

Location Aware Opportunistic Bandwidth Sharing between Static and Mobile Users with Stochastic Learning in Cellular Networks

Arpan Chattopadhyay, Bartłomiej Błaszczyszyn, Eitan Altman

Abstract—In this paper, we consider the problem of location-dependent opportunistic bandwidth sharing between static and mobile (i.e., moving) downlink users in a cellular network. Each cell of the network has some fixed number of static users. Mobile users enter the cell, move inside the cell for some time and then leave the cell. In order to provide higher data rate to the highly mobile users whose fast fading channel variation is difficult to track, we propose location dependent bandwidth sharing between the two classes of static and mobile users; the idea is to provide higher bandwidth to the mobile users at favourable locations, and provide higher bandwidth to the static users in other times. Our approach is agnostic to the way the bandwidth is further shared within the same class of users; it can be combined with any particular bandwidth allocation policy employed for one of these two classes of users. We formulate the problem as a long run average reward Markov decision process (MDP) where the per-step reward is a linear combination of instantaneous data volumes received by static and mobile users, and find the optimal policy. The optimal policy is binary in nature; it allocates the entire bandwidth either to the static users or to the mobile users at any given time. The reward structure of this MDP is not known in general, and it may change with time. To alleviate these issues, we propose a learning algorithm based on single timescale *stochastic approximation*. Also, noting that the MDP problem can be used to maximize the long run average data rate for mobile users subject to a constraint on the long run average data rate of static users, we provide a learning algorithm based on multi-timescale stochastic approximation. We prove asymptotic convergence of the bandwidth sharing policies under these learning algorithms to the optimal policy. The results are extended to address the issue of fair bandwidth sharing between the two classes of static and mobile users, where the notion of fairness is motivated by the popular notion of α -fairness in the literature. Numerical results exhibit significant performance improvement by our scheme, as well as fast convergence, and also demonstrate the trade-off between performance gain and fairness requirement.

Index Terms—Cellular network, mobility, dynamic bandwidth sharing, location-dependent bandwidth sharing, fair allocation, Markov decision process, stochastic approximation.



1 INTRODUCTION

In recent years, cellular traffic has shown unprecedented growth, due to the proliferation of high-specification handheld/mobile devices such as smartphones and tablets. It is speculated that increasing use of applications such as video streaming or downloading, image or media file transfer, social networking applications and cloud services (requested or run by these devices) will further increase this traffic demand in the coming years. In order to meet the enormous bandwidth demand for these applications, the use of macro-assisted small cell networks (see [1], [2], [3], [4], [5]) have recently become popular; the small cells (such as femtocells and picocells) can meet the bandwidth demand of the users, while the macro base stations are supposed to provide cellular coverage.

While small cell networks can provide high through-

put to the static users, the performance of mobile users (i.e., fast moving users) deteriorates due to frequent handoff at cell boundaries resulting in huge signaling overhead (see [6]) and temporary data outage for each handoff. As a solution to this problem, the use of heterogeneous network architecture has been proposed (see [7]), where only macro base stations can serve the mobile users; the relatively large cell size of the macro base stations result in a much smaller handoff rate for mobile users in this architecture. Static users can be served by either macro or micro base stations. However, this alone is not capable of meeting the growing traffic demand from mobile users, and hence new improvements in PHY and MAC techniques are essential.

In order to address the above issue, we propose opportunistic (and dynamic) sharing of the total allocated bandwidth to a base station, by the two classes of static and mobile *downlink* users, based on user locations. The transmission bandwidth available for a macro base station can be shared among its users in many ways. However, the interference field and downlink path-loss vary over various locations inside a macro cell, due to spatio-temporal variation in fast fading, distance and shadowing from various interfering base stations. Hence,

Arpan Chattopadhyay is with Electrical Engineering department, University of Southern California, Los Angeles, USA. Bartłomiej Błaszczyszyn are with Inria/ENS, Paris, France. Eitan Altman is with Inria, Sophia-Antipolis, France. Email: achattop@usc.edu, Bartek.Blaszczyszyn@ens.fr, eitan.altman@sophia.inria.fr .

This work was done when Arpan Chattopadhyay was a postdoctoral researcher in Inria/ENS, Paris, France.

All appendices are provided in the supplementary material.

a natural way to improve user throughput is to employ dynamic bandwidth sharing among the static and mobile users inside a macro cell, depending on their instantaneous location, direction of motion and speed; the idea is to provide more bandwidth to the mobile users opportunistically when they are at favourable locations, in a distributed fashion so that the base stations need not communicate among themselves to decide on bandwidth allocation. This approach also alleviates the problem of measuring fast fading channel variations from the base station to the highly mobile users. Our goal is to maximize the time average of a linear combination of the expected data rates of mobile and static users. We formulate the problem as an average reward Markov decision process (MDP), and establish the policy structure. However, the decision making requires information on the location of other base stations and the shadowing realizations from other base stations to various locations in the macro cell; these quantities might not be known to the macro base station, and some of them might even change over time. Hence, instantaneous data rate for a fast moving mobile user may not be computable in the presence of the spatially varying unknown interference field; only the cumulative amount of data downloaded by the mobile user over an interval will be available to the macro BS. Hence, we provide a learning algorithm using the theory of stochastic approximation, and prove its asymptotic convergence to the optimal dynamic bandwidth sharing policy. Next, we propose a learning algorithm based on multi-timescale stochastic approximation, which converges to the optimal policy for the problem of maximizing the time-average expected data rate of mobile users subject to a constraint on the time-average expected data rate of static users. Hence, the learning algorithms can be used by the macrocells to dynamically adapt the bandwidth sharing policy depending on mobile user locations. We also explain how to adapt the dynamic bandwidth sharing scheme when fair bandwidth sharing between the classes of static and mobile (i.e., moving) users is required. Finally, we demonstrate numerically that the proposed dynamic (opportunistic) bandwidth allocation scheme can improve user performance significantly, and also demonstrate the trade-off between performance improvement and a measure of the degree of fairness in allocation.

1.1 Related Work

There has been a vast literature on the impact of user mobility in wireless networks. The authors in [8] have shown that mobility increases the capacity. [9] has explored the trade-off between delay and throughput in ad-hoc networks in presence of mobility. The papers [10], [11], [12], [13], [14] study the impact of inter and intra cell mobility on capacity, and also the trade-off between throughput and fairness; these results show that mobility increases the capacity of cellular networks when base

stations cooperate among themselves.

However, in practice, base stations may not cooperate. Moreover, due to frequent handoff of fast moving mobile users, significant control bandwidth has to be dedicated for handoff management; it is often the case that handoff results in temporary data outage for mobile users. In order to optimize the performance of cellular networks under user mobility, we propose to use *optimal* dynamic bandwidth sharing between the two classes of static and mobile users (depending on user locations); this problem is formulated as an average reward MDP (where the reward is a time average linear combination of data rates of static and mobile users) and later learning algorithms for computation of the optimal policy are provided. There have been some work in the literature relevant to our paper. The authors of [15] also have evaluated gain in performance due to mobility by favouring users with good radio channel conditions, however they did not propose any optimal bandwidth allocation scheme among users under mobility when channel condition may not be measured accurately. The paper [16] deals with proportional fair scheduling algorithm for a *fixed* population of users with time-varying channel conditions due to mobility; this paper proposes a *single-timescale* stochastic approximation based algorithm to estimate throughput of each user, and analyses its convergence. The paper [17] essentially considers location based proportional fair scheduling over a finite time horizon to a fixed user population, but it does not provide any convergence analysis of the proposed algorithm. The authors of [17] allocate bandwidth among users opportunistically via the construction of a spatial radio map which depends on path-loss and slow fading but is averaged over fast fading; in other words, they pursue a more experimental and data-driven approach. Reference [18] solves the problem of bandwidth allocation among users as a static optimization problem that yields the fraction of bandwidth to be allocated to each user at a given state; this method maximizes the sum throughput of users, is easy to implement, but it requires the knowledge of user mobility statistics at the base station.

However, to the best of our knowledge, there has been no prior work that considers optimal dynamic bandwidth sharing depending on user location in order to maximize the sum data rate of mobile users subject to a constraint on the time-average sum data rate of static users, and proposes any learning algorithm based on multi-timescale stochastic approximation for this problem with provable convergence guarantee; in our current paper, we seek to address these problems.

1.2 Organization and Our Contribution

The rest of the paper is organized as follows:

- The system model is described in Section 2.
- In Section 3, we develop optimal bandwidth sharing strategy between the two classes of static and mobile users in a *single* cell, so as to maximize the time average of a linear combination of expected sum

throughputs of static and mobile users inside the cell. This unconstrained optimization problem is formulated as an average reward Markov decision process (MDP), and optimal policy structure is derived analytically. To the best of our knowledge, this model and mathematical contributions including the specific policy structures are new contributions to the literature and they can be used in practical wireless cellular networks.

- In Section 4, we provide a learning algorithm based on stochastic approximation, which converges asymptotically to the optimal bandwidth sharing policy, without using the transition and cost structure of the MDP.
- Noting that the unconstrained optimization problem can be used to solve the constrained problem of maximizing the time average sum rate of the mobile users subject to a constraint on the time-average sum rate of the static users, we provide, in Section 5, a learning algorithm based on multi-timescale stochastic approximation, which *provably* converges to the optimal policy for the constrained problem. This multi-timescale stochastic approximation based learning algorithm yields a randomized policy, and the randomization technique proposed in this paper is novel to the literature. This randomized bandwidth allocation technique can be used in practical cellular network where a precise radio map of the cell is not available.
- In Section 6, we show how the dynamic (and opportunistic) bandwidth sharing schemes developed in previous sections can be adapted to ensure a fair allocation between the two classes of static and mobile users; we specifically adapt the notion of α -fairness and extend our algorithms to this framework.
- In Section 7, we numerically demonstrate considerable performance gain due to opportunistic bandwidth sharing, and also explore the trade-off between performance gain and fairness in allocation. Fast convergence of the proposed learning algorithm is also demonstrated.
- In Section 8, we show the equivalence of the global problem of decentralized maximization of the time average of a linear combination of the expected sum throughputs of mobile and static users, with a problem where each base station seeks to maximize the time average of a linear combination of the expected sum rates of all mobile users and all static users via location aware opportunistic bandwidth sharing between the two classes of static and mobile users. We also explain how to modify our algorithms in case the location of users in a cell are not known perfectly, thereby extending the algorithms to a more practical regime. We also motivate the need for location based bandwidth sharing instead of channel estimation based bandwidth allocation. It has also been argued how to extend the proposed algorithms for more generalised system model.

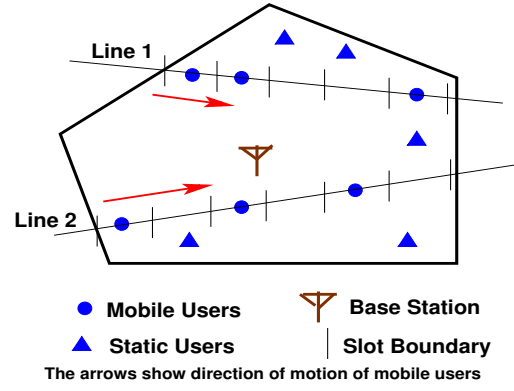


Fig. 1. A snapshot of one cell with base station, static users and mobile users. Two lines (line 1 and line 2) cross the cell; these lines have lengths $l_1 v \sigma$ and $l_2 v \sigma$ respectively inside the cell, where $l_1 = 5$, $l_2 = 6$, v is the velocity of mobile users and σ is the time slot duration. Mobile users entering the cell stay for l_1 or l_2 slots and then leave the cell. In each slot, depending on the instantaneous location of all users inside the cell, bandwidth is shared opportunistically between the classes of static and mobile users.

- Finally, we conclude in Section 9.
- All proofs are provided in the appendix.

2 SYSTEM MODEL

We consider a cellular network with multiple (possibly infinite and heterogeneous) base stations (BSs) on the two dimensional plane. Among these BSs, we consider one *single* BS and focus on the cell served by that BS (see Figure 1); this BS can be a macro BS if the network is heterogeneous. We consider two classes of *downlink* users served by this BS: *Static users* (SU) and *mobile users* (MU). We assume that there exist multiple directed lines/routes (e.g., roads) crossing the cell, and MUs are moving along these lines with constant speed v . This can be a model for the roads in urban or suburban areas where users sitting in fast moving cars download contents from the base stations. Given a realization of the line segments inside the cell, we assume that, MUs are entering a cell along each line according to a time-homogeneous process, and the arrival rate is potentially different along different routes. We assume that all the base stations transmit simultaneously (either on the same band or using frequency reuse), and these transmissions create interference at the SUs and MUs.

In order to mathematically formulate the dynamic bandwidth sharing problem, we make the following simplified modeling assumptions (also, see Figure 1 for a clear pictorial description):

- Time is discretized into slots of duration σ . Hence, a MU moves $v \sigma$ distance in one slot.
- The BS under consideration knows the locations of all static users associated with it.

- The BS knows the lines intersecting with its cell, and also the lengths of these corresponding line segments. Let us denote the number of line segments intersecting with the cell under consideration, by n . Let the i -th line segment have length $l_i \sigma v$, i.e., a MU can remain in the i -th line segment for l_i number of time slots. Thus, the model is as follows: any MU that enters the cell (after coming out of a handoff) along the i -th line segment spends a time of l_i slots, and finally enters another cell. The values $\{l_i\}_{1 \leq i \leq n}$ are known to the base station.
- We allow the MU arrival rates along the n line segments to be unequal. We denote the number of arrivals to the cell along line i at time slot t by a *bounded* random variable $A_{i,t}$; we assume that $A_{i,t}$ is i.i.d. across t and independent across i ; the arrival process will be correlated across cells, but that does not affect the resource allocation problem for a single cell.
- At the beginning of each slot τ , the BS under consideration decides the fraction η_τ of the available bandwidth to be dedicated for transmission to the mobile users. The remaining bandwidth is assigned to the set of static users. *In each slot, a base station can allocate equal bandwidth to all available mobile users, or possibly unequal (arbitrarily) bandwidth sharing among the mobile users is done.* Similarly, the $(1 - \eta_\tau)$ fraction of bandwidth can be shared arbitrarily among all static users. *In this paper, we assume equal bandwidth sharing within one user class for the sake of illustration.* It is important to note that, for fast moving users, traditional channel estimation may not be very accurate since the user might travel the fading coherence distance very fast; hence, dynamic bandwidth allocation among users based on instantaneous channel qualities may not be feasible. This necessitates location dependent bandwidth sharing which works on a slower timescale compared to variation in fading due to high speed of users. However, our scheme of sharing bandwidth between two classes of users can well accommodate any scheduling policy employed within the same class of users. Another reason for not considering location-dependent (resp., channel quality based) bandwidth allocation to individual users (instead of user classes) is that this will result in allocation of bandwidth to the user having the best location (resp., channel quality) at any given time slot, which might be unfair to all other users. *See Section 8.6 for detailed discussion on the necessity of location-dependent bandwidth sharing instead of channel quality measurement based bandwidth allocation.*
- At the beginning of each slot, the base station gets to know the number m of existing mobile users inside the cell (including the newly arrived MU), the index set z_1, z_2, \dots, z_m of lines on which each of these mobile users are moving, and also the remaining sojourn times (in terms of slots) t_1, t_2, \dots, t_m of those mobile users. This can be done via the GPS

connection of the mobile users. Otherwise, since the base station records the time and line of entry of a new MU into the cell, and since the velocity is known, the base station can always calculate the location of any mobile station inside the cell. We define $s := (\{t_i, z_i\}_{i=1}^m)$ to be the state of the system at the beginning of a slot.

We will explain in Section 8.2 how we can relax the assumption on availability of perfect information of the system state to the decision maker.

- At state s , if all available bandwidth (assumed to be equal to 1 unit) is allocated only to MUs, then, given a bandwidth sharing scheme among all SUs and a bandwidth sharing scheme among all MUs, and given the realization of shadowing and path-loss from each BS to each location in the cell, the amount of data each MU will be able to download over a slot is a random variable since the fading process seen by each user (from the serving BS and interfering BSs) over this slot is random. However, if these quantities and the fading distribution is known, the base station can calculate the expected data volume each user will be able to download until the beginning of the next slot.¹ Let us define $R_{mobile}(s)$ to be the (random) total amount of data the MUs download **per unit bandwidth** if the entire bandwidth is allocated to MUs, and similar meaning applies for $R_{static}(s)$ (i.e., this is the random amount of data the SUs can download in a slot in case the entire bandwidth is allocated to SUs). In presence of fading, the expectations of these two random variables (expectation taken over fading distribution) are denoted by $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$. Note that, $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ are dependent on the shadowing realizations from all base stations to the static and mobile users over various locations in the cell (since they will determine the signal to interference ratio for various users at different locations).

We will assume in Section 3 that $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ are known to the BS; this assumption will be relaxed in subsequent sections.

3 OPPORTUNISTIC BANDWIDTH ALLOCATION UNDER PERFECT KNOWLEDGE OF MEAN USER RATES $\bar{R}_{mobile}(s)$ AND $\bar{R}_{static}(s)$

3.1 Markov decision process formulation

We formulate the dynamic (i.e., opportunistic) bandwidth allocation problem for a BS as a Markov decision process. We assume in this section that the base station

1. Note that, for a given realization of the location of base stations and for a given realization of the spatially varying shadowing process, the amount of interference at any location is not a random variable if there is no fading. Even the interference averaged over random, time-varying fading is a deterministic quantity. But this quantity is unknown in general to a BS which does not possess global information about the base station locations and shadowing process.

knows $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ perfectly for each state s at each slot.

3.1.1 State Space

New mobile users arrive to the cell in each slot. The state of the system at the beginning of a slot is considered. The state at the beginning of a slot (after new arrivals in the previous slot) is of the form $s := (\{t_i, z_i\}_{i=1}^m)$, where m is the number of mobile users present in the cell, t_i is the residual sojourn time of the i -th user in the cell, and $z_i \in \{1, 2, \dots, n\}$ is the index of the line along which the i -th mobile user is moving; if $z_i = k$, then $t_i \in \{1, 2, \dots, l_k\}$. Note that the state space is finite since the number of arrivals in each slot is bounded and each mobile user stays inside the cell for a bounded number of slots.

3.1.2 Action Space

Given a state, the BS takes an action $x \in [0, 1]$; x is the fraction of bandwidth the BS decides to allocate to the mobile users. Hence, our action space is $[0, 1]$. In this work, we assume that, at any given time, all static users share the $(1 - x)$ fraction of bandwidth equally among themselves, and all mobile users share the x fraction of bandwidth equally among them.²

3.1.3 State Transition

For current state $s := (\{t_i, z_i\}_{i=1}^m)$, if p MUs arrive to the cell in a slot, then the next state will be $s' = (\{(t_i - 1)^+, z_i\}_{i=1}^m, \{t_i, z_i\}_{i=m+1}^{m+p})$, where $z_i \in \{1, 2, \dots, n\} \forall i \in \{m+1, m+2, \dots, m+p\}$ is the index of the line along which the i -th new arrival at the slot enters the cell, and $t_i = l_k$ if $z_i = k$ for $i \in \{m+1, m+2, \dots, m+p\}$. In course of this, if $(t_i - 1)^+ = 0$ for any $i \in \{1, 2, \dots, m\}$, then information of that user is removed from the state since he has already left the cell. We denote the state at time slot τ by $s(\tau)$.

3.1.4 Policy

A stationary policy $\eta(\cdot|\cdot)$ is a family of probability distributions $\eta(\cdot|s)$ on the action space $[0, 1]$ conditioned on the state s ; i.e., $\eta(\cdot|s)$ denotes the probability distribution of the action taken whenever the system reaches state s . If $\eta(\cdot|s)$ is such that for each state s , the policy chooses one action with probability 1, then the policy is called a stationary deterministic policy $\eta(\cdot)$; in this case, $\eta(s)$ denotes the action taken at state s . We denote by η_τ the action taken at time τ (i.e., the fraction of bandwidth allocated to the class of MUs in slot τ); this will be equal to $\eta(s(\tau))$ if a stationary deterministic policy $\eta(\cdot)$ is used in decision making. We denote by η_τ a number in $[0, 1]$, and by $\eta(\cdot)$ a function.

2. From the optimization point of view, it will always be better to allocate x fraction of bandwidth to the *best* mobile user at a given slot, and $(1 - x)$ fraction to the *best* static user for ever. But this will result in complete starvation for many static users, and short-term service unfairness among the mobile users; each mobile users will get high data rate in some slots, and very low (possibly zero) data rate in some other slots.

3.1.5 Single Stage Reward

If the system state is $s(\tau)$ at slot τ , and if an action $\eta_\tau \in [0, 1]$ is taken, the total (random) reward for the base station at decision epoch τ is defined as

$$R(\tau) := \eta_\tau R_{mobile}(s(\tau)) + \xi(1 - \eta_\tau) R_{static}(s(\tau)).$$

3.1.6 Objective Function

Let us denote the expectation under policy $\eta(\cdot|\cdot)$ by $\mathbf{E}_{\eta(\cdot|\cdot)}$; the expectation is over the randomness in the policy and over the randomness in state evolution. We seek to solve the following problem of maximizing the time average of the expected reward per slot:

$$\sup_{\eta(\cdot|\cdot)} \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta(\cdot|\cdot)} \left(\eta_\tau \bar{R}_{mobile}(s(\tau)) + \xi(1 - \eta_\tau) \bar{R}_{static}(s(\tau)) \right) \quad (1)$$

Here $\xi \geq 0$ can be considered as a Lagrange multiplier; it captures the emphasis we put on the time average sum throughput of SUs and MUs in the objective function. This problem is an unconstrained optimization problem.

Note that, there are two expectations in this objective function: one is over randomness in the fading process (which are captured by $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$), and the other one is over the randomness in the policy and over the randomness in the state evolution (captured by $\mathbf{E}_{\eta(\cdot|\cdot)}$).

The problem (1) has a stationary, deterministic optimal policy (by standard MDP theory), which we denote by $\eta_\xi^*(\cdot)$. Under the deterministic policy $\eta_\xi^*(\cdot)$, the optimal action at state s is denoted by $\eta_\xi^*(s)$ (parametrized by ξ) or simply by $\eta^*(s)$. The optimal value for the objective in (1) is denoted by $\lambda^*(\xi)$ or simply by λ^* .

It has to be noted that, under $\eta_\xi^*(\cdot)$, we have $\lim_{t \rightarrow \infty} \frac{\sum_{\tau=1}^t R(\tau)}{t} = \lambda^*(\xi)$ almost surely (by the ergodicity of the Markov chain $\{s(\tau)\}_{\tau \geq 1}$).

Later in Section 8.1, we relate (1) to a global optimization problem over multiple cells.

3.1.7 Connection Between the Unconstrained Problem and a Constrained Problem

The unconstrained optimization problem (1) can be used to solve the following constrained optimization problem of maximizing the time-average sum data rate for the mobile users while satisfying a minimum time-average sum data rate constraint R_0 for static users:

$$\begin{aligned} & \sup_{\eta(\cdot|\cdot)} \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta(\cdot|\cdot)} \left(\eta_\tau \bar{R}_{mobile}(s(\tau)) \right) \\ \text{s.t., } & \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta(\cdot|\cdot)} \left((1 - \eta_\tau) \bar{R}_{static}(s(\tau)) \right) \geq R_0 \end{aligned} \quad (2)$$

It is well-known that by choosing an appropriate value ξ^* for ξ and solving the optimization problem (1), one

can find an optimal policy $\eta_{\xi^*}^*(\cdot|\cdot)$ for the constrained problem (2) as well.

The following standard result tells us how to choose the optimal *Lagrange multiplier* ξ^* (see [19, Theorem 4.3]):

Theorem 1: Consider the constrained problem (2). If there exists a multiplier $\xi^* \geq 0$ and a policy $\eta_{\xi^*}^*(\cdot|\cdot)$ such that $\eta_{\xi^*}^*(\cdot|\cdot)$ is an optimal policy for the unconstrained problem (1) under ξ^* and the constraint in (2) is met with equality under policy $\eta_{\xi^*}^*(\cdot|\cdot)$, then $\eta_{\xi^*}^*(\cdot|\cdot)$ is an optimal policy for the constrained problem (2) also. \square

Remark: We will see in Section 5 that, in order to meet the constraint in (2) with equality, we will need randomization between two deterministic policies (contrary to the fact that (1) has a stationary, deterministic, optimal policy).

3.2 Optimal Policy Structure

In this section, we will only consider the unconstrained problem (1). We formulate the problem as a Markov decision process (MDP). The average reward optimality equation for this MDP is given by (see [20, Chapter 7, Section 4]):

$$h^*(s) = \max_{x \in [0,1]} \left(x \bar{R}_{mobile}(s) + \xi(1-x) \bar{R}_{static}(s) - \lambda^* + \mathbf{E}(h^*(S')) \right) \quad (3)$$

where λ^* is the optimal average reward per slot for the problem (1), $h^*(s)$ is the optimal differential cost at state s (see [20, Chapter 7, Section 4] for thorough interpretation of the differential cost $h^*(s)$), and S' is the (random) next state whose distribution depends on s and the realization of new arrivals. Note that, state transition is independent of the action taken in any slot; hence, the expectation in $\mathbf{E}(h^*(S'))$ is taken only over the randomness in the new arrivals of MUs to the BS in one slot.

Theorem 2: (Optimal policy $\eta_{\xi}^*(\cdot|\cdot)$) If the state s is such that, $\bar{R}_{mobile}(s) - \xi \bar{R}_{static}(s) > 0$, then optimal action is $\eta_{\xi}^*(s) = 1$. If $\bar{R}_{mobile}(s) - \xi \bar{R}_{static}(s) < 0$, then $\eta_{\xi}^*(s) = 0$. If $\bar{R}_{mobile}(s) - \xi \bar{R}_{static}(s) = 0$, then we can choose any action $\eta_{\xi}^*(s)$.

Proof: From (3), we can say that:

$$\eta_{\xi}^* = \arg \max_{x \in [0,1]} \left(x \bar{R}_{mobile}(s) + \xi(1-x) \bar{R}_{static}(s) - \lambda^* + \mathbf{E}(h^*(S')) \right),$$

i.e., η_{ξ}^* should be the maximizer in the average cost optimality equation. Since λ^* , ξ , $\bar{R}_{mobile}(s)$, $\bar{R}_{static}(s)$ and $\mathbf{E}(h^*(S'))$ are independent of x in this optimization problem, we have

$$\eta_{\xi}^* = \arg \max_{x \in [0,1]} x \left(\bar{R}_{mobile}(s) - \xi \bar{R}_{static}(s) \right)$$

This proves the theorem. \square

Remark: The binary nature of the optimal policy in Theorem 2 makes it very easy to use the policy for optimal bandwidth allocation in a practical cellular network.

Comments on Fairness: Note that, each static user will asymptotically receive positive throughput, since with positive probability a cell will have zero mobile user at a given time slot. On the other hand, a mobile user might get zero throughput in the current cell. In order to ensure a fair bandwidth sharing inside each cell, we describe in Section 6 how to share bandwidth between the two classes for a modified objective function which is motivated by the notion of α -fairness (see [21] for reference). The modified objective function ensures that both classes receive a positive throughput at the same time.

Let us denote the steady-state probability of occurrence of state s by $g(s)$, with $\sum_s g(s) = 1$. Under policy $\eta_{\xi}^*(\cdot)$, the optimal data rate for the mobile users per slot is given by:

$$\bar{R}_{mobile}^*(\xi) := \sum_s g(s) \bar{R}_{mobile}(s) \eta_{\xi}^*(s).$$

Similarly, we define the optimal data rate of static users per slot by

$$\bar{R}_{static}^*(\xi) := \sum_s g(s) \bar{R}_{static}(s) (1 - \eta_{\xi}^*(s)).$$

Lemma 1: $\bar{R}_{mobile}^*(\xi)$ decreases with ξ , and $\bar{R}_{static}^*(\xi)$ increases in ξ .

Proof: See Appendix A. \square

Error in estimating user location: This issue is addressed in Section 8.2 in detail.

4 LEARNING ALGORITHM FOR THE UNCONSTRAINED PROBLEM

In Section 3, we assumed that perfect knowledge of $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ is available to the BS. However, in practice, unknown path-loss factor (since path-loss exponent and location of interfering base stations are unknown to the BS), unknown shadowing variation over space and unknown fading distribution will make it impossible for the base station to compute $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$. Hence, the base station cannot use the simple policy structure given by Theorem 2. However, the base station can get a feedback from the users about how much data the users were able to download between two successive decision instants; this can happen if the base station keeps on sending data packets to the users, and the users measure packet error rate in the received data and send feedback to the base station before a new decision is made. In this section, we propose a sequential bandwidth allocation and learning algorithm, which maintains a running estimate of $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ for each state s , and updates these running estimates as new user feedbacks are gathered, so as to converge asymptotically to a stationary policy solving the unconstrained problem (1).

Assumption 1: The fading gain between any base station (serving or interfering) and a specific location in the cell comes from an ergodic Markov process (across time slots) taking values from a bounded subset of the nonnegative real line, and it is identically distributed across locations in the cell and across various BSs. \square

Note that, this assumption ensures that if we sample $R_{mobile}(s)$ infinitely often, we can essentially average over fading, and obtain a correct estimate of $\bar{R}_{mobile}(s)$, even though the slot duration σ might be smaller than the time required to average over all possible fading states by a mobile user.

Note that, by Theorem 2, we can restrict ourselves to the action space $\{0, 1\}$ instead of $[0, 1]$. With this reduced state space, we present our sequential bandwidth allocation and learning algorithm, which is motivated by the theory of stochastic approximation (see [22]).

4.1 The Learning Algorithm

Some notation: Let $\eta_\tau \in \{0, 1\}$ denote the decision to be taken at decision instant τ . Let $R_{mobile}(s)$ and $R_{static}(s)$ be the (random) realization of the total rates received between decision instant τ and decision instant $\tau + 1$ by the mobile (resp., static) users, provided that $\eta_\tau = 1$ (resp., $\eta_\tau = 0$).

Fix any small number $\epsilon > 0$. Suppose that at the decision instant τ , the Markov chain has reached state s , and let the current estimates of $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ be $R_{mobile}^{(\tau)}(s)$ and $R_{static}^{(\tau)}(s)$, respectively.

Let us define $\nu(s, 1, \tau) := \sum_{t=1}^{\tau} \mathbf{1}\{s_t = s, \eta_t = 1\}$ and $\nu(s, 0, \tau) := \sum_{t=1}^{\tau} \mathbf{1}\{s_t = s, \eta_t = 0\}$.

Let $\{a(t)\}_{t \geq 1}$ be a decreasing sequence of positive numbers with $\sum_{t=1}^{\infty} a(t) = \infty$ and $\sum_{t=1}^{\infty} a^2(t) < \infty$.

Algorithm 1: Start with arbitrary $R_{mobile}^{(0)}(s)$ and $R_{static}^{(0)}(s)$.

(Decision on bandwidth sharing:) At decision instant τ , with probabilities $\frac{\epsilon}{2}$ each, allocate the entire bandwidth to the static users (i.e., take $\eta_\tau = 0$) or to the mobile users (i.e., take $\eta_\tau = 1$). Else (with probability $(1 - \epsilon)$), allocate the entire bandwidth to mobile users (i.e., $\eta_\tau = 1$) if $R_{mobile}^{(\tau)}(s) - \xi R_{static}^{(\tau)}(s) > 0$, allocate the entire bandwidth to static users (i.e., $\eta_\tau = 0$) if $R_{mobile}^{(\tau)}(s) - \xi R_{static}^{(\tau)}(s) < 0$, and allocate the entire bandwidth arbitrarily either to SUs or to MUs if $R_{mobile}^{(\tau)}(s) - \xi R_{static}^{(\tau)}(s) = 0$.

(Updating/learning the estimates:) Just before the $(\tau + 1)$ -st decision instant, for each possible state s , make the following update:

$$\begin{aligned} R_{mobile}^{(\tau+1)}(s) &= R_{mobile}^{(\tau)}(s) + a(\nu(s, 1, \tau)) \mathbf{1}\{s(\tau) = s, \eta_\tau = 1\} \\ &\quad \times \left(R_{mobile}(s) - R_{mobile}^{(\tau)}(s) \right) \\ R_{static}^{(\tau+1)}(s) &= R_{static}^{(\tau)}(s) + a(\nu(s, 0, \tau)) \mathbf{1}\{s(\tau) = s, \eta_\tau = 0\} \\ &\quad \times \left(R_{static}(s) - R_{static}^{(\tau)}(s) \right) \end{aligned}$$

\square

4.2 Optimality of the Learning Algorithm

Let us denote the average expected reward per slot under Algorithm 1 by $\lambda_\epsilon^*(\xi)$.

Theorem 3: Under Assumption 1 and Algorithm 1, for each state s , we have $\lim_{\tau \rightarrow \infty} R_{mobile}^{(\tau)}(s) = \bar{R}_{mobile}(s)$ and $\lim_{\tau \rightarrow \infty} R_{static}^{(\tau)}(s) = \bar{R}_{static}(s)$ almost surely. Consequently, $\lim_{\epsilon \downarrow 0} \lambda_\epsilon^*(\xi) = \lambda^*(\xi)$ (note that, ϵ cannot be taken to be equal to 0).

Proof: See Appendix A. \square

4.3 Remarks

- Theorem 3 tells us that in a practical cellular network where the shadowing realizations at all locations and the location of interfering base stations are not known, one can still learn the asymptotically optimal bandwidth sharing policy by learning only $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$.
- At any state s , we randomize our decision with probabilities ϵ and $(1 - \epsilon)$ for the following reason. A sufficient condition for the convergence of $R_{mobile}^{(\tau)}(s)$ to $\bar{R}_{mobile}(s)$ and convergence of $R_{static}^{(\tau)}(s)$ to $\bar{R}_{static}(s)$ is $\liminf_{\tau \rightarrow \infty} \frac{\nu(s, 1, \tau)}{\tau} > 0$ and $\liminf_{\tau \rightarrow \infty} \frac{\nu(s, 0, \tau)}{\tau} > 0$ almost surely for each s ; i.e., all state-action pairs should occur comparatively often. We ensure this by the proposed randomized decision making and using the fact that the states come from an ergodic discrete-time finite state Markov chain. Very small or very large value of ϵ might lead to possibly sample-path dependent slow convergence rate.
- It is easy to see that:

$$|\lambda_\epsilon^*(\xi) - \lambda^*(\xi)| \leq \frac{\epsilon}{2} \sum_s g(s) \mathbf{E} |R_{mobile}(s) - R_{static}(s)|.$$

Hence, by choosing ϵ small, we can achieve a mean reward per slot which is arbitrarily close to the optimal value, but the convergence rate might be slow depending on the initial values of the iterates and the realization of the sample path.

- The above problem of yielding an average reward slightly different than $\lambda^*(\xi)$ can be solved in the following way. At the decision instant τ , instead of using the randomization with probability ϵ (as defined in Algorithm 1), one could randomize for

state s with a probability $\frac{\epsilon}{\nu(s, \tau)}$ where $\nu(s, \tau)$ is the number of occurrence of state s up to time τ . Since the Markov chain is finite state, positive recurrent, irreducible and independent of the actions taken by the base station, and since $\sum_{k=1}^{\infty} \frac{\epsilon}{k} = \infty$, by the second Borel-Cantelli lemma we can say that

$$\mathbf{P}(\lim_{\tau \rightarrow \infty} \nu(s, 1, \tau) = \infty) = 1;$$

this is sufficient to prove Theorem 3. However, we did not use this randomization probability because it will not ensure the conditions $\liminf_{\tau \rightarrow \infty} \frac{\nu(s, 1, \tau)}{\tau} > 0$ and $\liminf_{\tau \rightarrow \infty} \frac{\nu(s, 0, \tau)}{\tau} > 0$ almost surely for each s , which is necessary for the convergence proof of the multi-timescale learning algorithm (Algorithm 2 in Section 5.3) which is inspired by Algorithm 1.

- A special choice would be $a(t) = \frac{1}{t}$, which will lead to sample averaging of the iterates (of course, with the imperfection created by randomized sampling). But we use the general step size $a(t)$ here because it will help in developing multi-timescale learning algorithm for a constrained problem explained in Section 5.
- The rate of convergence is dependent on sample path (i.e., realization of arrival process and the fading process at various locations), and also on the size of the state space. However, convergence is guaranteed by Theorem 3 so long as the state space is finite.
- Speed of convergence will also depend on the choice of $a(t)$; however, choosing a suitable step size sequence is beyond the scope of this paper and we propose to leave it for future research work in this domain.

5 LEARNING ALGORITHM FOR THE CONSTRAINED PROBLEM

In Section 4, we had provided a learning algorithm that solves problem (1) for a given ξ . However, let us recall from Theorem 1 that, in order to solve the constrained problem (2), we need to choose an appropriate ξ^* . Since the transition structure of the MDP in Section 3 might not be known apriori (as discussed in Section 4), in this section we develop a sequential decision and learning algorithm for dynamic bandwidth sharing between the two classes of static and mobile users; this algorithm maintains an estimate of ξ^* and updates this estimate each time user is observed before a new MU enters the cell. We prove asymptotic convergence of the policy to the set of optimal policies.

5.1 Need for Randomization

Note that, while an optimal Lagrange multiplier ξ^* may exist for a feasible constraint R_0 , the optimal policy $\eta_{\xi^*}^*(\cdot|\cdot)$ solving the constrained problem (2) may not be a

deterministic policy. This can be explained in the following way. By Lemma 1, the optimal per-slot sum data rate for static users $\bar{R}_{static}^*(\xi)$ increases with ξ . However, since there are finite number of states and only two actions $\{0, 1\}$, there are finite number of deterministic policies in the class specified by Theorem 2. The mapping from state space to action space can only change a finite number of times as we increase ξ from 0 to ∞ . Hence, the plot of the optimal time-average sum rate of static users under policy $\eta_{\xi}^*(\cdot)$ (i.e., $\bar{R}_{static}^*(\xi)$), as a function of ξ , would look like an increasing staircase function where the discontinuities correspond to the values of ξ where, by increasing ξ^- to ξ^+ , the policy changes because the optimal action for exactly one state changes from 1 to 0. Let the set of ξ values where this plot is discontinuous, be denoted by S . Also, let \mathcal{D} denote the set of values of mean data rate per slot for static users, which can be achieved only via $\eta_{\xi}^*(\cdot)$ by varying ξ from 0 to ∞ .

In light of the above discussion, it is clear that a way to meet the constraint in (2) with equality (if $R_0 \notin \mathcal{D}$) is to randomize between the two policies $\eta_{\xi^*+}^*(\cdot)$ and $\eta_{\xi^*-}^*(\cdot)$ at each decision instant, with probabilities $1 - p$ and p respectively; these two deterministic policies differ in the action for exactly one state (if $R_0 \notin \mathcal{D}$).

5.2 A special randomization technique

In Algorithm 2 presented next, we implement this randomization in a slightly unconventional way in order to tackle certain technical issues. Let us recall the policy $\eta_{\xi}^*(\cdot)$ from Theorem 2. We choose a very small number $\delta > 0$ (choice of δ is explained in Algorithm 2 later in Section 5.3), and define a probability density function $f_p(\cdot)$ (parametrized by a probability p) as follows:

$$f_p(y) = \frac{p}{\delta} \text{ if } y \in [-\delta, 0], f_p(y) = \frac{1-p}{\delta} \text{ if } y \in (0, \delta], \text{ and } f_p(y) = 0 \text{ for all other values of } y.$$

For any given ξ , in each slot τ one can sample a random variable $\Delta_{\tau} \sim f_p$ ($\{\Delta_{\tau}\}_{\tau \geq 1}$ i.i.d. across τ) and use the policy $\eta_{\xi+\Delta_{\tau}}^*(\cdot)$ (i.e., take action $\eta_{\xi+\Delta_{\tau}}^*(s(\tau))$ in slot τ). If $\xi = \xi^*$ and R_0 does not belong to \mathcal{D} , then this scheme will correspond to randomizing between $\eta_{\xi^*+}^*(\cdot)$ and $\eta_{\xi^*-}^*(\cdot)$ with probabilities $1-p$ and p in each slot (but this randomization is applicable to all possible values of ξ).

Let the optimal value of p for a given value of ξ be denoted by $p^*(\xi)$; this is the optimal value of p under multiplier ξ so that the corresponding randomized algorithm (described just above using the probability density function $f_p(\cdot)$) meets the constraint with equality (if possible, given the value of ξ , as explained later in this section).

Definition 1: The set $\mathcal{K}(R_0) \subset [0, 1] \times [0, A]$ is defined to be the set of tuples $(p^*(\xi), \xi)$ under which the randomized policy described above meets the constraint in (2) with equality.

Assumption 2: There exists $\xi^* > 0$ and $p^*(\xi^*) \in [0, 1]$ such that the corresponding randomized policy with

these parameters is optimal for the constrained problem (2), while the constraint is satisfied with equality. In other words, the set $\mathcal{K}(R_0)$ is nonempty. \square

Note that, $\mathcal{K}(R_0)$ involves the function $p^*(\xi)$, and $p^*(\xi)$ can be 0 or 1 also, depending on the value of ξ . If ξ is such that $\sum_s g(s) \mathbf{P}(\eta(s) = 0 | \xi, p) \bar{R}_{static}(s) < R_0$ for all $p \in [0, 1]$, then we will have $p^*(\xi) = 0$. If ξ is such that $\sum_s g(s) \mathbf{P}(\eta(s) = 0 | \xi, p) \bar{R}_{static}(s) > R_0$ for all $p \in [0, 1]$, then we will have $p^*(\xi) = 1$. These two events happen if the value of ξ does not fall within a δ -neighbourhood of the element from \mathcal{S} for which the constraint can be met with equality, and R_0 does not belong to \mathcal{D} ; the constraint cannot be met with equality in this case under this ξ . If R_0 does not belong to \mathcal{D} but the value of ξ is within δ -neighbourhood of the value from \mathcal{S} which can achieve this R_0 , then $p^*(\xi)$ can be anything in the interval $[0, 1]$, depending on the value of R_0 , so that the constraint is met with equality (if possible).

It is easy to prove the following:

Lemma 2: $p^*(\xi)$ is Lipschitz continuous in ξ .

Remark: This lemma will be required to prove desired convergence of our learning Algorithm 2. Note that, if we only randomize between policies $\eta_{\xi-\delta}^*(\cdot)$ and $\eta_{\xi+\delta}^*(\cdot)$ with probabilities $p^*(\xi)$ and $1 - p^*(\xi)$ in each slot, then the result in this lemma will not hold. This is the specific reason that we consider this special form of randomization.

Definition 2: Let the sets \mathcal{S} and \mathcal{D} change to \mathcal{S}_ϵ and \mathcal{D}_ϵ when, in each slot τ , we decide $\eta_\tau = 1$ or $\eta_\tau = 0$ with probabilities $\frac{\epsilon}{2}$ each, and use the policy $\eta_\xi^*(\cdot)$ with probability $(1 - \epsilon)$. Similarly, let the analogue of $\mathcal{K}(R_0)$ be $\mathcal{K}_\epsilon(R_0)$, and the analogue of $p^*(\xi)$ be $p_\epsilon^*(\xi)$.

5.3 The Learning Algorithm Based on Two Timescale Stochastic Approximation

Now we present a sequential bandwidth allocation and learning algorithm in order to solve the constrained problem (2). The algorithm maintains running estimates $\{R_{mobile}^{(\tau)}(s), R_{static}^{(\tau)}(s)\}$ for all s , the Lagrange multiplier $\xi^{(\tau)}$, and the randomizing parameter $p^{(\tau)}$; this algorithm is motivated by two-timescale stochastic approximation (see [22]).

Suppose that at the decision instant τ , the Markov chain has reached state s , and let the current iterates be $R_{mobile}^{(\tau)}(s)$, $R_{static}^{(\tau)}(s)$, $\xi^{(\tau)}$ and $p^{(\tau)}$. Let us define \mathcal{R}_τ to be the collection of $\{R_{mobile}^{(\tau)}(s), R_{static}^{(\tau)}(s)\}$ for all s . We define $\eta_\xi^*(\cdot, \cdot)$ to be the same policy as $\eta_\xi^*(\cdot)$ given in Theorem 2, except that $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ in Theorem 2 are replaced by the currents estimates $R_{mobile}^{(\tau)}(s)$ and $R_{static}^{(\tau)}(s)$ in slot τ ; the action taken in slot τ is $\eta_\xi^*(s(\tau), \mathcal{R}_\tau)$.

Let $\eta_\tau \in \{0, 1\}$ denote the decision at decision instant τ . Let $R_{mobile}(s)$ and $R_{static}(s)$ be the (random) realization of the total rates received between decision instant τ and decision instant $\tau + 1$ by the mobile (resp., static) users, provided that $\eta_\tau = 1$ (resp., $\eta_\tau = 0$).

Let us define $\nu(s, 1, \tau) := \sum_{t=1}^\tau \mathbf{1}\{s_t = s, \eta_t = 1\}$ and $\nu(s, 0, \tau) := \sum_{t=1}^\tau \mathbf{1}\{s_t = s, \eta_t = 0\}$.

Let $\{a(t)\}_{t \geq 1}$ and $\{b(t)\}_{t \geq 1}$ be decreasing sequences of positive numbers with $\sum_{t=1}^\infty a(t) = \sum_{t=1}^\infty b(t) = \infty$, $\sum_{t=1}^\infty a^2(t) < \infty$, $\sum_{t=1}^\infty b^2(t) < \infty$ and $\lim_{t \rightarrow \infty} \frac{b(t)}{a(t)} = 0$. More specifically, we choose $a(t) = \frac{1}{t^{n_1}}$ and $b(t) = \frac{1}{t^{n_2}}$, with $\frac{1}{2} < n_1 < n_2 \leq 1$. Let $[x]_0^A$ denote the projection of x on the compact interval $[0, A]$, and let us choose the value of A is chosen so large that $\bar{R}_{static}^*(\xi = A) > R_0$.

Fix any small number $\epsilon > 0$. We choose $\delta > 0$ to be a very small number, smaller than $1/10$ -th of ϵ and $1/10$ -th of the smallest difference between two successive values of ξ from the set \mathcal{S}_ϵ .

Algorithm 2: Start with $R_{mobile}^{(0)}(s)$, $R_{static}^{(0)}(s)$, $p^{(0)}$, $\xi^{(0)}$.

(Decision on bandwidth sharing:) At decision instant τ , with probabilities $\frac{\epsilon}{2}$ each, allocate the entire bandwidth to the static users or to the mobile users. Else, (with probability $(1 - \epsilon)$) sample a random variable Δ_τ (independent across τ) from the distribution $f_{p^{(\tau)}}(\cdot)$ independent of all other random variables, and use the policy $\eta_{\xi^{(\tau)} + \Delta_\tau}^*(\cdot, \cdot)$ (i.e., take an action $\eta_\tau = \eta_{\xi^{(\tau)} + \Delta_\tau}^*(s(\tau), \mathcal{R}_\tau)$). In other words, choose $\eta_\tau = 1$ if $R_{mobile}^{(\tau)}(s(\tau)) - (\xi^{(\tau)} + \Delta_\tau) R_{static}^{(\tau)}(s(\tau)) > 0$, choose $\eta_\tau = 0$ if $R_{mobile}^{(\tau)}(s(\tau)) - (\xi^{(\tau)} + \Delta_\tau) R_{static}^{(\tau)}(s(\tau)) < 0$ and choose η_τ arbitrarily if $R_{mobile}^{(\tau)}(s(\tau)) - (\xi^{(\tau)} + \Delta_\tau) R_{static}^{(\tau)}(s(\tau)) = 0$.

(Updating/learning the estimates:) Just before the $(\tau + 1)$ -st decision instant, for each s , update as follows:

$$\begin{aligned} R_{mobile}^{(\tau+1)}(s) &= R_{mobile}^{(\tau)}(s) + a(\nu(s, 1, \tau)) \mathbf{1}\{s(\tau) = s, \eta_\tau = 1\} \\ &\quad \times \left(R_{mobile}(s) - R_{mobile}^{(\tau)}(s) \right) \\ R_{static}^{(\tau+1)}(s) &= R_{static}^{(\tau)}(s) + a(\nu(s, 0, \tau)) \mathbf{1}\{s(\tau) = s, \eta_\tau = 0\} \\ &\quad \times \left(R_{static}(s) - R_{static}^{(\tau)}(s) \right) \\ p^{(\tau+1)} &= \left[p^{(\tau)} + a(\tau) \left(\sum_s \mathbf{1}\{s(\tau) = s, \eta_\tau = 0\} \right. \right. \\ &\quad \left. \left. \times R_{static}(s) - R_0 \right) \right]_0^1 \\ \xi^{(\tau+1)} &= \left[\xi^{(\tau)} + b(\tau) (R_0 - \right. \\ &\quad \left. \sum_s \mathbf{1}\{s(\tau) = s, \eta_\tau = 0\} R_{static}(s)) \right]_0^A \end{aligned}$$

\square

5.4 Optimality of the Learning Algorithm for the Constrained Problem

Let us denote the nonstationary, randomized policy induced by Algorithm 2 by $\eta^{(\epsilon)}(\cdot | \cdot, \cdot, \cdot, \cdot)$; the quantity $\eta^{(\epsilon)}(\cdot | s, \mathcal{R}_\tau, \xi^{(\tau)}, p^{(\tau)})$ denotes the probability distribution on the set of actions conditioned on the current state and the current values of the iterates.

Theorem 4: Under Assumption 1, Assumption 2 and Algorithm 2, we have $\lim_{\tau \rightarrow \infty} R_{mobile}^{(\tau)}(s) = \bar{R}_{mobile}(s)$ and $\lim_{\tau \rightarrow \infty} R_{static}^{(\tau)}(s) = \bar{R}_{static}(s)$ for all s almost surely.

Also, for any $\epsilon > 0$, $(p^{(\tau)}, \xi^{(\tau)}) \rightarrow \mathcal{K}_\epsilon(R_0)$ almost surely as $\tau \rightarrow \infty$.

Proof: See Appendix A. \square

Let us denote

$$\bar{R}_{static}^{rand,\epsilon} := \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta^{(\epsilon)}(\cdot|\cdot, \cdot, \cdot, \cdot)} \left((1 - \eta_\tau) \bar{R}_{static}(s(\tau)) \right)$$

and

$$\bar{R}_{mobile}^{rand,\epsilon} := \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta^{(\epsilon)}(\cdot|\cdot, \cdot, \cdot, \cdot)} \left(\eta_\tau \bar{R}_{mobile}(s(\tau)) \right)$$

Corollary 1: $\lim_{\epsilon \downarrow 0} \bar{R}_{mobile}^{rand,\epsilon}$ and $\lim_{\epsilon \downarrow 0} \bar{R}_{static}^{rand,\epsilon}$ exist, and these limit values are equal to the optimal value of the objective in the constrained problem (2) and R_0 , respectively.

Proof: See Appendix A. \square

Remark: Corollary 1 implies that Algorithm 2 approximately solves the constrained problem (2) for arbitrarily small $\epsilon > 0$. This result, which is derived from Theorem 4, allows us to optimally assign bandwidth between the static and mobile user classes even when the transition probability structure of the MDP is not known apriori.

5.5 Remarks on Theorem 4:

- *Two timescales:* The update scheme is based on two timescale stochastic approximation (see [22, Chapter 6]). Note that, $\lim_{t \rightarrow \infty} \frac{b(t)}{a(t)} = 0$; ξ is adapted in the *slower* timescale, and R_{mobile} , R_{static} and p are updated in the *faster* timescale). The dynamics behaves as if the slower timescale update equation views the faster timescale iterates as quasi-static, while a faster timescale update equation views the slower timescale update equations as almost equilibrated; as if ξ is being varied in a slow outer loop, while the other iterates are being varied in an inner loop.
- *Structure of the iteration:* Note that, the value of ξ is increased whenever the sum data downloaded by static users between two successive decision instants is less than the target R_0 , so that more emphasis is given to the static user rate in the objective function. Under the same situation, the value of p is reduced for the same reason. The goal is to converge to a randomized policy $\eta^{(\epsilon)}(\cdot|\cdot, \cdot, \cdot, \cdot)$ so that the corresponding randomized policy satisfies the constraint in (2) with equality.
- Algorithm 2 induces a nonstationary policy. But, by Theorem 4 and Corollary 1, the sequence of policies generated by Algorithm 2 converges close to the set of optimal stationary, randomized policies for the constrained problem (2).

6 FAIR BANDWIDTH SHARING BETWEEN STATIC AND MOBILE USER CLASSES

In previous sections, the proposed dynamic bandwidth sharing schemes do not guarantee nonzero throughput to each user all the time. While such schemes are suitable for elastic traffic applications, they are not at all suitable for streaming applications such as online video watching or voice call. In fact, opportunistic bandwidth sharing depending on user location as described before will result in unfair sharing of bandwidth. In order to incorporate fairness constraint into the opportunistic bandwidth sharing problem, we modify the objective function presented in Section 3.1.

Let us denote $\bar{R}_{static}^\alpha(s) := \mathbf{E} R_{static}^\alpha(s)$ and $\bar{R}_{mobile}^\alpha(s) := \mathbf{E} R_{mobile}^\alpha(s)$ where α is a real number.

In this section, we consider the following unconstrained problem:

$$\sup_{\eta(\cdot|\cdot)} \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta(\cdot|\cdot)} \left(\eta_\tau^\alpha \bar{R}_{mobile}^\alpha(s(\tau)) + \xi(1 - \eta_\tau)^\alpha \bar{R}_{static}^\alpha(s(\tau)) \right) \quad (4)$$

and also the associated constrained problem as follows:

$$\begin{aligned} & \sup_{\eta(\cdot|\cdot)} \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta(\cdot|\cdot)} \left(\eta_\tau^\alpha \bar{R}_{mobile}^\alpha(s(\tau)) \right) \\ s.t., \quad & \liminf_{N \rightarrow \infty} \frac{1}{N} \sum_{\tau=1}^N \mathbf{E}_{\eta(\cdot|\cdot)} \left((1 - \eta_\tau)^\alpha \bar{R}_{static}^\alpha(s(\tau)) \right) \geq R_0 \end{aligned} \quad (5)$$

This objective function is motivated by the notion of α -fairness (see [21]). The intuition is that the degree of fairness in resource allocation between mobile and static user classes can be controlled by appropriately tuning α in (4) and (5).

Let us recall the proof of Theorem 2; Theorem 2 provides optimal allocation for the special case $\alpha = 1$. In general, the function $x^\alpha \bar{R}_{mobile}^\alpha(s) + \xi(1 - x)^\alpha \bar{R}_{static}^\alpha(s)$ is strictly convex in x for $\alpha > 1$ or for $\alpha < 0$, strictly concave in x for $\alpha \in (0, 1)$, independent of x for $\alpha = 0$, and linear in x for $\alpha = 1$. Hence, for $\alpha > 1$ or for $\alpha < 0$, the optimization $\max_{x \in [0, 1]} x^\alpha \bar{R}_{mobile}^\alpha(s) + \xi(1 - x)^\alpha \bar{R}_{static}^\alpha(s)$ will always have $x^* = \eta^*(s) \in \{0, 1\}$. On the other hand, for $\alpha \in (0, 1)$, the optimal value of x lies in $(0, 1)$. Hence, in this section, we focus only on $\alpha \in (0, 1)$. Within this interval, α close to 0 yields more egalitarian solution, whereas α close to 1 provides more opportunistic bandwidth sharing.

Note that, α in this current paper is slightly different

from α in [21]).³

6.1 Policy structure under perfect knowledge of $\bar{R}_{static}^\alpha(s)$ and $\bar{R}_{mobile}^\alpha(s)$

Let us assume the availability of perfect knowledge at the base station (as assumed in Section 3), but let us consider the objective function (4). Similar to Theorem 2, the optimal action at state s is given by

$$\eta_\xi^*(s) = \arg \max_{x \in [0,1]} \left(x^\alpha \bar{R}_{mobile}^\alpha(s) + \xi(1-x)^\alpha \bar{R}_{static}^\alpha(s) \right).$$

Differentiation w.r.t. x and setting the derivative equal to 0, we obtain:

$$\eta_\xi^*(s) = \frac{(\xi \bar{R}_{static}^\alpha(s))^{\frac{1}{\alpha-1}}}{(\xi \bar{R}_{static}^\alpha(s))^{\frac{1}{\alpha-1}} + (\bar{R}_{mobile}^\alpha(s))^{\frac{1}{\alpha-1}}} \quad (6)$$

The optimal policy is to allocate $\eta_\xi^*(s)$ fraction of bandwidth to mobile users and $1-\eta_\xi^*(s)$ fraction of bandwidth to static users whenever the system reaches state s . Let this optimal policy be denoted by $\eta_\xi^*(\cdot)$.

The following lemma is easy to prove:

Lemma 3: $\eta_\xi^*(s)$ in (6) is strictly decreasing and Lipschitz continuous in ξ for all s and for all $\xi > 0$.

6.2 Learning Algorithm for the Unconstrained Problem (4)

Let us now consider imperfect knowledge at the base station as assumed in Section 4, but with the modified objective function (4) with $\alpha \in (0, 1)$. We seek to propose learning algorithms as done in Algorithm 1.

Note that, since a strictly positive fraction of bandwidth is always allocated to static and mobile users at any given time, samples of $R_{mobile}^\alpha(s)$ and $R_{static}^\alpha(s)$ are always available whenever the system reaches state s . As a result of this and the positive recurrence of the Markov chain associated with state evolution, the estimates of $\bar{R}_{mobile}^\alpha(s)$ and $\bar{R}_{static}^\alpha(s)$ will be updated infinitely often for each state s , and there is no need to do the randomization with probability ϵ as described in Section 4.

Let $\eta_\tau \in [0, 1]$ denote the decision at decision instant τ . Let $R_{mobile}^\alpha(s)$ and $R_{static}^\alpha(s)$ denote the samples (of the α -th moment of the corresponding rates) obtained between decision instant τ and decision instant $\tau + 1$.

3. In [21], if the resource (e.g., rate) allocated to user i is r_i , then α -fair utility function is given by $\sum_i \frac{r_i^{1-\alpha}-1}{1-\alpha}$. For $\alpha > 1$, it requires minimization of $\sum_i r_i^{1-\alpha}$, and for $\alpha < 1$, it requires maximization of $\sum_i r_i^{1-\alpha}$. For $\alpha = 1$, this reduces to maximization of $\sum_i \log(r_i)$, which is called proportional fair allocation; we exclude this case because this, in our current problem, will result in bandwidth sharing independent of the value of s . In our current paper, we have used α in a slightly different sense than [21], though the broad concept of fairness is the same as [21]; the only difference is that, unlike [21], we consider fairness only between two classes of users, and any bandwidth sharing policy can be employed within a single class.

Suppose that at the decision instant τ , the Markov chain has reached state s , and let the current estimates of $\bar{R}_{mobile}^\alpha(s)$ and $\bar{R}_{static}^\alpha(s)$ be $R_{mobile,\alpha}^{(\tau)}(s)$ and $R_{static,\alpha}^{(\tau)}(s)$, respectively.

Let $\nu(s, \tau) := \sum_{t=1}^\tau \mathbf{1}\{s_t = s\}$. Let $\{a(t)\}_{t \geq 1}$ be a decreasing sequence of positive numbers with $\sum_{t=1}^\infty a(t) = \infty$ and $\sum_{t=1}^\infty a^2(t) < \infty$.

We propose the following algorithm to learn the optimal policy for problem (4).

Algorithm 3: Start with any arbitrary initial estimates $R_{mobile,\alpha}^{(0)}(s) > 0$ and $R_{static,\alpha}^{(0)}(s) > 0$.

At decision instant τ , if the system is at state s , allocate the following fraction of bandwidth to the mobile users:

$$\eta_\tau = \frac{(\xi R_{static,\alpha}^{(\tau)}(s))^{\frac{1}{\alpha-1}}}{(\xi R_{static,\alpha}^{(\tau)}(s))^{\frac{1}{\alpha-1}} + (R_{mobile,\alpha}^{(\tau)}(s))^{\frac{1}{\alpha-1}}} \quad (7)$$

Just before the $(\tau + 1)$ -st decision instant, for each possible state s , make the following update:

$$\begin{aligned} R_{mobile,\alpha}^{(\tau+1)}(s) &= R_{mobile,\alpha}^{(\tau)}(s) + a(\nu(s, \tau)) \mathbf{1}\{s(\tau) = s\} \\ &\quad \times \left(R_{mobile}^\alpha(s) - R_{mobile,\alpha}^{(\tau)}(s) \right) \\ R_{static,\alpha}^{(\tau+1)}(s) &= R_{static,\alpha}^{(\tau)}(s) + a(\nu(s, \tau)) \mathbf{1}\{s(\tau) = s\} \\ &\quad \times \left(R_{static}^\alpha(s) - R_{static,\alpha}^{(\tau)}(s) \right) \end{aligned}$$

□

Theorem 5: Under Assumption 1 and Algorithm 3, for each state s , we have $\lim_{\tau \rightarrow \infty} R_{mobile,\alpha}^{(\tau)}(s) = \bar{R}_{mobile}^\alpha(s)$ and $\lim_{\tau \rightarrow \infty} R_{static,\alpha}^{(\tau)}(s) = \bar{R}_{static}^\alpha(s)$ almost surely.

Proof: The proof is similar to the proof of Theorem 3. □

6.3 Learning Algorithm for the constrained problem (5)

In this subsection, we seek to propose learning algorithms for the constrained problem (5), in a way similar to Section 5.3. Let us define

$$\bar{R}_{mobile,\alpha}^*(\xi) := \sum_s g(s) \bar{R}_{mobile}^\alpha(s) (\eta_\xi^*(s))^\alpha$$

and

$$\bar{R}_{static,\alpha}^*(\xi) := \sum_s g(s) \bar{R}_{static}^\alpha(s) (1 - \eta_\xi^*(s))^\alpha,$$

where $\eta_\xi^*(s)$ is defined in (6). By Lemma 3, $\bar{R}_{static,\alpha}^*(\xi)$ is strictly increasing and continuous in ξ . Hence, if the constraint in (5) is feasible, then there exists one $\xi^* > 0$ such that the constraint is met with equality under the optimal policy given in Section 6.1 with $\xi = \xi^*$, i.e., $\bar{R}_{static,\alpha}^*(\xi^*) = R_0$ under $\eta_{\xi^*}^*(\cdot)$.

Now we propose a sequential bandwidth allocation and learning algorithm (based on *single timescale* stochastic approximation) that will solve problem (5).

Suppose that at the decision instant τ , the Markov chain has reached state s , and let the current estimates of $\bar{R}_{mobile}^\alpha(s)$, $\bar{R}_{static}^\alpha(s)$ and ξ^* be $R_{mobile,\alpha}^{(\tau)}(s)$, $R_{static,\alpha}^{(\tau)}(s)$ and $\xi^{(\tau)}$, respectively. Let $\nu(s, \tau) := \sum_{t=1}^{\tau} \mathbf{1}\{s_t = s\}$.

Let $\eta_\tau \in [0, 1]$ denote the decision at decision instant τ . Let $R_{mobile}^\alpha(s)$ and $R_{static}^\alpha(s)$ denote the samples obtained between decision instant τ and decision instant $\tau + 1$.

Let $\{a(t)\}_{t \geq 1}$ be a decreasing sequence of positive numbers with $\sum_{t=1}^{\infty} a(t) = \infty$ and $\sum_{t=1}^{\infty} a^2(t) < \infty$. The numbers $B > 0$ and $A > B$ are such that $\xi^* \in (B, A)$.

Algorithm 4: Start with any arbitrary initial estimates $R_{mobile,\alpha}^{(0)}(s) > 0$, $R_{static,\alpha}^{(0)}(s) > 0$ and $\xi^{(0)}$.

At decision instant τ , if the system is at state s , allocate the following fraction of bandwidth to the mobile users:

$$\eta_\tau = \frac{(\xi^{(\tau)} R_{static,\alpha}^{(\tau)}(s))^{\frac{1}{\alpha-1}}}{(\xi^{(\tau)} R_{static,\alpha}^{(\tau)}(s))^{\frac{1}{\alpha-1}} + (R_{mobile,\alpha}^{(\tau)}(s))^{\frac{1}{\alpha-1}}} \quad (8)$$

Just before the $(\tau + 1)$ -st decision instant, for each possible state s , make the following update:

$$\begin{aligned} R_{mobile,\alpha}^{(\tau+1)}(s) &= R_{mobile,\alpha}^{(\tau)}(s) + a(\nu(s, \tau)) \mathbf{1}\{s(\tau) = s\} \\ &\quad \times \left(R_{mobile}^\alpha(s) - R_{mobile,\alpha}^{(\tau)}(s) \right) \\ R_{static,\alpha}^{(\tau+1)}(s) &= R_{static,\alpha}^{(\tau)}(s) + a(\nu(s, \tau)) \mathbf{1}\{s(\tau) = s\} \\ &\quad \times \left(R_{static}^\alpha(s) - R_{static,\alpha}^{(\tau)}(s) \right) \\ \xi^{(\tau+1)} &= \left[\xi^{(\tau)} + a(\tau) (R_0 - \sum_s \mathbf{1}\{s(\tau) = s\} R_{static}^\alpha(s)) \right]_B^A \end{aligned}$$

□

Theorem 6: Under Assumption 1 and Algorithm 4, we have $\lim_{\tau \rightarrow \infty} R_{mobile,\alpha}^{(\tau)}(s) = \bar{R}_{mobile}^\alpha(s)$ and $\lim_{\tau \rightarrow \infty} R_{static,\alpha}^{(\tau)}(s) = \bar{R}_{static}^\alpha(s)$ for all s , and $\xi^{(\tau)} \rightarrow \xi^*$ (if there exists $\xi^* \geq 0$ such that $\bar{R}_{static,\alpha}^*(\xi^*) = R_0$) almost surely.

Proof: The proof is similar to the proof of Theorem 4. □

Remark: If $\xi^{(\tau)} = 0$ for some τ , the entire bandwidth is allocated to the class of mobile users at that decision instant. To avoid this, we always maintain $\xi^{(\tau)} \geq B > 0$.

7 PERFORMANCE IMPROVEMENT THROUGH OPPORTUNISTIC BANDWIDTH ALLOCATION: A NUMERICAL STUDY

In this section, we numerically explore the improvement in performance for static and mobile users via opportunistic bandwidth allocation.

7.1 Asymptotic performance Improvement for various combinations of α and θ

We consider the following simulation environment:

- The base stations are located on the corners of a regular grid; the set of locations of base stations

is given by $\{(1000i, 1000j) : -10 \leq i \leq 10, -10 \leq j \leq 10\}$, where the unit of distance in the xy plane is meter. Hence, the smallest distance between two base stations is 1000 m. We consider Voronoi cells under this realization of the base stations. The base station whose cell is under consideration is located at the origin, and its cell is a 1000 m \times 1000 m square with the origin at its center.

- Path loss at a distance r is $r^{-\beta}$ with the path-loss exponent $\beta = 4$. There is no shadowing and fading in the wireless propagation environment. However, later we will also demonstrate the convergence rate of Algorithm 2 in presence of shadowing and fading.
- All base stations are transmitting at the same power levels. Since we do not assume any thermal noise at the receiving nodes, and since the signal-to-interference-ratio (SIR) remains unchanged if the transmit power of each base station is multiplied by the same factor, we can safely assume that the transmit power of each base station is 1 unit.
- There are 500 static users inside the cell containing the origin, and their locations are chosen independently with uniform distribution from the cell.
- Two roads along the $x = 25$ line and $y = 50$ line intersect the cell under consideration. Each of these 1000 m long line segments are divided into 10 segments of length 100 m each (it can be segmented further to the shadowing decorrelation distance level). We assume that mobile users enter the cell along these lines at a velocity 50 m/sec, and the slot duration is 2 sec so that each MU covers one 100 m distance segment in one slot, i.e., each MU traverses the cell in 20 seconds (10 slots).
- The number of arrivals (of MUs) to the cell in each slot is 50 times a Bernoulli distributed random variable with mean $\frac{1}{\theta}$; this is a batch arrival process.
- We assume that, if the entire bandwidth (assumed to be 1 unit) is allocated to a single MU, then the total amount of data this MU can download over a 100 m long line segment is given by $\log_2(1 + SIR_{centre})$ where SIR_{centre} is the SIR value at the center of the segment. For example, the total data rate assigned to a MU when it is crossing the distance between $(-500, 50)$ and $(-400, 50)$ in a single slot is given by $\log_2(1 + SIR_{(-450, 50)})$ (provided that the entire bandwidth is allocated to this user). On the other hand, the amount of data downloaded by a static user when the entire bandwidth is allocated to this user is given by $\log_2(1 + SIR)$ when the SIR corresponds to the location of the static user.

Since the stochastic approximation algorithms presented in this paper asymptotically converge to the optimal value, we first consider perfect knowledge scenario where the MDP transition and cost structures are known to the decision maker. For opportunistic (i.e., location-dependent) bandwidth allocation, we assume that all

α	θ	ξ	$\bar{R}_{mobile,\alpha}^{equal}$	$\bar{R}_{static,\alpha}^{equal}$	$\bar{R}_{mobile,\alpha}^*$	$\bar{R}_{static,\alpha}^*$
0.1	0.1	2.3	0.4867	0.4655	0.5263	0.4654
0.1	0.9	1	0.7675	0.6161	0.7680	0.6212
0.9	0.1	1.5	0.0237	0.2539	0.0387	0.2780
0.9	0.9	1.9	0.0935	0.0275	0.0941	0.0289

TABLE 1

Comparison of equal bandwidth sharing among all users against opportunistic (dynamic) bandwidth sharing between the static and mobile user classes, for various combinations of α and θ (and correspondingly appropriate choice of ξ). Under dynamic bandwidth sharing (columns 6 and 7 in the table), it is assumed that all users in the same user class (static or mobile) share equally the bandwidth available to that class at any moment. The notation has been defined in the text.

users in the same class (i.e., static or mobile) share equal bandwidth among themselves all the time.

We have done extensive simulation over a range of parameter values, for various realizations of the location of static users. In this section, we only provide a few of them to illustrate the performance gains and trade-offs.

We first focus on the problem (5) for $\alpha \in (0, 1)$ for comparison. For each combination of α and θ , we first compute non-opportunistic performance metrics $\bar{R}_{mobile,\alpha}^{equal}$ and $\bar{R}_{static,\alpha}^{equal}$ which are analogous to $\bar{R}_{mobile,\alpha}^*$ and $\bar{R}_{static,\alpha}^*$ defined in Section 6.3 (with ξ dropped from the notation), except that $\bar{R}_{mobile,\alpha}^{equal}$ and $\bar{R}_{static,\alpha}^{equal}$ are calculated assuming equal bandwidth sharing among all static and mobile users at any point of time. Then we chose an appropriate value of ξ (for a given α and θ) so that, under the corresponding optimal policies given in Section 6.1 with this choice of ξ , the constraint in (5) is (approximately) met with equality; clearly, our objective is to solve the constrained problem (5). The quantities $\bar{R}_{mobile,\alpha}^*$ and $\bar{R}_{static,\alpha}^*$ under $\alpha = 1$ become \bar{R}_{mobile}^* and \bar{R}_{static}^* (defined in Section 3.2). Our goal is to see how much improvement is possible (via opportunistic bandwidth sharing) in the time-average sum data rate of mobile users which keeping the same quantity unchanged for static users. The results are summarized in Table 1. Note that, each row in Table 1 corresponds to an independent set of static user locations.

From Table 1, we observe that even 60% improvement is possible in the time-average throughput of mobile users, while keeping the time-average throughput of static users almost unchanged; this clearly shows that it is worth employing the proposed opportunistic bandwidth allocation algorithms in cellular networks. We also observe that the margin of performance improvement decreases as α becomes smaller. This happens because of two reasons: (i) choice of $\alpha \in (0, 1)$ allows more fair allocation at the cost of opportunistic gain, (ii) it is also an artifact of the choice of $\alpha \in (0, 1)$ since the derivative of the concave function x^α is decreasing in x . On the

other hand, performance gain in the data rate for mobile users becomes smaller if θ becomes close to 1. When θ is small, bandwidth is allocated only to the mobile users when they come close to the base station; however, when θ is large, there are a large number of mobile users inside the cell at any time with high probability, and hence equal bandwidth sharing among the mobile users results in significant bandwidth allocation to the mobile users which are either close or away from the base station.

One should also note that, the amount of gain will vary depending on the topology of a cell, location of interfering base stations, static user locations, shadowing realizations in various locations as well as fading process statistics; it is hard to quantify these effects but some intuitive conclusions can be drawn. For example, if static users are very close to the base station, then the performance gain in the throughput of mobile users will be less since opportunistic allocation will assign more bandwidth to static users. The numerical work presented in this section is only an illustration for possible performance gain by location-dependent dynamic bandwidth allocation.

7.2 Convergence of Algorithm 2 for $\alpha = 1$

Here we consider the same network model as in Section 7.1 except that (i) the shadowing between any base station and any static user location or road segment center is assumed to be independent lognormal random variable with standard deviation 8 dB, (ii) the fading gain between the origin and any location inside the cell is exponentially distributed with mean 1 (Rayleigh fading), but the fading in any interfering link is averaged out, (iii) $\theta = 0.2$, (iv) there is only one line $x = 50$ along which the mobile users traverse.

The convergence of Algorithm 2 is examined under this network setting. We first generate the network and compute the time-average expected data rate to static and mobile user classes when all users are allocated equal bandwidth in each slot; the mean data rate per slot for the static user class is then set as the target R_0 in (2), and Algorithm 2 is employed with stepsize sequences $a(t) = \frac{1}{t^{0.6}}$, $b(t) = \frac{1}{t}$ and initial estimates $R_{mobile}^{(0)}(s) = 1$, $R_{static}^{(0)}(s) = 1$ and $\xi^{(0)} = 2$. Under Algorithm 2, the evolution of $\frac{1}{N} \sum_{\tau=1}^N \eta_\tau R_{mobile}(s(\tau))$, $\frac{1}{N} \sum_{\tau=1}^N (1 - \eta_\tau) R_{static}(s(\tau))$ against N and $\xi^{(\tau)}$ against τ for a single sample path are shown in Figure 2. From Figure 2, we can see that all iterates converge asymptotically, and they are close to the respective limiting values within 20000 iterations. In practice, the convergence will be faster because the initial values $R_{mobile}^{(0)}(s)$, $R_{static}^{(0)}(s)$ and $\xi^{(0)}$ will be chosen based on prior experience in previous days, and will be chosen close to the target values. Also, the convergence speed will depend on the network parameters, network topology, wireless propagation model and the step size sequences $\{a(\tau), b(\tau)\}_{\tau \geq 0}$. From Figure 2, we can see that more than

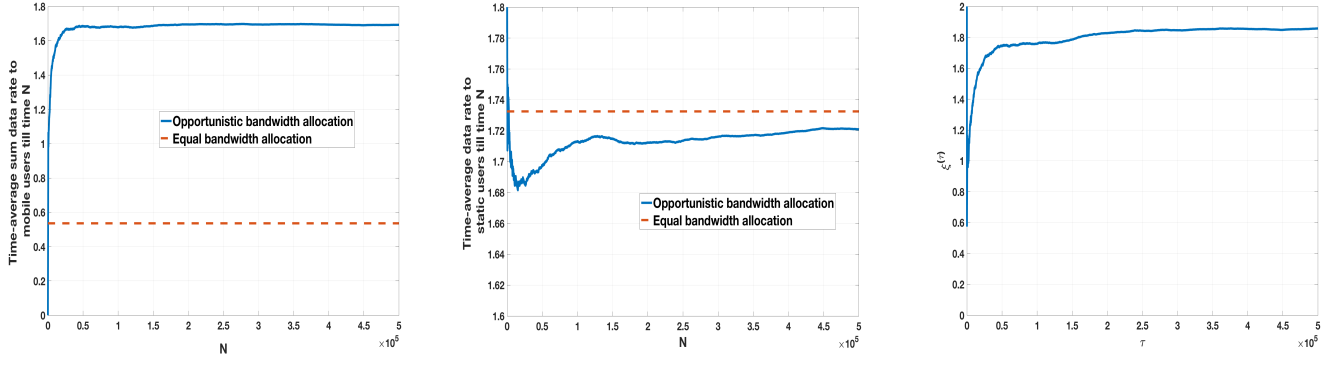


Fig. 2. Convergence of Algorithm 2.

150% improvement is possible in the sum throughput of mobile users while achieving the same sum throughput for static users.

8 ADDITIONAL DISCUSSION

8.1 Connection Between the Cell Level Problem and a Global Problem

Let us consider a heterogeneous network consisting of two tiers of base stations (BSs). The two tiers are modeled by two independent stationary, ergodic processes Φ_{macro} and Φ_{micro} (such as homogeneous Poisson point processes). SUs are assumed to be located on \mathbb{R}^2 according to a stationary, ergodic point process \mathcal{U}_{static} of intensity λ_{SU} . We assume that MUs are moving with constant speed v along a collection of directed routes; these routes are modeled by a stationary, ergodic line process (such as directed homogeneous Poisson line process, see [23, Chapter 8]). Given a realization of the line process, we assume that, at any time t , two successive MUs on any line of \mathcal{L} are separated by an exponentially distributed distance with mean $\frac{1}{\lambda_{MU}}$, i.e., the MUs on any line form a Poisson point process of intensity λ_{MU} at any time t . Hence, the crossing of any point of a line by the MUs form a time homogeneous Poisson process with intensity $\lambda_{MU}v$.

A static user is served by a macro or micro BS, and the association rule can be arbitrary (e.g., a SU can be associated with the BS that sends strongest signal to the SU). Each MU is served by the nearest macro BS. We call the Voronoi cells generated by Φ_{macro} as macro cells. From now on, unless specified, a cell will mean a macro cell.

Note that, the heterogeneous network model is used to illustrate our model in advanced cellular network context (e.g., for LTE). But the analysis presented in this paper will be valid even if the network is homogeneous and each BS is allowed to serve both SUs and MUs.

Let us consider the time-slotted simplification of the above system and the problem addressed in Section 3. The unconstrained optimization problem (1) can be used for performance optimization in a single macro cell. Let us enumerate the macro BSs on the plane

by $\{1, 2, \dots\}$. Since the base stations do not communicate for making the decision on bandwidth allocation, and since each macro BS has different number of line segments intersecting its cell and different number of SUs associated with it, the dynamic bandwidth sharing policy adopted by the network is $\underline{\eta} = \times_{k=1}^{\infty} \eta^{(k)}$ where $\eta^{(k)}$ is the policy used by the k -th macro BS. Let us denote the numerator in (1) for the k -th macro BS, i.e., $\sum_{\tau=1}^N \mathbb{E}_{\eta^{(k)}} \left(\eta_{\tau} \bar{R}_{mobile,k}(s(\tau)) + \xi(1 - \eta_{\tau}) \bar{R}_{static,k}(s(\tau)) \right)$ by $r(k, N)$. Let us consider the following problem:

$$\sup_{\underline{\eta}} \liminf_{M \rightarrow \infty} \liminf_{N \rightarrow \infty} \frac{\sum_{k=1}^M r(k, N)}{NM} \quad (9)$$

Now, since $\underline{\eta} = \times_{k=1}^{\infty} \eta^{(k)}$, the above problem can be rewritten as:

$$\liminf_{M \rightarrow \infty} \frac{1}{M} \sum_{k=1}^M \sup_{\eta^{(k)}} \liminf_{N \rightarrow \infty} \frac{r(k, N)}{N}$$

Let the optimal mean reward per slot for the problem (1) for cell k be λ_k . Now, for (9), $\liminf_{M \rightarrow \infty} \frac{1}{M} \sum_{k=1}^M \lambda_k$ is almost surely equal to the expected optimal time-average reward for the typical macro cell. Hence, by solving the problem (1) for each cell, we can maximize the expected optimal time-average reward for the typical macro cell.

8.2 Addressing the Possibility of Error in Location Estimation for MUs

Let us recall the framework in Section 3. It has to be noted that there can be error in estimating the location of a mobile user, and therefore an error in estimating the residual sojourn time of a mobile user inside a cell is possible. Let us assume that the error in estimation of states at any two different time slots are independent, and that we know the error statistics (i.e., given the observed state \hat{s} , we know the conditional distribution $p(s|\hat{s})$ of the true state s). Since the action in a slot does not affect the state transition, the best possible action one can take in a slot is to choose

$$\begin{aligned}
& \eta_{\xi}^*(\hat{s}) \\
&= \arg \max_{\eta \in [0,1]} \sum_s p(s|\hat{s}) \left(\eta \bar{R}_{mobile}(s) + \xi(1-\eta) \bar{R}_{static}(s) \right) \\
&:= \arg \max_{\eta \in [0,1]} \left(\eta \tilde{r}_{mobile}(\hat{s}) + \xi(1-\eta) \tilde{r}_{static}(\hat{s}) \right).
\end{aligned}$$

The structure of the optimal policy will be similar to Theorem 2. Similarly, it will be optimal to work with the observed state \hat{s} in case learning algorithms are employed.

8.3 Deviation from movement along a line

The analysis in this paper can be trivially extended to the case where the location of mobile users vary according to a positive recurrent discrete-time Markov chain over the cell divided into a finite number of area segments. The state in this case should include the location of all mobile users at a given time. However, in such cases, the location of each user needs to be tracked by the base station in each slot; this is not required if the users move along straight lines with known velocity, since one can easily predict the location of a user at a time once the initial location of that user at a given time instant is known. In this paper, we considered movement of users along a line because it stands for vehicle movements along roads and it is a simple but powerful example.

8.4 Unequal bandwidth sharing within a single class of users

From the analysis presented in Section 3, it is clear that the decision at any given state s depends only on $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$, and not on the specific bandwidth fraction allocated to each user within a user class. The percentage bandwidth allocation among mobile users within the class of mobile users determines $\bar{R}_{mobile}(s)$ which affects the decision η_{ξ}^* .

8.5 Multiple user classes

In case there are multiple user classes with different velocities, analogous to Theorem 2, the optimal policy will be to allocate the entire bandwidth to a single user class at any given time. Algorithm 2 can be extended for maximizing the time average data rate for one class subject to a minimum time-average data rate constraint on each of the other classes; in this case, one Lagrange multiplier needs to be updated for each constraint, and the optimal solution for the constrained problem will involve randomization among multiple policies.

8.6 Channel estimation versus location-based bandwidth sharing

For fast moving users, traditional channel estimation may not be very accurate since the user might travel the

fast fading coherence distance very fast; hence, dynamic bandwidth allocation among users based on instantaneous channel qualities may not be feasible. Also, even if channel measurement is accurate, for a user moving at a velocity 72 kmph, the channel coherence time will be less than 50 ms; hence, gathering channel state information from each user every 50 ms will require huge signaling overhead. Moreover, it may be difficult to estimate the interference at any given location, since the interference at any location depends on path-loss, shadowing and time-varying fast fading gains from all interfering base stations. As an alternative, we propose location-dependent bandwidth sharing. Section 3 deals with the situation where $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ are known; this amounts to assuming that the path-loss and shadowing from serving and interfering base stations are known, and the distribution of fast fading gains from the serving and interfering base stations to each location are also known. *We emphasize that this is an idealistic assumption, and, in practice, the serving base station has to learn $\bar{R}_{mobile}(s)$ and $\bar{R}_{static}(s)$ over time from the download data volume reported by the users; this is discussed in Section 4 and Section 5. Clearly, the learning algorithms do not need any propagation based model. While channel quality measurement, if done accurately, can result in superior user performance, our proposed algorithms for location-based bandwidth algorithms with learning are useful when accurate channel estimates and interference estimates are not available due to high velocity of users.*

Location dependent bandwidth sharing has also been discussed in [15], where a preference is given to the mobile users located close to the base station.

9 CONCLUSION

In this paper, we have proposed and analyzed opportunistic (dynamic) bandwidth sharing depending on user location and mobility, in order to improve the performance of cellular networks. Even though we have solved the basic problem in this paper, there are numerous issues to improve upon: (i) In practice, there can be multiple (possibly uncountable) values of user velocity. Hence, a dynamic bandwidth sharing scheme that allocates bandwidth depending on exact velocity of each user needs to be developed (this might require classification of user velocities into a finite set). (ii) For general motion of users, one reasonable approach would be to divide the cell into various zones (or locations), and assume a Markov evolution of user locations; similar learning techniques as in our paper can be applied in such situation. (iii) Testing and optimizing the proposed and subsequent algorithms in real data-traffic networks will be an important requirement. We propose to pursue these topics in our future research endeavours.

REFERENCES

- [1] J. Andrews, H. Claussen, M. Dohler, S. Rangan, and M. Reed. Femtocells: Past, present and future. *IEEE Journal on Selected Areas in Communications*, 30(3):497–508, 2012.

- [2] V. Pauli, J. Diego Naranjo, and E. Seidel. Heterogeneous LTE networks and inter-cell interference coordination. <http://nomor.de/home/technology/white-papers/lte-hetnet-and-icic>, 2010. Nomor Research White Paper.
- [3] O. Stanze and A. Weber. Heterogeneous networks with lte-advanced technologies. *Bell Labs Technical Journal*, 18(1):41–58, 2013.
- [4] T. Nakamura, S. Nagata, A. Benjebbour, Y. Kishiyama, T. Hai, S. Xiaodong, Y. Ning, and L. Nan. Trends in small cell enhancements in lte advanced. *IEEE Communications Magazine*, 51(2):98–105, 2013.
- [5] H. Ishii, Y. Kishiyama, and H. Takahashi. A novel architecture for lte-b c-plane/u-plane split and phantom cell concept. In *IEEE Globecom Workshops*, pages 624–630, 2012.
- [6] T. Camp, J. Boleng, and V. Davies. A survey of mobility models for ad hoc network research. *Wireless Communication and Mobile Computing (WCMC): Special Issue on Mobile Ad Hoc Networking: Research, Trends and Applications*, 2:483–502, 2002.
- [7] A. Chattopadhyay, B. Blaszczyszyn, and E. Altman. Cell planning for mobility management in heterogeneous cellular networks. *Technical report, available in <http://arxiv.org/abs/1605.07341>*.
- [8] M. Grossglauser and D.N.C. Tse. Mobility increases the capacity of ad hoc wireless networks. *IEEE/ACM Transactions on Networking*, 10(4):477–486, 2002.
- [9] N. Bansal and Z. Liu. Capacity, delay and mobility in wireless ad-hoc networks. In *Twenty-Second Annual Joint Conference of the IEEE Computer and Communications (INFOCOM)*, Vol. 2, pages 1553–1563. IEEE, 2003.
- [10] T. Bonald, S. Borst, N. Hegde, M. Jonckheere, and A. Proutiere. Flow-level performance and capacity of wireless networks with user mobility. *Queueing Systems: Theory and Applications*, 63:131–164, 2009.
- [11] T. Bonald, S.C. Borst, and A. Proutiere. How mobility impacts the flow-level performance of wireless data systems. In *Twenty-third Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM)*, Vol. 3, pages 1872–1881. IEEE, 2004.
- [12] S. Borst, A. Proutiere, and N. Hegde. Capacity of wireless data networks with intra- and inter-cell mobility. In *25th IEEE International Conference on Computer Communications (INFOCOM)*, pages 1–12. IEEE, 2006.
- [13] M.K. Karay. Users mobility effect on the performance of wireless cellular networks serving elastic traffic. *Wireless Networks*, 17(1):247–262, 2011.
- [14] P.V. Orlik and S.S. Rappaport. On the handoff arrival process in cellular communications. *Wireless Networks*, 7:147–157, 2001.
- [15] N. Abbas, T. Bonald, and B. Sayrac. Opportunistic gains of mobility in cellular data networks. In *International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt)*, pages 315–322, 2015.
- [16] H.J. Kushner and P.A. Whiting. Convergence of proportional-fair sharing algorithms under general conditions. *IEEE Transactions on Wireless Communications*, 3(4):1250–1259, 2004.
- [17] R. Margolies, A. Sridharan, V. Aggarwal, R. Jana, N.K. Shankaranarayanan, V.A. Vaishampayan, and G. Zussman. Exploiting mobility in proportional fair cellular scheduling: Measurements and algorithms. *IEEE/ACM Transactions on Networking*, 24(1):355–367, 2016.
- [18] S.H. Ali, V. Krishnamurthy, and V.C.M. Leung. Optimal and approximate mobility-assisted opportunistic scheduling in cellular networks. *IEEE Transactions on Mobile Computing*, 6(6):633–648, 2007.
- [19] Frederick J. Beutler and Keith W. Ross. Optimal policies for controlled markov chains with a constraint. *Journal of Mathematical Analysis and Applications*, 112:236–252, 1985.
- [20] D.P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. I*. Athena Scientific, 2007.
- [21] T. Lan, D. Kao, M. Chiang, and A. Sabharwal. An axiomatic theory of fairness in network resource allocation. In *IEEE Infocom*, pages 1–9, 2010.
- [22] Vivek S. Borkar. *Stochastic approximation: a dynamical systems viewpoint*. Cambridge University Press, 2008.
- [23] S.N. Chiu, D. Stoyan, W.S. Kendall, and J. Mecke. *Stochastic Geometry and its Applications*. Wiley, 2013.
- [24] Walter Rudin. *Principles of Mathematical Analysis, Third Edition*. McGraw-Hill International Editions, 1976.
- [25] A. Chattopadhyay, M. Coupechoux, and A. Kumar. Sequential decision algorithms for measurement-based impromptu deployment of a wireless relay network along a line. *accepted in IEEE/ACM Transactions on Networking, longer version available in <http://arxiv.org/abs/1502.06878>*.
- [26] H.J. Kushner and D.S. Clark. *Stochastic Approximation Methods for Constrained and Unconstrained Systems*. Springer-Verlag, 1978.



networks, cyber-physical systems, machine learning, and networked control.



major international journals and conferences, as well as a two-volume book on *Stochastic Geometry and Wireless Networks* NoW Publishers, jointly with F. Baccelli.



and in particular congestion control, wireless communications and networking games. He is in the editorial board of several scientific journals: JEDC, COMNET, DEDS and WICON. He has been the (co)chairman of the program committee of several international conferences and workshops on game theory, networking games and mobile networks.

Arpan Chattopadhyay obtained his B.E. in Electronics and Telecommunication Engineering from Jadavpur University, Kolkata, India in the year 2008, and M.E. and Ph.D in Telecommunication Engineering from Indian Institute of Science, Bangalore, India in the year 2010 and 2015, respectively. After Ph.D., he did his first postdoc in the group DYOGENE of INRIA, Paris, France. He is currently working in University of Southern California, Los Angeles as a postdoctoral researcher. His research interests include

Bartłomiej Blaszczyszyn received his PhD degree and Habilitation qualification in applied mathematics from University of Wrocław (Poland) in 1995 and 2008, respectively. He is now a Senior Researcher at Inria (France), and a member of the Computer Science Department of Ecole Normale Supérieure in Paris. His professional interests are in applied probability, in particular in stochastic modeling and performance evaluation of communication networks. He coauthored several publications on this subject in

Eitan Altman received the B.Sc. degree in electrical engineering (1984), the B.A. degree in physics (1984) and the Ph.D. degree in electrical engineering (1990), all from the Technion-Israel Institute, Haifa. In (1990) he further received his B.Mus. degree in music composition in Tel-Aviv University. Since 1990, he has been with INRIA (National research institute in informatics and control) in Sophia-Antipolis, France. His current research interests include performance evaluation and control of telecommunication networks

Supplementary Material for “Location Aware Opportunistic Bandwidth Sharing between Static and Mobile Users with Stochastic Learning in Cellular Networks”

APPENDIX A

A.1 Proof of Lemma 1

Let $\xi \geq 0$ and $\kappa > 0$. By the optimality of $\eta_\xi^*(\cdot)$ and $\eta_{\xi+\kappa}^*(\cdot)$, we can write:

$$\bar{R}_{mobile}^*(\xi) + \xi \bar{R}_{static}^*(\xi) \geq \bar{R}_{mobile}^*(\xi + \kappa) + \xi \bar{R}_{static}^*(\xi + \kappa)$$

and

$$\bar{R}_{mobile}^*(\xi + \kappa) + (\xi + \kappa) \bar{R}_{static}^*(\xi + \kappa) \geq \bar{R}_{mobile}^*(\xi) + (\xi + \kappa) \bar{R}_{static}^*(\xi)$$

By adding these two equations, we obtain:

$$\bar{R}_{static}^*(\xi + \kappa) \geq \bar{R}_{static}^*(\xi).$$

Similarly we can prove that $\bar{R}_{mobile}^*(\xi)$ decreases in ξ .

A.2 Proof of Theorem 3

Let us rewrite the update equation in Algorithm 1 as follows:

$$\begin{aligned} R_{mobile}^{(\tau+1)}(s) &= R_{mobile}^{(\tau)}(s) + a(\nu(s, 1, \tau)) \mathbf{1}\{s(\tau) = s\} \mathbf{1}\{\eta_\tau = 1\} \\ &\quad \times \left(\bar{R}_{mobile}(s) - R_{mobile}^{(\tau)}(s) + N^{(\tau+1)}(s, 1) \right) \\ R_{static}^{(\tau+1)}(s) &= R_{static}^{(\tau)}(s) + a(\nu(s, 0, \tau)) \mathbf{1}\{s(\tau) = s\} \mathbf{1}\{\eta_\tau = 0\} \\ &\quad \times \left(\bar{R}_{static}(s) - R_{static}^{(\tau)}(s) + N^{(\tau+1)}(s, 0) \right) \end{aligned}$$

where

$$N^{(\tau+1)}(s, 1) := R_{mobile}(s) - \bar{R}_{mobile}(s),$$

and

$$N^{(\tau+1)}(s, 0) := R_{static}(s) - \bar{R}_{static}(s).$$

This is an asynchronous stochastic approximation iteration as described in [22], with $N(s, 1)$ and $N(s, 0)$ as Martingale difference noise sequences. However, for each s ,

$$\liminf_{\tau \rightarrow \infty} \frac{\nu(s, 1, \tau)}{\tau} \geq \frac{g(s)\epsilon}{2} > 0,$$

where $g(s)$ has been defined in Section 3.2, and

$$\liminf_{\tau \rightarrow \infty} \frac{\nu(s, 0, \tau)}{\tau} \geq \frac{g(s)\epsilon}{2} > 0$$

almost surely.

Since each iterate is updated infinitely often, and since the iterations of various components of the iterates are uncoupled, for each s we can cast the iteration as an ordinary stochastic approximation as defined in [22, Chapter 2].

Now we will check some conditions from [22]. Let us denote $\underline{R} := \{R_{mobile}(s), R_{static}(s)\}_{\forall s}$.

Checking Assumption (A1) of [22, Chapter 2]: Clearly, $\bar{R}_{mobile}(s) - R_{mobile}^{(\tau)}(s)$ is Lipschitz in $R_{mobile}^{(\tau)}(s)$ and $\bar{R}_{static}(s) - R_{static}^{(\tau)}(s)$ is Lipschitz in $R_{static}^{(\tau)}(s)$ for each s ; hence, this assumption is satisfied.

Checking Assumption (A2) of [22, Chapter 2]: This assumption is satisfied by the choice of the step size sequence.

Checking Assumption (A3) of [22, Chapter 2]: It is easy to see that, $\{N^{(\tau+1)}(s, 1), N^{(\tau+1)}(s, 0)\}_{\tau \geq 1}$ for each s is a sequence of Martingale difference noise with zero mean, adapted to the sigma algebra generated by $\{N^{(k)}(s, 1), N^{(k)}(s, 0)\}_{0 \leq k \leq \tau, \forall s}$. Also, the conditional mean of $|N^{(\tau+1)}(s, 1)|^2$ given all the noise values up to time τ is uniformly upper bounded by some constant, since T is Geometrically distributed and fading process is bounded by Assumption 1. Hence, Assumption (A3) of [22, Chapter 2] is satisfied.

Checking Assumption (A5) of [22, Chapter 3]: Note that, $\lim_{c \rightarrow \infty} \frac{\bar{R}_{mobile}(s) - cx(s, 1)}{c} = -x(s, 1)$ is continuous in $x(s, 1)$ for all s . Also, $\frac{\bar{R}_{mobile}(s) - cx(s, 1)}{c}$ is decreasing in c . Hence, by Theorem 7.13 of [24], convergence of $\frac{\bar{R}_{mobile}(s) - cx(s, 1)}{c}$ over compacts is uniform. Also, the collection of ODEs of the form $\dot{x}(s, 1) = \lim_{c \rightarrow \infty} \frac{\bar{R}_{mobile}(s) - cx(s, 1)}{c} = -x(s, 1)$ has a unique globally asymptotically stable equilibrium $x(s, 1) = x(s, 0) = 0$ for all s . Hence, this assumption is satisfied.

Let us consider the following ODE for all s :

$$\begin{aligned} \dot{x}(s, 1) &= \bar{R}_{mobile}(s) - x(s, 1) \\ \dot{x}(s, 0) &= \bar{R}_{static}(s) - x(s, 0) \end{aligned}$$

The above ODE has a unique globally asymptotically stable equilibrium $x(s, 1) = \bar{R}_{mobile}(s)$ and $x(s, 0) = \bar{R}_{static}(s)$. Hence, by [22, Theorem 7, Chapter 3] and [22, Theorem 2, Chapter 2], convergence of Algorithm 1 follows.

Now, it is easy to see that,

$$|\lambda_\epsilon^*(\xi) - \lambda^*(\xi)| \leq \frac{\epsilon}{2} \sum_s g(s) \mathbf{E} |R_{mobile}(s) - R_{static}(s)|$$

The second part of the theorem follows from this. \square

A.3 Proof of Theorem 4

We will prove desired convergence in the two timescales one by one.

A.3.1 Convergence in the faster timescale

Lemma 4: Under Algorithm 2, we have $\lim_{\tau \rightarrow \infty} R_{mobile}^{(\tau)}(s) = \bar{R}_{mobile}(s)$ and $\lim_{\tau \rightarrow \infty} R_{static}^{(\tau)}(s) = \bar{R}_{static}(s)$ for all s almost surely.

Proof: The proof is similar to the proof of Theorem 3. \square

Lemma 5: Under Algorithm 2, we have $\lim_{\tau \rightarrow \infty} |p^{(\tau)} - p^*(\xi^{(\tau)})| = 0$.

Proof: Note that, we can rewrite the p iteration (using Taylor series expansion) as follows:

$$\begin{aligned} p^{(\tau+1)} &= \left[p^{(\tau)} + a(\tau) \left(\sum_s \mathbf{1}\{s(\tau) = s, \eta_\tau = 0\} \right. \right. \\ &\quad \left. \left. \times R_{static}(s) - R_0 \right) \right]_0^1 \\ &= p^{(\tau)} + o(a(\tau)) \\ &\quad + \lim_{\beta \rightarrow \infty} \frac{\left[p^{(\tau)} + \beta \left(\sum_s \mathbf{1}\{s(\tau) = s, \eta_\tau = 0\} R_{static}(s) - R_0 \right) \right]_0^1 - p^{(\tau)}}{\beta} \\ &= p^{(\tau)} + o(a(\tau)) + N^{(\tau)} \\ &\quad + a(\tau) \mathbf{E} \lim_{\beta \rightarrow \infty} \frac{\left[p^{(\tau)} + \beta \left(\sum_s \mathbf{1}\{s(\tau) = s, \eta_\tau = 0\} R_{static}(s) - R_0 \right) \right]_0^1 - p^{(\tau)}}{\beta} \end{aligned}$$

where $N^{(\tau)}$ is a Martingale difference noise sequence, $o(b(\tau))$ is the tail of the Taylor series expansion, and the expectation is under the randomized policy $\eta^{(\epsilon)}(\cdot|\cdot, \cdot, \cdot, \cdot)$.

Now,

$$\begin{aligned} &\mathbf{P}(\eta_\tau = 0 | s(\tau) = s, \mathcal{R}_\tau, \xi^{(\tau)}, p^{(\tau)}) \\ &= \frac{\epsilon}{2} + (1 - \epsilon) \mathbf{P}(R_{mobile}^{(\tau)}(s) - (\xi^{(\tau)} + \Delta_\tau) R_{static}^{(\tau)}(s) \leq 0) \end{aligned}$$

Note that, $\mathbf{P}(\eta_\tau = 0 | s(\tau) = s, \mathcal{R}_\tau, \xi^{(\tau)}, p^{(\tau)})$ is continuous in $(R_{mobile}^{(\tau)}(s), R_{static}^{(\tau)}(s), p^{(\tau)}, \xi^{(\tau)})$. As a result of this and Lemma 4, the difference of the above expectation under the policies $\eta^{(\epsilon)}(\cdot|\cdot, \cdot, \cdot, \cdot)$ and $\eta^{(\epsilon)}(\cdot|\cdot, \{\bar{R}_{mobile}(s), \bar{R}_{static}(s)\}_{\forall s}, \cdot, \cdot)$ go to 0 as $\tau \rightarrow \infty$.

Now we claim that $(p^{(\tau)}, \xi^{(\tau)})$ converges to the internally chain transitive invariant sets of the o.d.e.

$$\begin{aligned} \dot{p}(t) &= \mathbf{E} \lim_{\beta \rightarrow 0} \frac{[p(t) + \beta \left(\sum_s \mathbf{1}\{s(t) = s, \eta_t = 0\} R_{static}(s) - R_0 \right)]_0^1 - p(t)}{\beta}, \\ \xi(t) &= 0, \end{aligned}$$

where the expectation is under the policy $\eta^{(\epsilon)}(\cdot|\cdot, \{\bar{R}_{mobile}(s), \bar{R}_{static}(s)\}_{\forall s}, \cdot, \cdot)$ and the fading distribution.

Note that, this o.d.e becomes $\dot{p}(t) \geq 0$ at $p(t) = 0$, $\dot{p}(t) \leq 0$ at $p(t) = 1$, and else

$$\begin{aligned} \dot{p}(t) &= \mathbf{E} \left(\sum_s \mathbf{1}\{s(t) = s, \eta_t = 0\} R_{static}(s) - R_0 \right) \\ &= \sum_s g(s) \mathbf{P}(\eta_t(s) = 0) \bar{R}_{static}(s) - R_0. \end{aligned}$$

Since $\mathbf{P}(\eta(s) = 0)$ is decreasing in $p(t)$, the o.d.e.

$$\dot{p}(t) = \mathbf{E} \lim_{\beta \rightarrow 0} \frac{[p(t) + \beta \left(\sum_s \mathbf{1}\{s(t) = s, \eta_t = 0\} R_{static}(s) - R_0 \right)]_0^1 - p(t)}{\beta}$$

can have at most one limit point. This limit point, which we call $p_\epsilon^*(\xi)$, is either in $\{0, 1\}$ or it is a stationary

point of the above o.d.e.

Also, by an argument similar to Lemma 2, $p_\epsilon^*(\xi)$ is Lipschitz continuous in ξ .

Hence, using an argument similar to [22, Chapter 6, Lemma 1], we prove the lemma. \square

A.3.2 Convergence in the slower timescale

We first prove the following lemma.

Lemma 6: $\mathbf{P}(\eta(s) = 0)$ under the randomized policy $\eta^{(\epsilon)}(\cdot|s, \{\bar{R}_{mobile}(s), \bar{R}_{static}(s)\}_{\forall s}, p_\epsilon^*(\xi), \xi)$ is continuous in ξ .

Proof: We have:

$$\begin{aligned} &\mathbf{P}(\eta(s) = 0) \\ &= \frac{\epsilon}{2} + (1 - \epsilon) \mathbf{P}(\bar{R}_{mobile}(s) - (\xi + \Delta) \bar{R}_{static}(s) \leq 0) \\ &= \frac{\epsilon}{2} + (1 - \epsilon) \mathbf{P}\left(\Delta \geq \frac{\bar{R}_{mobile}(s)}{\bar{R}_{static}(s)} - \xi\right) \end{aligned}$$

This is continuous in ξ and $p_\epsilon^*(\xi)$, and by an argument similar to Lemma 2, $p_\epsilon^*(\xi)$ is continuous in ξ . This proves the lemma. \square

Now we state the final lemma.

Lemma 7: Almost surely, as $\tau \rightarrow \infty$, the iterates $\xi^{(\tau)}$ converges to the projection of $\mathcal{K}_\epsilon(R_0)$ onto the ξ axis.

Proof: The proof follows using similar arguments as [25, Appendix E.C.3, Appendix E.C.4 and Appendix E.C.5]. The arguments require the results in Lemma 6, Lemma 4 and Lemma 5. The choice of A should be sufficiently large, otherwise $\xi^{(\tau)}$ might converge to A .

The proof of the slowest timescale convergence in [25, Theorem 12] involves checking of five conditions required for [26, Theorem 5.3.1]. The constrained MDP associated with [25, Theorem 12] had two constraints and hence two slower timescale iterates, whereas we have only one slowest timescale iterate for which it is easier to check these five conditions. Since these conditions hold, we can claim that the ξ iteration converges almost surely to the set of stationary points of a suitable o.d.e. Again, since we have only one slower timescale iterate, using large enough A is sufficient to ensure that the stationary points of that o.d.e. lie in $(0, A)$; it was much more complicated in [25, Theorem 12] since there were two slower timescale iterates.

The proof of this lemma requires Lemma 6 and the fact that $R_{static}(s)$ has a bounded support (since by Assumption 1, fading gain distribution has bounded support). \square

In light of Lemma 4, Lemma 5 and Lemma 7, the theorem is proved. \square

A.4 Proof of Corollary 1

Using arguments similar to the proof of Lemma 6, under $\eta^{(\epsilon)}(\cdot|\cdot, \cdot, \cdot, \cdot)$, one can claim that

$$\mathbf{P}(\eta_\tau = 0 | s(\tau) = s, \mathcal{R}_\tau = \mathcal{R}, \xi^{(\tau)} = \xi, p^{(\tau)} = p, \epsilon)$$

is continuous in $(\mathcal{R}, \xi, p, \epsilon)$ for given s . Hence, by Theorem 4, we can claim that: $\mathbf{P}(\eta_\tau = 0 | s(\tau) = s, \epsilon)$ converges to the set $\{x : x = \eta^{(\epsilon)}(0 | s, \{\bar{R}_{mobile}(s), \bar{R}_{static}(s)\}_{\forall s}, \xi, p), (\xi, p) \in \mathcal{K}_\epsilon(R_0)\}$ almost surely as $\tau \rightarrow \infty$.

Now,

$$\begin{aligned} & \lim_{\epsilon \downarrow 0} \{x : x = \eta^{(\epsilon)}(0 | s, \{\bar{R}_{mobile}(s), \bar{R}_{static}(s)\}_{\forall s}, \xi, p), (\xi, p) \in \mathcal{K}_\epsilon(R_0)\} \\ &= \{x : x = \eta^{(0)}(0 | s, \{\bar{R}_{mobile}(s), \bar{R}_{static}(s)\}_{\forall s}, \xi, p), (\xi, p) \in \mathcal{K}(R_0)\}. \end{aligned}$$

The proof trivially follows from this. \square