

Low-Complexity Algorithm for Restless Bandits with Imperfect Observations

Keqin Liu^{1,2*}, Richard Weber³ and Chengzhong Zhang²

^{1*}Department of Mathematics, Nanjing University, Nanjing, 210093, China.

²National Center for Applied Mathematics, Nanjing, 210093, China.

³Department of Mathematics, University of Cambridge, Cambridge, CB3 0WB, UK.

*Corresponding author(s). E-mail(s): kqliu@nju.edu.cn;

Contributing authors: rrw1@cam.ac.uk; 171840780@smail.nju.edu.cn;

Abstract

We consider a class of restless bandit problems that finds a broad application area in reinforcement learning and stochastic optimization. We consider N independent discrete-time Markov processes, each of which had two possible states: 1 and 0 ('good' and 'bad'). Only if a process is both in state 1 and observed to be so does reward accrue. The aim is to maximize the expected discounted sum of returns over the infinite horizon subject to a constraint that only M ($< N$) processes may be observed at each step. Observation is error-prone: there are known probabilities that state 1 (0) will be observed as 0 (1). From this one knows, at any time t , a probability that process i is in state 1. The resulting system may be modeled as a restless multi-armed bandit problem with an information state space of uncountable cardinality. Restless bandit problems with even finite state spaces are PSPACE-HARD in general. We propose a novel approach for simplifying the dynamic programming equations of this class of restless bandits and develop a low-complexity algorithm that achieves a strong performance and is readily extensible to the general restless bandit model with observation errors. Under certain conditions, we establish the existence (indexability) of Whittle index and its equivalence to our algorithm. When those conditions do not hold, we show by numerical experiments the near-optimal performance of our algorithm in the general parametric space. Furthermore, we theoretically prove the optimality of our algorithm for homogeneous systems.

Keywords: restless bandits, continuous state space, observation errors, value functions, index policy

1 Introduction

The exploration-exploitation (EE) dilemma is well posed in optimization-over-time problems and mathematically modeled in various forms for reinforcement learning, to which a major category, *multi-armed bandits* (MAB) belongs. In the classical MAB model, a player chooses one out of N statistically independent arms to pull and possibly accrues reward determined by the state of the chosen arm which transits to a new state according to a known Markovian rule (Gittins et al., 2011). The states of other arms remain frozen. The objective is to maximize the expected total discounted reward summed over times $t = 1, 2, \dots$ to an infinite time horizon with discount factor $\beta \in (0, 1)$,

$$\max_{\pi \in \Pi} \mathbb{E}_{\pi} \left[\sum_{t=1}^{\infty} \beta^{t-1} R(t) \middle| S(t), A(t) \right]. \quad (1)$$

The expectation is taken under some policy π , chosen from the set of all feasible policies Π ; $S(t)$ is the joint state of all arms at time t , $A(t) \in \{1, 2, \dots, N\}$ is the arm pulled at t and $R(t)$ is the reward thus obtained. It follows from standard theory of Markov decision processes that there must exist an optimal stationary policy π^* , independent of time t . If each arm's state space has cardinality K , then the joint state space has size K^N . This means that a dynamic programming solution to the problem will have running time that grows geometrically as the number of arms increases. Gittins (1979) solved the problem by showing the optimality of an *index* policy, *i.e.*, for each state of each arm there exists an index depending solely on the parameters of that arm; it is then optimal at each time to choose the arm whose state has highest index. The running time of the Gittins index policy grows only linearly with the number of arms as they are decoupled when computing the index function (of the states of each arm). Whittle (1988) generalized Gittins index to the *restless* MAB model in which those arms that are not chosen may also change states and produce reward. Whittle's generalization has been shown to perform very well in theoretical and numerical studies (see, *e.g.*, Brown and Smith (2020), Chen et al. (2021), Gast et al. (2021), Hu and Frazier (2017), Liu and Zhao (2010), Weber and Weiss (1990, 1991), Zayas-Cabán et al. (2019)). In general, however, it is difficult to theoretically establish the condition that is necessary for the Whittle index to exist (so called *indexability*) and to solve for the closed-form Whittle index when it does exist due to the curse of dimensionality (see Sec. 2.1). For finite-state models, Gast et al. (2023) proposed an efficient algorithm to numerically test indexability and compute Whittle index. A natural question to ask is if one can transform a bandit with a continuous-state space into a finite-state one by discretization. Unfortunately, we found that how fine the discretization needs to be given the system parameters is itself a difficult problem. Furthermore, a finer discretization inevitably leads to a larger state transition matrix with increased algorithmic complexity and unpredictability of the steady-state system performance. This motivates us to consider alternative approaches to deal with bandits with infinite state spaces as addressed in this paper. In terms of searching for the general optimal policy, Papadimitriou and Tsitsiklis (1999) have shown that the

restless MAB with a finite state space is PSPACE-HARD in general. With an infinite state space, a restless bandit problem yields practical difficulties in implementing purely numerical methods as discussed above. In this paper, we show that for our particular Markovian model with an infinite state space, Whittle index policy can be efficiently implemented with satisfactory accuracy after theoretical analysis on the rich mathematical structure of the problem.

In this paper, we extend the work in [Liu and Zhao \(2010\)](#) (for a perfect observation model) and [Liu et al. \(2010\)](#) (for the myopic policy on stochastically identical arms) to build a near-optimal algorithm with low complexity for a class of restless bandits with an infinite state space and an imperfect observation model. Our model also belongs to the general framework of partially observable Markov decision processes (POMDP) ([Sondik, 1978](#)). Consider N processes each of which evolves on a 2-state Markov chain whose state is observed if and only if the process is chosen. Furthermore, the observation is *error-prone*: state 1 may be observed as 0 and vice versa. Each process is referred to as an arm. At time t , the player obtains reward of amount B_n if and only if arm n is currently chosen and accurately observed in state 1. Under resource constraints, the player’s goal is to select M ($M < N$) arms at each time and maximize the long-term reward. By formulating *the belief vector* as the system state for decision making, we show that the indexability is satisfied under certain conditions. Furthermore, we propose an efficient algorithm to compute an approximate Whittle index that achieves a near-optimal performance in general, even if the conditions for indexability do not hold.

[Rahul et al. \(2018\)](#), [Varun et al. \(2018\)](#) and [Kesav et al. \(2019\)](#) considered similar models (except for some nuances) and established indexability under much stricter conditions on the system parameters. For example, all the three papers require that the Markov transition probabilities of each arm have differences bounded by $1/3$ while we do not need any restriction on the transition probabilities. Furthermore, the three papers require that the discount factor β is less than $1/3$ while our condition on β is a more relaxed closed-form expression of the system parameters. In terms of the computation of the Whittle index, their algorithm is a direct application of general reinforcement learning while our algorithmic framework is based on the detailed analysis of the value functions with a quick convergence to the exact Whittle index function. In this paper, we also plot the performance of the optimal policy to demonstrate the near-optimality of our algorithm in addition to the comparison with the myopic policy. For homogeneous systems (stochastically identical arms), we show that our algorithm is equivalent to the myopic policy and theoretically prove its optimality under certain conditions. [Wang et al. \(2018\)](#) also considered our model and assumed the optimality of the threshold policy for a single arm while using a very coarse linear approximation to compute the Whittle index function (the key step (a) for the second equality in the proof of Lemma 6 in [Wang et al. \(2018\)](#) is incorrect). In this paper, we rigorously prove the optimality of the threshold policy for a single arm and establish the indexability under certain conditions and subsequently construct an efficient algorithm for computing the Whittle index function with arbitrary precision with its optimality numerically verified in general and formally proved for a class of homogeneous systems.

The rest of this paper is organized as follows: Sec. 2 presents our problem formulation and main results on the optimal threshold policy and indexability. Sec. 3 presents our algorithm with design details. Sec. 4 presents a theoretic proof of the optimality for homogeneous arms. Sec. 5 concludes this paper and provides some future research directions in this field.

2 Main Results

Consider a restless MAB having N internal 2-state Markov chains (arms) of potentially *different* transition probabilities. At each time t , the player chooses to observe the states of M ($< N$) arms. Let $S \in \{0 \text{ (bad)}, 1 \text{ (good)}\}$ denote the current state of an arm and let O denote its observation outcome (detection outcome). The error probabilities are $\delta = \Pr(O = 1 \mid S = 0)$ and $\epsilon = \Pr(O = 0 \mid S = 1)$, *i.e.*, the probabilities of miss detection and false alarm, respectively, in the observation model. In Levy (2008), it was shown that the error probabilities δ and ϵ follow the curve of receiver operating characteristics (ROC) under the optimal detector that makes $1 - \delta$ a concave increasing function from 0 to 1 over ϵ . This matches the intuition that making a detector more sensitive will reduce δ but increase ϵ . Since the optimal detector design is a solved problem and not the focus of this paper, we simply assume that δ and ϵ are given. If arm n in state $S = 1$ is observed in state 1 (*i.e.*, $S = O = 1$), then the player accrues B_n units of reward from this arm. One of many application examples of this observation model is to *cognitive radios*, where a secondary user aims to utilize a frequency band (channel/arm) currently unused by the primary users. Due to energy and policy constraints on the sensor of the secondary user, only a subset of channels can be sensed at each time and if any of them is sensed idle ($O = 1$), the user can send certain packets over it to its receiver and obtain an ACK (acknowledgement) in the end of the time slot if the channel is indeed idle ($S = 1$); otherwise no ACK from this channel would be received. Then the reward B_n is just the bandwidth of channel n . Clearly, the hard constraint here should be on the miss detection probability δ to guarantee the satisfaction of the primary users, *i.e.*, the disturbance (when a secondary user senses a busy channel as idle and subsequently sends data over it) to the primary users should be capped.

2.1 System Model and Belief Vector

At each discrete time t , the internal state (0/1) of an arm cannot be observed before deciding whether or not to observe the arm. Therefore, we cannot use the states of the Markov chains as the system state for decision making. Applying the general POMDP theory to our model the *belief state vector* consisting of probabilities that arms are in state 1 given all past observations is a sufficient statistics for making future decisions (Sondik, 1978):

$$\boldsymbol{\omega}(t) = (\omega_1(t), \omega_2(t), \dots, \omega_N(t)), \quad (2)$$

$$\omega_n(t) = \Pr(S_n(t) = 1 \mid \text{past observations on arm } n), \quad \forall n \in \{1, \dots, N\}, \quad (3)$$

where $\omega_n(t)$ is the *belief state* of arm n at time t and $S_n(t)$ its internal state. According to the Bayes' rule, the belief state (of any arm) itself evolves as a Markov chain with an infinite state space:

$$\omega_n(t+1) = \begin{cases} p_{11}^{(n)}, & n \in A(t), \text{ACK}_n(t) \ (S_n(t) = 1, O_n(t) = 1) \\ \mathcal{T}_n\left(\frac{\epsilon\omega_n(t)}{\epsilon\omega_n(t)+1-\omega_n(t)}\right), & n \in A(t), \text{no ACK}_n(t) \\ \mathcal{T}_n(\omega_n(t)), & n \notin A(t) \end{cases}, \quad (4)$$

$$\mathcal{T}_n(\omega_n(t)) = \omega_n(t)p_{11}^{(n)} + (1 - \omega_n(t))p_{01}^{(n)}, \quad (5)$$

where $A(t) \subset \{1, 2, \dots, N\}$ is the set of arms chosen at time t with $|A(t)| = M$, $S_n(t)$ and $O_n(t)$ are respectively the state and observation from arm n at time t if $n \in A(t)$, $\text{ACK}_n(t)$ the acknowledgement of successful utilization of arm n for slot t , $\mathcal{T}_n(\cdot)$ the one-step belief update operator without observation, and $\mathbf{P}^{(n)} = \{p_{ij}^{(n)}, i, j \in \{0, 1\}\}$ the transition matrix of the internal Markov chain of arm n . Note that when $\epsilon = 0$, all three expressions of the next belief state in (4) become linear functions and the problem is reduced to that in [Liu and Zhao \(2010\)](#) for perfect observations where the value functions for dynamic programming can be directly solved in closed-form. Later, we will see that when $\epsilon \neq 0$, the value functions are very difficult to analyze except for a few general properties. Without a fine-grain analysis on the value functions, the indexability and Whittle index cannot be established for our model. Our approach is to start with a finite time horizon and then utilize the backward induction (on time horizon) methodology until we take limits of corresponding functions as the time horizon goes to infinity. The key idea is to obtain analytic bounds on the value functions and their derivatives as functions of the system parameters instead of just numerical computations of these functions. We will show that these bounds, together with some detailed properties of the value functions, are sufficient for our purpose. From (5), the k -step belief update of an unobserved arm for k consecutive slots starting from any belief state ω is

$$\mathcal{T}_n^k(\omega) = \frac{p_{01}^{(n)} - (p_{11}^{(n)} - p_{01}^{(n)})^k(p_{01}^{(n)} - (1 + p_{01}^{(n)} - p_{11}^{(n)})\omega)}{1 + p_{01}^{(n)} - p_{11}^{(n)}}. \quad (6)$$

For simplicity of notations, we denote $\mathcal{T}_n^1(\cdot)$ by $\mathcal{T}_n(\cdot)$.

At time $t = 1$, the initial belief state $\omega_n(1)$ of arm n can be set as the stationary distribution $\omega_{n,o}$ of the internal Markov chain*:

$$\omega_n(1) = \omega_{n,o} = \lim_{k \rightarrow \infty} \mathcal{T}_n^k(\omega') = \frac{p_{01}^{(n)}}{p_{01}^{(n)} + p_{10}^{(n)}}, \quad (7)$$

*Here we assume the internal Markov chain with transition matrix $\mathbf{P}^{(n)}$ is irreducible and aperiodic.

where $\omega_{n,o}$ is the unique solution to $\mathcal{T}_n(\omega) = \omega$ and $\omega' \in [0, 1]$ an arbitrary probability. Given the initial belief vector $\boldsymbol{\omega}(1) = (\omega_1(1), \omega_2(1), \dots, \omega_N(1))$, we arrive at the following *constrained* optimization problem:

$$\max_{\pi: \boldsymbol{\omega}(t) \rightarrow A(t)} \mathbb{E}_\pi \left[\sum_{t=1}^{\infty} \beta^{t-1} R(t) \middle| \boldsymbol{\omega}(1) \right], \quad (8)$$

$$\text{subject to } |A(t)| = M, \quad t \geq 1. \quad (9)$$

Now the decision problem has a countable state space as modelled by the belief vector for a fixed initial $\boldsymbol{\omega}(1)$ and an uncountable state space for an arbitrarily chosen $\boldsymbol{\omega}(1)$. It is clear that fixing $\boldsymbol{\omega}(1)$, the action-dependent belief vector $\boldsymbol{\omega}(t)$ takes possible values growing geometrically with time t , leading to a high-complexity in solving the problem; this is the so-called *curse of dimensionality*. In the following, we adopt Whittle's original idea of Lagrangian relaxation to decouple arms for an index policy and show some crucial properties of the value functions of a single arm.

2.2 Arm Decoupling by Lagrangian Relaxation

$$\max_{\pi: \boldsymbol{\omega}(t) \rightarrow A(t)} \mathbb{E}_\pi \left[\sum_{t=1}^{\infty} \beta^{t-1} \sum_{n=1}^N \mathbb{1}_{(n \in A(t))} \cdot S_n(t) \cdot O_n(t) \cdot B_n \middle| \boldsymbol{\omega}(1) \right] \quad (10)$$

$$\text{subject to } \mathbb{E}_\pi \left[\sum_{t=1}^{\infty} \beta^{t-1} \sum_{n=1}^N \mathbb{1}_{(n \notin A(t))} \middle| \boldsymbol{\omega}(1) \right] = \frac{N - M}{1 - \beta}. \quad (11)$$

Clearly constraint (11) is a relaxation on the player's action $A(t)$ from (9). Applying the Lagrangian multiplier μ to constraint (11), we arrive at the following *unconstrained* optimization problem:

$$\max_{\pi: \boldsymbol{\omega}(t) \rightarrow A(t)} \mathbb{E}_\pi \left[\sum_{t=1}^{\infty} \beta^{t-1} \sum_{n=1}^N [\mathbb{1}_{(n \in A(t))} S_n(t) O_n(t) B_n + \mu \cdot \mathbb{1}_{(n \notin A(t))}] \middle| \boldsymbol{\omega}(1) \right]. \quad (12)$$

Fixing μ , the above optimization is equivalent to N *independent* unconstrained optimization problem as shown below: for each $n \in \{1, 2, \dots, N\}$,

$$\max_{\pi: \omega_n(t) \rightarrow \{0, 1\}} \mathbb{E}_\pi \left[\sum_{t=1}^{\infty} \beta^{t-1} [\mathbb{1}_{(n \in A(t))} S_n(t) O_n(t) B_n + \mu \cdot \mathbb{1}_{(n \notin A(t))}] \middle| \omega_n(1) \right]. \quad (13)$$

Here π is a single-arm policy that maps the belief state of the arm to the binary action $u = 1$ (chosen/activated) or $u = 0$ (unchosen/made passive). It is thus sufficient to consider a single arm for solving problem (12). For simplicity, we will drop the subscript n in consideration of a single-armed bandit problem without loss of generality.

Let $V_{\beta,m}(\omega)$ denote the value of (13) with $\mu = m$ and $\omega_n(1) = \omega$, it is straightforward to write out the dynamic equation of the single-armed bandit problem as follows:

$$V_{\beta,m}(\omega) = \max\{V_{\beta,m}(\omega; u = 1); V_{\beta,m}(\omega; u = 0)\}, \quad (14)$$

where $V_{\beta,m}(\omega; u = 1)$ and $V_{\beta,m}(\omega; u = 0)$ denote, respectively, the maximum expected total discounted reward that can be obtained if the arm is activated or made passive at the current belief state ω , followed by an *optimal* policy in subsequent slots. Since we consider the infinite-horizon problem, a stationary optimal policy can be chosen and the time index t is not needed in (14). Define the nonlinear operator $\phi(\cdot)$ as

$$\phi(\omega) = \frac{\epsilon\omega}{\epsilon\omega + 1 - \omega}.$$

It is easy to see that $\mathcal{T} \circ \phi(\cdot)$ is Lipschitz continuous on $[0, 1]$:

$$\left| \mathcal{T}\left(\frac{\epsilon\omega}{\epsilon\omega + 1 - \omega}\right) - \mathcal{T}\left(\frac{\epsilon\omega'}{\epsilon\omega' + 1 - \omega'}\right) \right| = \left| \frac{p_{11}\epsilon\omega + (1-\omega)p_{01}}{\epsilon\omega + 1 - \omega} - \frac{p_{11}\epsilon\omega' + (1-\omega')p_{01}}{\epsilon\omega' + 1 - \omega'} \right| \quad (15)$$

$$= \left| \frac{\epsilon(p_{11} - p_{01})(\omega - \omega')}{(1 - (1-\epsilon)\omega)(1 - (1-\epsilon)\omega')} \right| \leq \frac{|p_{11} - p_{01}|}{\epsilon} |\omega - \omega'|. \quad (16)$$

We assume that $\epsilon \neq 0$ (otherwise the problem is reduced to that considered in [Liu and Zhao, 2010](#)) and $p_{11} \neq p_{01}$ (otherwise the belief update is independent of observations or actions and the problem becomes trivial). Without loss of generality, set $B = 1$. We have

$$\begin{cases} V_{\beta,m}(\omega; u = 1) = (1 - \epsilon)\omega + \beta[(1 - \epsilon)\omega V_{\beta,m}(p_{11}) \\ \quad + (1 - (1 - \epsilon)\omega)V_{\beta,m}(\mathcal{T}(\phi(\omega)))], \\ V_{\beta,m}(\omega; u = 0) = m + \beta V_{\beta,m}(\mathcal{T}(\omega)). \end{cases} \quad (17)$$

Define *passive set* $P(m)$ as the set of all belief states in which taking the passive action $u = 0$ is optimal:

$$P(m) \triangleq \{\omega : V_{\beta,m}(\omega; u = 1) \leq V_{\beta,m}(\omega; u = 0)\}. \quad (18)$$

Notice that the immediate reward under the active action cannot exceed 1 so it is optimal to always make the arm passive if the immediate reward under the passive action exceeds 1 ($m \geq 1$). This is because we would obtain the maximum immediate reward equal to m at each time step regardless of the state transitions in this case. On the other hand, if $m < -1/(1 - \beta)$ then the optimal action must be to activate the arm. To see this, note that the total discounted reward by any policy consists of two parts: the reward under the active actions and the reward under the passive actions. If the optimal policy is to make the arm passive now when $m < -1/(1 - \beta)$, then the reward obtained under the future active actions by this policy must be greater than $1/(1 - \beta)$ otherwise the total discounted reward would be negative, contradicting the optimality of the policy since any policy that always activates the arm achieves a nonnegative total discounted reward so should the optimal policy. This is again

a contradiction since the total discounted reward under the active actions is upper bounded by $1/(1-\beta)$. Consequently, the passive set $P(m)$ changes from the empty set to the closed interval $[0, 1]$ as m increases from $-\infty$ to ∞ . However, such change may not be monotonic as m increases. But if $P(m)$ does increase monotonically with m , then for each value ω of the belief state, one can define the unique m that makes it join $P(m)$ and stay in the set forever. Intuitively, such m measures in a well-ordered manner the attractiveness of activating the arm in belief state ω compared to other belief states: the larger is the m that is required for it to be passive, the more is the incentive to activate the arm in belief state ω , even in the problem without m . This Lagrangian multiplier m is thus called ‘subsidy for passivity’ by Whittle who formalized the following definition of *indexability* and *Whittle index* (Whittle, 1988).

Definition 1. *A restless multi-armed bandit is indexable if for each single-armed bandit in a problem with subsidy m for passivity, the set of arm states $P(m)$ in which passivity is optimal increases monotonically from the empty set to the whole state space as m increases from $-\infty$ to $+\infty$. Under indexability, the Whittle index of an arm state is defined as the infimum subsidy m such that the state remains in the passive set.*

For our model in which the arm state is given by the belief vector, the indexability is equivalent to the following:

$$\begin{aligned} & \text{If } V_{\beta,m}(\omega; u = 1) \leq V_{\beta,m}(\omega; u = 0), \text{ then} \\ & \forall m' > m, \quad V_{\beta,m'}(\omega; u = 1) \leq V_{\beta,m'}(\omega; u = 0). \end{aligned} \quad (19)$$

Under indexability, the Whittle index $W(\omega)$ of arm state ω is defined as

$$W(\omega) \triangleq \inf\{m : V_{\beta,m}(\omega; u = 1) \leq V_{\beta,m}(\omega; u = 0)\}. \quad (20)$$

In the following we derive useful properties of the value functions $V_{\beta,m}(\omega; u = 1)$, $V_{\beta,m}(\omega; u = 0)$ and $V_{\beta,m}(\omega)$. Our strategy is to first establish those properties for finite horizons and then extend them to the infinite horizon by the uniform convergence of the value functions of the former to the latter. Define the T -horizon value function $V_{1,T,\beta,m}(\omega)$ as the maximum expected total discounted reward achievable over the next T time slots starting from the initial belief state ω . Then

$$V_{1,T,\beta,m}(\omega) = \max\{V_{1,T,\beta,m}(\omega; u = 1); V_{1,T,\beta,m}(\omega; u = 0)\}, \quad (21)$$

where $V_{1,T,\beta,m}(\omega; u = 1)$ and $V_{1,T,\beta,m}(\omega; u = 0)$ denote, respectively, the maximum expected total discounted reward achievable given the initial active and passive actions over the next T time slots starting from the initial belief state ω :

$$\begin{aligned} V_{1,T,\beta,m}(\omega; u = 1) &= (1-\epsilon)\omega + (1-\epsilon)\omega\beta V_{1,T-1,\beta,m}(p_{11}) \\ &\quad + (1-(1-\epsilon)\omega)\beta V_{1,T-1,\beta,m}\left(\mathcal{T}\left(\frac{\epsilon\omega}{1-(1-\epsilon)\omega}\right)\right), \end{aligned} \quad (22)$$

$$V_{1,T,\beta,m}(\omega; u = 0) = m + \beta V_{1,T-1,\beta,m}(\mathcal{T}(\omega)), \quad (23)$$

$$V_{1,0,\beta,m}(\cdot) \equiv 0. \quad (24)$$

From the above recursive equations, we can analyze $V_{1,T,\beta,m}(\omega)$ by backward induction on T . It is easy to see that for any ω ,

$$V_{1,1,\beta,m}(\omega; u = 1) = (1 - \epsilon)\omega, \quad V_{1,1,\beta,m}(\omega; u = 0) = m. \quad (25)$$

Therefore $V_{1,T,\beta,m}(\omega)$ is the maximum of two linear equations and thus piecewise linear and convex for $T = 1$ (in both ω and m). Assume that $V_{1,T-1,\beta,m}(\omega)$ is piecewise linear and convex. The Bayes' rule shows that the following term

$$(1 - (1 - \epsilon)\omega)\beta V_{1,T-1,\beta,m}\left(\mathcal{T}\left(\frac{\epsilon\omega}{1 - (1 - \epsilon)\omega}\right)\right) \quad (26)$$

is piecewise linear and convex since the leading coefficient $(1 - (1 - \epsilon)\omega)$ also appears as the denominator of the argument of the linear operator \mathcal{T} in $V_{1,T-1,\beta,m}(\cdot)$ assumed to be piecewise linear and convex by the induction hypothesis. Henceforth, the recursive equation set (22) and (23) shows that $V_{1,T,\beta,m}(\omega)$ is the maximum of two convex and piecewise linear functions and thus piecewise linear and convex for any $T > 1$ (in both ω and m). Motivated by the Lipschitz continuity of $\mathcal{T} \circ \phi$, we show in Lemma 2 that $V_{1,T,\beta,m}(\omega)$ is also Lipschitz continuous under certain conditions. In the following, we first establish a monotonic property of $V_{1,T,\beta,m}(\omega)$ in the case of $p_{11} > p_{01}$ (positively correlated Markov chain).

Lemma 1. *If $p_{11} > p_{01}$, then $V_{1,T,\beta,m}(\omega)$ is monotonically increasing with $\omega \in [0, 1]$ for any $T \geq 1$.*

Proof. Since $V_{1,T,\beta,m}(\omega)$ is piecewise linear, it is differentiable almost everywhere except on a null set (under the Lebesgue measure on \mathbb{R}) consisting of finite points among which both the left and right derivatives at any point exist but not equal. To prove that the continuous function $V_{1,T,\beta,m}(\omega)$ is monotonically increasing with ω , we only need to show

$$V'_{1,T,\beta,m}(\omega) \geq 0, \quad \forall \omega \in (0, 1), \quad (27)$$

where $V'_{1,T,\beta,m}(\omega)$ denotes the *right* derivative of $V_{1,T,\beta,m}(\cdot)$ as a function of the belief state with m fixed. From (25), the value function $V_{1,1,\beta,m}(\omega) = \max\{(1 - \epsilon)\omega, m\}$ is monotonically increasing with nonnegative right derivative $1 - \epsilon$ or 0. Assume (27) is true for $T \geq 1$, then for $T + 1$ we have $V_{1,T+1,\beta,m}(\omega) = \max\{f_T(\omega), g_T(\omega)\}$ with

$$\begin{aligned} f_T(\omega) &\triangleq (1 - \epsilon)\omega + (1 - \epsilon)\omega\beta V_{1,T,\beta,m}(p_{11}) \\ &\quad + (1 - (1 - \epsilon)\omega)\beta V_{1,T,\beta,m}(\mathcal{T} \circ \phi(\omega)), \\ g_T(\omega) &\triangleq m + \beta V_{1,T,\beta,m}(\mathcal{T}(\omega)). \end{aligned} \quad (28)$$

From the above, we have

$$\begin{aligned}
f'_T(\omega) &= (1 - \epsilon) + (1 - \epsilon)\beta V_{1,T,\beta,m}(p_{11}) - (1 - \epsilon)\beta V_{1,T,\beta,m}(\mathcal{T} \circ \phi(\omega)) \\
&\quad + V'_{1,T,\beta,m}(\mathcal{T} \circ \phi(\omega)) \frac{\epsilon\beta(p_{11}-p_{01})}{1-(1-\epsilon)\omega}, \\
g'_T(\omega) &= \beta(p_{11} - p_{01})V'_{1,T,\beta,m}(\mathcal{T}(\omega)),
\end{aligned} \tag{29}$$

where $f'_T(\cdot)$, $g'_T(\cdot)$ and $V'_{1,T,\beta,m}(\cdot)$ denote the right derivatives of the corresponding functions. We have used the fact that $\phi(\cdot)$ is monotonically increasing and when $p_{11} > p_{01}$, $\mathcal{T}(\cdot)$ is also monotonically increasing and that

$$\phi'(\omega) = \frac{\epsilon}{(1 - (1 - \epsilon)\omega)^2}. \tag{30}$$

By the induction hypothesis and (29), if $p_{11} > p_{01}$ then $g_T(\omega)$ is monotonically increasing (since $g'_T(\omega) \geq 0$) and

$$\begin{aligned}
f'_T(\omega) &= (1 - \epsilon) + (1 - \epsilon)\beta V_{1,T,\beta,m}(p_{11}) - (1 - \epsilon)\beta V_{1,T,\beta,m}(\mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega})) \\
&\quad + (1 - (1 - \epsilon)\omega)\beta V'_{1,T,\beta,m}(\mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega}))(\mathcal{T} \circ \phi)'(\omega) \\
&\geq (1 - \epsilon) + \frac{\epsilon\beta(p_{11}-p_{01})}{1-(1-\epsilon)\omega} V'_{1,T,\beta,m}(\mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega})) > 0,
\end{aligned} \tag{31}$$

where both the first and second inequalities are due to the monotonically increasing property of $V_{1,T,\beta,m}(\cdot)$ under the assumption that $p_{11} > p_{01}$ by our induction hypothesis and

$$p_{01} \leq \mathcal{T}(\omega) \leq p_{11}, \quad 0 \leq \phi(\omega) \leq 1, \quad \forall \omega \in [0, 1]. \tag{32}$$

This proves the monotonically increasing property of $f_T(\omega)$. Thus $V_{1,T+1,\beta,m}(\omega) = \max\{f_T(\omega), g_T(\omega)\}$ is also monotonically increasing and the proof by induction is finished. \square

Now we show that under a constraint on the discount factor $\beta \in (0, 1)$, the value function $V_{1,T,\beta,m}(\omega)$ is a Lipschitz function:

Lemma 2. *Suppose the discount factor $\beta \in (0, 1)$ satisfies*

$$\beta < \frac{1}{(2-\epsilon)|p_{11}-p_{01}|}. \tag{33}$$

Then $\forall T \geq 1$ and $\forall \omega, \omega' \in [0, 1]$,

$$|V_{1,T,\beta,m}(\omega) - V_{1,T,\beta,m}(\omega')| \leq C|\omega - \omega'|, \quad \text{where} \tag{34}$$

$$C = \frac{1-\epsilon}{1-(2-\epsilon)\beta|p_{11}-p_{01}|}. \tag{35}$$

Proof. We prove this by induction. Without loss of generality, assume $\omega < \omega'$. For the case of $T = 1$,

$$|V_{1,1,\beta,m}(\omega) - V_{1,1,\beta,m}(\omega')| = \begin{cases} 0, & m \geq (1-\epsilon)\omega' \\ (1-\epsilon)\omega' - m, & \text{if } (1-\epsilon)\omega \leq m < (1-\epsilon)\omega' \\ (1-\epsilon)|\omega - \omega'|, & \text{if } m < (1-\epsilon)\omega \end{cases}$$

Thus $|V_{1,1,\beta,m}(\omega) - V_{1,1,\beta,m}(\omega')| \leq (1-\epsilon)|\omega - \omega'| \leq C|\omega - \omega'|$, where the second inequality is due to (33).

Assume that for $T \geq 1$, $|V_{1,T,\beta,m}(\omega) - V_{1,T,\beta,m}(\omega')| \leq C|\omega - \omega'|$ holds, *i.e.*, neither the left nor the right derivative of $V_{1,T,\beta,m}(\cdot)$ can have an absolute value exceeding C . We have the following inequalities:

$$\begin{aligned} |V_{1,T,\beta,m}(p_{11}) - V_{1,T,\beta,m}(\mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega}))| &\leq C|p_{11} - \mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega})| \\ &\leq C|p_{11} - p_{01}|, \\ \frac{\epsilon}{1-(1-\epsilon)\omega} |V'_{1,T,\beta,m}(\mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega}))| &\leq \frac{\epsilon}{1-(1-\epsilon)\omega} C = C. \end{aligned} \quad (36)$$

To prove $|V_{1,T+1,\beta,m}(\omega) - V_{1,T+1,\beta,m}(\omega')| \leq C|\omega - \omega'|$, recall the definitions of $f_T(\omega)$ and $g_T(\omega)$ in (28) and their right derivatives $f'_T(\omega)$ and $g'_T(\omega)$ in (29). We have

$$\begin{aligned} |f'_T(\omega) - (1-\epsilon)| &\leq (1-\epsilon)\beta C|p_{11} - p_{01}| + C\beta|p_{11} - p_{01}| = (2-\epsilon)\beta C|p_{11} - p_{01}|, \\ |g'_T(\omega)| &\leq C\beta|p_{11} - p_{01}|, \end{aligned} \quad (37)$$

where the inequality in the first (or second) line above is due to the first (or second) line in (36). Thus we have the following lower and upper bounds on $f'_T(\omega)$ and $g'_T(\omega)$:

$$\begin{aligned} (1-\epsilon) - (2-\epsilon)\beta C|p_{11} - p_{01}| &\leq f'_T(\omega) \leq (1-\epsilon) + (2-\epsilon)\beta C|p_{11} - p_{01}|, \\ -C\beta|p_{11} - p_{01}| &\leq g'_T(\omega) \leq C\beta|p_{11} - p_{01}|. \end{aligned} \quad (38)$$

From (38) and that $V_{1,T+1,\beta,m}(\omega) = \max\{f_T(\omega), g_T(\omega)\}$, we have

$$\begin{aligned} |V'_{1,T+1,\beta,m}(\omega)| &\leq (1-\epsilon) + (2-\epsilon)\beta C|p_{11} - p_{01}| \\ &= (1-\epsilon) + (2-\epsilon)\beta|p_{11} - p_{01}| \frac{1-\epsilon}{1-(2-\epsilon)\beta|p_{11} - p_{01}|} \\ &= \frac{1-\epsilon}{1-(2-\epsilon)\beta|p_{11} - p_{01}|} = C, \end{aligned}$$

where we used the fact that $2-\epsilon \geq 1$ in the above inequality. Since $V_{1,T+1,\beta,m}(\omega)$ is absolutely continuous, the above implies that

$$|V_{1,T+1,\beta,m}(\omega) - V_{1,T+1,\beta,m}(\omega')| \leq C|\omega - \omega'|.$$

The proof is thus finished by the induction process. \square

Last, we give a lemma establishing the order of $V'_{1,T,\beta,m}(\cdot; u = 1)$ and $V'_{1,T,\beta,m}(\cdot; u = 0)$ under certain conditions which further leads to a threshold structure of the optimal single-arm policy as detailed in Section 2.3.

Lemma 3. *Suppose that $p_{11} > p_{01}$ and $\beta \leq 1/[(3 - \epsilon)(p_{11} - p_{01})]$, we have*

$$V'_{1,T,\beta,m}(\omega; u = 1) \geq V'_{1,T,\beta,m}(\omega; u = 0), \quad (39)$$

where $V'_{1,T,\beta,m}(\omega; u = i)$ denotes the right derivative of $V_{1,T,\beta,m}(\cdot; u = i)$ at ω for $i \in \{0, 1\}$. The above inequality is also true if $p_{01} > p_{11}$ and $\beta \leq 1/[(5 - 2\epsilon)(p_{01} - p_{11})]$.

Proof. Again, we prove by induction on the time horizon T . When $T = 1$, it is clear that $V_{1,1,\beta,m}(\omega; u = 1) = (1 - \epsilon)\omega$ and $V_{1,1,\beta,m}(\omega; u = 0) = m$:

$$V'_{1,1,\beta,m}(\omega; u = 1) = 1 - \epsilon > V'_{1,1,\beta,m}(\omega; u = 0) = 0. \quad (40)$$

Assume that $V'_{1,T,\beta,m}(\omega; u = 1) \geq V'_{1,T,\beta,m}(\omega; u = 0)$ for $T \geq 1$. From (38), we have, in case of $p_{01} > p_{11}$ and $\beta \leq \frac{1}{(5-2\epsilon)(p_{01}-p_{11})}$, the inequality $C\beta(p_{01} - p_{11}) \leq (1 - \epsilon) - (2 - \epsilon)\beta C(p_{01} - p_{11})$, which shows that $f'_T(\omega) \geq g'_T(\omega)$. When $p_{11} > p_{01}$, $V_{1,T,\beta,m}(\omega)$ is increasing with ω therefore has nonnegative right derivatives by Lemma 1. We can thus obtain tighter bounds on $f'_T(\omega)$ and $g'_T(\omega)$ by (29):

$$\begin{aligned} (1 - \epsilon) &\leq f'_T(\omega) \leq (1 - \epsilon) + (2 - \epsilon)\beta C(p_{11} - p_{01}), \\ 0 &\leq g'_T(\omega) \leq C\beta(p_{11} - p_{01}). \end{aligned}$$

If $\beta \leq \frac{1}{(3-\epsilon)(p_{11}-p_{01})}$, we have $C\beta(p_{11}-p_{01}) \leq (1-\epsilon)$, which shows that $f'_T(\omega) \geq g'_T(\omega)$. The proof is thus complete. \square

2.3 Threshold Policy and Indexability

In this section, we show that the optimal single-arm policy is a threshold policy under the constraints on the discount factor β specified in Section 2.2 and analyze the conditions for indexability. First, for a finite-horizon single-armed bandit, a threshold policy π is defined by a time-dependent real number $\omega_{T,\beta}(m)$ such that

$$u_{T,m}(\omega) = \begin{cases} 1, & \text{if } \omega > \omega_{T,\beta}(m); \\ 0, & \text{if } \omega \leq \omega_{T,\beta}(m). \end{cases} \quad (41)$$

In the above $u_{T,m}(\omega) \in \{0, 1\}$ is the action taken under π at the current state ω with T slots remaining. Intuitively, the larger ω is, the larger expected immediate reward to accrue and thus more attractive to activate the arm. We formalize this intuition under certain conditions in the following theorem.

Theorem 1. *Suppose that $p_{11} > p_{01}$ and $\beta \leq \frac{1}{(3-\epsilon)(p_{11}-p_{01})}$. For any $T \geq 1$, the optimal single-arm policy π^* is a threshold policy, i.e., there exists $\omega_{T,\beta}^*(m) \in \mathbb{R}$ such*

that under π^* , the optimal action is

$$u_{T,m}^*(\omega) = \begin{cases} 1, & \text{if } \omega > \omega_{T,\beta}^*(m); \\ 0, & \text{if } \omega \leq \omega_{T,\beta}^*(m). \end{cases}$$

Furthermore, at the threshold $\omega_{T,\beta}^*(m)$,

$$V_{1,T,\beta,m}(\omega_{T,\beta}^*(m); u = 0) = V_{1,T,\beta,m}(\omega_{T,\beta}^*(m); u = 1). \quad (42)$$

The conclusion is also true for the case of $p_{01} > p_{11}$ and $\beta \leq \frac{1}{(5-2\epsilon)(p_{01}-p_{11})}$.

Proof. At $T = 1$, $V_{1,1,\beta,m}(\omega; u = 1) = (1 - \epsilon)\omega$, $V_{1,1,\beta,m}(\omega; u = 0) = m$. Thus we can choose $\omega_{1,\beta}^*(m)$ as follows:

$$\omega_{1,\beta}^*(m) = \begin{cases} c, & \text{if } m \geq 1 - \epsilon; \\ \frac{m}{1-\epsilon}, & \text{if } 0 \leq m < 1 - \epsilon; \\ b, & \text{if } m < 0, \end{cases}$$

where $b < 0, c > 1$ are arbitrary constants.

For $T \geq 1$, when the condition on β is satisfied, Lemma 3 shows that

$$h_T(\omega) \triangleq V_{1,T,\beta,m}(\omega; u = 1) - V_{1,T,\beta,m}(\omega; u = 0), \quad (43)$$

$$h_T'(\omega) \geq 0, \quad \forall \omega \in (0, 1). \quad (44)$$

This shows that $h_T(\cdot)$ is monotonically increasing and either has no zeros in the interval $[0, 1]$ or intersects with it over a closed interval (which can be a single point) only. Specially,

$$\begin{aligned} V_{1,T,\beta,m}(0; u = 1) &= \beta V_{1,T-1,\beta,m}(p_{01}), \\ V_{1,T,\beta,m}(0; u = 0) &= m + \beta V_{1,T-1,\beta,m}(p_{01}), \\ V_{1,T,\beta,m}(1; u = 1) &= (1 - \epsilon) + \beta V_{1,T-1,\beta,m}(p_{11}), \\ V_{1,T,\beta,m}(1; u = 0) &= m + \beta V_{1,T-1,\beta,m}(p_{11}). \end{aligned}$$

Consider the following three regions of m .

- (i) $0 \leq m < 1 - \epsilon$. In this case, $V_{1,T,\beta,m}(0; u = 1) \leq V_{1,T,\beta,m}(0; u = 0)$ and $V_{1,T,\beta,m}(1; u = 1) > V_{1,T,\beta,m}(1; u = 0)$. Therefore $h_T(\cdot)$ intersects over (at least) one point in $[0, 1]$. This point can thus be chosen as $\omega_{T,\beta}^*(m)$.
- (ii) $m < 0$. In this case, $V_{1,T,\beta,m}(0; u = 1) > V_{1,T,\beta,m}(0; u = 0)$ and $V_{1,T,\beta,m}(1; u = 1) > V_{1,T,\beta,m}(1; u = 0)$. So $h_T(\cdot)$ is strictly positive over $[0, 1]$ and we can choose $\omega_{T,\beta}^*(m) = b$ with any $b < 0$.

- (iii) $m \geq (1 - \epsilon)$. In this case, always choosing the passive action is clearly optimal as the expected immediate reward is uniformly upper-bounded by m over the whole belief state space. We can thus choose $\omega_{T,\beta}^*(m) = c$ with any $c > 1$.

In conclusion, when the conditions in the theorem are satisfied, the optimal finite-horizon single-arm policy is a threshold policy for any horizon length $T \geq 1$. \square

In the next theorem, we show that the optimal single-arm policy over the *infinite* horizon is also a threshold policy under the same conditions.

Theorem 2. *Fix the subsidy m . The finite-horizon value functions $V_{1,T,\beta,m}(\cdot)$, $V_{1,T,\beta,m}(\cdot; u = 1)$ and $V_{1,T,\beta,m}(\cdot; u = 0)$ uniformly converge to the infinite-horizon value functions $V_{\beta,m}(\cdot)$, $V_{\beta,m}(\cdot; u = 1)$ and $V_{\beta,m}(\cdot; u = 0)$ which are consequently obedient to the same properties established in Lemmas 1 and 2 and Theorem 1.*

Proof. The uniform convergence is obvious since $\beta < 1$ and the rest can be easily proved by contradiction following the uniform convergence. \square

Thus far we have established the threshold structure of the optimal single-arm policy with subsidy based on the analysis of $V_{\beta,m}(\omega)$ as a function of the belief state ω with m fixed. To study the indexability condition, we now analyze the properties of $V_{\beta,m}(\omega)$ as a function of the subsidy m with the starting belief ω fixed. From Definition 1 and the threshold structure of the optimal policy, the indexability of our model is reduced to requiring that the threshold $\omega_{\beta}^*(m)$ is monotonically increasing with m (if the threshold is a closed interval then the right end is selected). Note that for the infinite-horizon problem, the threshold $\omega_{\beta}^*(m)$ is independent of time. Furthermore, $V_{\beta,m}(\omega)$ is also convex in m as for any $m_1, m_2 \in \mathbb{R}$ and $\theta \in (0, 1)$ the optimal policy $\pi_{\beta}^*(\theta m_1 + (1 - \theta)m_2)$ achieving $V_{\beta,\theta m_1 + (1 - \theta)m_2}(\omega)$ applied respectively on the problem with subsidies m_1 and m_2 cannot outperform those achieving $V_{\beta,m_1}(\omega)$ and $V_{\beta,m_2}(\omega)$. Specifically, let r_a be the expected total discounted reward from the active action and $r_p(m)$ that from the passive action under $\pi_{\beta}^*(\theta m_1 + (1 - \theta)m_2)$ applied to the problem with subsidy m , then

$$\begin{aligned} \theta V_{\beta,m_1}(\omega) + (1 - \theta)V_{\beta,m_2}(\omega) &\geq r_a + \theta r_p(m_1) + (1 - \theta)r_p(m_2) \\ &= r_a + r_p(\theta m_1 + (1 - \theta)m_2) \\ &= V_{\beta,\theta m_1 + (1 - \theta)m_2}(\omega). \end{aligned}$$

Since $V_{\beta,m}(\omega)$ is convex in m , its left and right derivatives with m exist at every point $m_0 \in \mathbb{R}$. Furthermore, consider two policies $\pi_{\beta}^*(m_1)$ and $\pi_{\beta}^*(m_2)$ achieving $V_{\beta,m_1}(\omega)$ and $V_{\beta,m_2}(\omega)$ for any $m_1, m_2 \in \mathbb{R}$, respectively. With a similar interchange argument of $\pi_{\beta}^*(m_1)$ and $\pi_{\beta}^*(m_2)$ as above, we have $|V_{\beta,m_1}(\omega) - V_{\beta,m_2}(\omega)| \leq \frac{1}{1 - \beta}|m_1 - m_2|$ and $V_{\beta,m}(\omega)$ is Lipschitz continuous in m . By the Rademacher theorem (see [Heinonen, 2005](#)), $V_{\beta,m}(\omega)$ is differentiable almost everywhere in m . For a small increase of m , the rate at which $V_{\beta,m}(\omega)$ increases is at least the expected total discounted passive time under any optimal policy for the problem with subsidy m starting from the belief state ω . In the following theorem, we formalize this relation between the value function and the passive time as well as a sufficient condition for the indexability of our model.

Theorem 3. Let $\Pi_\beta^*(m)$ denote the set of all optimal single-arm policies achieving $V_{\beta,m}(\omega)$ with initial belief state ω . Define the passive time

$$D_{\beta,m}(\omega) \triangleq \max_{\pi_\beta^*(m) \in \Pi_\beta^*(m)} \mathbb{E}_{\pi_\beta^*(m)} \left[\sum_{t=1}^{\infty} \beta^{t-1} \mathbb{1}_{(u(t)=0)} \mid \omega(1) = \omega \right]. \quad (45)$$

The right derivative of the value function $V_{\beta,m}(\omega)$ with m , denoted by $\frac{dV_{\beta,m}(\omega)}{(dm)^+}$, exists at every value of m and

$$\frac{dV_{\beta,m}(\omega)}{(dm)^+} \Big|_{m=m_0} = D_{\beta,m_0}(\omega). \quad (46)$$

Furthermore, the single-armed bandit is indexable if at least one of the following condition is satisfied:

- i. for any $m_0 \in [0, 1 - \epsilon)$ the optimal policy is a threshold policy with threshold $\omega_\beta^*(m_0) \in [0, 1)$ (if the threshold is a closed interval then the right end is selected) and

$$\frac{dV_{\beta,m}(\omega_\beta^*(m_0); u=0)}{(dm)^+} \Big|_{m=m_0} > \frac{dV_{\beta,m}(\omega_\beta^*(m_0); u=1)}{(dm)^+} \Big|_{m=m_0}. \quad (47)$$

- ii. for any $m_0 \in \mathbb{R}$ and $\omega \in P(m_0)$, we have

$$\frac{dV_{\beta,m}(\omega; u=0)}{(dm)^+} \Big|_{m=m_0} \geq \frac{dV_{\beta,m}(\omega; u=1)}{(dm)^+} \Big|_{m=m_0}. \quad (48)$$

Proof. The proof of (46) follows directly from the argument in Theorem 1 in Liu (2021) and is omitted here. To prove the sufficiency of (47), we note that if it is true then there exists a $\Delta m > 0$ such that $\forall m \in (m_0, m_0 + \Delta m)$,

$$\begin{aligned} & V_{\beta,m}(\omega_\beta^*(m_0); u=0) - V_{\beta,m_0}(\omega_\beta^*(m_0); u=0) \\ & > V_{\beta,m}(\omega_\beta^*(m_0); u=1) - V_{\beta,m_0}(\omega_\beta^*(m_0); u=1). \end{aligned}$$

Since $V_{\beta,m_0}(\omega_\beta^*(m_0); u=0) = V_{\beta,m_0}(\omega_\beta^*(m_0); u=1)$, we have $V_{\beta,m}(\omega_\beta^*(m_0); u=0) > V_{\beta,m}(\omega_\beta^*(m_0); u=1)$ which implies that the threshold $\omega_\beta^*(m_0)$ remains in the passive set as m continuously increases so $P(m)$ is monotonically increasing with m . This conclusion is clearly true for the trivial case of $m < 0$ or $m \geq 1 - \epsilon$. The sufficiency of (48) is obvious because then it is impossible for any $\omega \in P(m)$ to escape from $P(m)$ as m increases due to the nondecreasing property of $V_{\beta,m}(\omega; u=0) - V_{\beta,m}(\omega; u=1)$ enforced by (48). \square

Theorem 3 essentially provides a way for checking the indexability condition in terms of the passive times. For example, equation (47) is equivalent to for any $m \in [0, 1 - \epsilon)$,

$$\beta \left[(1 - \epsilon) \omega_\beta^*(m) D_{\beta,m}(p_{11}) + (1 - (1 - \epsilon) \omega_\beta^*(m)) D_{\beta,m} \left(\mathcal{T} \left(\frac{\epsilon \omega_\beta^*(m)}{\epsilon \omega_\beta^*(m) + 1 - \omega_\beta^*(m)} \right) \right) \right]$$

$$< 1 + \beta D_{\beta,m}(\mathcal{T}(\omega_{\beta}^*(m))). \quad (49)$$

The above strict inequality clearly holds if $\beta < 0.5$ since $D_{\beta,m}(\cdot) \in [0, \frac{1}{1-\beta}]$ for any $m \in \mathbb{R}$. When $\beta = 0.5$, we prove by contradiction that the strict inequality (49) must hold under the threshold structure of the optimal policy. If $\omega_{\beta}^*(m) = 0$ then (49) is clearly true. Assume that the left and right sides of (49) are equal and $\omega_{\beta}^*(m) \neq 0$. In this case, we have

$$D_{\beta,m}(p_{11}) = \frac{1}{1-\beta}, \quad (50)$$

$$D_{\beta,m}(\mathcal{T}(\omega_{\beta}^*(m))) = 0. \quad (51)$$

Equation (51) implies that starting from $\mathcal{T}(\omega_{\beta}^*(m))$, always activating the arm is strictly optimal. This means that the threshold $\omega_{\beta}^*(m)$ is strictly below p_{11} and we have a contradiction to (50). Another easier way to see that the bandit is indexable if $\beta \leq 0.5$ is that (48) would be satisfied where no strict inequality is required. However, condition (49) provides a convenient way for approximately computing the passive times as well as the value functions which leads to an efficient algorithm for evaluating the indexability and solving for the Whittle index function for any $\beta \in (0, 1)$, as detailed in the next section.

Corollary 1. *The restless bandit is indexable if $\beta \leq 0.5$.*

3 The Whittle Index Policy

In this section, we design an efficient algorithm by approximating the Whittle index.

3.1 The Approximated Whittle Index

The threshold structure of the optimal single-arm policy under certain conditions yields the following iterative nature of the dynamic equations for both $D_{\beta,m}(\omega)$ and $V_{\beta,m}(\omega)$. Define *the first crossing time*

$$L(\omega, \omega') = \min_{0 \leq k < \infty} \{k : \mathcal{T}^k(\omega) > \omega'\}. \quad (52)$$

In the above $\mathcal{T}^0(\omega) \triangleq \omega$ and we set $L(\omega, \omega') = +\infty$ if $\mathcal{T}^k(\omega) \leq \omega'$ for all $k \geq 0$. Clearly $L(\omega, \omega')$ is the minimum time slots required for a belief state ω to stay in the passive set $P(m)$ before the arm is activated given a threshold $\omega' \in [0, 1)$. Consider the nontrivial case where $p_{01}, p_{11} \in (0, 1)$ and $p_{01} \neq p_{11}$ such that the Markov chain of the internal arm states is aperiodic and irreducible and that the belief update is action-dependent. From (6), if $p_{11} > p_{01}$ then

$$L(\omega, \omega') = \begin{cases} 0, & \omega > \omega' \\ \left\lceil \log_{\frac{p_{01}-\omega'(1-p_{11}+p_{01})}{p_{01}-\omega(1-p_{11}+p_{01})}} \right\rceil + 1, & \omega \leq \omega' < \omega_o \\ \infty, & \omega \leq \omega', \omega' \geq \omega_o \end{cases} \quad (53)$$

or if $p_{11} < p_{01}$ then

$$L(\omega, \omega') = \begin{cases} 0, & \omega > \omega' \\ 1, & \omega \leq \omega', \mathcal{T}(\omega) > \omega' \\ \infty, & \omega \leq \omega', \mathcal{T}(\omega) \leq \omega' \end{cases}. \quad (54)$$

To illustrate (53) and (54), note that a belief state will converge to the stationary belief value given by (7) monotonically for $p_{11} > p_{01}$ or in an oscillating manner for $p_{11} < p_{01}$ under passive actions (see Fig. 3 and Fig. 4 in Liu and Zhao (2010)). Suppose that the following conditions are satisfied such that the optimal single-arm policy is a threshold policy and the indexability holds:

$$\beta \leq \begin{cases} \min \left\{ \frac{1}{(3-\epsilon)(p_{11}-p_{01})}, 0.5 \right\}, & \text{if } p_{11} > p_{01} \\ \min \left\{ \frac{1}{(5-2\epsilon)(p_{01}-p_{11})}, 0.5 \right\}, & \text{if } p_{11} < p_{01} \end{cases}. \quad (55)$$

To solve for the Whittle index function $W(\omega)$, given the current arm state ω , we aim to find out the minimum subsidy m that makes it as a threshold:

$$V_{\beta,m}(\omega) = V_{\beta,m}(\omega; u = 1) = V_{\beta,m}(\omega; u = 0), \quad (56)$$

$$V_{\beta,m}(\omega; u = 1) = (1 - \epsilon)\omega + \beta[(1 - \epsilon)\omega V_{\beta,m}(p_{11}) + (1 - (1 - \epsilon)\omega)V_{\beta,m}(\mathcal{T} \circ \phi(\omega))], \quad (57)$$

$$V_{\beta,m}(\omega; u = 0) = m + \beta V_{\beta,m}(\mathcal{T}(\omega)). \quad (58)$$

Given a threshold $\omega_{\beta}^*(m) \in [0, 1]$ and any $\omega \in [0, 1]$, the value function $V_{\beta,m}(\omega)$ can be expanded by the first crossing time as

$$\begin{aligned} V_{\beta,m}(\omega) &= \frac{1 - \beta \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}}{1 - \beta} m + \beta \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))} V_{\beta,m}(\mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega); u = 1) \\ &= \frac{1 - \beta \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}}{1 - \beta} m + \beta \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))} \left\{ (1 - \epsilon) \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega) \right. \\ &\quad \left. + \beta \left[(1 - \epsilon) \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega) V_{\beta,m}(p_{11}) + (1 - (1 - \epsilon) \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega)) \right. \right. \\ &\quad \left. \left. \times V_{\beta,m} \left(\mathcal{T} \left(\frac{\epsilon \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega)}{\epsilon \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega) + 1 - \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega)} \right) \right) \right] \right\}. \end{aligned} \quad (59)$$

There is no doubt that the last item of the above equation has caused us trouble in solving for $V_{\beta,m}(\omega)$. However, if we let

$$\begin{aligned} f(\omega, \omega_{\beta}^*(m)) &= \mathcal{T} \left(\frac{\epsilon \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega)}{\epsilon \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega) + 1 - \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega)} \right) \\ &= \frac{p_{11} \epsilon \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega) + p_{01} (1 - \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega))}{\epsilon \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega) + 1 - \mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega)} \end{aligned} \quad (60)$$

and construct iteratively the sequence $\{k_n\}$ as $k_{n+1} = f(k_n, \omega_\beta^*(m))$ with $k_0 = \omega$. We then get the following sequence of equations:

$$\begin{aligned}
V_{\beta,m}(k_0) &= \frac{1 - \beta^{L(k_0, \omega_\beta^*(m))}}{1 - \beta} m + \beta^{L(k_0, \omega_\beta^*(m))} \left\{ (1 - \epsilon) \mathcal{T}^{L(k_0, \omega_\beta^*(m))}(k_0) + \beta [(1 - \epsilon) \right. \\
&\quad \left. \times \mathcal{T}^{L(k_0, \omega_\beta^*(m))}(k_0) V_{\beta,m}(p_{11}) + (1 - (1 - \epsilon) \mathcal{T}^{L(k_0, \omega_\beta^*(m))}(k_0)) V_{\beta,m}(k_1)] \right\} \\
V_{\beta,m}(k_1) &= \frac{1 - \beta^{L(k_1, \omega_\beta^*(m))}}{1 - \beta} m + \beta^{L(k_1, \omega_\beta^*(m))} \left\{ (1 - \epsilon) \mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1) + \beta [(1 - \epsilon) \right. \\
&\quad \left. \times \mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1) V_{\beta,m}(p_{11}) + (1 - (1 - \epsilon) \mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1)) V_{\beta,m}(k_2)] \right\} \\
&\dots \\
V_{\beta,m}(k_n) &= \frac{1 - \beta^{L(k_n, \omega_\beta^*(m))}}{1 - \beta} m + \beta^{L(k_n, \omega_\beta^*(m))} \left\{ (1 - \epsilon) \mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n) \right. \\
&\quad \left. + \beta [(1 - \epsilon) \cdot \mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n) V_{\beta,m}(p_{11}) + (1 - (1 - \epsilon) \right. \\
&\quad \left. \times \mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n)) V_{\beta,m}(k_{n+1})] \right\} \\
&\dots
\end{aligned}$$

For sufficiently large n , we can get an estimation of $V_{\beta,m}(\omega) = V_{\beta,m}(k_0)$ with an arbitrarily small error by setting $V_{\beta,m}(k_{n+1}) = 0$ whose error is discounted by β in computing $V_{\beta,m}(k_n)$ thus causing a geometrically decreasing error propagation in the backward computation process for $V_{\beta,m}(k_0)$. Note that we first compute $V_{\beta,m}(p_{11})$ in the same way by setting $k_0 = p_{11}$ in the above equation set. Therefore we can have an estimation of $V_{\beta,m}(\omega)$ with arbitrarily high precision for any $\omega \in [0, 1]$. Interestingly, extensive numerical results found that $\{k_n\}$ quickly converges to a limit belief state k^* (independent of k_0) (see Fig. 1 for an example). Specifically, after 4 iterations, the difference $|k_4 - k|$ becomes too small to affect the performance of our algorithm as discussed in Sec. 5.1. So we can set $V_{\beta,m}(k_5) = V_{\beta,m}(k_4)$ and efficiently solve the *finite* linear equation set (up to $V_{\beta,m}(k_4)$). In general, the n -iteration Whittle index is based on the solution of the following equations:

$$\begin{aligned}
V_{\beta,m}(p_{11}) &= \frac{1 - \beta^{L(p_{11}, \omega_\beta^*(m))}}{1 - \beta} m + \beta^{L(p_{11}, \omega_\beta^*(m))} \left\{ (1 - \epsilon) \mathcal{T}^{L(p_{11}, \omega_\beta^*(m))}(p_{11}) \right. \\
&\quad \left. + \beta [(1 - \epsilon) \mathcal{T}^{L(p_{11}, \omega_\beta^*(m))}(p_{11}) V_{\beta,m}(p_{11}) + (1 - (1 - \epsilon) \right. \\
&\quad \left. \times \mathcal{T}^{L(p_{11}, \omega_\beta^*(m))}(p_{11})) V_{\beta,m}(k_1)] \right\} \\
V_{\beta,m}(k_1) &= \frac{1 - \beta^{L(k_1, \omega_\beta^*(m))}}{1 - \beta} m + \beta^{L(k_1, \omega_\beta^*(m))} \left\{ (1 - \epsilon) \mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1) + \beta [(1 - \epsilon) \right. \\
&\quad \left. \times \mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1) V_{\beta,m}(p_{11}) + (1 - (1 - \epsilon) \mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1)) V_{\beta,m}(k_2)] \right\} \\
&\vdots \\
V_{\beta,m}(k_n) &= \frac{1 - \beta^{L(k_n, \omega_\beta^*(m))}}{1 - \beta} m + \beta^{L(k_n, \omega_\beta^*(m))} \left\{ (1 - \epsilon) \mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n) \right.
\end{aligned}$$

$$+\beta[(1-\epsilon)\mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n)V_{\beta,m}(p_{11}) + (1-(1-\epsilon)) \\ \times \mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n)V_{\beta,m}(k_n)]\}.$$

After the value functions are (approximately) solved, we can plot the active value function $V_{\beta,m}(\omega; u = 1)$ and the passive one $V_{\beta,m}(\omega; u = 0)$ to see that they intersect at one single point (the threshold) verifying the optimality of the threshold policy proven by Theorem 2 (See Fig. 2 for an example). According to Theorem 3, the passive time $D_{\beta,m}(\omega)$ can also be approximately solved based on the following equations:

$$\begin{aligned} D_{\beta,m}(p_{11}) &= \frac{1 - \beta^{L(p_{11}, \omega_\beta^*(m))}}{1 - \beta} + \beta^{L(p_{11}, \omega_\beta^*(m))+1} \{(1 - \epsilon)\mathcal{T}^{L(p_{11}, \omega_\beta^*(m))}(p_{11}) \\ &\quad \times D_{\beta,m}(p_{11}) + (1 - (1 - \epsilon))\mathcal{T}^{L(p_{11}, \omega_\beta^*(m))}(p_{11})D_{\beta,m}(k_1)\} \\ D_{\beta,m}(k_1) &= \frac{1 - \beta^{L(k_1, \omega_\beta^*(m))}}{1 - \beta} + \beta^{L(k_1, \omega_\beta^*(m))+1} \{(1 - \epsilon)\mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1) \\ &\quad \times D_{\beta,m}(p_{11}) + (1 - (1 - \epsilon))\mathcal{T}^{L(k_1, \omega_\beta^*(m))}(k_1)D_{\beta,m}(k_2)\} \\ &\quad \vdots \\ D_{\beta,m}(k_n) &= \frac{1 - \beta^{L(k_n, \omega_\beta^*(m))}}{1 - \beta} + \beta^{L(k_n, \omega_\beta^*(m))+1} \{(1 - \epsilon)\mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n) \\ &\quad \times D_{\beta,m}(p_{11}) + (1 - (1 - \epsilon))\mathcal{T}^{L(k_n, \omega_\beta^*(m))}(k_n)D_{\beta,m}(k_n)\} \end{aligned}$$

Substituting ω for $\omega_\beta^*(m)$ in the above $n + 1$ linear equations with $n + 1$ unknowns (first solving for $\omega = p_{11}$), we can obtain $V_{\beta,m}(\omega')$ and $D_{\beta,m}(\omega')$ for any $\omega' \in [0, 1]$ according to the linear equation sets. The indexability condition (47) in Theorem 3 can be checked *online*: for the original multi-armed bandit problem and for each arm at state $\omega(t)$ at time t , we compute its approximated Whittle index $W(\omega(t))$ by solving a set of linear equations, which has a polynomial complexity of the iteration number n , independent of the decision time t . At time t , for each arm, if $W(\cdot)$ is found to be nondecreasing with the arm states $(\omega(1), \omega(2), \dots, \omega(t))$ appeared so far starting from the initial belief vector $\boldsymbol{\omega}(1)$ defined in (2), then the indexability has not been violated. Interestingly, extensive numerical studies have shown that the indexability is always satisfied as illustrated in Figs. 3 and 4 in Sec. 5.1.

For large $\beta \in (0, 1)$ where the threshold structure of the optimal policy or the indexability may not hold (*i.e.*, condition (55) is not satisfied), we can still use the above process to solve for the subsidy m that makes (56) true if it exists. Note that after computing the value functions appeared in (56) in terms of m , both $V_{\beta,m}(\omega; u = 1)$ and $V_{\beta,m}(\omega; u = 0)$ are linear (affine) in m and their equality gives a *unique* solution of m if their linear coefficients are not equal. This m , if exists, can thus be used as the approximated Whittle index $W(\omega)$ *without* requiring indexability or threshold-based optimal policy. If it does not exist, we can simply set $W(\omega) = \omega B$. The existence of such an m is defined as *the relaxed indexability* in Liu (2021). Note that extensive numerical studies have shown that the relaxed indexability of our model with imperfect state observations is always satisfied as well. Before summarizing our general algorithm for

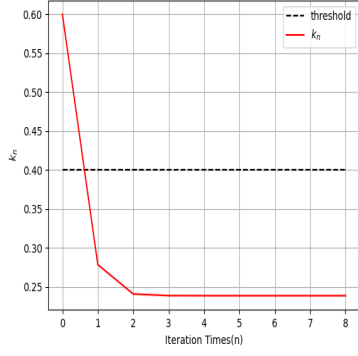


Fig. 1 The Convergence of k_n

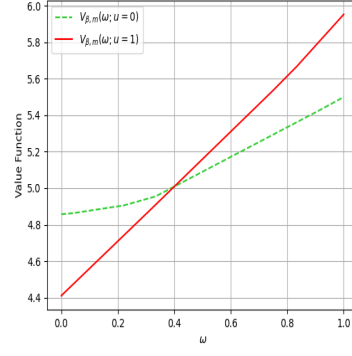


Fig. 2 The Optimality of Threshold Policy

all $\beta \in (0, 1)$ in Section 3.2, we solve for the approximated Whittle index function in closed-form for the simplest case of 0-iteration, which is referred to as *the imperfect Whittle index*. Note that if $\epsilon \rightarrow 0$ then $\mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega}) \rightarrow p_{01}$. Thus when ϵ is sufficiently small, we can approximate $V_{\beta,m}(\mathcal{T}(\frac{\epsilon\omega}{\epsilon\omega+1-\omega}))$ by $V_{\beta,m}(p_{01})$. Under this approximation, we have, for any $\omega \in [0, 1]$,

$$\begin{aligned} V_{\beta,m}(\omega; u=1) &= (1-\epsilon)\omega + \beta[(1-\epsilon)\omega V_{\beta,m}(p_{11}) + (1-(1-\epsilon)\omega)V_{\beta,m}(p_{01})] \\ V_{\beta,m}(\omega; u=0) &= m + \beta V_{\beta,m}(\mathcal{T}(\omega)) \\ V_{\beta,m}(\omega) &= \frac{1 - \beta^{L(\omega, \omega_{\beta}^*(m))}}{1 - \beta} m + \beta^{L(\omega, \omega_{\beta}^*(m))} \{ (1-\epsilon)\mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega) \\ &\quad + \beta[(1-\epsilon)\mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega)V_{\beta,m}(p_{11}) \\ &\quad + (1-(1-\epsilon)\mathcal{T}^{L(\omega, \omega_{\beta}^*(m))}(\omega))V_{\beta,m}(p_{01})] \} \end{aligned}$$

By using the above three equations, we can directly solve for $V_{\beta,m}(p_{01})$ and $V_{\beta,m}(p_{11})$ in closed-form.

When $p_{11} > p_{01}$, $V_{\beta,m}(p_{01}) =$

$$\begin{cases} \frac{(1-\epsilon)p_{01}}{(1-\beta)(1-\beta(1-\epsilon)p_{11} + \beta(1-\epsilon)p_{01})}, & \text{if } \omega_{\beta}^*(m) < p_{01} \\ \frac{(1-\beta(1-\epsilon)p_{11})(1-\beta^{L(p_{01}, \omega_{\beta}^*(m))})m + (1-\epsilon)(1-\beta)\beta^{L(p_{01}, \omega_{\beta}^*(m))}\mathcal{T}^{L(p_{01}, \omega_{\beta}^*(m))}(p_{01})}{(1-\beta(1-\epsilon)p_{11})(1-\beta)(1-\beta^{L(p_{01}, \omega_{\beta}^*(m))+1}) + (1-\epsilon)(1-\beta)2\beta^{L(p_{01}, \omega_{\beta}^*(m))+1}\mathcal{T}^{L(p_{01}, \omega_{\beta}^*(m))}(p_{01})}, & \text{if } p_{01} \leq \omega_{\beta}^*(m) < \omega_0 \\ \frac{m}{1-\beta}, & \text{if } \omega_{\beta}^*(m) \geq \omega_0 \end{cases},$$

$$V_{\beta,m}(p_{11}) = \begin{cases} \frac{(1-\epsilon)p_{11} + \beta(1-(1-\epsilon)p_{11})V_{\beta,m}(p_{01})}{1-\beta(1-\epsilon)p_{11}}, & \text{if } \omega_{\beta}^*(m) < p_{11} \\ \frac{m}{1-\beta}, & \text{if } \omega_{\beta}^*(m) \geq p_{11} \end{cases}.$$

The approximate Whittle index is given by

$$W(\omega) = \begin{cases} \frac{\omega(1-\epsilon)(1-\beta p_{11} + \beta p_{01})}{1-\beta(1-\epsilon)p_{11} + \beta(1-\epsilon)p_{01}}, & \text{if } \omega \leq p_{01} \\ \frac{(1-\epsilon)(\omega - \beta \mathcal{T}(\omega)) + C_2(1-\beta)\beta[(1-\beta(1-\epsilon)p_{11}) - (1-\epsilon)(\omega - \beta \mathcal{T}(\omega))]}{1-\beta(1-\epsilon)p_{11} - C_1\beta[(1-\beta(1-\epsilon)p_{11}) - (1-\epsilon)(\omega - \beta \mathcal{T}(\omega))]}, & \text{if } p_{01} < \omega \leq \omega_0 \\ \frac{(1-\epsilon)\omega}{1-\beta(1-\epsilon)p_{11} + \beta(1-\epsilon)\omega}, & \text{if } \omega_0 < \omega \leq p_{11} \\ (1-\epsilon)\omega, & \text{if } \omega > p_{11} \end{cases},$$

where

$$C_1 = \frac{(1-\beta(1-\epsilon)p_{11})(1-\beta^{L(p_{01}, \omega)})}{(1-\beta(1-\epsilon)p_{11})(1-\beta^{L(p_{01}, \omega)+1}) + (1-\epsilon)(1-\beta)\beta^{L(p_{01}, \omega)+1}\mathcal{T}^{L(p_{01}, \omega)}(p_{01})},$$

$$C_2 = \frac{(1-\epsilon)\beta^{L(p_{01}, \omega)}\mathcal{T}^{L(p_{01}, \omega)}(p_{01})}{(1-\beta(1-\epsilon)p_{11})(1-\beta^{L(p_{01}, \omega)+1}) + (1-\epsilon)(1-\beta)\beta^{L(p_{01}, \omega)+1}\mathcal{T}^{L(p_{01}, \omega)}(p_{01})}.$$

Similarly, when $p_{01} > p_{11}$, we have

$$V_{\beta, m}(p_{11}) = \begin{cases} \frac{(1-\epsilon)(p_{11}(1-\beta) + \beta p_{01})}{(1-\beta)(1-\beta(1-\epsilon)p_{11} + \beta(1-\epsilon)p_{01})}, & \text{if } \omega_{\beta}^*(m) < p_{11} \\ \frac{(1-\beta(1-(1-\epsilon)p_{01}))m + \beta(1-\epsilon)\mathcal{T}(p_{11})(1-\beta) + \beta^2(1-\epsilon)p_{01}}{(1-\beta)(1+\beta(1+\beta)(1-\epsilon)p_{01} - \beta^2(1-\epsilon)\mathcal{T}(p_{11}))}, & \text{if } p_{11} \leq \omega_{\beta}^*(m) < \mathcal{T}(p_{11}) \\ \frac{m}{1-\beta}, & \text{if } \omega_{\beta}^*(m) \geq \mathcal{T}(p_{11}) \end{cases},$$

$$V_{\beta, m}(p_{01}) = \begin{cases} \frac{(1-\epsilon)p_{01} + \beta(1-\epsilon)p_{01}V_{\beta, m}(p_{11})}{1-\beta(1-(1-\epsilon)p_{01})}, & \text{if } \omega_{\beta}^*(m) < p_{01} \\ \frac{m}{1-\beta}, & \text{if } \omega_{\beta}^*(m) \geq p_{01} \end{cases}.$$

The approximate Whittle index is given by

$$W(\omega) = \begin{cases} \frac{\omega(1-\epsilon)(1-\beta p_{11} + \beta p_{01})}{1-\beta(1-\epsilon)p_{11} + \beta(1-\epsilon)p_{01}}, & \text{if } \omega \leq p_{11} \\ \frac{(1-\epsilon)(1-\beta + C_4\beta)(\beta p_{01} + \omega - \beta \mathcal{T}(\omega))}{1-\beta(1-(1-\epsilon)p_{01}) + (1-\epsilon)C_3\beta(\beta \mathcal{T}(\omega) - \beta p_{01} - \omega)}, & \text{if } p_{11} < \omega < \omega_0 \\ \frac{(1-\epsilon)(1-\beta + \beta C_4)(\beta p_{01} + \omega(1-\beta))}{1-\beta(1-(1-\epsilon)p_{01}) - (1-\epsilon)\beta C_3(\beta p_{01} + \omega - \beta \omega)}, & \text{if } \omega_0 \leq \omega < \mathcal{T}(p_{11}) \\ \frac{(1-\epsilon)(\beta p_{01} + (1-\beta)\omega)}{1+(1-\epsilon)\beta(p_{01} - \omega)}, & \text{if } \mathcal{T}(p_{11}) \leq \omega < p_{01} \\ (1-\epsilon)\omega, & \text{if } \omega \geq p_{01} \end{cases},$$

where

$$C_3 = \frac{1-\beta(1-(1-\epsilon)p_{01})}{1+\beta(1+\beta)(1-\epsilon)p_{01} - \beta^2(1-\epsilon)\mathcal{T}(p_{11})},$$

$$C_4 = \frac{\beta(1-\epsilon)\mathcal{T}(p_{11})(1-\beta) + \beta^2(1-\epsilon)p_{01}}{1+\beta(1+\beta)(1-\epsilon)p_{01} - \beta^2(1-\epsilon)\mathcal{T}(p_{11})}.$$

3.2 Algorithm

Our analysis leads to the algorithm for the RMAB model with imperfect observations in Algorithm 1 for all $\beta \in (0, 1)$.

Algorithm 1 Whittle Index Policy

Input: $\beta \in (0, 1)$, $T \geq 1, N \geq 2$, $1 \leq M < N$, iteration number k

Input: initial belief state $\omega_n(1)$, $\mathbf{P}^{(n)}$, B_n , $n = 1, \dots, N$

```
1: for  $t = 1, 2, \dots, T$  do
2:   for  $n = 1, \dots, N$  do
3:     Set the threshold  $\omega_\beta^*(m) = \omega_n(t)$  in (59)
4:     Compute  $L(p_{11}^{(n)}, \omega_n(t))$  and set  $\omega = p_{11}^{(n)}$  in (59)
5:     Expand (59) to the  $k$ th step and solve for  $V_{\beta, m}^{(n)}(p_{11}^{(n)})$ 
6:     Compute  $L(\mathcal{T} \circ \phi(\omega), \omega_n(t))$  and set  $\omega = \mathcal{T} \circ \phi(\omega)$  in (59)
7:     Expand (59) to the  $k$ th step and solve for  $V_{\beta, m}^{(n)}(\mathcal{T} \circ \phi(\omega))$  from  $V_{\beta, m}^{(n)}(p_{11}^{(n)})$ 
8:     Compute  $L(\mathcal{T}(\omega), \omega_n(t))$  and set  $\omega = \mathcal{T}(\omega)$  in (59)
9:     Expand (59) to the  $k$ th step and solve for  $V_{\beta, m}^{(n)}(\mathcal{T}(\omega))$  from  $V_{\beta, m}^{(n)}(p_{11}^{(n)})$ 
10:    Solve for  $V_{\beta, m}^{(n)}(\omega; u = 1)$  by  $V_{\beta, m}^{(n)}(p_{11}^{(n)})$  and  $V_{\beta, m}^{(n)}(\mathcal{T} \circ \phi(\omega))$  as in (57)
11:    Solve for  $V_{\beta, m}^{(n)}(\omega; u = 0)$  by  $V_{\beta, m}^{(n)}(\mathcal{T}(\omega))$  as in (58)
12:    Evaluate the solvability of the linear equation of  $m$ :  $V_{\beta, m}^{(n)}(\omega; u = 1)$ 
       $= V_{\beta, m}^{(n)}(\omega; u = 0)$ 
13:    Set  $W(\omega_n(t)) = \omega_n(t)B_n$  and skip Step 14 if the above is unsolvable
14:    Compute  $W(\omega_n(t))$  as the solution to  $V_{\beta, m}^{(n)}(\omega; u = 1) = V_{\beta, m}^{(n)}(\omega; u = 0)$ 
15:  end for
16:  Choose the top  $M$  arms with the largest Whittle Indices  $W(\omega_n(t))$ 
17:  Observe the selected  $M$  arms and accrue reward  $O_n(t)S_n(t)B_n$  from each
18:  observed arm
19:  for  $n = 1, \dots, N$  do
20:    Update the belief state  $\omega_n(t)$  according to (4)
21:  end for
22: end for
```

4 Optimality for Homogeneous Systems

A space-wise homogeneous system for a restless bandit is defined as the system with N stochastically identical arms, *i.e.*, the parameters $\mathbf{P}^{(n)}$ and B_n do not depend on n . In this case, our algorithm is equivalent to the myopic policy that chooses the arms with the largest belief values and is optimal.

Theorem 4. *Consider a space-wise homogeneous model with positively correlated arms ($p_{11} \geq p_{01}$) and ϵ satisfying*

$$\epsilon \leq \frac{p_{01}(1 - p_{11})}{p_{11}(1 - p_{01})} = \frac{p_{01}p_{10}}{p_{11}p_{00}}, \quad (61)$$

the myopic policy is optimal over both finite and infinite horizons.

Proof. We adopt notations similar to that in Liu et al. (2011) for the case of perfect observation ($\epsilon = \delta = 0$) but need several non-trivial differences due to the additional complexity introduced by observation errors. Consider N arms in total and we choose K arms to active at each step. Let $W_s(\omega_1, \dots, \omega_N)$ denote the expected total discounted reward over s steps when all arms are ordered so the probabilities that the underlying random processes are in state 1 are $\omega_1 \geq \dots \geq \omega_N$. In Liu et al. (2010), it has been proved that the myopic policy has a dynamic queuing structure if the error probability ϵ satisfies (61). Then we have

$$W_{s+1}(\omega_1, \dots, \omega_N) = (1 - \epsilon) \sum_{i=1}^K \omega_i + \beta \mathbb{E}[W_s(p_{11}, \dots, p_{11}, \tau(\omega_{K+1}), \dots, \tau(\omega_N), \sigma(\cdot), \dots, \sigma(\cdot))],$$

where $W_0(\cdot) = 0$, $\tau(\omega) = p_{11}\omega + p_{01}(1 - \omega)$, $\sigma(\omega) = \tau(\frac{\epsilon\omega}{\epsilon\omega+1-\omega})$, and the expectation is taken over possible outcomes that can occur when the K arms that are observed are those at the left end (*i.e.*, having probabilities $\omega_1, \dots, \omega_K$ that the underlying random processes are in state 1). We will describe it more specifically. For a belief state sequence $(\omega_1, \dots, \omega_K)$, we call $(\omega_{i_1}, \dots, \omega_{i_s})$ and $(\omega_{j_1}, \dots, \omega_{j_t})$ a partition of $(\omega_1, \dots, \omega_K)$ if they satisfy:

- (i) $(\omega_{i_1}, \dots, \omega_{i_s}, \omega_{j_1}, \dots, \omega_{j_t})$ is a rearrangement of $(\omega_1, \dots, \omega_K)$;
- (ii) $i_1 < \dots < i_s$ and $j_1 < \dots < j_t$.

Let \mathcal{P} be the collection of all partitions of $(\omega_1, \dots, \omega_K)$, the expectation above can be written as follows:

$$\begin{aligned} & \mathbb{E}[W_s(p_{11}, \dots, p_{11}, \tau(\omega_{K+1}), \dots, \tau(\omega_N), \sigma(\cdot), \dots, \sigma(\cdot))] \\ &= \sum_{\mathcal{P}} \left(\prod_{m=1}^s (1 - \epsilon)\omega_{i_m} \right) \left(\prod_{n=1}^t (1 - (1 - \epsilon)\omega_{j_n}) \right) \\ & \quad \times W_s(p_{11}, \dots, p_{11}, \tau(\omega_{K+1}), \dots, \tau(\omega_N), \sigma(\omega_{j_1}), \dots, \sigma(\omega_{j_t})). \end{aligned}$$

We can see that when $\omega_1 \geq \omega_2 \geq \dots \geq \omega_N$, $W_s(\omega_1, \dots, \omega_N)$ is the value function for the myopic policy.

To prove the theorem, we first prove that $\forall s$, $W_s(\omega_1, \dots, \omega_N)$ is linear in ω_i ($1 \leq i \leq N$). We will prove it by induction. It is obvious that $W_0 = 0$, $W_1(\omega_1, \dots, \omega_N) = (1 - \epsilon) \sum_{i=1}^K \omega_i$ are linear in $\omega_1, \dots, \omega_N$. Assume it is true for s , *i.e.*, $W_s(\omega_l) = W_s(\omega_1, \dots, \omega_l, \dots, \omega_N) = a_l \omega_l + b_l$ ($\forall 1 \leq l \leq N$), where a_l and b_l are constants independent of ω_l . Now consider $W_{s+1}(\omega_l) = W_{s+1}(\omega_1, \dots, \omega_l, \dots, \omega_N)$. When $l > K$,

$$W_{s+1}(\omega_1, \dots, \omega_l, \dots, \omega_N) = (1 - \epsilon) \sum_{i=1}^K \omega_i + \beta \mathbb{E}[W_s(p_{11}, \dots, p_{11}, \tau(\omega_{K+1}), \dots, \tau(\omega_l), \dots, \tau(\omega_N), \sigma(\cdot), \dots, \sigma(\cdot))]$$

In this case, probability terms in expectation are only related to $\omega_1, \dots, \omega_K, \tau(\omega_l)$ is linear in ω_l and W_s is linear in $\tau(\omega_l)$, thus W_{s+1} is also linear in ω_l . When $l \leq k$, for any partition $(\omega_{i_1}, \dots, \omega_{i_s})$ and $(\omega_{j_1}, \dots, \omega_{j_t})$ of $(\omega_1, \dots, \omega_K)$, if $\omega_l \in \{\omega_{i_1}, \dots, \omega_{i_s}\}$, the corresponding term

$$\left(\prod_{m=1}^s (1 - \epsilon)\omega_{i_m} \right) \cdot \left(\prod_{n=1}^t (1 - (1 - \epsilon)\omega_{j_n}) \right) \\ \times W_s(p_{11}, \dots, p_{11}, \tau(\omega_{K+1}), \dots, \tau(\omega_N), \sigma(\omega_{j_1}), \dots, \sigma(\omega_{j_t}))$$

in the expectation is linear in ω_l . If $\omega_l \in \{\omega_{j_1}, \dots, \omega_{j_t}\}$, by inductive hypothesis there exists \tilde{a}, \tilde{b} ,

$$(1 - (1 - \epsilon)\omega_l)W_s(p_{11}, \dots, p_{11}, \tau(\omega_{K+1}), \dots, \tau(\omega_N), \sigma(\omega_{j_1}), \dots, \sigma(\omega_l), \dots, \sigma(\omega_{j_t})) \\ = (1 - (1 - \epsilon)\omega_l)(\tilde{a}\sigma(\omega_l) + \tilde{b}) = \tilde{a}(\epsilon p_{11}\omega_l + p_{01}(1 - \omega_l)) + \tilde{b}(1 - (1 - \epsilon)\omega_l).$$

The equation above shows that W_{s+1} is linear in ω_l , thus the proposition is proved.

From above, we can assume $W_s(\omega_1, \dots, x, \dots, y, \dots, \omega_N) = ax + by + cxy + d$, where a, b, c, d are constants. If we swap the positions of x and y and make differences between the two, we have

$$W_s(\omega_1, \dots, x, \dots, y, \dots, \omega_N) - W_s(\omega_1, \dots, y, \dots, x, \dots, \omega_N) \\ = (x - y)[W_s(\omega_1, \dots, 1, \dots, 0, \dots, \omega_N) - W_s(\omega_1, \dots, 0, \dots, 1, \dots, \omega_N)].$$

Next we will prove two important properties of W_s . We let $\bar{\omega}_i$ denote any sequence of ω_i s, possibly empty. We still adopt induction to prove next two properties:

- (A) $1 - \epsilon + W_s(\bar{\omega}_1, y, \bar{\omega}_2, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_2, y, \bar{\omega}_3) \geq 0$.
- (B) $\forall y > x, W_s(\bar{\omega}_1, y, \bar{\omega}_2, x, \bar{\omega}_3) - W_s(\bar{\omega}_1, x, \bar{\omega}_2, y, \bar{\omega}_3) \geq 0$.

These are clearly true for $s = 1$. We will begin by proving an induction step for (B). As above, the expression in (B) is equal to $(y - x)[W_s(\bar{\omega}_1, 1, \bar{\omega}_2, 0, \bar{\omega}_3) - W_s(\bar{\omega}_1, 0, \bar{\omega}_2, 1, \bar{\omega}_3)]$. Suppose the position exchange occur in the i th and j th place, $i < j$. If $i, j \leq K$, for some $\bar{\omega}_1', \bar{\omega}_2', \bar{\omega}_3'$ (which are stochastically determined by the observations from the top K arms in the queue), by inductive hypothesis,

$$W_s(\bar{\omega}_1, 1, \bar{\omega}_2, 0, \bar{\omega}_3) - W_s(\bar{\omega}_1, 0, \bar{\omega}_2, 1, \bar{\omega}_3) \\ = \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{01}, \bar{\omega}_3') - W_{s-1}(\bar{\omega}_1', p_{01}, \bar{\omega}_2', \bar{\omega}_3')] \geq 0.$$

Similarly if $i, j > K$,

$$W_s(\bar{\omega}_1, 1, \bar{\omega}_2, 0, \bar{\omega}_3) - W_s(\bar{\omega}_1, 0, \bar{\omega}_2, 1, \bar{\omega}_3) \\ = \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', p_{11}, \bar{\omega}_2', p_{01}, \bar{\omega}_3') - W_{s-1}(\bar{\omega}_1', p_{01}, \bar{\omega}_2', p_{11}, \bar{\omega}_3')] \geq 0.$$

The interesting case is $i \leq K < j$. In this case,

$$\begin{aligned}
& W_s(\bar{\omega}_1, 1, \bar{\omega}_2, 0, \bar{\omega}_3) - W_s(\bar{\omega}_1, 0, \bar{\omega}_2, 1, \bar{\omega}_3) \\
&= 1 - \epsilon + \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', p_{11}, \bar{\omega}_2', p_{01}, \bar{\omega}_3', \bar{\omega}_4') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, \bar{\omega}_3', p_{01}, \bar{\omega}_4')] \\
&\geq 1 - \epsilon + \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, p_{01}, \bar{\omega}_3', \bar{\omega}_4') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, \bar{\omega}_3', p_{01}, \bar{\omega}_4')] \\
&= (1 - \epsilon)(1 - \beta) + \beta \mathbb{E}[(1 - \epsilon) \\
&\quad + W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, p_{01}, \bar{\omega}_3', \bar{\omega}_4') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, \bar{\omega}_3', p_{01}, \bar{\omega}_4')] \\
&\geq 0,
\end{aligned}$$

where the first inequality follows from the inductive hypothesis for (B) and the second follows from (A).

Next we will prove (A) by induction. Suppose that y occurs within the two expressions in the i th and j th place, $i < j$. If $i, j \leq K$, similarly for some $\bar{\omega}_1', \bar{\omega}_2', \bar{\omega}_3'$ (they depend on the observations),

$$\begin{aligned}
& 1 - \epsilon + W_s(\bar{\omega}_1, y, \bar{\omega}_2, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_2, y, \bar{\omega}_3) \\
&= 1 - \epsilon + \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', \sigma(y), \bar{\omega}_2', \bar{\omega}_3') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', \sigma(y), \bar{\omega}_3')] \\
&= (1 - \epsilon)(1 - \beta) + \beta \mathbb{E}[(1 - \epsilon) + W_{s-1}(\bar{\omega}_1', \sigma(y), \bar{\omega}_2', \bar{\omega}_3') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', \sigma(y), \bar{\omega}_3')] \\
&\geq 0.
\end{aligned}$$

If $i, j > K$, we have

$$\begin{aligned}
& 1 - \epsilon + W_s(\bar{\omega}_1, y, \bar{\omega}_2, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_2, y, \bar{\omega}_3) \\
&= (1 - \epsilon)(1 - \beta) + \beta \mathbb{E}[(1 - \epsilon) + W_{s-1}(\bar{\omega}_1', \tau(y), \bar{\omega}_2', \bar{\omega}_3') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', \tau(y), \bar{\omega}_3')] \\
&\geq 0.
\end{aligned}$$

The interesting case is $i \leq K < j$. Let $\bar{\omega}_2 = (\bar{\omega}_{21}, x, \bar{\omega}_{22})$, where $\bar{\omega}_1$ and $\bar{\omega}_{21}$ represent $K - 1$ states in total. Then

$$\begin{aligned}
& 1 - \epsilon + W_s(\bar{\omega}_1, y, \bar{\omega}_2, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_2, y, \bar{\omega}_3) \\
&= 1 - \epsilon + W_s(\bar{\omega}_1, y, \bar{\omega}_{21}, x, \bar{\omega}_{22}, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_{21}, x, \bar{\omega}_{22}, y, \bar{\omega}_3).
\end{aligned}$$

The above expression is a function of x and y , of the form $ax + by + cxy + d$. To prove the expression above is nonnegative for all $x, y \in [0, 1]$, we just need to check out that it is true for $(x, y) \in \{(0, 0), (0, 1), (1, 0), (1, 1)\}$.

If $x = y = 0$, then

$$\begin{aligned}
& 1 - \epsilon + W_s(\bar{\omega}_1, 0, \bar{\omega}_{21}, 0, \bar{\omega}_{22}, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_{21}, 0, \bar{\omega}_{22}, 0, \bar{\omega}_3) \\
&= 1 - \epsilon + \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', p_{01}, \tau(\bar{\omega}_{22}), \bar{\omega}_3', p_{01}, \bar{\omega}_4') \\
&\quad - W_{s-1}(\bar{\omega}_1', \tau(\bar{\omega}_{22}), p_{01}, \bar{\omega}_3', \bar{\omega}_4', p_{01})] \\
&\geq (1 - \epsilon)(1 - \beta) + \beta \mathbb{E}[(1 - \epsilon) \\
&\quad + W_{s-1}(\bar{\omega}_1', p_{01}, \tau(\bar{\omega}_{22}), \bar{\omega}_3', p_{01}, \bar{\omega}_4') - W_{s-1}(\bar{\omega}_1', \tau(\bar{\omega}_{22}), \bar{\omega}_3', p_{01}, \bar{\omega}_4', p_{01})]
\end{aligned}$$

$$\geq 0.$$

If $x = y = 1$, then

$$\begin{aligned} & 1 - \epsilon + W_s(\bar{\omega}_1, 1, \bar{\omega}_{21}, 1, \bar{\omega}_{22}, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_{21}, 1, \bar{\omega}_{22}, 1, \bar{\omega}_3) \\ &= 1 - \epsilon + \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', p_{11}, \bar{\omega}_2', p_{11}, \tau(\bar{\omega}_{22}), \bar{\omega}_3') \\ &\quad - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, \tau(\bar{\omega}_{22}), p_{11}, \bar{\omega}_3')] \\ &\geq (1 - \epsilon)(1 - \beta) + \beta \mathbb{E}[(1 - \epsilon) \\ &\quad + W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, p_{11}, \tau(\bar{\omega}_{22}), \bar{\omega}_3') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_2', p_{11}, \tau(\bar{\omega}_{22}), p_{11}, \bar{\omega}_3')] \\ &\geq 0. \end{aligned}$$

If $x = 0, y = 1$, then

$$\begin{aligned} & 1 - \epsilon + W_s(\bar{\omega}_1, 1, \bar{\omega}_{21}, 0, \bar{\omega}_{22}, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_{21}, 0, \bar{\omega}_{22}, 1, \bar{\omega}_3) \\ &= 2(1 - \epsilon) + \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', p_{11}, \bar{\omega}_3', p_{01}, \tau(\bar{\omega}_{22}), \bar{\omega}_4') \\ &\quad - W_{s-1}(\bar{\omega}_1', \bar{\omega}_3', \tau(\bar{\omega}_{22}), p_{11}, \bar{\omega}_4', p_{01})] \\ &\geq 2(1 - \epsilon)(1 - \beta) + \beta \mathbb{E}[2 - 2\epsilon \\ &\quad + W_{s-1}(\bar{\omega}_1', \bar{\omega}_3', p_{01}, \tau(\bar{\omega}_{22}), p_{11}, \bar{\omega}_4') - W_{s-1}(\bar{\omega}_1', \bar{\omega}_3', \tau(\bar{\omega}_{22}), p_{11}, \bar{\omega}_4', p_{01})] \\ &> 0. \end{aligned}$$

If $x = 1, y = 0$, then

$$\begin{aligned} & 1 - \epsilon + W_s(\bar{\omega}_1, 0, \bar{\omega}_{21}, 1, \bar{\omega}_{22}, \bar{\omega}_3) - W_s(\bar{\omega}_1, \bar{\omega}_{21}, 1, \bar{\omega}_{22}, 0, \bar{\omega}_3) \\ &= \beta \mathbb{E}[W_{s-1}(\bar{\omega}_1', p_{11}, \tau(\bar{\omega}_{22}), \bar{\omega}_3', p_{01}, \bar{\omega}_4') - W_{s-1}(\bar{\omega}_1', p_{11}, \tau(\bar{\omega}_{22}), p_{01}, \bar{\omega}_3', \bar{\omega}_4')] \\ &\geq 0. \end{aligned}$$

Thus (A) is true. In fact, (B) shows that the myopic policy is optimal over finite horizons. By contradiction, it is easy to show that the myopic policy also maximizes the expected total discounted reward and the expected average reward over the infinite horizon. Furthermore, our proof does not depend on the time-homogeneous nature of the system so the optimality result holds even if the system parameters are time-varying as long as $p_{11}(t) \geq p_{01}(t)$ and ϵ satisfies (61). \square

5 Numerical Analysis and Conclusion

In this section, we illustrate the near-optimality and efficiency of our approximated Whittle index policy for non-homogeneous arms through simulation examples. After the discussions and illustrations on these numerical results, we conclude this paper and propose several future research directions on relevant problems.

5.1 Numerical Examples

We will show that the 4-iteration approximation algorithm is sufficient to yield the same performance as the exact Whittle index policy and use it to plot the

approximated Whittle index function $W(\omega)$ for the following parameters: in Fig. 3 $p_{11} = 0.2, p_{01} = 0.9, \beta = 0.9, \epsilon = 0.1, B = 1$; in Fig. 4 $p_{11} = 0.6, p_{01} = 0.3, \beta = 0.9, \epsilon = 0.1, B = 1$. Note that the monotonic increasing property of $W(\omega)$ implies the indexability numerically while the nonlinearity of $W(\omega)$ illustrates its difference to the myopic policy (with index ωB as a linear function in ω).

We now compare the performance of Whittle index policy with the optimal policy which is computed by dynamic programming over a finite horizon of length T . In other words, we recursively call the following equation with terminating state $V_{1,0,\beta,m}(\cdot) = 0$:

$$V_{1,T,\beta,m}(\omega) = \max\{V_{1,T,\beta,m}(\omega; u = 1); V_{1,T,\beta,m}(\omega; u = 0)\}, \quad (62)$$

where $V_{1,T,\beta,m}(\omega; u = 1)$ and $V_{1,T,\beta,m}(\omega; u = 0)$ are respectively given by (22) and (23) in terms of $V_{1,T-1,\beta,m}(\cdot)$. Clearly, the number of observed belief states grows *exponentially* with both the number of arms (as arms are not decoupled by any relaxation) and the time horizon T due to the tree-expansion type of the belief update given in (4). In contrast, our approximate Whittle index has a *linear* complexity in both T and the number of arms. In Fig. 5 and Fig. 6, we compare the real running times between the optimal policy and our algorithm to illustrate the efficiency of the latter. Note that the optimal policy for any finite time horizon T provides an upper bound on the total discounted reward over T achieved by the infinite-horizon optimal policy. This is because one can definitely apply the infinite-horizon optimal policy to the finite-horizon problem up to time T . Henceforth, the near-optimality of our algorithm is well demonstrated by comparing to the finite-horizon optimal policy over the first T steps as shown in Figures 7-22 (see Table 1 and Table 2 for system parameters). Furthermore, all numerical experiments with randomly generated system parameters showed that setting the iteration number $k = 4$ is sufficient for the Whittle indices of all arms to converge such that their rank remains the same as k increases at each time step t . In other words, setting $k = 4$ makes the approximated Whittle index have the same action path as the exact Whittle index, leading to the same performance. When k becomes smaller, the approximation error will be larger and cause more performance loss as shown in Fig. 23.

We also illustrate the performance of the myopic policy that chooses the M arms with the largest $\omega_n B_n$ for comparison. From Figures 7-22, we observe that Algorithm 1 outperforms the myopic policy. Interestingly, the Whittle index policy may have some performance loss in the middle but eventually catches up with the optimal policy as time goes. This is consistent with the conjecture that Whittle index policy is asymptotically optimal as time goes to infinity as the Lagrangian relaxation should not fundamentally alter the state and action paths of the optimal policy for the original problem from the perspective of large deviation theory (Weber and Weiss, 1990). On the contrary, the myopic policy is unable to follow the optimal action path and never catches up! Since the myopic policy only cares about maximizing the immediate reward, its performance for $T = 1$ is optimal (thus better than any other policy) because the state transitions do not matter in this case. Definitely, the myopic policy has the lowest complexity but this advantage is negligible given its significant performance loss and the efficiency of our algorithm compared to the optimal policy as shown in Fig. 5 and Fig. 6.

5.2 Conclusion and Future Work

In this paper, we proposed a low-complexity algorithm based on the idea of value function approximations arisen in solving for the Whittle index policy of a class of RMAB with an infinite state space and an imperfect observation model. By exploring and exploiting the rich mathematical structure of this RMAB model, our algorithm was designed to be implemented online to control the approximation error such that it becomes equivalent to the original Whittle index policy. Extensive numerical examples showed that our algorithm achieves a near-optimal performance with a complexity linear in the key system parameters such as the time horizon T and the number of arms. From Figures 7-22, we observe that in some instances the four-iteration policy is closer to optimality than in other instances. Unfortunately, it is still unknown how to theoretically quantify the performance gap of Whittle index policy to optimality in finite regime: only some asymptotic results were obtained for restless bandits with finite states as both the number of arms and the time horizon go to infinity under the time-average reward criterion and some conditions only verified for bandits with 2 or 3 states (Weber and Weiss, 1990, 1991).

Future work includes the theoretic study of the performance loss of Whittle index policy by more in-depth analysis on the value functions to further improve our algorithm. Another research direction is to consider more complex system models such as non-Markovian state models (see, *e.g.*, Liu et al., 2011) or high-dimensional state models (see, *e.g.*, Liu, 2021) and study the convergence patterns of the state path to approximate the value functions. Future work also includes the generalization of constructing a finite set of linear equations to solve for dynamic programming problems, *e.g.*, the simplification of non-linearity by threshold-based first crossing time, the error-control process by backward induction, and further complexity reduction by minimizing the number of linear equations required.

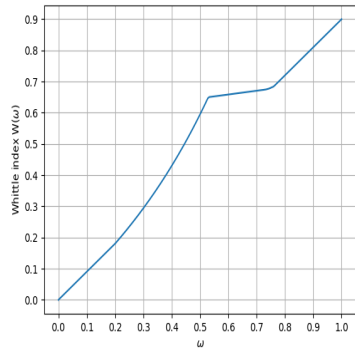


Fig. 3 Approximated $W(\omega)$ ($p_{11} < p_{01}$)

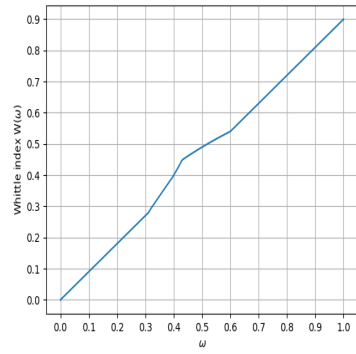


Fig. 4 Approximated $W(\omega)$ ($p_{11} > p_{01}$)

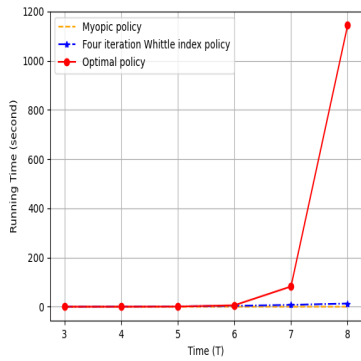


Fig. 5 The Complexity with T

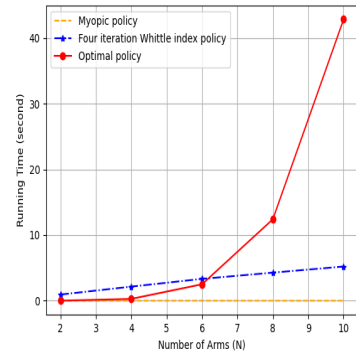


Fig. 6 The Complexity with N

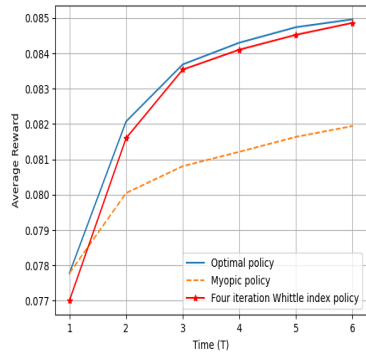


Fig. 7 Example-1

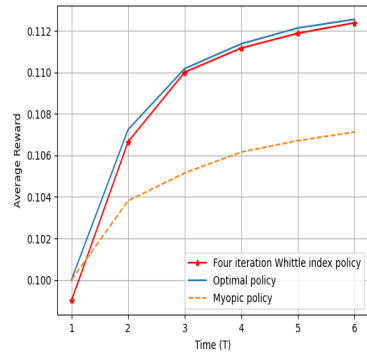


Fig. 8 Example-2

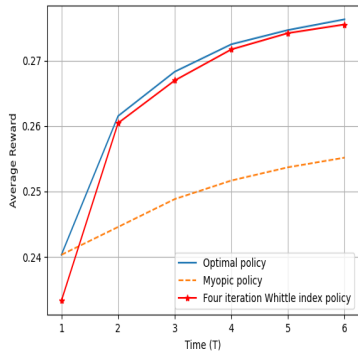


Fig. 9 Example-3

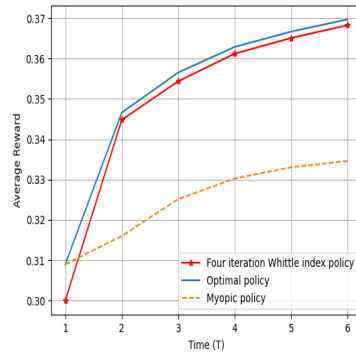


Fig. 10 Example-4

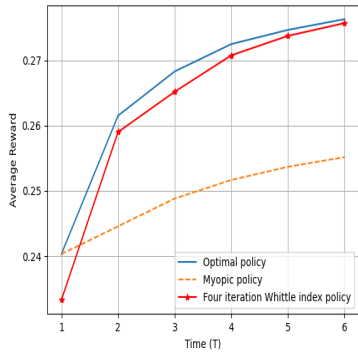


Fig. 11 Example-5

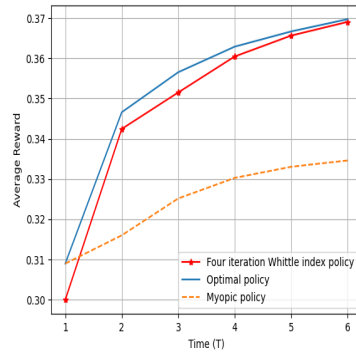


Fig. 12 Example-6

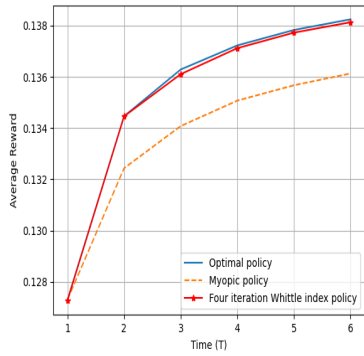


Fig. 13 Example-7

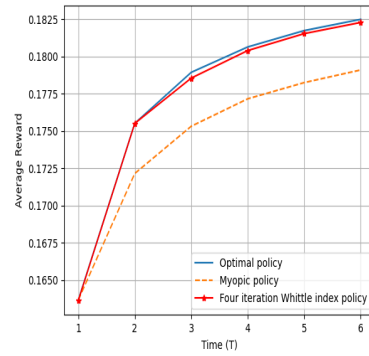


Fig. 14 Example-8

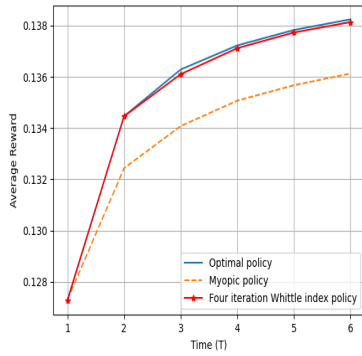


Fig. 15 Example-9

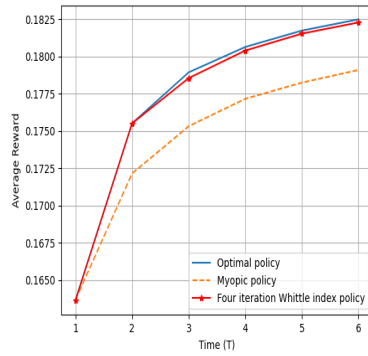


Fig. 16 Example-10

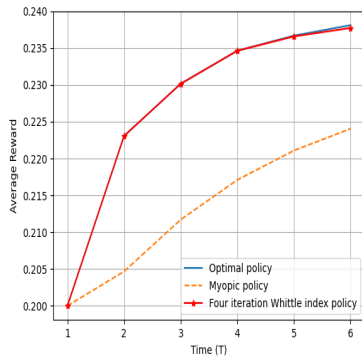


Fig. 17 Example-11

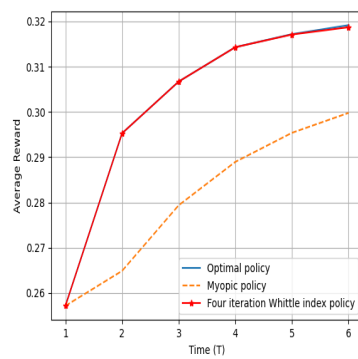


Fig. 18 Example-12

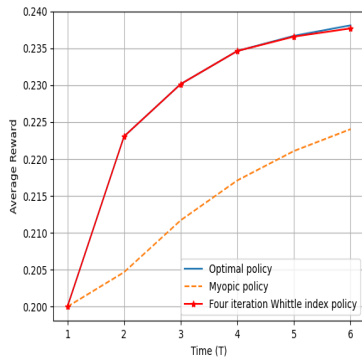


Fig. 19 Example-13

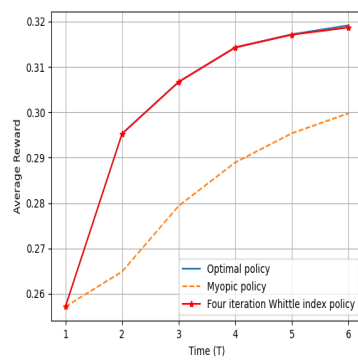


Fig. 20 Example-14

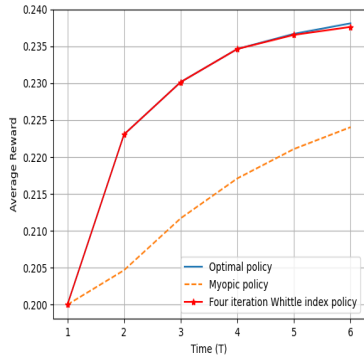


Fig. 21 Example-15

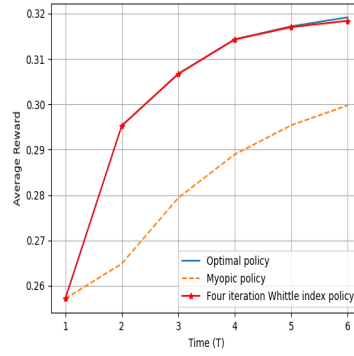


Fig. 22 Example-16

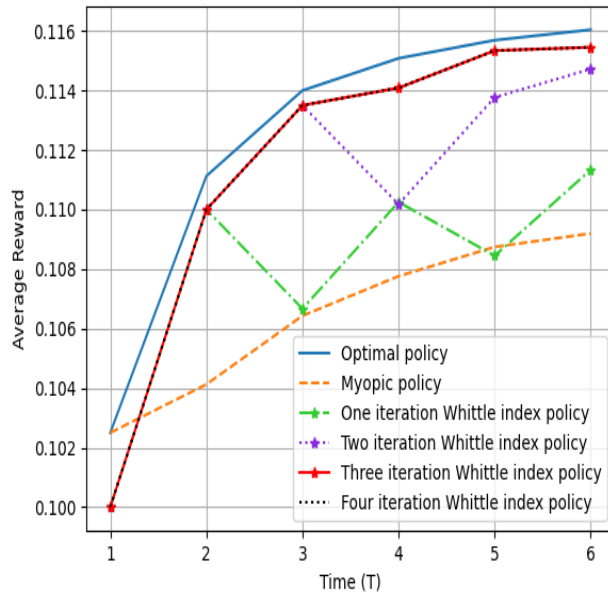


Fig. 23 The Performance Comparison for k

Table 1 Experiment Setting

System-1	$\{p_{11}^{(i)}\}_{i=1}^7$	{0.3, 0.6, 0.4, 0.7, 0.2, 0.6, 0.8}
	$\{p_{01}^{(i)}\}_{i=1}^7$	{0.1, 0.4, 0.3, 0.4, 0.1, 0.3, 0.5}
	$\{B_i\}_{i=1}^7$	{0.8800, 0.2200, 0.3300, 0.1930, 1.0000, 0.2558, 0.1549}
System-2	$\{p_{11}^{(i)}\}_{i=1}^7$	{0.6, 0.4, 0.2, 0.2, 0.4, 0.1, 0.3}
	$\{p_{01}^{(i)}\}_{i=1}^7$	{0.8, 0.6, 0.4, 0.9, 0.8, 0.6, 0.7}
	$\{B_i\}_{i=1}^7$	{0.5150, 0.6666, 1.0000, 0.6296, 0.5833, 0.8100, 0.6700}
System-3	$\{p_{11}^{(i)}\}_{i=1}^7$	{0.1, 0.4, 0.3, 0.4, 0.1, 0.3, 0.5}
	$\{p_{01}^{(i)}\}_{i=1}^7$	{0.3, 0.6, 0.4, 0.7, 0.2, 0.6, 0.8}
	$\{B_i\}_{i=1}^7$	{0.7273, 0.3636, 0.5000, 0.3377, 1.0000, 0.3939, 0.2955}
System-4	$\{p_{11}^{(i)}\}_{i=1}^7$	{0.6, 0.7, 0.2, 0.6, 0.4, 0.5, 0.3}
	$\{p_{01}^{(i)}\}_{i=1}^7$	{0.8, 0.4, 0.9, 0.5, 0.7, 0.2, 0.6}
	$\{B_i\}_{i=1}^7$	{0.4286, 0.5000, 0.5397, 0.5143, 0.5306, 1.0000, 0.6190}

Table 2 Experiment Setting (continued)

System	Example	ϵ	β	Meet threshold conditions?	Meet indexability conditions?
System-1	1	0.3	0.999	yes	no
	2	0.1	0.999	yes	no
System-2	3	0.3	0.29	yes	yes
	4	0.1	0.29	yes	yes
	5	0.3	0.48	no	yes
	6	0.1	0.48	no	yes
System-3	7	0.3	0.69	yes	no
	8	0.1	0.69	yes	no
	9	0.3	0.48	yes	yes
	10	0.1	0.48	yes	yes
System-4	11	0.3	0.29	yes	yes
	12	0.1	0.29	yes	yes
	13	0.3	0.48	no	yes
	14	0.1	0.48	no	yes
	15	0.3	0.999	no	no
	16	0.1	0.999	no	no

Declarations

- **Funding** Not applicable
- **Conflict of interest/Competing interests** Not applicable
- **Ethics approval** Not applicable
- **Consent to participate** Not applicable
- **Consent for publication** Yes
- **Availability of data and materials** Available upon request
- **Code availability** Available upon request

- **Authors' contributions** Keqin Liu constructed the proof sketch for each theorem and the main algorithm and contributed to the writing of the paper. Richard Weber outlined the proof strategy for the optimality of the myopic policy in homogeneous systems and contributed to the verification and writing of the paper. Chengzhong Zhang filled out the details of the proofs and conducted the numerical simulations.

References

- Brown DB, Simth JE (2020) Index policies and performance bounds for dynamic selection problems. *Manage Sci* 66(7):3029–3050.
- Chen M, Wu K, Song L (2021) A Whittle index approach to minimizing age of multi-packet information in IoT network. *IEEE Access* 9: 31467 - 31480.
- Gast N, Gaujal B, Khun K (2023) Testing indexability and computing Whittle and Gittins index in subcubic time. *Math Meth Oper Res* 97:391–436.
- Gast N, Gaujal B, Yan C (2021, working paper) (Close to) Optimal policies for finite horizon restless bandits. https://hal.inria.fr/hal-03262307/file/LP_paper.pdf.
- Gittins JC (1979) Bandit processes and dynamic allocation indices. *J R Stat Soc* 41(2):148–177.
- Gittins JC, Glazebrook KD, Weber RR (2011) Multi-Armed Bandit Allocation Indices. Wiley, Chichester.
- Hu W, Frazier PI (2017, working paper) An asymptotically optimal index policy for finite-horizon restless bandits. <https://arxiv.org/abs/1707.00205>.
- Heinonen J (2005) *Lectures on Lipschitz analysis*. www.math.jyu.fi/research/reports/rep100.pdf.
- Kesav K, Rahul M, Varun M, Shabbir NM (2019). Sequential decision making with limited observation capability: Application to wireless networks. *IEEE Trans Cognit Commun Network* 5(2):237–251.
- Levy BC (2008) Principles of Signal Detection and Parameter Estimation. Springer, Verlag.
- Liu K (2021) Index policy for a class of partially observable Markov decision processes. <https://arxiv.org/abs/2107.11939>.
- Liu K, Weber RR, Zhao Q (2011) Indexability and Whittle index for restless bandit problems involving reset processes. *Proc. of the 50th IEEE Conference on Decision and Control* 7690–7696.
- Liu K, Zhao Q (2010) Indexability of restless bandit problems and optimality of Whittle index for dynamic multichannel access. *IEEE Trans Inf Theory*

56(11):5547–5567.

- Liu K, Zhao Q, Krishnamachari B (2010) Dynamic multichannel access with imperfect channel state detection. *IEEE Trans Signal Process* 2795–2808.
- Papadimitriou CH, Tsitsiklis JN (1999) The complexity of optimal queueing network control. *Math Oper Res* 24(2):293–305.
- Rahul, M., D. M., & Aditya, G. (2018). On the whittle index for restless multiarmed hidden markov bandits. *IEEE Trans Autom Control* 63(9):3046–3053.
- Sondik EJ (1978) The optimal control of partially observable Markov processes over the infinite horizon: discounted costs. *Oper Res* 26(2):282–304.
- Varun, M., Rahul, M., Kesav, K., Shabbir, N.M. & Uday, B.D. (2018). Rested and restless bandits with constrained arms and hidden states: Applications in social networks and 5g networks. *IEEE Access* 6:56782–56799.
- Wang K, Chen L, Yu J, Win M (2018) Opportunistic Multichannel Access with Imperfect Observation: A Fixed Point Analysis on Indexability and Index-based Policy. *Proceedings of IEEE INFOCOM*.
- Weber RR, Weiss G (1990) On an index policy for restless bandits. *J Appl Probab* 27:637–648.
- Weber RR, Weiss G (1991) Addendum to ‘On an index policy for restless bandits’. *Adv Appl Prob* 23:429–430.
- Whittle P (1988) Restless bandits: Activity allocation in a changing world. *J Appl Probab* 25:287–298.
- Zayas-Cabán G, Jasin S, Wang G (2019) An asymptotically optimal heuristic for general nonstationary finite-horizon restless multi-armed, multi-action bandits. *Adv Appl Probab* 51:745–772.