


Article

A Fast-Pivoting Algorithm for Whittle's Restless Bandit Index

José Niño-Mora 

Department of Statistics, Carlos III University of Madrid, 28903 Getafe, Spain; jose.nino@uc3m.es

Received: 20 November 2020; Accepted: 10 December 2020; Published: 15 December 2020

Abstract: The Whittle index for restless bandits (two-action semi-Markov decision processes) provides an intuitively appealing optimal policy for controlling a single generic project that can be active (engaged) or passive (rested) at each decision epoch, and which can change state while passive. It further provides a practical heuristic priority-index policy for the computationally intractable multi-armed restless bandit problem, which has been widely applied over the last three decades in multifarious settings, yet mostly restricted to project models with a one-dimensional state. This is due in part to the difficulty of establishing indexability (existence of the index) and of computing the index for projects with large state spaces. This paper draws on the author's prior results on sufficient indexability conditions and an adaptive-greedy algorithmic scheme for restless bandits to obtain a new fast-pivoting algorithm that computes the n Whittle index values of an n -state restless bandit by performing, after an initialization stage, n steps that entail $(2/3)n^3 + O(n^2)$ arithmetic operations. This algorithm also draws on the parametric simplex method, and is based on elucidating the pattern of parametric simplex tableaux, which allows to exploit special structure to substantially simplify and reduce the complexity of simplex pivoting steps. A numerical study demonstrates substantial runtime speed-ups versus alternative algorithms.

Keywords: restless bandits; Whittle index; stochastic scheduling; index policies; indexability; index algorithm; Markov decision processes

1. Introduction

We consider a general two-action (1: engaged/active; 0: rested/passive) semi-Markov decision process (SMDP) restless bandit model (see, for example, ([1], Ch. 11) and [2]) of a dynamic and stochastic *project*, whose state $X(t)$ moves over continuous time $t \in [0, \infty)$ across the *state space* \mathcal{N} , which is assumed finite consisting of $n \triangleq |\mathcal{N}|$ states. At each of an increasing sequence of *decision epochs* $\{t_k\}_{k=0}^{\infty}$ starting at $t_0 = 0$ and with $t_k \nearrow \infty$ as $k \nearrow \infty$, the *embedded state* $X_k \triangleq X(t_k)$ is observed and then an *action* $A_k \triangleq A(t_k) \in \{0, 1\}$ is chosen, which remains fixed over the subsequent *stage* $[t_k, t_{k+1})$, so $A(t) = A(t_k)$ for $t \in [t_k, t_{k+1})$. Note that the project state $X(t)$ can change within such stages, and that processes $X(t)$ and $A(t)$ are piecewise constant and right-continuous.

When the project is in state $X(t) = i$ under action $A(t) = a$, it accrues rewards and consumes a generic resource at rates R_i^a and Q_i^a per unit time, respectively, which are exponentially discounted with rate $\alpha > 0$. Resource consumption rates are assumed to satisfy the natural requirement that $0 < Q_i^1 \geq Q_i^0 \geq 0$, so when the project is active, it consumes resources at a positive rate, which is not lower than when it is passive.

To select actions, a *policy* π is adopted from the class Π of history-dependent randomized policies, which base the choice of action at each decision epoch t_k on the history $\mathcal{H}_k \triangleq \{(X_{k'}, A_{k'}) : 0 \leq k' < k\} \cup \{X_k\}$ of embedded states and actions.

Suppose that the amount of resource consumed by the project must be paid for at the unit price λ , so, writing as $\mathbb{E}_i^\pi[\cdot]$ the expectation starting from $X(0) = i$ under π ,

$$V_i^\pi(\lambda) \triangleq \mathbb{E}_i^\pi \left[\int_0^\infty (R_{X(t)}^A - \lambda Q_{X(t)}^A) e^{-\alpha t} dt \right]$$

is the (expected total discounted) project value starting from i under π , and $V_i^*(\lambda) \triangleq \sup_{\pi \in \Pi} V_i^\pi(\lambda)$ is the optimal project value starting from state i . Consider now the λ -price problem

$$(P_\lambda): \quad \text{find } \pi_\lambda^* \in \Pi \text{ such that } V_i^{\pi_\lambda^*}(\lambda) = V_i^*(\lambda) \text{ for every state } i \in \mathcal{N}, \tag{1}$$

which entails finding a policy that maximizes the project value for every initial state. We will call P_λ -optimal a policy π_λ^* solving the λ -price problem (P_λ) in (1).

Under mild assumptions, it is well-known from the theory of finite-state and -action SMDPs (see ([1], Ch. 11)) that, to solve problem (1), it suffices to consider the class Π^{SD} of stationary deterministic policies, which prescribe a fixed action at each decision epoch based on the current project state. For any given resource price λ , a P_λ -optimal policy in Π^{SD} can be computed by classic algorithms, such as value or policy iteration.

Instead, we will treat the resource price λ as a parameter, and consider a solution approach for the parametric collection \mathcal{P} of all λ -price problems (P_λ) as $\lambda \in \mathbb{R}$. Thus, let us say that the project, or more precisely, the parametric problem collection \mathcal{P} , is *indexable* if, for every project state i , there exists a finite critical price λ_i^* such that, for any problem (P_λ) and at a decision epoch in which the project is in state i : (i) it is optimal to *engage* the project if, and only if, $\lambda_i^* \geq \lambda$; and (ii) it is optimal to *rest* the project if, and only if, $\lambda_i^* \leq \lambda$. Hence, both actions will be optimal in state i if, and only if, $\lambda_i^* = \lambda$.

We will refer to the mapping $i \mapsto \lambda_i^*$ as the project's *Whittle index*, since it was Whittle who introduced such a concept in [2], in a Markovian setting with resource consumption $Q_i^a \triangleq a$. As for the extension to general resource consumption Q_i^a , it was introduced in [3]. In the case, $Q_i^a \triangleq a$, the Whittle index extends the *Gittins index* for classic (non-restless) bandit projects, which do not change state while passive. Thus, strictly speaking, the index considered herein for general Q_i^a is a *generalized Whittle index*, which reduces to a *generalized Gittins index* in the non-restless case.

Considering a semi-Markov instead of a purely Markovian setting significantly expands the modeling power of the resulting restless bandits, as in some applications, the time during which a project remains engaged before the next decision can be made may follow a general distribution. As an example (see [4]), imagine that a “project” represents a queue of jobs with Poisson arrivals and generally distributed service times, where serving a job (active action) is non-preemptive, i.e., it cannot be interrupted once started. In [5], semi-Markov restless bandits are used to model dynamic job assignment for energy-efficient server farms. Semi-Markov restless bandits can also be used to model classic (non-restless) bandits where changing from engaging one project to another entails project-dependent switching times, as shown in [6].

Besides its intrinsic interest for solving the aforementioned single-project parametric problem collection \mathcal{P} , Whittle proposed in [2] to use the index λ_i^* as the basis of a widely popular heuristic for the *multi-armed restless bandit problem* (MARBP), in which M out of $N > M$ restless bandit projects must be selected to be engaged at each time to maximize the value (under a discounted or long-run average criterion) earned from the N projects over an infinite horizon. For a sample of recent applications, see, for example, [5,7–18]. While the MARBP is computationally intractable (PSPACE-hard; see [19]), the *Whittle index policy* is an intuitively appealing heuristic where, at each time, M projects with the highest current indices are engaged, so the Whittle index plays the role of a *priority index* for a project to be engaged. This policy is convenient for practical implementation, as it avoids the *curse of dimensionality*, since each project has its own Whittle index. Furthermore, an extensive body of numerical evidence accumulated over the last three decades has shown that the policy is often nearly

optimal. Though its exact analysis remains elusive, Whittle's conjecture that his proposed policy enjoys a form of asymptotic optimality has been established under certain conditions. See [20–23].

Yet, unlike the Gittins index, which is well-defined for any classic bandit, the Whittle index exists only for a limited type of restless bandits, which are called *indexable*. Typically, researchers use ad hoc analyses to prove indexability and calculate the Whittle index in particular models. In contrast, the author has introduced and developed in [3,24] a methodology to establish indexability and compute the Whittle index for general finite-state restless bandits, extended to the semi-Markov denumerable-state case in [4] and to the continuous-state case in [25]. The effectiveness of such an approach, based on verification of so-called *PCL-indexability conditions*—as they are grounded on satisfaction by project performance metrics of *partial conservation laws* (PCLs)—has been demonstrated in diverse models. See, for example, [14,26–35].

In the case of finite-state restless bandits, [3,24] introduced the concept of satisfaction of PCLs by project performance metrics, along with a related one-pass *adaptive-greedy index algorithm*, which calculates the n index values of an n -state *PCL-indexable* project in n steps. This has its early roots in Klimov's algorithm in [36] for calculating the optimal priority indices for scheduling a multiclass queue with Bernoulli feedback, which was extended in [37] to a framework of stochastic scheduling systems satisfying so-called *generalized conservation laws*, such as the classic (non-restless) multi-armed bandit problem and branching bandits. See also [38,39]. Yet, the aforementioned work has not addressed the efficient computational implementation of such an algorithm, which is necessary for its actual deployment and widespread application.

This paper develops an efficient computational scheme for calculating the Whittle index of a general finite-state PCL-indexable restless bandit, by extending the approach in [40] from discrete-time classic bandits to semi-Markov restless bandits. The main contribution is a new fast-pivoting block implementation of the adaptive-greedy Whittle index algorithm in [3,24] that performs—after an initialization stage that entails solving a block $n \times n$ linear equation system— $(2/3)n^3 + O(n^2)$ arithmetic operations. The complexity of the resulting algorithm can be further reduced in particular models, by exploiting the special structure of the underlying state space and transition matrices. However, it appears unlikely that such an operation count can be further reduced for general restless bandits, since as shown in [3,24] computing the Whittle index, even if the ordering of the project states was known in advance, entails the solution of an $n \times n$ linear equation system, which would be solved in $(2/3)n^3 + O(n^2)$ operations by Gaussian elimination.

Note that this is the fastest operation count for a general Whittle index algorithm presented to date. In contrast, the Whittle index algorithm in [41], which applies only to the average criterion, has a complexity of $O(n^4 2^n)$ operations, which reduces to $O(n^5)$ for projects with a one-dimensional state that are known beforehand to be both indexable and solved optimally by threshold policies.

Structure of the Paper

The rest of the paper proceeds as follows. Section 2 presents a review of the related literature. Section 3 reviews previous results on Whittle indexation for finite-state restless bandits via PCL-indexability. Section 4 lays the groundwork for an efficient implementation of the adaptive-greedy index algorithm, drawing on *dynamic programming* (DP) and *linear programming* (LP) methods. Section 5 applies such results to develop a fast-pivoting computational implementation of the adaptive-greedy index algorithm. Section 6 outlines how to extend the previous results to the average criterion. Section 7 presents the results of a numerical study testing the runtimes of several index algorithms and demonstrating that the proposed fast-pivoting algorithm has significantly lower runtimes than alternative algorithms. Section 8 discusses the results presented in the paper. Section 9 concludes the paper.

2. Review of Related Literature

In this section, we review related literature, focusing on relatively recent work. We refer the reader to the recent monograph [42] on multi-armed bandits, which discusses both the MABP and the MARBP and their widespread applications.

When projects are classic (non-restless), the Whittle index reduces to the Gittins index. The efficient computation of the Gittins index has been extensively investigated in the literature. See, for example, [43–48]. While some algorithms proposed in the aforementioned papers perform $O(n^3)$ arithmetic operations for a general n -state bandit, [40] presents a *fast-pivoting* implementation of the adaptive-greedy Gittins index algorithm in [37] that performs $(2/3)n^3 + O(n^2)$ arithmetic operations, thus achieving lower complexity than alternative algorithms and matching that of Gaussian elimination for solving an $n \times n$ linear equation system. It is unlikely that such a complexity count can be improved, since it is shown in [37] that computing the Gittins index reduces precisely to solving an $n \times n$ linear equation system whose solution is subject to certain inequalities. The algorithm in [40] is based on an elucidation of the pattern of parametric simplex tableaux and exploitation of special structure to lower the computational effort of simplex pivoting steps with respect to standard pivoting.

In contrast, efficient computation of the Whittle index for restless bandits has received relatively scant research attention. Three approaches can be distinguished in the literature: deriving the Whittle index in closed form, iterative index approximation, and exact numerical computation. The first approach is to derive the Whittle index in closed form, as, for example, in [5,13,15–18]. This is to be preferred whenever possible, as the resulting analytical expressions for the Whittle index, besides facilitating its numerical evaluation, may provide valuable insight on the index dependence on model parameters. However, obtaining closed-form Whittle index formulae is only possible in relative simple models, typically with a one-dimensional state.

In more complex models in which the Whittle index cannot be evaluated in closed form, the most widespread approach, which has its roots in the calibration method for the Gittins index in [49], is to apply an iterative procedure for approximately computing the index. This is done, for example, in [7–12]. Besides the drawback that the resulting index is only an approximation to the true Whittle index, this approach is typically computationally expensive.

As for the third approach, efficient exact numerical computation of the Whittle index, this has received the least research attention to date. In [3,24], the author introduced an *adaptive-greedy algorithm* for computing the Whittle index in general finite-state restless bandits, provided that they satisfy the PCL-indexability conditions also introduced in such work. The algorithmic computes the n Whittle index values of an n -state bandit in n steps, proceeding in a greedy fashion at each step. However, as given in [3,24], such a method provides only an *algorithmic scheme*, as it is not specified how certain metrics that arise in its description are to be computed in practice. The effectiveness of such an approach is demonstrated in [3] in the setting of a broad birth–death restless bandit model, both under the discounted and long-run average criteria, motivated by queueing admission control and routing problems, where a specific implementation of the adaptive-greedy algorithm for such a case is obtained that computes the first n Whittle indices in $O(n)$ operations, under mild conditions on model parameters. An alternative approach to long-run average birth–death restless bandits is developed in [50]. Yet, to prove indexability, Proposition 2 in that paper assumes both that threshold policies are optimal (which in the birth–death model in [3] is obtained as a byproduct of the PCL-indexability conditions) and that a certain function of steady-state probabilities is strictly increasing, which is nontrivial to verify. Further, the Reference [50] does not give an index-computing algorithm, but an expression for the Whittle index in terms of steady-state metrics, which also appears in [3].

A Whittle index computing algorithm for general continuous-time finite-state restless bandits has been recently proposed in [41], focusing on the average criterion, as it is stated there that the approach does not apply to the discounted criterion. As for the arithmetic-operation complexity of the Whittle index algorithm in [41], which also checks for indexability, Remark 3.5 states that it performs $O(n^{42^n})$ operations. Even in the case that indexability is known to hold and that threshold policies are optimal,

it is stated in Remark 4.1 of that paper that the complexity of that index algorithm reduces to $O(n^5)$ operations. This is to be contrasted with the Whittle index algorithm presented herein, which has a cubic operation complexity in the number n of bandit states.

3. Review of Finite-State Restless Bandit Whittle Indexation via PCL-Indexability

This section reviews key results of the author’s approach to restless bandit indexation, as it applies to a single finite-state semi-Markov restless bandit project, which we will simply refer to as a “project” in the sequel.

3.1. SMDP Restless Bandits and Their Discrete-Stage Reformulation

Consider a SMDP restless bandit project, as outlined in Section 1. We next describe its standard discrete-stage reformulation (see, for example, ([1], Ch. 11)). If, at decision epoch t_k the project lies in state $X_k = i$ and action $A_k = a$ is chosen, the joint distribution of the length $t_{k+1} - t_k$ of the following (i, a) -stage and embedded state X_{k+1} is characterized by the transition distribution function

$$H_{ij}^a(t) \triangleq \mathbb{P}_i^a \{t_{k+1} - t_k \leq t, X_{k+1} = j\}, \quad t \geq 0,$$

having Laplace–Stieltjes transform (LST)

$$\psi_{ij}^a(\alpha) \triangleq \int_0^\infty e^{-\alpha t} dH_{ij}^a(t) = \mathbb{E}_i^a [1_{\{X_{k+1}=j\}} e^{-\alpha(t_{k+1}-t_k)}], \quad \alpha > 0,$$

where \mathbb{P}_i^a and \mathbb{E}_i^a denote probability and expectation conditional on starting from state i with action a ($X_k = i, A_k = a$).

From $H_{ij}^a(t)$, we have that the distribution of the length of an (i, a) -stage is

$$H_i^a(t) \triangleq \mathbb{P}_i^a \{t_{k+1} - t_k \leq t\} = \sum_{j \in \mathcal{N}} H_{ij}^a(t),$$

having LST

$$\psi_i^a(\alpha) \triangleq \mathbb{E}_i^a [e^{-\alpha(t_{k+1}-t_k)}] = \sum_{j \in \mathcal{N}} \psi_{ij}^a(\alpha), \tag{2}$$

and mean

$$m_i^a \triangleq \mathbb{E}_i^a [t_{k+1} - t_k \mid X_k = i, A_k = a] = \int_0^\infty t dH_i^a(t).$$

The one-stage transition probabilities for the embedded process are

$$p_{ij}^a \triangleq \mathbb{P}_i^a \{X_{k+1} = j \mid X_k = i, A_k = a\} = \lim_{t \rightarrow \infty} H_{ij}^a(t) = \lim_{\alpha \searrow 0} \psi_{ij}^a(\alpha).$$

Recall that the process $X(t)$ may change state between successive decision epochs. Its dynamics within an (i, a) -stage are characterized by

$$\tilde{p}_{ij}^a(s) \triangleq \mathbb{P}_i^a \{X(t_k + s) = j \mid t_{k+1} - t_k > s\}, \quad s \geq 0,$$

the conditional probability that, s time units after the start of an (i, a) -stage, and given that this is still ongoing, the project occupies state j . We can thus formulate the expected discounted amount of resource consumed and reward obtained in an (i, a) -stage, respectively, as

$$q_i^a \triangleq \mathbb{E}_i^a \left[\int_{t_k}^{t_{k+1}} Q_{X(t)}^A e^{-\alpha(t-t_k)} dt \right] = \sum_{j \in \mathcal{N}} Q_j^a \int_0^\infty \tilde{p}_{ij}^a(s) \bar{H}_i^a(s) e^{-\alpha s} ds \tag{3}$$

and

$$r_i^a \triangleq \mathbb{E}_i^a \left[\int_{t_k}^{t_{k+1}} R_{X(t)}^{A_k} e^{-\alpha(t-t_k)} dt \right] = \sum_{j \in \mathcal{N}} R_j^a \int_0^\infty \tilde{p}_{ij}^a(s) \bar{H}_i^a(s) e^{-\alpha s} ds, \tag{4}$$

where $\bar{H}_i^a(s) \triangleq 1 - H_i^a(s)$ is the tail distribution of the length of an (i, a) -stage.

Recall that we denote by $n \triangleq |\mathcal{N}|$ the number of project states.

3.2. Indexability, Whittle Index, and the Achievable Resource–Reward Performance Region

We start by considering the discounted criterion with rate $\alpha > 0$, deferring discussion of the average criterion to Section 6. We consider the following project performance metrics to evaluate a policy $\pi \in \Pi$ starting from a state i : the *reward metric*

$$F_i^\pi \triangleq \mathbb{E}_i^\pi \left[\int_0^\infty R_{X(t)}^{A(t)} e^{-\alpha t} dt \right] = \mathbb{E}_i^\pi \left[\sum_{k=0}^\infty r_{X_k}^{A_k} e^{-\alpha t_k} \right], \tag{5}$$

measuring the expected discounted value of rewards obtained, and the *resource metric*

$$G_i^\pi \triangleq \mathbb{E}_i^\pi \left[\int_0^\infty Q_{X(t)}^{A(t)} e^{-\alpha t} dt \right] = \mathbb{E}_i^\pi \left[\sum_{k=0}^\infty q_{X_k}^{A_k} e^{-\alpha t_k} \right], \tag{6}$$

giving the expected discounted quantity of resource consumed. Note that the right-hand side expectations in (5) and (6) are formulated in terms of the aforementioned discrete-stage reformulation.

We will also refer to the metrics obtained by drawing the initial state from a probability mass vector $\mathbf{p} = (p_i)_{i \in \mathcal{N}}$, given by

$$F_{\mathbf{p}}^\pi \triangleq \sum_{i \in \mathcal{N}} p_i F_i^\pi \quad \text{and} \quad G_{\mathbf{p}}^\pi \triangleq \sum_{i \in \mathcal{N}} p_i G_i^\pi. \tag{7}$$

To calibrate the marginal value of engaging the project in each state, we introduce a parameter $\lambda \in \mathbb{R}$ representing the *resource (unit) price*, and define the (*net*) *project value metric* $V_{\mathbf{p}}^\pi(\lambda) \triangleq F_{\mathbf{p}}^\pi - \lambda G_{\mathbf{p}}^\pi$, so the *optimal project value* is $V_{\mathbf{p}}^*(\lambda) \triangleq \sup_{\pi \in \Pi} V_{\mathbf{p}}^\pi(\lambda)$. We also write $V_i^\pi(\lambda)$ and $V_i^*(\lambda)$ as in Section 1. Consider the parametric family \mathcal{P} of λ -price problems (P_λ) defined in (1) for $\lambda \in \mathbb{R}$. Thus, for given λ , problem (1) is to find an admissible policy that maximizes the project value.

The *dynamic programming* (DP) optimality equations for λ -price problem (P_λ) are

$$V_i^*(\lambda) = \max_{a \in \{0,1\}} r_i^a - \lambda q_i^a + \sum_{j \in \mathcal{N}} \psi_{ij}^a V_j^*(\lambda), \quad i \in \mathcal{N}, \tag{8}$$

where we write $\psi_{ij}^a = \psi_{ij}^a(\alpha)$. Classic results of the SMDP theory ensure existence of a P_λ -optimal policy π_λ^* solving (1) that is stationary deterministic ($\pi_\lambda^* \in \Pi^{\text{SD}}$).

We will also consider the optimal project value starting from state i with initial action a , given by

$$V_i^{(a,*)}(\lambda) \triangleq r_i^a - \lambda q_i^a + \sum_{j \in \mathcal{N}} \psi_{ij}^a V_j^*(\lambda), \tag{9}$$

and say that *action a is P_λ -optimal in state i* if $V_i^{(a,*)}(\lambda) \geq V_i^{(1-a,*)}(\lambda)$. It will be convenient to reformulate such a definition in terms of the sign of $\Delta_{a=0}^{a=1} V_i^{(a,*)}(\lambda) \triangleq V_i^{(1,*)}(\lambda) - V_i^{(0,*)}(\lambda)$, the *marginal value of engaging the project in state i* . Thus, action $a = 1$ is P_λ -optimal in i if $\Delta_{a=0}^{a=1} V_i^{(a,*)}(\lambda) \geq 0$; $a = 0$ is P_λ -optimal in i if $\Delta_{a=0}^{a=1} V_i^{(a,*)}(\lambda) \leq 0$; and, hence, both actions are P_λ -optimal in i if, and only if,

$$\Delta_{a=0}^{a=1} V_i^{(a,*)}(\lambda) = 0. \tag{10}$$

Since actions are binary, it is convenient to represent a stationary deterministic policy by its *active set* $S \subseteq \mathcal{N}$, which is the subset of states where, at a decision epoch, the policy chooses the active action. We will refer to the *S-active policy* and write F_i^S, G_i^S , and $V_i^S(\lambda)$.

Thus, we can reformulate λ -price problem (1) as the *discrete optimization problem*

$$(P_\lambda): \quad \text{find } S^* \subseteq \mathcal{N} \text{ such that } V_i^{S^*}(\lambda) = V_i^*(\lambda) \text{ for every } i \in \mathcal{N}. \tag{11}$$

For a given resource price λ , the P_λ -optimal active and passive sets are given by

$$S^{*,1}(\lambda) \triangleq \{i \in \mathcal{N} : \Delta_{a=0}^{a=1} V_i^{(a,*)}(\lambda) \geq 0\} \quad \text{and} \quad S^{*,0} \triangleq \{i \in \mathcal{N} : \Delta_{a=0}^{a=1} V_i^{(a,*)}(\lambda) \leq 0\}. \tag{12}$$

3.3. Indexability

We will address the parametric problem collection \mathcal{P} in (1) through the concept of *indexability*, extended in [2] from classic to restless bandits with resource consumption $Q_i^a \triangleq a$, and further extended in [3] to general Q_i^a .

Under indexability, the P_λ -optimal active and passive sets $S^{*,1}(\lambda)$ and $S^{*,0}(\lambda)$ are characterized by an *index* attached to project states. Note that the definition below refers to the *sign* function $\text{sgn}: \mathbb{R} \rightarrow \{-1, 0, 1\}$, and follows the formulation of indexability in ([25], Definition 1).

Definition 1 (Indexability and Whittle index). *The project is indexable if there exists a mapping $i \mapsto \lambda_i^*$ for $i \in \mathcal{N}$, such that*

$$\text{sgn} \Delta_{a=0}^{a=1} V_i^{(a,*)}(\lambda) = \text{sgn}(\lambda_i^* - \lambda), \quad \text{for every state } i \in \mathcal{N} \text{ and price } \lambda \in \mathbb{R}. \tag{13}$$

We call λ_i^* the project's (generalized) Whittle index (Gittins index if the project is non-restless).

Thus, under indexability, the P_λ -optimal active and passive sets $S^{*,1}(\lambda)$ and $S^{*,0}(\lambda)$ are characterized as

$$S^{*,1}(\lambda) \triangleq \{i \in \mathcal{N} : \lambda_i^* \geq \lambda\} \quad \text{and} \quad S^{*,0} \triangleq \{i \in \mathcal{N} : \lambda_i^* \leq \lambda\}. \tag{14}$$

As for the economic interpretation of the Whittle index, it is shown in [3,4,27] that λ_i^* measures the *marginal productivity of the resource at state i* .

3.4. PCL-Indexability and Adaptive-Greedy Algorithm

In applications of Whittle indexation to given restless bandit models, researchers are concerned with establishing analytically their indexability and efficiently computing the index. For that purpose, the author introduced the *PCL-indexability conditions* in [3,4,24] because they are grounded on the satisfaction of *partial conservation laws* (PCLs). See [38] for a survey on conservation laws in stochastic scheduling and multi-class queueing systems.

Such conditions can be used to establish indexability relative to a nonempty family $\mathcal{F} \subseteq 2^{\mathcal{N}}$ of active sets that one should posit a priori (based on insight on the model's structure).

Definition 2 (\mathcal{F} -indexability). *We say that the project, or more precisely, the parametric collection \mathcal{P} of λ -price problems (P_λ) in (11), is \mathcal{F} -indexable if: (i) it is indexable; and (ii) for every price λ , there exists an optimal active set $S^*(\lambda) \in \mathcal{F}$, such that the $S^*(\lambda)$ -active policy is P_λ -optimal.*

Thus, for an \mathcal{F} -indexable project, problem (11) can be reduced to

$$(P_\lambda): \quad \text{find } S^* \in \mathcal{F} \text{ such that } V_i^{S^*}(\lambda) = V_i^*(\lambda) \text{ for every } i \in \mathcal{N}, \tag{15}$$

since, in such a case \mathcal{F} -policies (those with active sets $S \in \mathcal{F}$) are optimal for (11).

Algorithmic considerations lead us to require that active-set family \mathcal{F} satisfies the connectedness conditions stated next. We will refer to the *inner boundary of S with respect to \mathcal{F}* , given by

$$\partial_{\mathcal{F}}^{\text{in}} S \triangleq \{i \in S : S \setminus \{i\} \in \mathcal{F}\},$$

and to the *outer boundary of S with respect to \mathcal{F}* , given by

$$\partial_{\mathcal{F}}^{\text{out}} S \triangleq \{j \in S^c : S \cup \{j\} \in \mathcal{F}\}.$$

- Assumption 1.** (i) For every nonempty $S \in \mathcal{F}$, $\partial_{\mathcal{F}}^{\text{in}} S$ is nonempty.
 (ii) For every $S \in \mathcal{F} \setminus \{\mathcal{N}\}$, $\partial_{\mathcal{F}}^{\text{out}} S$ is nonempty.

Note that from the requirement that \mathcal{F} be nonempty, along with Assumption 1, it follows that $\emptyset, \mathcal{N} \in \mathcal{F}$.

To formulate the following result, we need to consider certain *marginal metrics*. For an action $a \in \{0, 1\}$ and an active set $S \subseteq \mathcal{N}$, denote by $\langle a, S \rangle$ the policy taking initially action a , and thereafter following the S -active policy. For a state i and an active set S , define the *marginal resource (consumption) metric* by

$$g_i^S \triangleq G_i^{\langle 1, S \rangle} - G_i^{\langle 0, S \rangle}, \tag{16}$$

that is, as the marginal increase in resource consumed resulting from taking first the active rather than the passive action at state i , provided that the S -active policy is followed thereafter.

Define also the *marginal reward metric* by

$$f_i^S \triangleq F_i^{\langle 1, S \rangle} - F_i^{\langle 0, S \rangle}, \tag{17}$$

that is, as the marginal increase in the value of rewards gained. Finally, for $g_i^S \neq 0$, define the *marginal productivity rate metric* by

$$\lambda_i^S \triangleq \frac{f_i^S}{g_i^S}. \tag{18}$$

The following definition refers to the *adaptive-greedy index algorithm*, which is given in Algorithms 1 and 2 in its top-down and bottom-up versions, respectively. In the former, index values are computed from largest to smallest, whereas in the latter they are computed in the opposite order.

Algorithm 1: Adaptive-greedy algorithm: top-down version $AG_{\mathcal{F}}^{\text{TD}}$.

Output: $\{j_k, \lambda_{j_k}^*\}_{k=1}^n$

$S_0 := \emptyset$

for $k := 1$ **to** n **do**

pick $j_k \in \arg \max \{\lambda_j^{S_{k-1}} : j \in \partial_{\mathcal{F}}^{\text{out}} S_{k-1}\}$

$\lambda_{j_k}^* := \lambda_{j_k}^{S_{k-1}}; S_k := S_{k-1} \cup \{j_k\}$

end { **for** }

Algorithm 2: Adaptive-greedy algorithm: bottom-up version $AG_{\mathcal{F}}^{\text{BU}}$.

Output: $\{i_k, \lambda_{i_k}^*\}_{k=1}^n$

$S'_0 := \mathcal{N}$

for $k := 1$ **to** n **do**

pick $i_k \in \arg \min \{\lambda_i^{S'_{k-1}} : i \in \partial_{\mathcal{F}}^{\text{in}} S'_{k-1}\}$

$\lambda_{i_k}^* := \lambda_{i_k}^{S'_{k-1}}; S'_k := S'_{k-1} \setminus \{i_k\}$

end { **for** }

Definition 3 (PCL(\mathcal{F})-indexability). We say that the project is PCL(\mathcal{F})-indexable if:

- (i) for every active set $S \in \mathcal{F}$, $g_i^S > 0$ for each $i \in \mathcal{N}$; and
- (ii) algorithm $AG_{\mathcal{F}}^{\text{TD}}$ computes a monotone non-increasing index sequence $(\lambda_{j_1}^* \geq \dots \geq \lambda_{j_n}^*)$; or algorithm $AG_{\mathcal{F}}^{\text{BU}}$ computes a monotone non-decreasing index sequence $(\lambda_{i_1}^* \leq \dots \leq \lambda_{i_n}^*)$.

Theorem 1. Suppose that the project is PCL(\mathcal{F})-indexable. Then, it is \mathcal{F} -indexable and the index λ_i^* computed by either adaptive-greedy index algorithm is its Whittle index.

Theorem 1 was introduced and proven, in increasingly general settings, in ([24], Corollary 2), (Ref. [3], Theorem 6.3) and ([4], Theorem 4.1). The author’s work reviewed in [27] demonstrates the applicability of such a result to a number of relevant restless bandit models.

Regarding the interpretation of positive marginal resource condition (i) in Definition 3, it is shown in ([3], Proposition 6.2) and in ([4], Lemma 4.3) that it can be reformulated in terms of the resource metric as follows: for $\mathbf{p} > \mathbf{0}$,

$$G_{\mathbf{p}}^{S \setminus \{i\}} < G_{\mathbf{p}}^S < G_{\mathbf{p}}^{S \cup \{j\}}, \quad S \in \mathcal{F}, i \in S, j \in S^c. \tag{19}$$

Thus, condition (i) represents a strong form of monotone increasingness of resource metric $G_{\mathbf{p}}^S$ in its active set S relative to inclusion. Additionally, recalling the definition of $G_{\mathbf{p}}^{\pi}$ in (7) in terms of metrics G_i^{π} , (19) can, in turn, be reformulated in terms of the latter metrics as follows: for $S \in \mathcal{F}$, $i \in S$ and $j \in S^c$,

$$\mathbf{G}^{S \setminus \{i\}} \preceq \mathbf{G}^S \preceq \mathbf{G}^{S \cup \{j\}}, \quad S \in \mathcal{F}, i \in S, j \in S^c \tag{20}$$

where $\mathbf{G}^{\pi} = (G_{i_0}^{\pi})_{i_0 \in \mathcal{N}}$ and “ \preceq ” means “less than or equal to componentwise, but not equal.”

Along with condition (i), condition (ii) in Definition 3 is interpreted in ([3], Section 6.4) in terms of satisfaction of the law of diminishing marginal returns (to resource usage) in economics, relative to \mathcal{F} -policies. Such an interpretation is based on the result in ([3], Proposition 6.4(a)) that the marginal productivity rates λ_i^S in (18) for $S \in \mathcal{F}$ can be recast in terms of resource and reward metrics as

$$\lambda_i^S = \begin{cases} \frac{F_{\mathbf{p}}^S - F_{\mathbf{p}}^{S \setminus \{i\}}}{G_{\mathbf{p}}^S - G_{\mathbf{p}}^{S \setminus \{i\}}}, & i \in S \\ \frac{F_{\mathbf{p}}^{S \cup \{i\}} - F_{\mathbf{p}}^S}{G_{\mathbf{p}}^{S \cup \{i\}} - G_{\mathbf{p}}^S}, & i \in S^c. \end{cases} \tag{21}$$

Such a reformulation allows us, in turn, to reformulate the adaptive-greedy algorithms above in a geometrically intuitive form in the resource–reward (γ, ϕ) plane, as shown in Algorithms 3 and 4. We thus see that the top-down algorithm $AG_{\mathcal{F}}^{\text{TD}}$ aims to traverse the upper boundary $\bar{\partial}\mathcal{R}_{\mathbf{p}}$ of the achievable resource–reward performance region $\mathcal{R}_{\mathbf{p}}$ from left to right using only active sets in \mathcal{F} , whereas the bottom-up version $AG_{\mathcal{F}}^{\text{BU}}$ seeks to traverse such an upper boundary from right to left, proceeding in a greedy fashion at each step. In such a setting, condition (ii) in Definition 3 means that the function obtained by linear interpolation on the points $(G_{\mathbf{p}}^S, F_{\mathbf{p}}^S)$ produced by either algorithm is concave. The remarkable result in Theorem 1 is that this, along with condition (i), suffices to ensure that such a function characterizes the upper boundary $\bar{\partial}\mathcal{R}_{\mathbf{p}}$ and the Whittle index.

In [4], such results are extended to restless projects on the denumerable state space of nonnegative integers, for which the resource metric is increasing along the nested family of active sets induced by the corresponding ordering. Further, in Section 3 of that paper the \mathcal{F} -indexability of such projects is characterized in terms of satisfaction of the law of diminishing marginal returns relative to \mathcal{F} -policies. In [25], such results are further extended to continuous-state projects.

Algorithm 3: Reformulated adaptive-greedy index algorithm: top-down version $AG_{\mathcal{F}}^{\text{TD}}$.

Output: $\{j_k, \lambda_{j_k}^*\}_{k=1}^n$

$S_0 := \emptyset$

for $k := 1$ **to** n **do**

pick $j_k \in \arg \max \left\{ \frac{F_{\mathbf{p}}^{S_{k-1} \cup \{j\}} - F_{\mathbf{p}}^{S_{k-1}}}{G_{\mathbf{p}}^{S_{k-1} \cup \{j\}} - G_{\mathbf{p}}^{S_{k-1}}} : j \in \partial_{\mathcal{F}}^{\text{out}} S_{k-1} \right\}$

$\lambda_{j_k}^* := \lambda_{j_k}^{S_{k-1}}; S_k := S_{k-1} \cup \{j_k\}$

end { for }

Algorithm 4: Reformulated adaptive-greedy index algorithm: bottom-up version $AG_{\mathcal{F}}^{\text{BU}}$.

Output: $\{i_k, \lambda_{i_k}^*\}_{k=1}^n$

$S'_0 := \mathcal{N}$

for $k := 1$ **to** n **do**

pick $i_k \in \arg \min \left\{ \frac{F_{\mathbf{p}}^{S'_{k-1}} - F_{\mathbf{p}}^{S'_{k-1} \setminus \{i\}}}{G_{\mathbf{p}}^{S'_{k-1}} - G_{\mathbf{p}}^{S'_{k-1} \setminus \{i\}}} : i \in \partial_{\mathcal{F}}^{\text{in}} S'_{k-1} \right\}$

$\lambda_{i_k}^* := \lambda_{i_k}^{S'_{k-1}}; S'_k := S'_{k-1} \setminus \{i_k\}$

end { for }

4. Laying the Groundwork for an Efficient Implementation of the Adaptive-Greedy Algorithm

This section lays the groundwork for developing an efficient implementation of the adaptive-greedy index algorithm. It draws on LP methods, based on formulating the λ -price problem (1) as a parametric LP problem, and elucidating the structure of its simplex tableaux.

4.1. Optimality Equations and Parametric LP Formulation

While the LP formulation below is well-known in SMDP theory (cf. [51]), for convenience we next outline it, starting from the optimality Equations (8) for (1). The primal LP formulation of such optimality equations is

$$V_{\mathbf{p}}^*(\lambda) = \min \sum_{j \in \mathcal{N}} p_j v_j$$

subject to

$$x_i^a : v_i - \sum_{j \in \mathcal{N}} \psi_{ij}^a v_j \geq r_i^a - \lambda q_i^a, \quad (i, a) \in \mathcal{N} \times \{0, 1\}.$$

Our analyses will instead use the dual problem,

$$V_{\mathbf{p}}^*(\lambda) = \max \sum_{(j,a) \in \mathcal{N} \times \{0,1\}} (r_j^a - \lambda q_j^a) x_j^a$$

subject to

$$v_j : \sum_{a \in \{0,1\}} (x_j^a - \sum_{i \in \mathcal{N}} \psi_{ij}^a x_i^a) = p_j, \quad j \in \mathcal{N}$$

$$x_j^a \geq 0, \quad (j, a) \in \mathcal{N} \times \{0, 1\}.$$

It will be convenient to deal with the latter in matrix notation, as

$$\begin{aligned}
 V_{\mathbf{p}}^*(\lambda) &= \max (\mathbf{r}^0 - \lambda \mathbf{c}^0) \mathbf{x}^0 + (\mathbf{r}^1 - \lambda \mathbf{c}^1) \mathbf{x}^1 \\
 &\text{subject to} \\
 &\left[(\mathbf{I} - \Psi^0)^\top \quad (\mathbf{I} - \Psi^1)^\top \right] \begin{pmatrix} \mathbf{x}^0 \\ \mathbf{x}^1 \end{pmatrix} = \mathbf{p} \\
 &\mathbf{x}^0, \mathbf{x}^1 \geq \mathbf{0},
 \end{aligned} \tag{22}$$

where $\mathbf{x}^a = (x_j^a)$ is a column vector, $\mathbf{r}^a = (r_j^a)$ and $\mathbf{c}^a = (q_j^a)$ are row vectors, and $^\top$ is the transposition operator.

Dual variables x_j^a correspond to the project’s *discounted state-action occupation measures*. For an initial state i , policy π , action a and state j , let

$$x_{ij}^{a,\pi} \triangleq \mathbb{E}_i^\pi \left[\sum_{k=0}^{\infty} \mathbf{1}_{\{X(t_k)=j, A(t_k)=a\}} e^{-\alpha t_k} \right]$$

be the expected discounted number of (j, a) -stages under π starting from i . If the initial state is drawn from \mathbf{p} , x_j^a corresponds to occupation measure $x_{\mathbf{p}j}^{a,\pi} \triangleq \sum_i p_i x_{ij}^{a,\pi}$. Note that reward and resource metrics are linear functions of such occupation measures: writing $\mathbf{x}_{\mathbf{p}}^{a,\pi} = (x_{\mathbf{p}j}^{a,\pi})_{j \in \mathcal{N}}$,

$$F_{\mathbf{p}}^\pi = \sum_{(j,a) \in \{0,1\} \times \mathcal{N}} r_j^a x_{\mathbf{p}j}^{a,\pi} = \mathbf{r}^0 \mathbf{x}_{\mathbf{p}}^{0,\pi} + \mathbf{r}^1 \mathbf{x}_{\mathbf{p}}^{1,\pi} \text{ and } G_{\mathbf{p}}^\pi = \sum_{(j,a) \in \{0,1\} \times \mathcal{N}} q_j^a x_{\mathbf{p}j}^{a,\pi} = \mathbf{c}^0 \mathbf{x}_{\mathbf{p}}^{0,\pi} + \mathbf{c}^1 \mathbf{x}_{\mathbf{p}}^{1,\pi}. \tag{23}$$

4.2. Bases, Basic Feasible Solutions, and Reduced Costs

We next analyze the LP problem (22), starting by elucidating its *basic feasible solutions* (BFS). Each BFS is obtained from a *basis* corresponding to an active set $S \subseteq \mathcal{N}$, and hence we will refer to the *S-active BFS*. Yet, note that different S ’s might yield the same BFS. For given S , we decompose the above vectors and matrices as

$$\mathbf{x}^a = \begin{pmatrix} \mathbf{x}_S^a \\ \mathbf{x}_{S^c}^a \end{pmatrix}, \quad \mathbf{p} = \begin{pmatrix} \mathbf{p}_S \\ \mathbf{p}_{S^c} \end{pmatrix}, \quad \Psi^a = \begin{pmatrix} \Psi_{SS}^a & \Psi_{SS^c}^a \\ \Psi_{S^cS}^a & \Psi_{S^cS^c}^a \end{pmatrix}, \quad \mathbf{I} = \begin{pmatrix} \mathbf{I}_{SS} & \mathbf{0}_{SS^c} \\ \mathbf{0}_{S^cS} & \mathbf{I}_{S^cS^c} \end{pmatrix},$$

and introduce the matrices

$$\begin{aligned}
 \Psi^S &\triangleq \begin{pmatrix} \Psi_{SS}^1 & \Psi_{SS^c}^1 \\ \Psi_{S^cS}^0 & \Psi_{S^cS^c}^0 \end{pmatrix}, \quad \Psi^{S^c} \triangleq \begin{pmatrix} \Psi_{SS}^0 & \Psi_{SS^c}^0 \\ \Psi_{S^cS}^1 & \Psi_{S^cS^c}^1 \end{pmatrix}, \\
 \mathbf{B}^S &\triangleq (\mathbf{I} - \Psi^S)^\top, \quad \mathbf{N}^S \triangleq (\mathbf{I} - \Psi^{S^c})^\top, \quad \mathbf{H}^S \triangleq (\mathbf{B}^S)^{-1}, \quad \mathbf{A}^S \triangleq \mathbf{H}^S \mathbf{N}^S.
 \end{aligned} \tag{24}$$

Note that Ψ^S is the transition-transform matrix for the S -active policy. Additionally, \mathbf{B}^S is the *basis matrix* in LP problem (22) for the S -active BFS, which has as *basic variables*

$$\begin{pmatrix} \mathbf{x}_S^1 \\ \mathbf{x}_{S^c}^0 \end{pmatrix},$$

and \mathbf{N}^S is the non-basic matrix of LP problem (22) corresponding to *non-basic variables*

$$\begin{pmatrix} \mathbf{x}_S^0 \\ \mathbf{x}_{S^c}^1 \end{pmatrix}.$$

We can thus rearrange the constraints in (22), decomposing them as

$$\mathbf{B}^S \begin{pmatrix} \mathbf{x}_S^1 \\ \mathbf{x}_{S^c}^0 \end{pmatrix} + \mathbf{N}^S \begin{pmatrix} \mathbf{x}_S^0 \\ \mathbf{x}_{S^c}^1 \end{pmatrix} = \mathbf{p}.$$

We next evaluate performance metrics under the S -active policy. The notation $x_{pj}^{a,S}$ below refers to occupation measure $x_{pj}^{a,\pi}$ under such a policy. Further, $\mathbf{F}^S = (F_j^S)_{j \in \mathcal{N}}$, $\mathbf{G}^S = (G_j^S)_{j \in \mathcal{N}}$, $\mathbf{f}^S = (f_j^S)_{j \in \mathcal{N}}$ and $\mathbf{g}^S = (g_j^S)_{j \in \mathcal{N}}$ are row vectors.

Lemma 1.

- (a) $\begin{pmatrix} \mathbf{x}_{\mathbf{p}^S}^{0,S} \\ \mathbf{x}_{\mathbf{p}^{S^c}}^{1,S} \end{pmatrix} = \mathbf{0}$ and $\begin{pmatrix} \mathbf{x}_{\mathbf{p}^S}^{1,S} \\ \mathbf{x}_{\mathbf{p}^{S^c}}^{0,S} \end{pmatrix} = \mathbf{H}^S \mathbf{p}$.
- (b) $\mathbf{G}^S = \begin{pmatrix} \mathbf{c}_S^1 & \mathbf{c}_{S^c}^0 \end{pmatrix} \mathbf{H}^S$.
- (c) $\mathbf{F}^S = \begin{pmatrix} \mathbf{r}_S^1 & \mathbf{r}_{S^c}^0 \end{pmatrix} \mathbf{H}^S$.
- (d) $\begin{pmatrix} \mathbf{g}_S^S & -\mathbf{g}_{S^c}^S \end{pmatrix} = \begin{pmatrix} \mathbf{c}_S^1 & \mathbf{c}_{S^c}^0 \end{pmatrix} \mathbf{A}^S - \begin{pmatrix} \mathbf{c}_S^0 & \mathbf{c}_{S^c}^1 \end{pmatrix}$.
- (e) $\begin{pmatrix} \mathbf{f}_S^S & -\mathbf{f}_{S^c}^S \end{pmatrix} = \begin{pmatrix} \mathbf{r}_S^1 & \mathbf{r}_{S^c}^0 \end{pmatrix} \mathbf{A}^S - \begin{pmatrix} \mathbf{r}_S^0 & \mathbf{r}_{S^c}^1 \end{pmatrix}$.

Proof. (a) Set to zero non-basic variables: $\mathbf{x}_{\mathbf{p}^S}^{0,S} = \mathbf{0}$ and $\mathbf{x}_{\mathbf{p}^{S^c}}^{1,S} = \mathbf{0}$. To calculate the values of basic variables, note that

$$\mathbf{B}^S \begin{pmatrix} \mathbf{x}_{\mathbf{p}^S}^1 \\ \mathbf{x}_{\mathbf{p}^{S^c}}^0 \end{pmatrix} = \mathbf{p}, \text{ and hence } \begin{pmatrix} \mathbf{x}_{\mathbf{p}^S}^{1,S} \\ \mathbf{x}_{\mathbf{p}^{S^c}}^{0,S} \end{pmatrix} = \mathbf{H}^S \mathbf{p}.$$

(b) Use part (a) with $\mathbf{p} = \mathbf{e}_j$ (the j th unit coordinate vector) to formulate resource metrics as

$$\mathbf{G}_j^S = \begin{pmatrix} \mathbf{c}_S^1 & \mathbf{c}_{S^c}^0 \end{pmatrix} \begin{pmatrix} \mathbf{x}_{j^S}^{1,S} \\ \mathbf{x}_{j^{S^c}}^{0,S} \end{pmatrix} = \begin{pmatrix} \mathbf{c}_S^1 & \mathbf{c}_{S^c}^0 \end{pmatrix} \mathbf{H}^S \mathbf{e}_j \implies \mathbf{G}^S = \begin{pmatrix} \mathbf{c}_S^1 & \mathbf{c}_{S^c}^0 \end{pmatrix} \mathbf{H}^S.$$

(c) Proceed as in (b) to formulate reward metrics as

$$\mathbf{F}_j^S = \begin{pmatrix} \mathbf{r}_S^1 & \mathbf{r}_{S^c}^0 \end{pmatrix} \begin{pmatrix} \mathbf{x}_{j^S}^{1,S} \\ \mathbf{x}_{j^{S^c}}^{0,S} \end{pmatrix} = \begin{pmatrix} \mathbf{r}_S^1 & \mathbf{r}_{S^c}^0 \end{pmatrix} \mathbf{H}^S \mathbf{e}_j \implies \mathbf{F}^S = \begin{pmatrix} \mathbf{r}_S^1 & \mathbf{r}_{S^c}^0 \end{pmatrix} \mathbf{H}^S.$$

(d) Represent marginal resource metrics (cf. (16)) as

$$\mathbf{g}_S^S = \mathbf{G}_S^S - \mathbf{c}_S^0 - \mathbf{G}^S (\mathbf{\Psi}_{S^c N}^0)^\top \quad \text{and} \quad \mathbf{g}_{S^c}^S = \mathbf{c}_{S^c}^1 + \mathbf{G}^S (\mathbf{\Psi}_{S^c N}^1)^\top - \mathbf{G}_{S^c}^S. \tag{25}$$

Recast now the equalities in (25), using (b), as

$$\begin{aligned} \begin{pmatrix} \mathbf{g}_S^S & -\mathbf{g}_{S^c}^S \end{pmatrix} &= \mathbf{G}^S \mathbf{N}^S - \begin{pmatrix} \mathbf{c}_S^0 & \mathbf{c}_{S^c}^1 \end{pmatrix} = \begin{pmatrix} \mathbf{c}_S^1 & \mathbf{c}_{S^c}^0 \end{pmatrix} \mathbf{H}^S \mathbf{N}^S - \begin{pmatrix} \mathbf{c}_S^0 & \mathbf{c}_{S^c}^1 \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{c}_S^1 & \mathbf{c}_{S^c}^0 \end{pmatrix} \mathbf{A}^S - \begin{pmatrix} \mathbf{c}_S^0 & \mathbf{c}_{S^c}^1 \end{pmatrix}. \end{aligned}$$

(e) This part follows as part (d). \square

The following result gives the *reduced costs* of the LP problem (22) in terms of the marginal resource and reward metrics in (16) and (17). It further uses such a result to represent such an LP problem’s objective in terms of reduced costs.

Lemma 2. *The reduced costs for non-basic variables corresponding to the S-active BFS for LP problem (22) are given by*

$$\left(\mathbf{f}_S^S - \lambda \mathbf{g}_S^S \quad -\mathbf{f}_{S^c}^S + \lambda \mathbf{g}_{S^c}^S \right), \tag{26}$$

and hence, the LP problem’s objective can be formulated as

$$\sum_{(j,a) \in \mathcal{N} \times \{0,1\}} (r_j^a - \lambda q_j^a) x_j^a = F_{\mathbf{p}}^S - \lambda G_{\mathbf{p}}^S - \sum_{j \in S} (f_j^S - \lambda g_j^S) x_j^0 + \sum_{j \in S^c} (f_j^S - \lambda g_j^S) x_j^1. \tag{27}$$

Proof. The results follow from the well-known representation of reduced costs in LP theory, as given by Lemma 1 (d,e), along with the well-known representation of the LP problem’s objective in terms of the value of the current BFS and reduced costs. \square

The following result, which follows from Lemma 2, represents metrics $G_{\mathbf{p}}^\pi, F_{\mathbf{p}}^\pi$ and objective $F_{\mathbf{p}}^\pi - \lambda G_{\mathbf{p}}^\pi$ in terms of their values under the S-active policy. These *decomposition identities* were first obtained in ([24], Theorem 3) and ([3], Proposition 6.1) via ad hoc arguments.

Lemma 3. *For any policy $\pi \in \Pi$:*

- (a) $G_{\mathbf{p}}^\pi = G_{\mathbf{p}}^S - \sum_{j \in S} g_j^S x_{\mathbf{p}j}^{0,\pi} + \sum_{j \in S^c} g_j^S x_{\mathbf{p}j}^{1,\pi}.$
- (b) $F_{\mathbf{p}}^\pi = F_{\mathbf{p}}^S - \sum_{j \in S} f_j^S x_{\mathbf{p}j}^{0,\pi} + \sum_{j \in S^c} f_j^S x_{\mathbf{p}j}^{1,\pi}.$
- (c) $F_{\mathbf{p}}^\pi - \lambda G_{\mathbf{p}}^\pi = F_{\mathbf{p}}^S - \lambda G_{\mathbf{p}}^S - \sum_{j \in S} (f_j^S - \lambda g_j^S) x_{\mathbf{p}j}^{0,\pi} + \sum_{j \in S^c} (f_j^S - \lambda g_j^S) x_{\mathbf{p}j}^{1,\pi}.$

The following result, first established in ([3], Corollary 6.1), elucidates the relationship between resource and reward metrics and the corresponding marginal metrics.

Lemma 4.

- (a) *For $j \in S^c, G_{\mathbf{p}}^{S \cup \{j\}} = G_{\mathbf{p}}^S + g_j^S x_{\mathbf{p}j}^{1, S \cup \{j\}}$ and $F_{\mathbf{p}}^{S \cup \{j\}} = F_{\mathbf{p}}^S + f_j^S x_{\mathbf{p}j}^{1, S \cup \{j\}}.$*
- (b) *For $j \in S, G_{\mathbf{p}}^{S \setminus \{j\}} = G_{\mathbf{p}}^S - g_j^S x_{\mathbf{p}j}^{1, S \setminus \{j\}}$ and $F_{\mathbf{p}}^{S \setminus \{j\}} = F_{\mathbf{p}}^S - f_j^S x_{\mathbf{p}j}^{1, S \setminus \{j\}}.$*

Proof. To obtain (a) use $\pi = S \cup \{j\}$ in Lemma 3 (a). Additionally, similarly with (b) with $\pi = S \setminus \{j\}$ and Lemma 3 (b). \square

5. A Fast-Pivoting Index Algorithm for PCL(\mathcal{F})-Indexable Projects

This section develops an efficient implementation of the adaptive-greedy index algorithm above, focusing on its top-down version, for a project that is PCL(\mathcal{F})-indexable.

We start by noting that the index of such a project can be computed by deploying in the LP formulation (22) of the λ -price problem the classic *parametric-objective simplex algorithm* in [52]. In the present setting, the *parametric simplex tableau* for the S-active BFS is shown in Table 1. This tableau has basic variables x_S^1 and $x_{S^c}^0$ in rows and non-basic variables x_S^0 and $x_{S^c}^1$ in columns, and further, has two rows of reduced-costs for non-basic variables. The tableau does not include right-hand sides or objectives, as they are not required in this context. The tableau is displayed just before *pivoting* on a_{jj}^S , where $j \in S^c$. That is, just before moving variable x_j^0 out of the basis, and carrying x_j^1 into the basis, which corresponds to changing from the S- to the $S \cup \{j\}$ -active BFS. After the pivot step, the updated tableau is shown in Table 2.

One can readily use such tableaux to implement the top-down adaptive-greedy algorithm $AG_{\mathcal{F}}^{\text{TD}}$ in Algorithm 3, by first constructing the initial tableau for the \emptyset -active policy, and then carrying out $n = |\mathcal{N}|$ pivot steps, at each of which the former BFS active set is augmented by a state. Such a direct approach results in an implementation that we call the *Conventional-Pivoting Index* (CP) algorithm.

An immediate counting argument shows that the n pivot steps of that algorithm perform $2n^3 + O(n^2)$ arithmetic operations—without considering computation of the initial tableau, which will be addressed in Section 5.1 below.

Note that this algorithm can also be applied to a restless bandit instance to test numerically whether it is indexable. The project will be indexable if, and only if, the successive pivot steps, as the price parameter λ decreases from ∞ to $-\infty$ in the parametric-objective simplex algorithm of [52], can be performed augmenting the current BFS by adding a state, thus producing a nested active-set family $\mathcal{F}_0 = \{S_0, S_1, \dots, S_n\}$ with $S_0 = \emptyset \subset \dots \subset S_n = \mathcal{N}$.

Table 1. Simplex tableau (parametric) for S -activeBFS, with pivot a_{jj}^S .

	$(\mathbf{x}_S^0)^\top$	x_j^1	$(\mathbf{x}_{S^c \setminus \{j\}}^1)^\top$
\mathbf{x}_S^1	\mathbf{A}_{SS}^S	\mathbf{A}_{Sj}^S	$\mathbf{A}_{SS^c \setminus \{j\}}^S$
x_j^0	\mathbf{A}_{jS}^S	a_{jj}^S	$\mathbf{A}_{jS^c \setminus \{j\}}^S$
$\mathbf{x}_{S^c \setminus \{j\}}^0$	$\mathbf{A}_{S^c \setminus \{j\}, S}^S$	$\mathbf{A}_{S^c \setminus \{j\}, j}^S$	$\mathbf{A}_{S^c \setminus \{j\}, S^c \setminus \{j\}}^S$
	\mathbf{g}_S^S	$-g_j^S$	$-\mathbf{g}_{S^c \setminus \{j\}}^S$
	\mathbf{f}_S^S	$-f_j^S$	$-\mathbf{r}_{S^c \setminus \{j\}}^S$

Table 2. Tableau for $S \cup \{j\}$ -active BFS, after pivoting on a_{jj}^S .

	$(\mathbf{x}_S^0)^\top$	x_j^0	$(\mathbf{x}_{S^c \setminus \{j\}}^1)^\top$
\mathbf{x}_S^1	$\mathbf{A}_{SS}^S - (a_{jj}^S)^{-1} \mathbf{A}_{Sj}^S \mathbf{A}_{jS}^S$	$-(a_{jj}^S)^{-1} \mathbf{A}_{Sj}^S$	$\mathbf{A}_{SS^c \setminus \{j\}}^S - (a_{jj}^S)^{-1} \mathbf{A}_{Sj}^S \mathbf{A}_{jS^c \setminus \{j\}}^S$
x_j^1	$(a_{jj}^S)^{-1} \mathbf{A}_{jS}^S$	$(a_{jj}^S)^{-1}$	$(a_{jj}^S)^{-1} \mathbf{A}_{jS^c \setminus \{j\}}^S$
$\mathbf{x}_{S^c \setminus \{j\}}^0$	$\mathbf{A}_{S^c \setminus \{j\}, S}^S - (a_{jj}^S)^{-1} \mathbf{A}_{S^c \setminus \{j\}, j}^S \mathbf{A}_{jS}^S$	$-(a_{jj}^S)^{-1} \mathbf{A}_{S^c \setminus \{j\}, j}^S$	$\mathbf{A}_{S^c \setminus \{j\}, S^c \setminus \{j\}}^S - (a_{jj}^S)^{-1} \mathbf{A}_{S^c \setminus \{j\}, j}^S \mathbf{A}_{jS^c \setminus \{j\}}^S$
	$\mathbf{g}_S^S + g_j^S (a_{jj}^S)^{-1} \mathbf{A}_{jS}^S$	$g_j^S (a_{jj}^S)^{-1}$	$-\mathbf{g}_{S^c \setminus \{j\}}^S + g_j^S (a_{jj}^S)^{-1} \mathbf{A}_{jS^c \setminus \{j\}}^S$
	$\mathbf{f}_S^S + f_j^S (a_{jj}^S)^{-1} \mathbf{A}_{jS}^S$	$f_j^S (a_{jj}^S)^{-1}$	$-\mathbf{f}_{S^c \setminus \{j\}}^S + f_j^S (a_{jj}^S)^{-1} \mathbf{A}_{jS^c \setminus \{j\}}^S$

We now consider how to improve the efficiency of the CP algorithm, drawing on the observation that the tableau’s rows for basic variables \mathbf{x}_S^1 are not used to update the reduced costs. Thus, it is enough to update and store the *reduced tableaux* shown in Table 3. Table 2 demonstrates that a reduced tableau can be updated without using the rows that have been deleted. The resulting simplification of the CP algorithm yields the *Reduced-Pivoting* (RP) algorithm. By an elementary counting argument, it is readily seen that the RP algorithm carries out the n pivot steps in $n^3 + O(n^2)$ operations, thus improving by a factor of 2 the arithmetic operation complexity of the CP algorithm.

Table 3. Reduced tableau for S -active BFS, with pivot a_{jj}^S .

	$(\mathbf{x}_S^0)^\top$	x_j^1	$(\mathbf{x}_{S^c \setminus \{j\}}^1)^\top$
x_j^0	\mathbf{A}_{jS}^S	a_{jj}^S	$\mathbf{A}_{jS^c \setminus \{j\}}^S$
$\mathbf{x}_{S^c \setminus \{j\}}^0$	$\mathbf{A}_{S^c \setminus \{j\}, S}^S$	$\mathbf{A}_{S^c \setminus \{j\}, j}^S$	$\mathbf{A}_{S^c \setminus \{j\}, S^c \setminus \{j\}}^S$
	\mathbf{g}_S^S	$-g_j^S$	$-\mathbf{g}_{S^c \setminus \{j\}}^S$
	\mathbf{f}_S^S	$-f_j^S$	$-\mathbf{f}_{S^c \setminus \{j\}}^S$

We can exploit the assumption of PCL(\mathcal{F})-indexability to reduce even further the operation count, by updating and storing only the *minimal tableaux* in Table 4. Such a tableau for the $S \cup \{j\}$ -active

BFS is computed from that for the S -active BFS in Table 4, as displayed in Table 5. This results in the *fast-pivoting (FP) adaptive-greedy* index algorithm $FP_{\mathcal{F}}$ in Algorithm 5.

Table 4. Minimal tableau for S -active BFS.

$$x_{S^c}^0 \begin{array}{|c|} \hline (x_{S^c}^1)^\top \\ \hline \mathbf{A}_{S^c S^c}^S \\ \hline \mathbf{g}_{S^c}^S \\ \hline \mathbf{f}_{S^c}^S \\ \hline \end{array}$$

Table 5. Minimal tableau for $S \cup \{j\}$ -active BFS, obtained after pivoting on a_{jj}^S .

$$x_{S^c \setminus \{j\}}^0 \begin{array}{|c|} \hline (x_{S^c \setminus \{j\}}^1)^\top \\ \hline \mathbf{A}_{S^c \setminus \{j\} S^c \setminus \{j\}}^S - (a_{jj}^S)^{-1} \mathbf{A}_{S^c \setminus \{j\} j}^S \mathbf{A}_{j S^c \setminus \{j\}}^S \\ \hline \mathbf{g}_{S^c \setminus \{j\}}^S - g_j^S (a_{jj}^S)^{-1} \mathbf{A}_{S^c \setminus \{j\} j}^S \\ \hline \mathbf{f}_{S^c \setminus \{j\}}^S - f_j^S (a_{jj}^S)^{-1} \mathbf{A}_{S^c \setminus \{j\} j}^S \\ \hline \end{array}$$

Algorithm 5: The fast-pivoting adaptive-greedy index algorithm $FP_{\mathcal{F}}$.

Output: $\{j_k, \lambda_{j_k}^*\}_{k=1}^n$

solve $\mathbf{A}^{(0)} \begin{pmatrix} \mathbf{I}_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}} - \Psi_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}}^0 & \tilde{\mathbf{m}}^0 \end{pmatrix} = \begin{pmatrix} \mathbf{I}_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}} - \Psi_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}}^1 & \tilde{\mathbf{m}}^1 \end{pmatrix}$
 $\begin{pmatrix} \mathbf{g}^{(0)} \\ \mathbf{f}^{(0)} \end{pmatrix} := \begin{pmatrix} \mathbf{c}^1 \\ \mathbf{r}^1 \end{pmatrix} - \begin{pmatrix} \mathbf{c}^0 \\ \mathbf{r}^0 \end{pmatrix} \mathbf{A}^{(0)}; S_0 := \emptyset$

for $k := 1$ **to** n **do**

$\lambda_j^{(k-1)} := f_j^{(k-1)} / g_j^{(k-1)}, j \in \partial_{\mathcal{F}}^{\text{out}} S_{k-1}$

pick $j_k \in \arg \max \{\lambda_i^{(k-1)} : j \in \partial_{\mathcal{F}}^{\text{out}} S_{k-1}\}; \lambda_{j_k}^* := \lambda_{j_k}^{(k-1)}; S_k := S_{k-1} \cup \{j_k\}$

if $k < n$ **then**

$\mathbf{A}_{S_k^c j_k}^{(k)} := (1/a_{j_k j_k}^{(k-1)}) \mathbf{A}_{S_k^c j_k}^{(k-1)}; \mathbf{A}_{S_k^c S_k^c}^{(k)} := \mathbf{A}_{S_k^c S_k^c}^{(k-1)} - \mathbf{A}_{S_k^c j_k}^{(k)} \mathbf{A}_{j_k S_k^c}^{(k-1)}$

end { if }

$\mathbf{g}_{S_k^c}^{(k)} := \mathbf{g}_{S_k^c}^{(k-1)} - g_{j_k}^{(k-1)} \mathbf{A}_{S_k^c j_k}^{(k)}; \mathbf{f}_{S_k^c}^{(k)} := \mathbf{f}_{S_k^c}^{(k-1)} - f_{j_k}^{(k-1)} \mathbf{A}_{S_k^c j_k}^{(k)}$

end { for }

The next result evaluates the operation count of algorithm $FP_{\mathcal{F}}$, showing that it outperforms significantly that of algorithm RP. Note that the complexity of its loop—which performs the n pivot steps—matches that of Gaussian elimination for solving an $n \times n$ system of linear equations, and is hence unlikely that such complexity can be improved for general restless bandits.

Proposition 1. *The loop of algorithm $FP_{\mathcal{F}}$ entails $(2/3)n^3 + O(n^2)$ arithmetic operations.*

Proof. The operation count in the loop is dominated by the update of matrix $\mathbf{A}_{S_k^c S_k^c}^{S_k}$ at each step k , taking $2(n - k)^2$ arithmetic operations. This yields the total operation count as stated. \square

5.1. Computing the Initial Tableau

We next address computation of the initial tableau, which corresponds to the \emptyset -active BFS, in a form that is numerically stable and that applies both to the discounted and the average criterion to be addressed in Section 6. The tableaux for the average criterion arise as limits letting the discount rate α vanish in the discounted tableaux.

Note that (cf. (24))

$$\mathbf{B}^\varnothing = (\mathbf{I} - \Psi^0)^\top, \quad \mathbf{N}^\varnothing = (\mathbf{I} - \Psi^1)^\top, \quad \mathbf{H}^\varnothing = (\mathbf{B}^\varnothing)^{-1}, \quad \mathbf{A}^\varnothing = \mathbf{H}^\varnothing \mathbf{N}^\varnothing. \tag{28}$$

Thus, a straightforward approach to computing \mathbf{A}^\varnothing is to solve the linear equation system

$$(\mathbf{A}^\varnothing)^\top (\mathbf{I} - \Psi^0) = (\mathbf{I} - \Psi^1). \tag{29}$$

However, this has a disadvantage: as α vanishes, the matrices $\mathbf{I} - \Psi^a$ become increasingly ill-conditioned, as they are singular for $\alpha = 0$ —because they converge to $\mathbf{I} - \mathbf{P}^a$ where $\mathbf{P}^a \triangleq (p_{ij}^a)$.

To overcome such a drawback, we draw on the identity $(\mathbf{I} - \Psi^a)\mathbf{1} = \mathbf{1} - \Psi^a$, which is a consequence from (2). From this and (28), we have

$$(\mathbf{A}^\varnothing)^\top (\mathbf{1} - \Psi^0) = \mathbf{1} - \Psi^1.$$

The latter identity gives a useful counterpart as $\alpha \searrow 0$. Thus, writing as ζ_i^a the length of an (i, a) -stage (cf. Section 3.1), and applying the MacLaurin series

$$\psi_i^a = \mathbb{E}[e^{-\alpha \zeta_i^a}] = 1 - \alpha m_i^a + O(\alpha^2) \quad \text{as } \alpha \searrow 0,$$

one obtains in the limit

$$(\mathbf{A}^\varnothing)^\top \mathbf{m}^0 = \mathbf{m}^1,$$

with m_i^a the mean length of an (i, a) -stage and $\mathbf{m} = (m_i^a)_{i \in \mathcal{N}}$.

We thus obtain the following approach to computing the initial tableau, for $\alpha \geq 0$. Letting

$$\tilde{m}_i^a \triangleq \begin{cases} m_i^a & \text{if } \alpha = 0 \\ (1 - \psi_i^a)/\alpha & \text{if } \alpha > 0, \end{cases}$$

and $\tilde{\mathbf{m}}^a = (\tilde{m}_i^a)_{i \in \mathcal{N}}$, pick any state $j^* \in \mathcal{N}$ and solve the (block) linear equation system

$$\left(\mathbf{I}_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}} - \Psi_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}}^0 \quad \tilde{\mathbf{m}}^0 \right)^\top \mathbf{A}^\varnothing = \left(\mathbf{I}_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}} - \Psi_{\mathcal{N}, \mathcal{N} \setminus \{j^*\}}^1 \quad \tilde{\mathbf{m}}^1 \right)^\top \tag{30}$$

to obtain \mathbf{A}^\varnothing . Then, calculate initial reduced costs using (28) and Lemma 1 (d,e) as

$$\begin{aligned} \mathbf{g}^\varnothing &= \mathbf{c}^1 - \mathbf{c}^0 \mathbf{A}^\varnothing \\ \mathbf{f}^\varnothing &= \mathbf{r}^1 - \mathbf{r}^0 \mathbf{A}^\varnothing. \end{aligned} \tag{31}$$

6. Extension to the Average Criterion

In applications where the (long-run) average criterion is employed, one must consider the version of λ -price problem (1) based on the average reward and resource metrics given by

$$F_i^\pi \triangleq \liminf_{T \nearrow \infty} \frac{1}{T} \mathbb{E}_i^\pi \left[\int_0^T R_{X(t)}^A dt \right] = \liminf_{K \nearrow \infty} \frac{1}{K} \mathbb{E}_i^\pi \left[\sum_{k=0}^K r_{X_k}^{a_k} \right], \tag{32}$$

and

$$G_i^\pi \triangleq \limsup_{T \nearrow \infty} \frac{1}{T} \mathbb{E}_i^\pi \left[\int_0^T Q_{X(t)}^A dt \right] = \limsup_{K \nearrow \infty} \frac{1}{K} \mathbb{E}_i^\pi \left[\sum_{k=0}^K q_{X_k}^{a_k} \right]. \tag{33}$$

As in ([3], Section 6.5), we assume that the embedded process X_k is *communicating*, so each state can be reached from every other state under some stationary policy. Under this assumption, the above metrics do not depend on the initial state under stationary deterministic policies, so one can write F^S

and G^S for active sets $S \subseteq N$. This allows a straightforward extension of the above indexability theory to the average criterion.

Regarding the algorithms discussed above, they apply without change to the average criterion, as the results presented in Section 5.1 show that the corresponding tableaux are simply the limits of the discounted tableaux as the discount rate goes to zero, and also outline how to evaluate the initial tableau. To properly extend the results, one must further assume that the active-set family \mathcal{F} satisfies that, for every $S \in \mathcal{F}$, the S -active policy is *unichain*, so it induces a single recurrent class on the embedded process X_k plus a class of transient states, which may be empty.

7. Numerical Experiments

This section discusses results of numerical experiments, based on implementations by the author of the aforementioned algorithms.

Comparing Runtimes of Index Algorithms

The runtime performance of an algorithm depends not only on its arithmetic operation count, but also on its memory-access patterns, which can actually be the dominant factor. Hence, to compare the algorithms considered herein, a numerical study has been conducted, using MATLAB implementations developed by the author. The experiments were run on a PC with an Intel Core i7-8700 CPU at 3.2 GHz with 16 GB of RAM using MATLAB 2020b under Windows 10 Enterprise. For the state space sizes $n = 1000, 2000, \dots, 15,000$, a discrete-time restless bandit instance was randomly constructed. Transition matrices were generated from random matrices with Uniform[0, 1] entries, dividing each row by its sum. Active rewards were randomly generated with Uniform[0, 1] entries, and passive rewards were zero. The discount factor was $\beta = 0.8$.

For each generated instance, the CP algorithm was first used to test for indexability and for PCL-indexability (by checking the signs of marginal resource metrics for the nested active-set family obtained). Since such tests were positive in each instance, the Whittle index was calculated using the CP, RP, and FP algorithms (taking $\mathcal{F} = 2^N$ in the latter).

Table 6 records the runtimes for the loop of each algorithm, without counting the initialization stage of computing the initial tableau, while Figure 1 plots them, along with cubic least-squares fitting curves. These results highlight that the FP algorithm, whose loop operation count is of $(2/3)n^3 + O(n^2)$, is indeed the fastest algorithm, followed by the CP and RP algorithms. Recall that $2n^3 + O(n^2)$ and $n^3 + O(n^2)$ are the loop operation counts for the CP and the RP algorithms.

Table 6. Runtimes (s) of index algorithms.

n	FP	RP	CP
1000	1.8	2.6	2.2
2000	14.8	23.6	19.6
3000	53.3	77.1	67.8
4000	121.4	193.4	156.3
5000	227.6	342.8	295.9
6000	433.2	647.9	541.8
7000	699.6	1034.8	862.3
8000	1118.3	1531.5	1280.9
9000	1530.8	2173.6	1822.3
10,000	2100.7	2919.3	2500.9
11,000	2687.2	3580.0	3277.7
12,000	3575.9	5055.7	4300.6
13,000	4747.4	6610.4	5539.2
14,000	5629.5	7923.1	6829.9
15,000	7254.1	8871.0	8250.5

The observed discrepancies between operation count complexity and actual runtime performance are explained by taking into account the memory-access patterns of the algorithms. Algorithm CP, which carries out conventional pivot steps, has efficient memory-access patterns, because the coefficient matrix \mathbf{A} is always updated as a memory block of contiguous storage. Yet, both algorithms RP and FP achieve a reduction in operations by operating on submatrices of \mathbf{A} , with resulting noncontiguous memory-access patterns that are time-consuming. However, for the FP algorithm, as the decrease in arithmetic operations is substantial, it compensates such inefficiencies, and comes out in practice as the fastest algorithm.

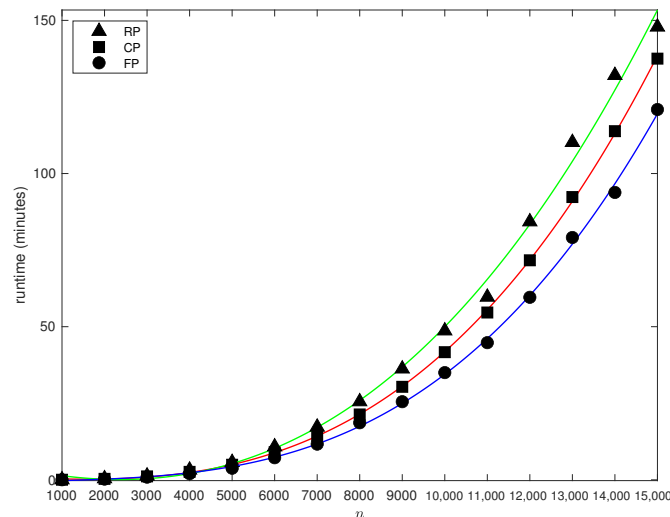


Figure 1. Runtimes of index algorithms vs. number of project states with cubic least-squares fits.

8. Discussion

This paper has presented a new algorithm for computing the Whittle index of a general finite-state semi-Markov restless bandit, based on an efficient implementation of the adaptive-greedy algorithmic scheme introduced in [3,24] for restless bandits, in which it was not specified how to evaluate certain metrics arising in the algorithm description. The algorithm extends to restless bandits the fast-pivoting implementation developed in [40] for classic (non-restless) bandits, and results from a similar approach, exploiting the structure of parametric simplex tableaux to reduce the operation count down to $(2/3)n^3 + O(n^2)$ —apart from the computational effort to compute the initial tableau. It is unlikely that such a complexity count can be improved for general restless bandits, as $(2/3)n^3 + O(n^2)$ is the complexity of Gaussian elimination for solving an $n \times n$ linear equation system, and it is shown in [3,24] that computing the Whittle index entails at least solving an equation system with such dimensions (albeit with an ordering of the states that generally is not known in advance).

The complexity of the proposed fast-pivoting algorithm is the best that has been reported in the literature. In contrast, the Whittle index algorithm recently presented in [41]—for the average criterion—has a complexity of $O(n^4 2^n)$, which reduces to $O(n^5)$ for one-dimensional indexable bandits provided it is known that threshold policies are optimal for them.

The new algorithm presented herein, whose implementation is straightforward, will be most useful for computing the Whittle index in complex models with large-scale, multi-dimensional state spaces for which closed index formulae cannot be derived, and an efficient computational approach is needed. Given the explosion of research interest in restless bandit models in the last decade, the proposed algorithm thus has the potential of becoming a useful tool, allowing researchers to expand the scope of Whittle's index policy to large-scale complex models.

Future research directions include developing efficient implementations with substantially lower complexity by exploiting the special structure of relevant model classes arising in applications, and testing the algorithm in large-scale real world models with real data. Another avenue of research

is to develop software implementations that improve the efficiency of the computationally costly block matrix operations required by the fast-pivoting algorithm.

9. Conclusions

To conclude, the findings of this paper can be summarized as follows:

- A new algorithm to compute the Whittle index of a general n -state semi-Markov restless bandit is presented, which can also be used to test numerically for indexability. After an initialization step, the algorithm computes the n index values in an n -step loop with a complexity of $(2/3)n^3 + O(n^2)$ arithmetic operations.
- The algorithm extends to Whittle's index the fast-pivoting $(2/3)n^3 + O(n^2)$ algorithm introduced by the author in [30] for the Gittins index of classic (non-restless) bandits, which also has the lowest complexity to date.
- The proposed algorithm has substantially better complexity than alternative algorithms proposed in the literature.
- The algorithm will be especially useful for computing the Whittle index in large-scale multi-dimensional models where the index cannot be derived in closed form and alternative algorithms will result in prohibitive computation times.

Funding: This research has been developed over a number of years, and has been funded by the Spanish Government under grants MEC MTM2004-02334 and PID2019-109196GB-I00/AEI/10.13039/501100011033. This research was also funded in part by the Comunidad de Madrid in the setting of the multi-year agreement with Universidad Carlos III de Madrid within the line of activity "Excelencia para el Profesorado Universitario", in the framework of the V Regional Plan of Scientific Research and Technological Innovation 2016–2020.

Acknowledgments: Preliminary early versions of this work were published by the author in the abridged proceedings [53] of the Second International Workshop on Tools for Solving Structured Markov Chains (SMCtools 2007), Nantes, France, 2007, and in the working paper [54].

Conflicts of Interest: The author declares no conflict of interest.

References

1. Puterman, M.L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*; Wiley: New York, NY, USA, 1994.
2. Whittle, P. Restless bandits: Activity allocation in a changing world. *J. Appl. Probab.* **1988**, *25A*, 287–298. [[CrossRef](#)]
3. Niño-Mora, J. Dynamic allocation indices for restless projects and queueing admission control: A polyhedral approach. *Math. Program.* **2002**, *93*, 361–413. [[CrossRef](#)]
4. Niño-Mora, J. Restless bandit marginal productivity indices, diminishing returns and optimal control of make-to-order/make-to-stock M/G/1 queues. *Math. Oper. Res.* **2006**, *31*, 50–84. [[CrossRef](#)]
5. Fu, J.; Moran, B.; Guo, J.; Wong, E.W.M.; Zukerman, M. Asymptotically optimal job assignment for energy-efficient processor-sharing server farms. *IEEE J. Sel. Areas Commun.* **2016**, *34*, 4008–4023. [[CrossRef](#)]
6. Niño-Mora, J. Computing an index policy for bandits with switching penalties. In Proceedings of the 2nd International Conference on Performance Evaluation Methodologies and Tools, Nantes, France, 22–27 October 2007. [[CrossRef](#)]
7. Qian, Y.; Zhang, C.; Krishnamachari, B.; Tambe, M. Restless poachers: Handling exploration-exploitation tradeoffs in security domains. In Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems, Singapore, 9–13 May 2016; pp. 123–131.
8. Borkar, V.S.; Pattathil, S. Whittle indexability in egalitarian processor sharing systems. *Ann. Oper. Res.* **2017**, *1*–21. [[CrossRef](#)]
9. Borkar, V.S.; Ravikumar, K.; Saboo, K. An index policy for dynamic pricing in cloud computing under price commitments. *Appl. Math.* **2017**, *44*, 215–245. [[CrossRef](#)]
10. Borkar, V.S.; Kasbekar, G.S.; Pattathil, S.; Shetty, P.Y. Opportunistic scheduling as restless bandits. *IEEE Trans. Control Netw. Syst.* **2018**, *5*, 1952–1961. [[CrossRef](#)]

11. Gerum, P.C.L.; Altay, A.; Baykal-Gursoy, M. Data-driven predictive maintenance scheduling policies for railways. *Transp. Res. Part C Emerg. Technol.* **2019**, *107*, 137–154. [[CrossRef](#)]
12. Abbou, A.; Makis, V. Group maintenance: A restless bandits approach. *INFORMS J. Comput.* **2019**, *31*, 719–731. [[CrossRef](#)]
13. Ayer, T.; Zhang, C.; Bonifonte, A.; Spaulding, A.C.; Chhatwal, J. Prioritizing hepatitis C treatment in US prisons. *Oper. Res.* **2019**, *67*, 853–873. [[CrossRef](#)]
14. Niño-Mora, J. Resource allocation and routing in parallel multi-server queues with abandonments for cloud profit maximization. *Comput. Oper. Res.* **2019**, *103*, 221–236. [[CrossRef](#)]
15. Fu, J.; Moran, B. Energy-efficient job-assignment policy with asymptotically guaranteed performance deviation. *IEEE/ACM Trans. Netw.* **2020**, *28*, 1325–1338. [[CrossRef](#)]
16. Hsu, Y.P.; Modiano, E.; Duan, L.J. Scheduling algorithms for minimizing age of information in wireless broadcast networks with random arrivals. *IEEE Trans. Mob. Comput.* **2020**, *19*, 2903–2915. [[CrossRef](#)]
17. Sun, J.Z.; Jiang, Z.Y.; Krishnamachari, B.; Zhou, S.; Niu, Z.S. Closed-form Whittle’s index-enabled random access for timely status update. *IEEE Trans. Commun.* **2020**, *68*, 1538–1551. [[CrossRef](#)]
18. Li, D.; Ding, L.; Connor, S. When to switch? Index policies for resource scheduling in emergency response. *Prod. Oper. Manag.* **2020**, *29*, 241–262. [[CrossRef](#)]
19. Papadimitriou, C.H.; Tsitsiklis, J.N. The complexity of optimal queuing network control. *Math. Oper. Res.* **1999**, *24*, 293–305. [[CrossRef](#)]
20. Weber, R.R.; Weiss, G. On an index policy for restless bandits. *J. Appl. Probab.* **1990**, *27*, 637–648. [[CrossRef](#)]
21. Weber, R.R.; Weiss, G. On an index policy for restless bandits. *Adv. Appl. Probab.* **1991**, *23*, 429–430. [[CrossRef](#)]
22. Ouyang, W.Z.; Eryilmaz, A.; Shroff, N.B. Downlink scheduling over Markovian fading channels. *IEEE/ACM Trans. Netw.* **2016**, *24*, 1801–1812. [[CrossRef](#)]
23. Verloop, I.M. Asymptotically optimal priority policies for indexable and nonindexable restless bandits. *Ann. Appl. Probab.* **2016**, *26*, 1947–1995. [[CrossRef](#)]
24. Niño-Mora, J. Restless bandits, partial conservation laws and indexability. *Adv. Appl. Probab.* **2001**, *33*, 76–98. [[CrossRef](#)]
25. Niño-Mora, J. A verification theorem for threshold-indexability of real-state discounted restless bandits. *Math. Oper. Res.* **2020**, *45*, 465–496. [[CrossRef](#)]
26. Niño-Mora, J. Marginal productivity index policies for scheduling a multiclass delay-/loss-sensitive queue. *Queueing Syst.* **2006**, *54*, 281–312. [[CrossRef](#)]
27. Niño-Mora, J. Dynamic priority allocation via restless bandit marginal productivity indices. *Top* **2007**, *15*, 161–198. [[CrossRef](#)]
28. Cao, J.; Nyberg, C. Linear programming relaxations and marginal productivity index policies for the buffer sharing problem. *Queueing Syst.* **2008**, *60*, 247–269. [[CrossRef](#)]
29. Huberman, B.A.; Wu, F. The economics of attention: Maximizing user value in information-rich environments. *Adv. Complex Syst.* **2008**, *11*, 487–496. [[CrossRef](#)]
30. Niño-Mora, J. A faster index algorithm and a computational study for bandits with switching costs. *INFORMS J. Comput.* **2008**, *20*, 255–269. [[CrossRef](#)]
31. Niño-Mora, J. Admission and routing of soft real-time jobs to multiclusters: Design and comparison of index policies. *Comput. Oper. Res.* **2012**, *39*, 3431–3444. [[CrossRef](#)]
32. Niño-Mora, J. Towards minimum loss job routing to parallel heterogeneous multiserver queues via index policies. *Eur. J. Oper. Res.* **2012**, *220*, 705–715. [[CrossRef](#)]
33. He, T.; Chen, S.; Kim, H.; Tong, L.; Lee, K.W. Scheduling parallel tasks onto opportunistically available cloud resources. In Proceedings of the 2012 IEEE Fifth International Conference on Cloud Computing, Honolulu, HI, USA, 24–29 June 2012; pp. 180–187.
34. Menner, M.; Zeilinger, M.N. A user comfort model and index policy for personalizing discrete controller decisions. In Proceedings of the 2018 European Control Conference (ECC), Limassol, Cyprus, 12–15 June 2018; pp. 1759–1765.
35. Dance, C.R.; Silander, T. Optimal policies for observing time series and related restless bandit problems. *J. Mach. Learn. Res.* **2019**, *20*, 35.
36. Klimov, G.P. Time-sharing service systems. I. *Theory Probab. Appl.* **1974**, *19*, 532–551. [[CrossRef](#)]
37. Bertsimas, D.; Niño-Mora, J. Conservation laws, extended polymatroids and multiarmed bandit problems; a polyhedral approach to indexable systems. *Math. Oper. Res.* **1996**, *21*, 257–306. [[CrossRef](#)]

38. Niño-Mora, J. Conservation laws and related applications. In *Wiley Encyclopedia of Operations Research and Management Science*; Wiley: New York, NY, USA, 2011. [[CrossRef](#)]
39. Niño-Mora, J. Klimov's model. In *Wiley Encyclopedia of Operations Research and Management Science*; Cochran, J.J., Ed.; Wiley: New York, NY, USA, 2011. [[CrossRef](#)]
40. Niño-Mora, J. A $(2/3)n^3$ fast-pivoting algorithm for the Gittins index and optimal stopping of a Markov chain. *INFORMS J. Comput.* **2007**, *19*, 596–606. [[CrossRef](#)]
41. Ayesta, U.; Gupta, M.K.; Verloop, I.M. On the computation of Whittle's index for Markovian restless bandits. *Math. Methods Oper. Res.* **2020**, 1–30. [[CrossRef](#)]
42. Zhao, Q. *Multi-Armed Bandits: Theory and Applications to Online Learning in Networks*; Morgan & Claypool: San Rafael, CA, USA, 2020.
43. Varaiya, P.P.; Walrand, J.C.; Buyukkoc, C. Extensions of the multiarmed bandit problem: The discounted case. *IEEE Trans. Automat. Control* **1985**, *30*, 426–439. [[CrossRef](#)]
44. Chen, Y.R.; Katehakis, M.N. Linear programming for finite state multi-armed bandit problems. *Math. Oper. Res.* **1986**, *11*, 180–183. [[CrossRef](#)]
45. Kallenberg, L.C.M. Computation of the Gittins index. *Math. Oper. Res.* **1986**, *11*, 184–186. [[CrossRef](#)]
46. Katehakis, M.N.; Veinott, A.F., Jr. The multi-armed bandit problem: Decomposition and computation. *Math. Oper. Res.* **1987**, *12*, 262–268. [[CrossRef](#)]
47. Katta, A.K.; Sethuraman, J. The multi-armed bandit problem: Decomposition and computation. *SIAM J. Discret. Math.* **2004**, *18*, 110–113. [[CrossRef](#)]
48. Sonin, I.M. A generalized Gittins index for a Markov chain and its recursive calculation. *Stat. Probab. Lett.* **2008**, *78*, 1526–1533. [[CrossRef](#)]
49. Gittins, J.C. Bandit processes and dynamic allocation indices. *J. R. Stat. Soc. Ser. B* **1979**, *41*, 148–177. [[CrossRef](#)]
50. Larrañaga, M.; Ayesta, U.; Verloop, I.M. Dynamic control of birth-and-death restless bandits: Application to resource-allocation problems. *IEEE/ACM Trans. Netw.* **2016**, *24*, 3812–3825. [[CrossRef](#)]
51. Kallenberg, L.C.M. Finite state and action MDPs. In *Handbook of Markov Decision Processes*; Kluwer: Boston, MA, USA, 2002; pp. 21–87.
52. Gass, S.; Saaty, T. The computational algorithm for the parametric objective function. *Nav. Res. Logist. Q.* **1955**, *2*, 39–45. [[CrossRef](#)]
53. Niño-Mora, J. Characterization and computation of restless bandit marginal productivity indices. In *SMCtools '07: Proceedings from the Second International Workshop on Tools for Solving Structured Markov Chains*; ICST: Brussels, Belgium, 2007.
54. Niño-Mora, J. Characterization and computation of restless bandit marginal productivity indices. Universidad Carlos III de Madrid (UC3M) Working Papers, Statistics and Econometrics 07–11. 2007. Available online: <http://hdl.handle.net/10016/796> (accessed on 15 November 2020).

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



© 2020 by the author. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).