interesting when n is large enough (say $n \geq 2000$).

To give a visual idea of how the various variants of the algorithms compare, we plot in Figure 6a the runtime of the four variants along with the runtime reported [29], that we believe was the fastest algorithm to compute Whittle index so far for discounted restless bandit. We plot in Figure 6b the runtime of each variant divided by the runtime from [29]. For large n , our best implementation is about 10 to 20 times faster than the one of [29]. For instance, for $n = 15000$, our implementation takes about 7 minutes to compute the index (or 9 minutes when checking indexability) whereas [29] reports 2 hours. Recall that the algorithm of [29] cannot check indexability on the fly. In our implementation, not testing indexability reduces the computation time of 15 – 20% for Subroutine 4 or 25 – 30% for Subroutine 3 compared to the version that tests indexability.

It should be clear that the comparison between our implementation and the numbers reported in [29] has its limits. First, there is no public implementation of the algorithm of [29] and we did not re-implement it. Hence, the computation was not done on the same machine, and not implemented in the same programming language. The characteristic of the machine described in [29] seems comparable if not better to ours, and our implementation is in python/numpy whereas the one of [29] is in matlab. Second, it is written in [29] that the reported numbers do not “count the initialization stage of computing the initial tableau”. There is no mention of the cost of the initialization stage in [29]. On the contrary, the numbers that we report in this paper include all initialization stages, which puts our running time at disadvantage.

6.2 Statistics of indexable problems

To the best of our knowledge, our algorithm provides the first indexability test that scales well with the dimension n . We used this to answer a very natural question: given a randomly generated arm, how likely is it to be indexable? This question was partially answered in [32] that shows that when generating dense arms, the probability of generating a non-indexable arm is close to 10^{-n} for $n \in \{3, \dots, 7\}$. This suggests that most arms are indexable. Below,



(a) Runtime of our implementations as a function of the state size n .

(b) Speedup of each implementation over the data from [29] as a function of the state size n .

Fig. 6: Numerical result over 7 simulations taking around 21 hours in total: in each simulation, we run the algorithm over randomly generated RBs with the state size ranging over $\{1000, \dots, 15000\}$. We plot the average runtime over 7 simulations. The solid lines represent the result of Subroutine 3 and the dashed-dot lines represent the one of Subroutine 4. The marker “+” indicates that algorithms test indexability and the triangles indicates that algorithms do not test indexability.

we answer two questions: what happens for larger values of n , and more importantly, what happens when the state transition matrices are not dense?

To answer these questions, we consider randomly generated arms where the matrices \mathbf{P}^0 and \mathbf{P}^1 are b -diagonal matrices with b non-null diagonals. In particular, $b = 3$ corresponds to tridiagonal matrices, $b = 5$ corresponds to pentadiagonal matrices and $b = 7$ corresponds to septadiagonal matrices. We also compare with the classical case of dense matrices (which corresponds to $b = 2n - 1$). For each model, the entries are generated from the exponential distribution for each row and we divide the row by the sum of generated entries for this row. We vary n from 3 to 50 and for each case, we generate 100000 arms. We report in Table 2 the number of indexable arms for each case. Note that pentadiagonal matrices are dense matrices for $n = 3$ and septadiagonal matrices are dense matrices for $n = 4$ and do not make sense for $n = 3$, which is why no numbers are reported.

Table 2: Number of indexable problems among 100 000 randomly generated problems.

Problem size n	Tridiagonal	5-diagonal	7-diagonal	Dense
3	98 731	–	–	99 883
4	95 067	99 655	–	99 931
5	89 198	99 309	99 902	99 969
10	54 129	90 377	98 914	100 000
30	7 094	29 699	66 143	100 000
50	1 823	9 332	32 069	100 000

Based on these results, we can assert that dense models are essentially always indexable which conforms with the data reported in [32]. The situation is, however, radically different for sparse models: the number of indexable problems decreases quickly with the number of states. For instance, there are only 1 823 indexable 50-state problems among 100 000 generated tridiagonal models (*i.e.* around 1.8% are indexable). Note that a tridiagonal model is a birth-death Markov chain which is frequently used for queueing systems. Hence, it is very important to check the indexability of the problem because it is not a prevalent property for sparse models. This also calls for new efficient policies in restless multi-arm bandit problems that are not based on Whittle indices.

7 Extension to the discounted case

The model described in Section 2 corresponds to the definition of Whittle index for a time-average criterion, for which Whittle index is known to be asymptotically optimal [6] for restless bandits. Yet, Whittle index can also be defined for the discounted case [29, 33]. Notably, the discounted Whittle index simplifies into Gittins index when the bandit is rested (*i.e.*, when $\mathbf{P}^0 = \mathbf{I}$ and $\mathbf{r}^0 = \mathbf{0}$). In this section, we show how to adapt our algorithm to the discounted case. As a by product, we obtain the first subcubic algorithm to compute Gittins index rested bandit.

7.1 Discounted Whittle index

We now consider a λ -penalized MDP in which the instantaneous reward received at time $t \geq 0$ is discounted by a factor β^t , where $\beta \in (0, 1)$ is called the discount factor: when executing action a in state i at time $t \geq 0$, the decision maker earns a reward $\beta^t(r_i^a - \lambda a)$. For a given policy π , we denote by $u_i^\pi(\lambda)$ the expected sum of discounted rewards earned by the decision maker when the MDP starts in state i at time 0. The vector $\mathbf{u}^\pi(\lambda) = [u_1^\pi(\lambda) \dots u_n^\pi(\lambda)]^\top$ is called the value function of the policy π . From [39], it satisfies Bellman's equation, that is, for all state i we have:

$$u_i^\pi(\lambda) = r_i^{\pi_i} - \lambda \pi_i + \beta \sum_{j=1}^n P_{ij}^{\pi_i} u_j^\pi(\lambda). \quad (21)$$

The above equation is a linear equation, whose solution is unique because $\beta < 1$. It is given by:

$$\mathbf{u}^\pi(\lambda) = (\mathbf{I} - \beta \mathbf{P}^\pi)^{-1}(\mathbf{r}^\pi - \lambda \boldsymbol{\pi}). \quad (22)$$

For a given penalty λ and a state i , we denote by $u_i^*(\lambda) := \max_{\pi} u_i^\pi(\lambda)$ be the optimal value of state i . A policy π is optimal for the penalty λ , *i.e.*, $\pi \in \Pi^*(\lambda)$, if for all state $i \in [n]$, $u_i^*(\lambda) = u_i^\pi(\lambda)$. By [39] such a policy exists, $|\Pi^*(\lambda)| > 0$. As mentioned in Section 2.2.3, the distinction between Bellman optimal and

gain optimal disappears in discounted MDP in which we are concerned with maximizing the value function. Similarly to the time-average criterion studied before, a β -discounted RB is called indexable if for all penalty $\lambda < \lambda'$, all $\pi \in \Pi^*(\lambda)$ and $\pi' \in \Pi^*(\lambda')$, one has $\pi \supseteq \pi'$.

7.2 Analogy between the time-average and the discounted versions

Let π be a policy and $i \in \pi$ be an active state. Similarly to average reward model studied before, the advantage of action activate over action rest in state i right before following policy π is given by, $\alpha_i^\pi(\lambda) := r_i^1 - r_i^0 - \lambda + \beta \sum_{j=1}^n (P_{ij}^1 - P_{ij}^0) u_j^\pi(\lambda)$. Then $\alpha_i^\pi(\lambda) = 0$ if and only if $\lambda = \delta_i + \sum_{j=1}^n \tilde{\Delta}_{ij} u_j^\pi(\lambda)$. where δ_i is defined as for the average reward model and $\tilde{\Delta}$ is such that for all⁷ states $i, j \in [n]$: $\tilde{\Delta}_{ij} := \beta(P_{ij}^1 - P_{ij}^0)$.

To finish the derivation of the algorithm, one should note that the value function $\mathbf{u}(\lambda)$ plays the same role as the vector $\mathbf{v}(\lambda)$ defined for the average reward model. In particular, the definition of \mathbf{u} in Equation (22) is the analogue of the definition of \mathbf{v} in (9) up to the replacement of the matrix \mathbf{A}^π in (9) by the matrix $\mathbf{I} - \beta\mathbf{P}^\pi$. This means that similarly to \mathbf{v} , the value function $\mathbf{u}(\lambda)$ is affine in λ .

Hence, following the same development in Section 4, we can modify Algorithm 2 to compute the discounted Whittle index by modifying only the initialization phase:

$$\tilde{\Delta} := \beta(\mathbf{P}^1 - \mathbf{P}^0) \text{ and } \mathbf{X}^1 := \tilde{\Delta}(\mathbf{I} - \beta\mathbf{P}^{\pi^1})^{-1}.$$

Note that we still have $\mathbf{y}^1 = \mathbf{0}$ because $(\mathbf{I} - \beta\mathbf{P}^{\pi^1})^{-1}\boldsymbol{\pi}^1 = \frac{1}{1-\beta}\mathbf{1}$ (value function in a β -discounted Markov reward process with reward equals to 1 in all states) and $\tilde{\Delta}\mathbf{1} = \mathbf{0}$. Also, if Subroutine 4 is used, Line 6 should be changed to $\mathbf{X}^{k_0} := \tilde{\Delta}(\mathbf{I} - \beta\mathbf{P}^{\pi^{k_0}})^{-1}$. Last but not least, in the discounted case, we no longer needed the arm to be weakly communicating nor the optimal policies to be unichain because the matrix $(\mathbf{I} - \beta\mathbf{P}^\pi)$ is invertible for any policy π as long as $\beta < 1$ (from Perron-Frobenius' Theorem).

7.3 Gittins index

The notion of “restless” bandit comes from the fact that even when the action “rest” is taken, the Markov chain can still change state and generate rewards. When this is not the case (*i.e.*, when $\mathbf{P}^0 = \mathbf{I}$ and $\mathbf{r}^0 = \mathbf{0}$), an arm is no longer restless and is simply called a Markovian bandit (or a *rested* Markovian bandit if one wants to emphasize that it is not restless).

In a discounted rested bandit, the notion of Whittle index coincides with the notion of Gittins index (In fact, Whittle index was first introduced as a

⁷Note that the definition of $\tilde{\Delta}$ is identical to Δ except when $j = 1$, for which $\Delta_{i1} := 0$ but $\tilde{\Delta}_{i1} := \beta(P_{i1}^1 - P_{i1}^0)$.

generalization of Gittins index to restless bandit in [3]). In such a case, there is no notion of indexability: a discounted rested bandit is always indexable. Its index can be computed by Algorithm 2 without testing indexability. The best known algorithms to compute Gittins index runs in $(2/3)n^3 + O(n^2)$ [28]. When using fast multiplication, our algorithm computes Gittins index in $O(n^{2.5286})$ which makes it the first algorithm to compute Gittins index in subcubic time.

Note that when $\mathbf{P}^1 = \mathbf{I}$, it is possible to compute $\mathbf{X}^1 = \hat{\Delta}/(1-\beta)$ without having to solve a linear system which means that, when $\mathbf{P}^1 = \mathbf{I}$, our algorithm with the variant Subroutine 3 has complexity $(2/3)n^3 + O(n^2)$. This shows that our cubic algorithm can compute the Gittins index also in $(2/3)n^3 + O(n^2)$ instead of $(2/3)n^3 + o(n^3)$.

8 Conclusion

In this paper, we propose a univocal definition of indexability and present an algorithm that is efficient for detecting the non-indexability and computing the Whittle index of all indexable finite-state restless bandits whose arms are all unichain. For weaker assumption where the arms are just weakly communicating, this algorithm can still test the indexability and compute the Whittle index of some arms that are multichain and remains efficient if it is able to do so. Our algorithm is based on the efficient application of the Sherman-Morrison formula. This is a unified algorithm that works for both discounted and non-discounted restless bandits, and can be used for Gittins index computation. We present a first version of our algorithm that runs in $n^3 + o(n^3)$ arithmetic operations (or in $(2/3)n^3 + o(n^3)$ if we do not test the indexability). So, we conclude that Whittle index is not harder to compute than Gittins index. The second version of our algorithm uses the fastest matrix multiplication method and has a complexity of $O(n^{2.5286})$. This makes it the first subcubic algorithm to compute Whittle index or Gittins index. We provide numerical simulations that show that our algorithm is very efficient in practice: it can test indexability and compute the index of a n -state restless bandit arm in less than one second for $n = 1000$, and in a few minutes for $n = 15000$. These numbers are provided for dense matrices. One might expect to have more efficient algorithms if the arm has a sparse structure. We leave this question for future work.

Acknowledgments. This work is supported by the ANR project REFINO (ANR-19-CE23-0015).

Appendix A Examples and counterexamples

In this section, we provide a few examples to illustrate the ambiguities in the classical definition of indexability, and to illustrate what can happen for some multichain arms. We also provide the parameters of arms presented in Figure 2.

A.1 Discussion on the definition of indexability

The classical notion of indexability used in the literature is to say that the optimal policy $\pi^*(\lambda)$ should be non-increasing in λ . Yet, we argue that this definition has two problems:

1. What does “increasing” mean when $\pi^*(\lambda)$ is not unique? Two possibilities are: for all penalties $\lambda < \lambda'$:
 - (\exists) there exist policies π, π' with π optimal for λ and π' optimal for λ' such that $\pi \supseteq \pi'$;
 - (\forall) for all policies π, π' such that π is optimal for λ and π' is optimal for λ' , we have $\pi \supseteq \pi'$.
2. What notion of “optimality” should be used? Two possibilities are:
 - (GO) “optimal” means gain optimal.
 - (BO) “optimal” means Bellman optimal.

The most problematic choice is the notion of increasingness: Interpretation (\exists) is more permissive: For instance, consider an arm with two states and assume that the optimal policy is $\{1, 2\}$ for $\lambda < 0$ and is either $\{1\}$ or \emptyset for $\lambda > 0$. Interpretation (\exists) says that the arm is indexable while interpretation (\forall) says that this arm is not indexable. If the arm is indexable, what should the index of state 1 be? Any choice $\lambda_1 \in [0, +\infty]$ seems reasonable. Saying that the arm is not indexable clarifies the situation. This is why we choose interpretation (\forall) in our paper.

In our paper, we choose the combination (\forall -BO) because we believe that, for a problem that has transient state, the notion of Bellman optimality is more meaningful than the notion of gain optimality. Also, our combination (\forall -BO) allows for more problems to be indexable compared to (\forall -GO) and is easier to characterize.

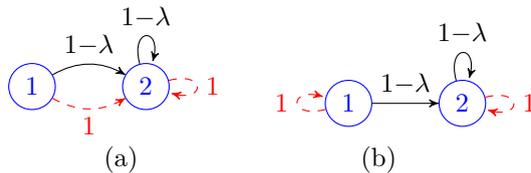


Fig. A1: Ambiguous examples: All transitions are deterministic and labels on transitions indicate rewards. Solid black arrows correspond to the action “activate” and dashed red arrows to the action “rest”.

We illustrate these different definitions in Figure A1. For example (a):

- The gain optimal policies are $\{1, 2\}$ and $\{2\}$ for $\lambda < 0$, and $\{1\}$ and \emptyset for $\lambda > 0$: According to the interpretation (\exists), the problem should be indexable but the index for state 1 is unclear. According to the interpretation (\forall), the problem should not be indexable.

- The Bellman optimal policy is $\{1, 2\}$ for $\lambda < 0$, and \emptyset for $\lambda > 0$. According to our definition, (\forall -BO), the problem is indexable and the indices are $\lambda_1 = \lambda_2 = 0$.

For example (b), the Bellman optimal and gain optimal policies are identical and equal to the gain optimal policies of example (a). Hence, example (b) is not indexable according to our definition. The output of our algorithm for this problem is "multi-chain".

Note that if the distinction between (BO) and (GO) disappears for discounted problems, the distinction between (\forall) and (\exists) remains.

A.2 Multichain arms and infinite indices

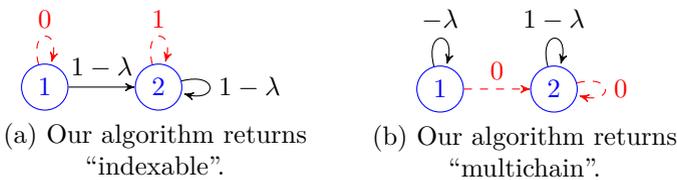


Fig. A2: Example of an indexable multichain problem with infinite Whittle index. Transitions are deterministic and labels on edges indicate rewards (for the λ -penalized arm). Solid black transitions correspond to action "activate" and dashed red to the action "rest".

Consider the two examples of Figure A2. The examples are multichain because policy \emptyset has two irreducible classes for example (a) and policy $\{1, 2\}$ has two irreducible classes for example (b). These two problems are indexable:

- For (a), the Bellman optimal policy for $\lambda < 0$ is $\{1, 2\}$ and $\{1\}$ for $\lambda > 0$. The indices are $\lambda_2 = 0$ and $\lambda_1 = +\infty$.
- For (b), the Bellman optimal policy is $\{2\}$ for $\lambda < 0$ and $\{1, 2\}$ for $\lambda > 0$. The indices are $\lambda_1 = 0$ and $\lambda_2 = -\infty$.

For the first example, our algorithm returns the correct indices because the constructed policies are $\pi^1 := \{1, 2\} \supseteq \pi_2 := \{1\}$ and they are both unichain. For the second example, our algorithm will start with the policy $\{1, 2\}$ and will stop by saying that this example is multichain.

A.3 Parameters for the example of Figure 2

The numerical data of the indexable arm presented in Figure 2a is

$$P^0 = \begin{bmatrix} 0.363 & 0.503 & 0.134 \\ 0.082 & 0.754 & 0.164 \\ 0.246 & 0.029 & 0.724 \end{bmatrix} \quad P^1 = \begin{bmatrix} 0.172 & 0.175 & 0.653 \\ 0.055 & 0.931 & 0.014 \\ 0.155 & 0.627 & 0.218 \end{bmatrix} \quad r^1 = \begin{bmatrix} 0.441 \\ 0.803 \\ 0.426 \end{bmatrix} \quad r^0 = \mathbf{0}$$

The numerical data of the non-indexable arm presented in Figure 2b is

$$\mathbf{P}^0 = \begin{bmatrix} 0.005 & 0.793 & 0.202 \\ 0.027 & 0.558 & 0.415 \\ 0.736 & 0.249 & 0.015 \end{bmatrix} \mathbf{P}^1 = \begin{bmatrix} 0.718 & 0.254 & 0.028 \\ 0.347 & 0.097 & 0.556 \\ 0.015 & 0.956 & 0.029 \end{bmatrix} \mathbf{r}^1 = \begin{bmatrix} 0.699 \\ 0.362 \\ 0.715 \end{bmatrix} \mathbf{r}^0 = \mathbf{0}$$

Appendix B Technical lemmas

B.1 Unicity of Bellman optimal policy

B.1.1 Definition and notation

We are given a MDP $\langle \mathcal{S}, \mathcal{A}, r, P \rangle$ with finite state and action spaces. As shown in [39, Chapter 9], for such a MDP, the optimal gain \mathbf{g}^* is a vector that satisfies the *multichain optimality equations*: for each $i \in \mathcal{S}$:

$$\max_{a \in \mathcal{A}} \left(\sum_{j \in \mathcal{S}} P_{ij}^a g_j^* - g_i^* \right) = 0 \quad (\text{B1})$$

$$\max_{a \in \mathcal{A}} \left(r_i^a - g_i^* + \sum_{j \in \mathcal{S}} P_{ij}^a h_j - h_i \right) = 0. \quad (\text{B2})$$

This system uniquely determines the optimal gain that we denote by \mathbf{g}^* . However, vector \mathbf{h} is not uniquely determined by the system. In the following, we denote by H the set of bias vector \mathbf{h} such that $(\mathbf{g}^*, \mathbf{h})$ is a solution of the optimality equations (B1)–(B2). A policy π is Bellman optimal if there exists a bias $\mathbf{h} \in H$ such that policy π attains the maximum (B2), *i.e.*, for all $i \in \mathcal{S}$:

$$\pi_i \in \arg \max_{a \in \mathcal{A}} \left(r_i^a - g_i + \sum_{j \in \mathcal{S}} P_{ij}^a h_j - h_i \right) = \arg \max_{a \in \mathcal{A}} \left(r_i^a + \sum_{j \in \mathcal{S}} P_{ij}^a h_j \right).$$

For a given policy $\pi : \mathcal{S} \mapsto \mathcal{A}$, we denote by \mathbf{r}^π and \mathbf{P}^π the reward vector and state transition matrix under policy π : $r_i^\pi = r_i^{\pi_i}$ and $P_{ij}^\pi = P_{ij}^{\pi_i}$. Let $\mathbf{g}^\pi, \mathbf{h}^\pi \in \mathbb{R}^{|\mathcal{S}|}$ be a solution of the following system:

$$\mathbf{g}^\pi - \mathbf{P}^\pi \mathbf{g}^\pi = \mathbf{0} \quad (\text{B3})$$

$$\mathbf{r}^\pi - \mathbf{g}^\pi + \mathbf{P}^\pi \mathbf{h}^\pi - \mathbf{h}^\pi = \mathbf{0}. \quad (\text{B4})$$

The vector \mathbf{g}^π is uniquely determined by this system of equations and is called the long-term average reward or gain of policy π . The vector \mathbf{h}^π is unique up to an element of the null space of $(\mathbf{I} - \mathbf{P}^\pi)$. Such a vector \mathbf{h}^π is called a bias of policy π .

If $(\mathbf{g}^\pi, \mathbf{h}^\pi)$ is a solution of (B3) and (B4), then the advantage of action a over the action π_i when the MDP is in state i is given by:

$$\begin{aligned} B_i^a(\mathbf{h}^\pi) &:= r_i^a + \sum_{j \in \mathcal{S}} P_{ij}^a h_j^\pi - g_i^\pi - h_i^\pi \\ &= r_i^a - r_i^{\pi_i} + \sum_{j \in \mathcal{S}} (P_{ij}^a - P_{ij}^{\pi_i}) h_j^\pi \end{aligned}$$

We recall the two notions of optimality:

- A policy π is gain optimal if its gain \mathbf{g}^π equals \mathbf{g}^* .
- A policy π is Bellman optimal if $\mathbf{g}^* = \mathbf{P}^\pi \mathbf{g}^*$ and if there exists a bias \mathbf{h}^π that is a solution of (B4) and also (B2).

Recall that a Bellman optimal policy is also gain optimal.

Note that by the definition the advantage $B(\cdot)$, a policy is Bellman optimal if $\mathbf{g}^* = \mathbf{P}^\pi \mathbf{g}^*$ and if there exists $\mathbf{h} \in H$ such that $B_i^{\pi_i}(\mathbf{h}) = 0$.

Useful notations

Let $\pi : \mathcal{S} \rightarrow \mathcal{A}$ be a policy. We say that a state is recurrent for π if it is recurrent for the Markov chain whose transition matrix is \mathbf{P}^π . In other words, a state $i \in \mathcal{S}$ is recurrent if when the chain starts in i at time 0, it almost surely visits the state i at some time $t \geq 1$. We denote by \mathcal{R}^π the set of recurrent states of policy π .

We also define $\bar{\mathbf{P}}^\pi$ as the Cesaro limit of the sequence $\{(\mathbf{P}^\pi)^t\}_{t=1}$:

$$\bar{\mathbf{P}}^\pi := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T (\mathbf{P}^\pi)^{t-1}.$$

From [39, Section A.4 of Appendix A], the matrix $\bar{\mathbf{P}}^\pi$ exists and has the following properties:

- $\bar{\mathbf{P}}^\pi$ is a stochastic matrix, and satisfies $\mathbf{P}^\pi \bar{\mathbf{P}}^\pi = \bar{\mathbf{P}}^\pi \mathbf{P}^\pi = \bar{\mathbf{P}}^\pi$.
- For all state i, j , if $j \notin \mathcal{R}^\pi$, then $\bar{P}_{ij}^\pi = 0$.
- If π is unichain, then the rows of $\bar{\mathbf{P}}^\pi$ are identical.

B.1.2 Characterization of gain optimal policies

The following lemma characterizes gain optimal policy by showing that the policy must satisfies (B2) on their recurrent states.

Lemma 4. *Let $\pi : \mathcal{S} \mapsto \mathcal{A}$ be a policy and recall that \mathcal{R}^π is the set of recurrent states of policy π . The three properties below are equivalent.*

- $\mathbf{P}^\pi \mathbf{g}^* = \mathbf{g}^*$ and for all $\mathbf{h} \in H$, $B_i^{\pi_i}(\mathbf{h}) = 0$ for all $i \in \mathcal{R}^\pi$
- $\mathbf{P}^\pi \mathbf{g}^* = \mathbf{g}^*$ and for some $\mathbf{h} \in H$, $B_i^{\pi_i}(\mathbf{h}) = 0$ for all $i \in \mathcal{R}^\pi$
- π is gain optimal.

Proof (i) \Rightarrow (ii) is trivial.

(ii) \Rightarrow (iii): By definition of $B_i^a(\mathbf{h})$, we have $r_i^\pi - g_i^* = h_i - \sum_{j \in \mathcal{S}} P_{ij}^\pi h_j$ for any recurrent state i of π . Multiply this with \bar{P}_{ki}^π and sum over $i \in \mathcal{S}$ (if i is not recurrent, then $\bar{P}_{ki}^\pi = 0$) gives

$$\begin{aligned} \sum_{i \in \mathcal{S}} \bar{P}_{ki}^\pi (r_i^\pi - g_i^*) &= \sum_{i \in \mathcal{S}} \bar{P}_{ki}^\pi h_i - \underbrace{\sum_{i \in \mathcal{S}} \bar{P}_{ki}^\pi \sum_{j \in \mathcal{S}} P_{ij}^\pi h_j}_{= \sum_{j \in \mathcal{S}} \bar{P}_{kj}^\pi h_j \text{ since } \bar{\mathbf{P}}^\pi \mathbf{P}^\pi = \bar{\mathbf{P}}^\pi} = \mathbf{0}. \end{aligned}$$

By Theorem 8.2.6 of [39], the average reward of π is $\bar{\mathbf{P}}^\pi \mathbf{r}^\pi$. The above equation shows that $\bar{\mathbf{P}}^\pi \mathbf{r}^\pi = \bar{\mathbf{P}}^\pi \mathbf{g}^*$. Moreover, the assumption $\mathbf{P}^\pi \mathbf{g}^* = \mathbf{g}^*$ implies that $\bar{\mathbf{P}}^\pi \mathbf{g}^* = \mathbf{g}^*$ which in turn implies that $\bar{\mathbf{P}}^\pi \mathbf{r}^\pi = \mathbf{g}^*$. This shows that the average reward of π is \mathbf{g}^* and therefore π is gain optimal.

(iii) \Rightarrow (i): If π is gain optimal, then $\mathbf{P}^\pi \mathbf{g}^* = \mathbf{g}^*$ and $\bar{\mathbf{P}}^\pi (\mathbf{r}^\pi - \mathbf{g}^*) = \mathbf{0}$. The latter rewrites as $\sum_{i \in \mathcal{S}} \bar{P}_{ki}^\pi (r_i^\pi - g_i^*) = 0$ for all state k . Let $\mathbf{h} \in H$ be an optimal bias. For all state k , we have

$$\sum_{i \in \mathcal{S}} \bar{P}_{ki}^\pi B_i^{\pi_i}(\mathbf{h}) = \sum_{i \in \mathcal{S}} \bar{P}_{ki}^\pi (r_i^\pi - g_i^* + \sum_{j \in \mathcal{S}} P_{ij}^\pi h_j - h_i) = 0. \quad (\text{B5})$$

As \mathbf{h} satisfied (B2), for all action a , we have $B_i^a(\mathbf{h}) \leq 0$ for all states $i \in \mathcal{S}$ and in particular $B_i^{\pi_i}(\mathbf{h}) \leq 0$. This shows that for any state i such that $\bar{P}_{ki}^\pi > 0$, one must have $B_i^{\pi_i}(\mathbf{h}) = 0$. Such state i are the recurrent states of π . This shows that $B_i^{\pi_i}(\mathbf{h}) = 0$ for all $i \in \mathcal{R}^\pi$. \square

B.1.3 Characterization of Bellman optimal policies

The previous lemma shows that a policy is gain optimal if and only if the actions for the recurrent states of the policy satisfy (B2). The following lemma shows the relationship between two policies that are unichain and satisfy (B2) on all states.

Lemma 5. *Suppose that two policies π and θ are Bellman optimal, unichain and have at least one common recurrent state: $\mathcal{R}^\pi \cap \mathcal{R}^\theta \neq \emptyset$.*

Then for any \mathbf{h}^π and \mathbf{h}^θ solutions of (B4) for π and θ respectively, there exists a constant c such that for all state i : $h_i^\pi - h_i^\theta = c$. Moreover, in this case, $B_i^{\theta_i}(\mathbf{h}^\pi) = B_i^{\pi_i}(\mathbf{h}^\theta) = 0$ for all i .

Proof Since π and θ are Bellman optimal, $\mathbf{h}^\pi, \mathbf{h}^\theta \in H$. In consequence, we have

$$\mathbf{h}^\pi \geq \mathbf{r}^\theta - \mathbf{g}^* + \mathbf{P}^\theta \mathbf{h}^\pi.$$

By Lemma 4 (i), the above inequality is an equality for all $i \in \mathcal{R}^\theta$ because θ is gain optimal.

As \mathbf{h}^θ satisfies (B4), we have

$$\mathbf{h}^\theta - \mathbf{h}^\pi \leq \mathbf{r}^\theta - \mathbf{g}^* + \mathbf{P}^\theta \mathbf{h}^\theta - (\mathbf{r}^\theta - \mathbf{g}^* + \mathbf{P}^\theta \mathbf{h}^\pi) = \mathbf{P}^\theta (\mathbf{h}^\theta - \mathbf{h}^\pi),$$

with equality for all state $i \in \mathcal{R}^\theta$. This shows that for all t , $\mathbf{h}^\theta - \mathbf{h}^\pi \leq (\mathbf{P}^\theta)^t (\mathbf{h}^\theta - \mathbf{h}^\pi)$ which implies that $\mathbf{h}^\theta - \mathbf{h}^\pi \leq \bar{\mathbf{P}}^\theta (\mathbf{h}^\theta - \mathbf{h}^\pi)$ with equality for all states $i \in \mathcal{R}^\theta$. Similarly, $\mathbf{h}^\pi - \mathbf{h}^\theta \leq \bar{\mathbf{P}}^\pi (\mathbf{h}^\pi - \mathbf{h}^\theta)$ with equality for any state $i \in \mathcal{R}^\pi$.

Let $c_i^\pi = \sum_{j \in \mathcal{S}} \bar{P}_{ij}^\pi (h_j^\pi - h_j^\theta)$ and $c_i^\theta = \sum_{j \in \mathcal{S}} \bar{P}_{ij}^\theta (h_j^\pi - h_j^\theta)$. By what we have just shown, for all state i , we have

$$c_i^\theta \underbrace{\leq}_{\text{equality if } i \in \mathcal{R}^\theta} h_i^\pi - h_i^\theta \underbrace{\leq}_{\text{equality if } i \in \mathcal{R}^\pi} c_i^\pi$$

As both policies are unichain, c_i^π and c_i^θ do not depend on i . Moreover, if there exists $i \in \mathcal{R}^\theta \cap \mathcal{R}^\pi$, we have $c_i^\pi = c_i^\theta =: c$. In consequence, $h_i^\pi - h_i^\theta = c$ for all state i . \square

B.1.4 Unicity of Bellman optimal policy

Lemma 6. *Let π be a Bellman optimal policy that is unichain. If π is not the unique Bellman optimal policy, then there exists a state i and an action $a \neq \pi_i$ such that $B_i^a(\mathbf{h}^\pi) = 0$.*

Proof Let $\theta \neq \pi$ be another Bellman optimal policy. Since θ is gain optimal and $\mathbf{h}^\pi \in H$, Lemma 4 implies that $B_i^{\theta_i}(\mathbf{h}^\pi) = 0$ for all $i \in \mathcal{R}^\theta$. If there exists $i \in \mathcal{R}^\theta$ such that $\theta_i \neq \pi_i$, then the proof is concluded. Otherwise, $\theta_i = \pi_i$ for all $i \in \mathcal{R}^\theta$. This shows that π and θ coincide for all recurrent states of θ and that $\mathcal{R}^\theta = \mathcal{R}^\pi$. Moreover, as π is unichain, θ is also unichain. Hence, Lemma 5 implies that $B_i^{\theta_i}(\mathbf{h}^\pi) = 0$ for all i . Since $\theta \neq \pi$, there exists at least one state $i \in \mathcal{S}$ such that $\theta_i \neq \pi_i$. \square

B.2 Unichain property

Lemma 7. *Given a two-action MDP $\langle [n], \{0, 1\}, r, P \rangle$, let P^π be the transition matrix under a Bellman optimal policy π . Policy π is unichain if and only if the matrix*

$$\mathbf{A}^\pi = \begin{bmatrix} 1 & -P_{12}^\pi & \dots & -P_{1n}^\pi \\ 1 & 1 - P_{22}^\pi & \dots & -P_{2n}^\pi \\ \vdots & & & \\ 1 & -P_{n2}^\pi & \dots & 1 - P_{nn}^\pi \end{bmatrix}$$

is invertible.

Proof \mathbf{A}^π is not invertible if there exists a column vector $\mathbf{u} \neq \mathbf{0}$ such that $\mathbf{u}^\top \mathbf{A}^\pi = \mathbf{0}$. We prove that such \mathbf{u} does not exist when policy π is unichain. Let $\mathbf{u} \in \mathbb{R}^n$ be an arbitrary vector such that $\mathbf{u}^\top \mathbf{A}^\pi = \mathbf{0}$. Then, we have

$$\begin{cases} \sum_{i=1}^n u_i & = 0 \\ u_i - \sum_{j=1}^n u_j P_{ji}^\pi & = 0, \text{ for } 2 \leq i \leq n \end{cases}$$

Combining the above equation with $\sum_j P_{ji}^\pi = 1$, we get:

$$u_1 = - \sum_{i=2}^n u_i = - \sum_{i=2}^n \sum_{j=1}^n u_j P_{ji}^\pi = - \sum_{j=1}^n u_j (1 - P_{j1}^\pi) = \sum_{j=1}^n u_j P_{j1}^\pi,$$

where we used that $\sum_{i=1}^n u_i = 0$ to obtain the last equality. This shows that

$$\begin{cases} \sum_{i=1}^n u_i &= 0 \\ \mathbf{u}^\top \mathbf{P}^\pi &= \mathbf{u}^\top \end{cases}$$

The set of vector \mathbf{u} such that $\mathbf{u}^\top \mathbf{P}^\pi$ is a vector space. It is of dimension 1 if and only if π is unichain, in which case the vector \mathbf{u} verifying $\mathbf{u}^\top \mathbf{P}^\pi = \mathbf{u}^\top$ are multiples of a stationary distribution under policy π [39]. Thus, if the policy π induces a unichain Markov chain, then $\sum_{i=1}^n u_i = 0$ implies $\mathbf{u} = \mathbf{0}$. If policy π is not unichain, there exists $\mathbf{u} \neq \mathbf{0}$ such that $\sum_{i=1}^n u_i = 0$. \square

Appendix C Implementations

C.1 Arithmetic complexity of Subroutine 3 and memory usage

Recall that in Subroutine 3, we compute the values X_{ij}^ℓ by doing the update (for all iteration k , for all $\ell = 1$ to k and for all $i \in [n]$ or all $i \in \pi^{\ell+1}$ if we do not test indexability):

$$X_{i\sigma^k}^{\ell+1} = X_{i\sigma^k}^\ell - \frac{X_{i\sigma^\ell}^\ell}{1 + X_{\sigma^\ell\sigma^\ell}^\ell} X_{\sigma^\ell\sigma^k}^\ell \quad (\text{C6})$$

If we test indexability, there are $\sum_{k=1}^n kn = n^3/2 + O(n^2)$ such updates. If we do not test indexability, there are $\sum_{k=1}^n \sum_{\ell=1}^k (n - \ell) = n^3/3 + O(n^2)$ such updates. Below, we show each update of Equation (C6) can be done in two arithmetic operations (one addition and one multiplication), which leads to the complexity of $n^3 + O(n^2)$ (or $(2/3)n^3 + O(n^2)$) arithmetic operations for the computation of all the needed X_{ij}^k . We also show how to reduce the memory size to $O(n^2)$.

Let $W_{i\ell} := X_{i\sigma^\ell}^\ell / (1 + X_{\sigma^\ell\sigma^\ell}^\ell)$ and $V_i := X_{i\sigma^k}^\ell$. Using this, Equation (C6) can be rewritten as:

$$V_i = V_i - W_{i\ell} V_{\sigma^\ell}. \quad (\text{C7})$$

This results in the following loop at iteration k :

- Initialize V_i from $X_{i\sigma^k}^1$.
- For all $\ell \in \{1, \dots, k-1\}$, and all $i \in [n]$ (or $i \in \pi^{\ell+1}$), apply (C7).
- Compute $W_{ik} = V_i / (1 + V_{\sigma^k})$

Note that the value of \mathbf{V} is not necessary for iteration k (only the values of $W_{i\ell}$ are needed). This shows that the algorithm can be implemented with a memory $O(n^2)$.

C.2 Speedup when not checking the indexability: First found go last

When the indexability is not tested, the update (C7) is computed for all $i \in \pi^{\ell+1}$. This creates inefficiencies (due to inefficient cache usage) because the elements V_i are not accessed sequentially.

To speedup the memory accesses, our solution is to sort the items during the execution of the algorithm. At iteration k , the algorithm computes σ^k . When this is done, our implementation switches all quantities in positions σ^k and $n - k + 1$. These quantities are δ, y, z, \mathbf{W} and \mathbf{X} . For instance, once σ^1 is found, we know that the state at position n is state n and we do the following switches:

$$\begin{aligned} \delta_{\sigma^1}, \delta_n &\rightarrow \delta_n, \delta_{\sigma^1} \\ y_{\sigma^1}^1, y_n^1 &\rightarrow y_n^1, y_{\sigma^1}^1 \\ z_{\sigma^1}^1, z_n^1 &\rightarrow z_n^1, z_{\sigma^1}^1 \\ \mathbf{W}_{\sigma^1:}, \mathbf{W}_n &\rightarrow \mathbf{W}_n, \mathbf{W}_{\sigma^1:} \\ \text{and } \mathbf{X}_{\sigma^1:}, \mathbf{X}_n &\rightarrow \mathbf{X}_n, \mathbf{X}_{\sigma^1:}. \end{aligned}$$

To do so, we need an array to store all states such that at iteration k , the first $n - k$ states of the array are the active states. We will need to track the position of each state in such array.

Appendix D Analysis of the experimental time to solve a linear system

In this section, we report in Figure D3 the time taken by the default implementation to solve a linear system of the form $\mathbf{A}\mathbf{X} = \mathbf{B}$ where \mathbf{A} and \mathbf{B} are two square matrices. To obtain this figure, we generated random (full) matrices where each entry is between 0 and 1 and use the function `scipy.linalg.solve` from the library `scipy`. The reported numbers suggest that the complexity of the solver is closer to $O(n^{2.8})$ than to $O(n^3)$, although we agree that the difference between the $O(n^{2.8})$ and the $O(n^3)$ curves is small. Note that this is in accordance with the papers [44, 45] that claim that the fastest implementations of matrix multiplication and inversion are based on Strassen's algorithm and should therefore be in $O(n^{2.8})$.

Appendix E Detailed comparison with [33] and [29]

In this section, we compare our algorithm with two main related works for finite-state restless bandits problem.

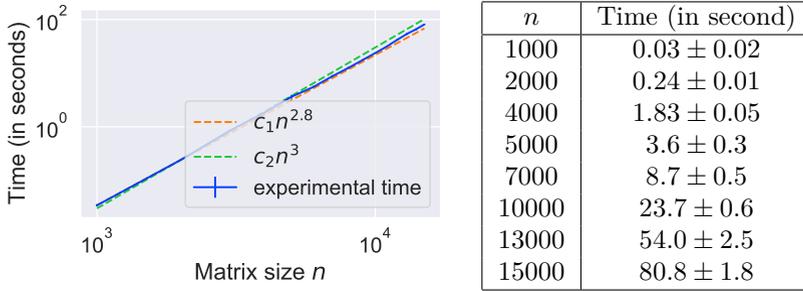


Fig. D3: Time taken of the default implementations `scipy.linalg.solve` of `scipy` to solve a linear system $\mathbf{AX} = \mathbf{B}$ where \mathbf{A} and \mathbf{B} are two square $n \times n$ matrices.

E.1 Comparison with [33]

The paper presents an algorithm that computes Whittle indices in $O(n^3)$ (no explicit constant before n^3 is given) for all indexable problems. Despite following a different approach, our algorithm for computing Whittle index can be viewed as a refinement of this work. Let us recall once again that our approach also allows one to check the indexability of general restless bandits.

In the following, we show how we can refine the work of [33] to obtain an algorithm that is exactly the same as ours. Let D^π and N^π be two vectors defined as in [33] (we use the same notation, D^π and N^π , as the cited paper),

$$D^\pi = (1 - \beta)(\mathbf{I} - \beta\mathbf{P}^\pi)^{-1}\mathbf{r}^\pi, \quad \text{and} \quad N^\pi = (1 - \beta)(\mathbf{I} - \beta\mathbf{P}^\pi)^{-1}\boldsymbol{\pi}.$$

Then, we have $D^\pi - \lambda N^\pi = (1 - \beta)\mathbf{u}^\pi(\lambda)$ where $\mathbf{u}^\pi(\lambda)$ is defined as in (22). In our proposition, at each iteration k , we compute μ_i^k by Line 13. Instead, it is defined in [33] by two steps:

1. for all state $j \in [n]$ such that $N_j^{\pi^k \setminus \{i\}} \neq N_j^{\pi^k}$, one needs to compute

$$\mu_{ij}^k = \frac{D_j^{\pi^k \setminus \{i\}} - D_j^{\pi^k}}{N_j^{\pi^k \setminus \{i\}} - N_j^{\pi^k}}$$

2. compute $\mu_i^k = \arg \min_{j \in [n]: N_j^{\pi^k \setminus \{i\}} \neq N_j^{\pi^k}} \mu_{ij}^k$.

From [33, Theorem 2], in an indexable problem, for state σ^k , there exists a state $j \in [n]$ such that $N_j^{\pi^k \setminus \{\sigma^k\}} \neq N_j^{\pi^k}$. Now, suppose that for any active state $i \in \pi^k$, there exists $j \in [n]$ such that $N_j^{\pi^k \setminus \{i\}} \neq N_j^{\pi^k}$. Using the Sherman-Morrison formula, we have⁸

$$D^{\pi^k \setminus \{i\}} - D^{\pi^k} = -\frac{(1 - \beta)\delta_i + \tilde{\Delta}_i D^{\pi^k}}{1 + \tilde{\Delta}_i [(\mathbf{I} - \beta\mathbf{P}^{\pi^k})^{-1}]_{:i}},$$

⁸the expression of $D^{\pi \setminus \{i\}}$ and $N^{\pi \setminus \{i\}}$ given by Equation (18) in [33] are erroneous.

$$\text{and } N^{\pi^k \setminus \{i\}} - N^{\pi^k} = -\frac{(1-\beta) + \tilde{\Delta}_i N^{\pi^k}}{1 + \tilde{\Delta}_i [(\mathbf{I} - \beta \mathbf{P}^{\pi^k})^{-1}]_{:i}} [(\mathbf{I} - \beta \mathbf{P}^{\pi^k})^{-1}]_{:i}.$$

Then, for any $j \in [n]$ such that $N_j^{\pi^k \setminus \{i\}} \neq N_j^{\pi^k}$, $\mu_{ij}^k = \frac{(1-\beta)\delta_i + \tilde{\Delta}_i D^{\pi^k}}{(1-\beta) + \tilde{\Delta}_i N^{\pi^k}}$

which does not depend on j . Then, we simply have $\mu_i^k = \frac{(1-\beta)\delta_i + \tilde{\Delta}_i D^{\pi^k}}{(1-\beta) + \tilde{\Delta}_i N^{\pi^k}}$.

Also, we have

$$\tilde{\Delta}_i N^{\pi^k} = (1-\beta)\tilde{\Delta}_i (\mathbf{I} - \beta \mathbf{P}^{\pi^k})^{-1} \boldsymbol{\pi}^k = -(1-\beta)y_i^k \quad \text{and}$$

$$\tilde{\Delta}_i D^{\pi^k} = (1-\beta)\tilde{\Delta}_i \mathbf{u}^{\pi^k} (\mu_{\min}^{k-1}) + \mu_{\min}^{k-1} \tilde{\Delta}_i N^{\pi^k} = (1-\beta)z_i^{k-1} - (1-\beta)\mu_{\min}^{k-1} y_i^k.$$

So, replacing these terms in μ_i^k , we get the formula in Equation (12) of our work.

Note that the algorithm of [33] was only developed for the discounted case. Our approach for the time-average reward case is different because we use the active advantage function defined in (3) instead of working with the expected discounted total reward D^{π^k} and total number of activations N^{π^k} under policy π^k . Note that the counterpart of D^{π^k} in undiscounted MDP is the average reward g^{π^k} and as we have seen in Appendix A.1, utilizing average reward optimality is not rich enough for undiscounted MDPs with transient states. In addition, our code is also optimized to avoid unnecessary computation and to reduce memory usage. Finally, the way we do the update of our matrix \mathbf{X} makes it possible to obtain a subcubic algorithm whereas their approach does not (see also below).

E.2 Comparison with the algorithm of [29]

The algorithm [29] has the best complexity up to date for discounted restless bandit. There is a square matrix \mathbf{A} that plays a similar role as the square matrix \mathbf{X} in our proposed algorithm. The most costly operations in the algorithm of [29] is to update their matrix \mathbf{A} at each iteration and it is done by Equation (19) that we recall here (using the same notation \mathbf{A} as the cited paper):

$$\text{for } i, j \in \pi^k, \mathbf{A}_{ij}^{k+1} = \mathbf{A}_{ij}^k - \frac{\mathbf{A}_{i\sigma^k}^k}{\mathbf{A}_{\sigma^k\sigma^k}^k} \mathbf{A}_{\sigma^k j}^k. \quad (\text{E8})$$

This incurs a total complexity of $(2/3)n^3 + O(n^2)$ arithmetic operations. As mentioned in Section 4.3, if we updated \mathbf{X}^{k+1} as given by (19), our algorithm would also have a $(2/3)n^3 + O(n^2)$ complexity but this version of update cannot be optimized by using fast matrix multiplication.

E.3 Their approach cannot be directly transformed into a subcubic algorithm

In addition to all the previously cited differences, one of the major contribution of our algorithm with respect to [29, 33] is that the most advanced version of our algorithm runs in a subcubic time. The approach⁹ used in [29, 33] is to update the full matrix \mathbf{X}^{k+1} at iteration k , by using (E8). This idea is represented in Figure E4(a): for a given ℓ , their algorithm compute $\mathbf{X}_{:\sigma^k}^\ell$ for all k (i.e., the full vertical lines represented by arrows). Our first Subroutine 3 uses an horizontal approach based on (20), which we recall here:

$$X_{i\sigma^k}^{\ell+1} = X_{i\sigma^k}^\ell - \frac{X_{i\sigma^\ell}^\ell}{1 + X_{\sigma^\ell\sigma^\ell}^\ell} X_{\sigma^\ell\sigma^k}^\ell.$$

At iteration $k+1$, we use $\mathbf{X}_{:\sigma^k}^1$ to compute all values of $\mathbf{X}_{:\sigma^k}^\ell$ up to $\ell = k+1$. This is represented in Figure E4(b). Our approach can be used to obtain the subcubic algorithm illustrated in Figure E4(c) by using subcubic algorithms for multiplication.

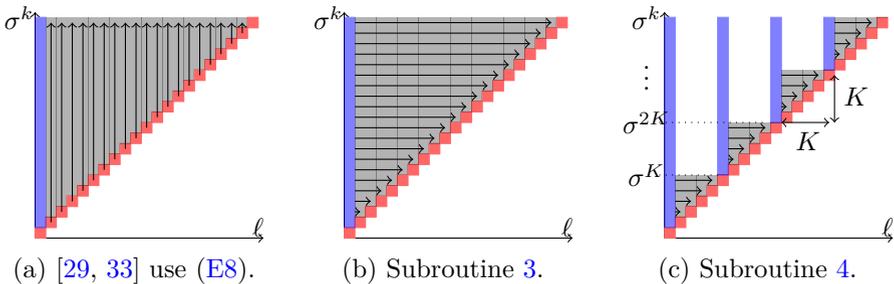


Fig. E4: Comparison of the computation load of (E8) used in [29, 33] with the one of Subroutine 3 and Subroutine 4.

This leads to the next fundamental question: why should the computation of Whittle index be harder than matrix inversion (or multiplication)? To us, the main difference is that when computing Whittle indices, the permutation σ is not known a priori but discovered as the algorithm progresses: σ^k is only known at iteration k . Hence, while all terms of the matrices X_{ij}^k are not needed, it is difficult to know a priori which ones are needed and which ones are not. Hence, a simple divide and conquer algorithm cannot be used. This is why when recomputing \mathbf{X}^{k+1} in Subroutine 4, we recompute the whole matrix (the vertical blue line) and not just the part that will be used to compute the gray zone: we do not know *a priori* what part of X_{ij}^{k+1} will be useful or not.

⁹Equation (18) of [33], which is central to their algorithm is the same as the above equation (E8).

References

- [1] Gittins, J.C.: Bandit Processes and Dynamic Allocation Indices. *Journal of the Royal Statistical Society: Series B (Methodological)* **41**(2), 148–164 (1979)
- [2] Papadimitriou, C.H., Tsitsiklis, J.N.: The complexity of optimal queueing network control. In: *Proceedings of IEEE 9th Annual Conference on Structure in Complexity Theory*, pp. 318–322 (1994). IEEE
- [3] Whittle, P.: Restless bandits: Activity allocation in a changing world. *Journal of applied probability* **25**(A), 287–298 (1988)
- [4] Verloop, I.M.: Asymptotically optimal priority policies for indexable and nonindexable restless bandits. *The Annals of Applied Probability* **26**(4), 1947–1995 (2016)
- [5] Lott, C., Teneketzis, D.: On the optimality of an index rule in multi-channel allocation for single-hop mobile networks with multiple service classes. *Probability in the Engineering and Informational Sciences* **14**(3), 259–297 (2000)
- [6] Weber, R.R., Weiss, G.: On an index policy for restless bandits. *Journal of applied probability*, 637–648 (1990)
- [7] Glazebrook, K.D., Ruiz-Hernandez, D., Kirkbride, C.: Some indexable families of restless bandit problems. *Advances in Applied Probability* **38**(3), 643–672 (2006)
- [8] Ansell, P., Glazebrook, K.D., Nino-Mora, J., O’Keeffe, M.: Whittle’s index policy for a multi-class queueing system with convex holding costs. *Mathematical Methods of Operations Research* **57**(1), 21–39 (2003)
- [9] Glazebrook, K., Mitchell, H.: An index policy for a stochastic scheduling model with improving/deteriorating jobs. *Naval Research Logistics (NRL)* **49**(7), 706–721 (2002)
- [10] Aalto, S., Lassila, P., Taboada, I.: Whittle index approach to opportunistic scheduling with partial channel information. *Performance Evaluation* **136**, 102052 (2019)
- [11] Liu, K., Zhao, Q.: Indexability of restless bandit problems and optimality of whittle index for dynamic multichannel access. *IEEE Transactions on Information Theory* **56**(11), 5547–5567 (2010)
- [12] Avrachenkov, K.E., Borkar, V.S.: Whittle index based q-learning for restless bandits with average reward. *Automatica* **139**, 110186 (2022)

- [13] Niño-Mora, J.: A dynamic page-refresh index policy for web crawlers. In: International Conference on Analytical and Stochastic Modeling Techniques and Applications, pp. 46–60 (2014). Springer
- [14] Avrachenkov, K., Ayesta, U., Doncel, J., Jacko, P.: Congestion control of tcp flows in internet routers by means of index policy. *Computer Networks* **57**(17), 3463–3478 (2013)
- [15] Avrachenkov, K., Piunovskiy, A., Zhang, Y.: Impulsive control for g-aimd dynamics with relaxed and hard constraints. In: 2018 IEEE Conference on Decision and Control (CDC), pp. 880–887 (2018). IEEE
- [16] Scully, Z., Harchol-Balter, M., Scheller-Wolf, A.: Soap: One clean analysis of all age-based scheduling policies. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* **2**(1), 1–30 (2018)
- [17] Aalto, S., Ayesta, U., Righter, R.: Properties of the gittins index with application to optimal scheduling. *Probability in the Engineering and Informational Sciences* **25**(3), 269–288 (2011)
- [18] Aalto, S., Ayesta, U., Righter, R.: On the gittins index in the m/g/1 queue. *Queueing Systems* **63**(1-4), 437 (2009)
- [19] Borkar, V.S., Pattathil, S.: Whittle indexability in egalitarian processor sharing systems. *Annals of Operations Research*, 1–21 (2017)
- [20] Larrañaga, M., Ayesta, U., Verloop, I.M.: Asymptotically optimal index policies for an abandonment queue with convex holding cost. *Queueing systems* **81**(2), 99–169 (2015)
- [21] Archibald, T.W., Black, D., Glazebrook, K.D.: Indexability and index heuristics for a simple class of inventory routing problems. *Operations research* **57**(2), 314–326 (2009)
- [22] Glazebrook, K.D., Kirkbride, C., Ouenniche, J.: Index policies for the admission control and routing of impatient customers to heterogeneous service stations. *Operations Research* **57**(4), 975–989 (2009)
- [23] Villar, S.S., Bowden, J., Wason, J.: Multi-armed bandit models for the optimal design of clinical trials: benefits and challenges. *Statistical science: a review journal of the Institute of Mathematical Statistics* **30**(2), 199 (2015)
- [24] Chen, Y.R., Katehakis, M.N.: Linear programming for finite state multi-armed bandit problems. *Mathematics of Operations Research* **11**(1), 180–183 (1986)

- [25] Katehakis, M.N., Veinott Jr, A.F.: The multi-armed bandit problem: decomposition and computation. *Mathematics of Operations Research* **12**(2), 262–268 (1987)
- [26] Niño-Mora, J.: A $(2/3) \ln 3$ fast-pivoting algorithm for the gittins index and optimal stopping of a markov chain. *INFORMS Journal on Computing* **19**(4), 596–606 (2007)
- [27] Sonin, I.M.: A generalized gittins index for a markov chain and its recursive calculation. *Statistics & Probability Letters* **78**(12), 1526–1533 (2008)
- [28] Chakravorty, J., Mahajan, A.: Multi-armed bandits, gittins index, and its calculation. *Methods and applications of statistics in clinical trials: Planning, analysis, and inferential methods* **2**(416-435), 455 (2014)
- [29] Niño-Mora, J.: A fast-pivoting algorithm for whittle’s restless bandit index. *Mathematics* **8**(12), 2226 (2020)
- [30] Akbarzadeh, N., Mahajan, A.: Restless bandits with controlled restarts: Indexability and computation of whittle index. In: 2019 IEEE 58th Conference on Decision and Control (CDC), pp. 7294–7300 (2019). IEEE
- [31] Akbarzadeh, N., Mahajan, A.: Maintenance of a collection of machines under partial observability: Indexability and computation of whittle index. *arXiv preprint arXiv:2104.05151* (2021)
- [32] Niño-Mora, J.: Dynamic priority allocation via restless bandit marginal productivity indices. *Top* **15**(2), 161–198 (2007)
- [33] Akbarzadeh, N., Mahajan, A.: Conditions for indexability of restless bandits and an $O(K^3)$ algorithm to compute Whittle index. *arXiv* (2020)
- [34] Gibson, L.J., Jacko, P., Nazarathy, Y.: A novel implementation of q-learning for the whittle index. In: *EAI International Conference on Performance Evaluation Methodologies and Tools*, pp. 154–170 (2021). Springer
- [35] Fu, J., Nazarathy, Y., Moka, S., Taylor, P.G.: Towards q-learning the whittle index for restless bandits. In: 2019 Australian & New Zealand Control Conference (ANZCC), pp. 249–254 (2019). IEEE
- [36] Nakhleh, K., Ganji, S., Hsieh, P.-C., Hou, I., Shakkottai, S., et al.: Neurwin: Neural whittle index network for restless bandits via deep rl. *Advances in Neural Information Processing Systems* **34** (2021)
- [37] Ayesta, U., Gupta, M.K., Verloop, I.M.: On the computation of whittle’s index for markovian restless bandits. *Mathematical Methods of*

Operations Research **93**(1), 179–208 (2021)

- [38] Strassen, V.: Gaussian elimination is not optimal. *Numerische mathematik* **13**(4), 354–356 (1969)
- [39] Puterman, M.L.: *Markov Decision Processes: Discrete Stochastic Dynamic Programming*, 1st edn. John Wiley & Sons, Inc., USA (1994)
- [40] Tsitsiklis, J.N.: Np-hardness of checking the unichain condition in average cost mdps. *Operations research letters* **35**(3), 319–323 (2007)
- [41] Schweitzer, P.J., Federgruen, A.: The functional equations of undiscounted markov renewal programming. *Mathematics of Operations Research* **3**(4), 308–321 (1978)
- [42] Gall, F.L., Urrutia, F.: Improved rectangular matrix multiplication using powers of the coppersmith-winograd tensor. In: *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1029–1046 (2018). SIAM
- [43] Woodbury, M.A.: *Inverting modified matrices*. Statistical Research Group (1950)
- [44] Huang, J., Smith, T.M., Henry, G.M., Van De Geijn, R.A.: Strassen’s algorithm reloaded. In: *SC’16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 690–701 (2016). IEEE
- [45] Huang, J., et al.: *Practical fast matrix multiplication algorithms*. PhD thesis (2018)
- [46] Gast, N., Gaujal, B., Yan, C.: Exponential convergence rate for the asymptotic optimality of whittle index policy. *arXiv preprint arXiv:2012.09064* (2020)