# Field Study in Deploying Restless Multi-Armed Bandits: Assisting Non-Profits in Improving Maternal and Child Health

Aditya Mate, Shresth Verma, Gargi Singh, Aparna Taneja and Milind Tambe

Google Research India

## 1 Introduction

<sup>1</sup> The wide-spread availability of cell phones has allowed non-profits to deliver targeted health information via voice or text messages to beneficiaries in underserved communities, often with significant demonstrated benefits to those communities [17, 29]. We focus in particular on non-profits that target improving maternal and infant health in low-resource communities in the global south. These non-profits deliver ante- and post-natal care information via voice and text to prevent adverse health outcomes [2, 14, 15].

Unfortunately, such information delivery programs are often faced with a key shortcoming: a large fraction of beneficiaries who enroll may drop out or reduce engagement with the information program. Yet non-profits often have limited health-worker time available on a periodic (weekly) basis to help prevent engagement drops. More specifically, there is limited availability of health-worker time where they can place crucial service calls (phone calls) to a limited number of beneficiaries, to encourage beneficiaries' participation, address complaints and thus prevent engagement drops.

Optimizing limited health worker resources to prevent engagement drops requires that we prioritize beneficiaries who would benefit most from service calls on a periodic basis. We model this resource optimization problem using Restless Multi-Armed Bandits (RMABs), with each beneficiary modeled as an RMAB arm. RMABs have been well studied for allocation of limited resources motivated by a myriad of application domains including preventive interventions for healthcare [25], planning anti-poaching patrols [32], machine repair and sensor maintenance [13] and communication systems [34]. However, RMABs have rarely seen real-world deployment, and to the best of our knowledge, never been deployed in the context of large-scale public health applications.

This chapter presents first results of an RMAB system in real-world public health settings. Based on available health worker time, RMABs choose m out of N total beneficiaries on a periodic basis for service calls, where the m are chosen to optimize prevention of engagement drops. The chapter presents two main contributions. First, previous work often assumes RMAB parameters as either known or easily learned over long periods of deployment. We show that neither

<sup>&</sup>lt;sup>1</sup> This work was pursued when Aditya was an intern at Google Research.

assumption holds in our real-world contexts; instead, we present clustering of offline historical data as a novel approach to infer unknown RMAB parameters.

Our second contribution is a real-world evaluation showing the benefit of our RMAB system, conducted in partnership with ARMMAN<sup>2</sup>, an NGO in India focused on maternal and child care. ARMMAN conducts a large-scale health information program, with concrete evidence of health benefits, which has so far served over a million mothers. As part of this program, an automated voice message is delivered to an expecting or new mother (beneficiary) over her cell phone on a weekly basis throughout pregnancy and for a year post birth in a language and time slot of her preference.

Unfortunately, ARMMAN's information delivery program also suffers from engagement drops. Therefore, in collaboration with ARMMAN we conducted a service quality improvement study to maximize the effectiveness of their service calls to ensure beneficiaries do not drop off from the program or stop listening to weekly voice messages. More specifically, the current standard of care in AR-MMAN's program is that any beneficiary may initiate a service call by placing a so called "missed call". This beneficiary-initiated service call is intended to help address beneficiaries' complaints and requests, thus encouraging engagement. However, given the overall decreasing engagement numbers in the current setup, key questions for our study are to investigate an approach for effectively conducting additional ARMMAN-initiated service calls (these are limited in number) to reduce engagement drops. To that end, our service quality improvement study comprised of 23,003 real-world beneficiaries spanning 7 weeks. Beneficiaries were divided into 3 groups, each adding to the current standard of care. The first group exercised ARMMAN's current standard of care (CSOC) without additional ARMMAN-initiated calls. In the second, the RMAB group, ARMMAN staff added to the CSOC by initiating service calls to 225 beneficiaries on average per week chosen by RMAB. The third was the Round-Robin group, where the exact same number of beneficiaries as the RMAB group were called every week based on a systematic sequential basis.

Results from our study demonstrate that RMAB provides statistically significant improvement over CSOC and round-robin groups. This improvement is also practically significant — the RMAB group achieves a ~ 30% reduction in engagement drops over the other groups. Moreover, the round-robin group does not achieve statistically significant improvement over the CSOC group, i.e., RMAB's optimization of service calls is crucial. To the best of our knowledge, this is the first large-scale empirical validation of use of RMABs in a public health context. Based on these results, the RMAB system is currently being transitioned to ARMMAN to optimize service calls to their ever growing set of beneficiaries. Additionally, this methodology can be useful in assisting engagement in many other awareness or adherence programs, e.g., [6, 36].

<sup>&</sup>lt;sup>2</sup> https://armman.org/

# 2 Related Work

Patient adherence monitoring in healthcare has been shown to be an important problem [24], and is closely related to the churn prediction problem, studied extensively in the context of industries like telecom [8], finance [33, 42], etc. The healthcare domain has seen several studies on patient adherence for diseases like HIV [38], cardiac problems [7, 35], Tuberculosis [19, 30], etc. These studies use a combination of patient background information and past adherence data, and build machine learning models to predict future adherence to prescribed medication <sup>3</sup>. However, such models treat adherence monitoring as a singleshot problem and are unable to appropriately handle the sequential resource allocation problem at hand. Additionally, the pool of beneficiaries flagged as high risk can itself be large, and the model can not be used to prioritize calls on a periodic basis, as required in our settings.

Campaign optimization (via phone outreach) has also been studied previously. Most existing works [9,20] however, rely on the availability of a customer social network based on preferences, behavior or demographics, to help identify the set of key customers who will increase the reach of the campaign. In our domains of interest, there is no evidence of a social network among the beneficiaries, so such campaign optimization techniques are inapplicable. Furthermore, campaign optimization relies on single-shot interventions for optimization, whereas, our problem requires tracking progress of beneficiaries over multiple timesteps.

The Restless Multi-Armed Bandit (RMAB) framework has been popularly adopted to tackle such sequential resource allocation problems [16, 40]. Computing the optimal solution for RMAB problems is shown to be PSPACE-hard. Whittle proposed an index-based heuristic [40], that can be solved in polynomial time and is now the dominant technique used for solving RMABs. It has been shown to be asymptotically optimal for the time average reward problem [39]. and other families of RMABs arising from stochastic scheduling problems [13]. Several works as listed in Section 1, show applicability of RMABs in different domains but these unrealistically assume perfect knowledge of the RMAB parameters, and have not been tested in real-world contexts. [3,5], present a Whittle Index-based Q-learning approach for unknown RMAB parameters. However, their techniques either assume identical arms or rely on receiving thousands of samples from each arm, which is unrealistic in our setting, given limited overall stay of a beneficiary in an information program — a beneficiary may drop out or stop engaging with the program few weeks post enrollment unless a service call convinces them to do otherwise. Instead, we present a novel approach that applies clustering to the available historical data to infer model parameters.

Clustering in the context of Multi-Armed Bandit and Contextual Bandits has received significant attention in the past [11, 21, 22, 43], but these settings do not consider restless bandit problems. [27] tackles a non-stationary setup

<sup>&</sup>lt;sup>3</sup> Similarly, in our previous preliminary study [28] published in a non-archival setting, we used demographic and message features to build models for predicting beneficiaries likely to drop-off from ARMMAN's information program.

with stochastic rewards, while [4] infers model parameters from independent studies in absence of historic data. In contrast, we focus on learning RMAB parameters using clustered historic beneficiary data. [23, 44] propose building predictive models per beneficiary in an online fashion, which is infeasible in our setup given the short stay of the beneficiaries.

## 3 Background: Restless Multi-Armed Bandits

An RMAB instance consists of N independent 2-action Markov Decision Processes (MDP) [31], where each MDP is defined by the tuple  $\{S, \mathcal{A}, R, \mathcal{P}\}$ . S denotes the state space,  $\mathcal{A}$  is the set of possible actions, R is the reward function  $R: S \times \mathcal{A} \times S \to \mathbb{R}$  and  $\mathcal{P}$  represents the transition function. We use  $P_{s,s'}^{\alpha}$  to denote the probability of transitioning from state s to state s' under the action  $\alpha$ . The policy  $\pi$ , is a mapping  $\pi: S \to \mathcal{A}$  that selects the action to be taken at a given state. The total reward accrued can be measured using either the discounted or average reward criteria to sum up the immediate rewards accrued by the MDP at each time step. Our formulation is amenable to both, although we use the discounted reward criterion in our study.

The expected discounted reward starting from state  $s_0$  is defined as  $V_{\beta}^{\pi}(s_0) = \mathbb{E}\left[\sum_{t=0}^{\infty} \beta^t R(s_t, \pi(s_t), s_{t+1} | \pi, s_0)\right]$  where the next state is drawn according to  $s_{t+1} \sim P_{s_t, s_{t+1}}^{\pi(s_t)}, \beta \in [0, 1)$  is the discount factor and actions are selected according to the policy mapping  $\pi$ . The planner's goal is to maximize the total reward.

## 4 Problem Statement

We model the engagement behavior of each beneficiary by an MDP corresponding to an arm of the RMAB. Pulling an arm corresponds to an active action, i.e., making a service call (denoted by  $\alpha = a$ ), while  $\alpha = p$  denotes the passive action of abstaining from a call. The state space S consists of binary valued states, s, that account for the recent engagement behavior of the beneficiary;  $s \in [NE, E]$ (or equivalently,  $s \in [0, 1]$ ) where E and NE denote the 'Engaging' and 'Not Engaging' states respectively. For example, in our domain, ARMMAN considers that if a beneficiary stays on the automated voice message for more than 30 seconds (average message length is 1 minute), then the beneficiary has engaged. If a beneficiary engages at least once with the automated voice messages sent during a week, they are assigned the engaging (E) state for that time step and nonengaging (NE) state otherwise. For each action  $\alpha \in \mathcal{A}$ , the beneficiary states follow a Markov chain represented by the 2-state Gilbert-Elliot model [12] with transition parameters given by  $P_{ss'}^{\alpha}$ , as shown in Figure 1. With slight abuse of notation, the reward function  $R(\cdot)$  of  $n^{th}$  MDP is simply given by  $R_n(s) = s$  for  $s \in \{0, 1\}$ .

We adopt the Whittle solution approach for solving the RMAB. It hinges around the key idea of a "passive subsidy", which is a hypothetical reward offered to the planner, in addition to the original reward function for choosing the passive



Fig. 1: The beneficiary transitions from a current state s to a next state s' under action  $\alpha$ , with probability  $P_{sc'}^{\alpha}$ .

action. The Whittle Index is then defined as the infimum subsidy that makes the planner indifferent between the 'active' and the 'passive' actions, i.e.,:

$$W(s) = inf_{\lambda}\{\lambda : Q_{\lambda}(s, p) = Q_{\lambda}(s, a)\}$$
(1)

We assume the planner has access to an offline historical data set of beneficiaries,  $\mathcal{D}_{train}$ . Each beneficiary data point  $\mathcal{D}_{train}[i]$  consists of a tuple,  $\langle f, \mathcal{E} \rangle$ , where f is beneficiary i's feature vector of static features, and  $\mathcal{E}$  is an episode storing the trajectory of  $(s, \alpha, s')$  pairs for that beneficiary, where s denotes the start state,  $\alpha$  denotes the action taken (passive v/s active), and s' denotes the next state that the beneficiary lands in after executing  $\alpha$  in state s. We assume that these  $(s, \alpha, s')$  samples are drawn according to fixed, latent transition matrices  $P_{ss'}^a[i]$  and  $P_{ss'}^p[i]$  (corresponding to the active and passive actions respectively), unknown to the planner, and potentially unique to each beneficiary.

Given  $D_{train}$ , we now consider a new beneficiary cohort  $\mathcal{D}_{test}$ , consisting of N beneficiaries, marked  $\{1, 2, \ldots, N\}$ , that the planner must plan service calls for. The MDP transition parameters corresponding to beneficiaries in  $\mathcal{D}_{test}$  are unknown to the planner, but assumed to be drawn at random from a distribution similar to the joint distribution of features and transition parameters of beneficiaries in the historical data distribution. We assume the planner has access to the feature vector f for each beneficiary in  $\mathcal{D}_{test}$ .

We now define the service call planning problem as follows. The planner has upto m resources available per round, which the planner may spend towards delivering service calls to beneficiaries. Beneficiaries are represented by N arms of the RMAB, of which the planner may pull upto m arms (i.e., m service calls) at each time step. We consider a round or timestep of one week which allows planning based on the most recent engagement patterns of the beneficiaries.

## 5 Method

Figure 2 shows our overall solution methodology. We use clustering techniques that exploit historical data  $D_{train}$  to estimate an offline RMAB problem instance relying solely on the beneficiaries' static features and state transition data. This enables overcoming the challenge of limited samples (time-steps) per beneficiary. Based on this estimation, we use the Whittle Index approach to prioritize service calls.



Fig. 2: RMAB Training and Testing pipelines proposed

#### 5.1 Clustering Methods

We use historical data  $\mathcal{D}_{train}$  to learn the impact of service calls on transition probabilities. While there is limited service call data (active transition samples) for any single beneficiary, clustering on the beneficiaries allows us to combine their data to infer transition probabilities for the entire group. Clustering offers the added advantage of reducing computational cost for resource limited NGOs; since all beneficiaries within a cluster share identical transition probability values we can compute their Whittle index all at once. We present four such clustering techniques below:

(i). Features-only Clustering (FO): This method relies on the correlation between the beneficiary feature vector f and their corresponding engagement behavior. We employ k-means clustering on the feature vector f of all beneficiaries in the historic dataset  $D_{train}$ , and then derive the representative transition probabilities for each cluster by pooling all the  $(s, \alpha, s')$  tuples of beneficiaries assigned to that cluster. At test time, the features f of a new, previously unseen beneficiary in  $D_{test}$  map the beneficiary to their corresponding cluster and estimated transition probabilities.

(ii). Feature + All Probabilities (FAP) In this 2-level hierarchical clustering technique, the first level uses a rule-based method, using features to divide beneficiaries into a large number of pre-defined buckets, B. Transition probabilities are then computed by pooling the  $(s, \alpha, s')$  samples from all the beneficiaries in each bucket. Finally, we perform a k-means clustering on the transition probabilities of these B buckets to reduce them to k clusters  $(k \ll B)$ . However, this method suffers from several smaller buckets missing or having very few active transition samples.

*(iii). Feature + Passive Probabilities (FPP):* This method builds on the FAP method, but only considers the passive action probabilities to preclude the issue of missing active transition samples.

(iv). Passive Transition-Probability based Clustering (PPF): The key motivation here is to group together beneficiaries with similar transition behaviors, irrespective of their features. To this end, we use k-means clustering on passive transition probabilities (to avoid issues with missing active data) of beneficiaries in  $D_{train}$  and identify cluster centers. We then learn a map  $\phi$  from the feature vector f to the cluster assignment of the beneficiaries that can be used to infer the cluster assignments of new beneficiaries at test-time solely from f. We use a random forest model as  $\phi$ .

The rule-based clustering on features involved in FPP and FAP methods can be thought of as using one specific, hand-tuned mapping function  $\phi$ . In contrast, the PPF method *learns* such a map  $\phi$  from data, eliminating the need to manually define accurate and reliable feature buckets.

#### 5.2 Evaluation of Clustering Methods

We use a historical dataset,  $\mathcal{D}_{train}$  from ARMMAN consisting of 4238 beneficiaries in total, who enrolled into the program between May-July 2020. We compare the clustering methods empirically, based on the criteria described below.

1. Representation: Cluster centers that are representative of the underlying data distribution better resemble the ground truth transition probabilities. This is of prime importance to the planner, who must rely on these values to plan actions. Fig 3 plots the ground truth transition probabilities and the resulting cluster centers determined using the proposed methods. Visual inspection reveals that the PPF method represents the ground truth well, as is corroborated by the quantitative metrics of Table 1 that compares the RMSE error across different clustering methods.

2. Balanced cluster sizes: A low imbalance across cluster sizes is desirable to preclude the possibility of arriving at few, gigantic clusters which will assign identical whittle indices to a large groups of beneficiaries. Working with smaller clusters also aggravates the missing data problem in estimation of active transition probabilities. Considering the variance in cluster sizes and RMSE error for the different clustering methods with  $k = \{20, 40\}$  as shown in Table 1, *PPF* outperforms the other clustering methods and was chosen for the pilot study.

Table 1: Average RMSE and cluster size variance over all beneficiaries for different methods. Total Beneficiaries = 4238,  $\mu_{20} = 211.9$ ,  $\mu_{40} = 105.95$  ( $\mu = average$  beneficiaries per cluster)

Clustering	Average	e RMSE	Standar	d Deviation
Method	k = 20	$\mathbf{k} = 40$	$\mathbf{k} = 20$	k = 40
FO	0.229	0.228	143.30	74.22
FPP	0.223	0.222	596.19	295.01
FAP	0.224	0.223	318.46	218.37
PPF	0.041	0.027	145.59	77.50

8



Fig. 3: Comparison of passive transition probabilities obtained from different clustering methods with cluster sizes  $k = \{20, 40\}$  with the ground truth transition probabilities. Blue dots represent the true passive transition probabilities for every beneficiary while red or green dots represent estimated cluster centres.

Next we turn to choosing k, the number of clusters: as k grows, the clusters become sparse in number of active samples aggravating the missing data problem while a smaller k suffers from a higher RMSE. We found k = 40 to be optimal and chose it for the pilot study.

Finally, we adopt the Whittle solution approach for RMABs to plan actions and pre-compute all of the possible 2 \* k index values that beneficiaries can take (corresponding to combinations of k possible clusters and 2 states). The indices can then be looked up at all future time steps in constant time, making this an optimal solution for large scale deployment with limited compute resources.

As we got this RMAB system ready for real-world use, there was an important observation for social impact settings: real-world use also required us to carefully handle several domain specific challenges, which were time consuming. For example, despite careful clustering, a few clusters may still be missing active probability values, which required employing a data imputation heuristic . Moreover, there were other constraints specific to ARMMAN, such as a beneficiary should receive only one service call every  $\eta$  weeks, which was addressed by introducing "sleeping states" for beneficiaries who receive a service call .

#### 6 **Resource Requirements**

#### 6.1 Data Collected by ARMMAN

Beneficiaries enroll into ARMMAN's information program with the help of health workers, who collect the beneficiary's demographic data such as age, education level, income bracket, phone owner in the family, gestation age, number of children, preferred language and preferred slots for the automated voice messages during enrolment. These features are referred to as Beneficiary Registration Features in rest of the chapter. Beneficiaries provided both written and digital consent for receiving automated voice messages and service calls. ARMMAN also stores listenership information regarding the automated voice messages together with the registration data in an anonymized fashion.

## 7 Field Evaluation

In this section, we discuss a real-world quality improvement study. We also simulate the expected outcome in other synthetically constructed situations and demonstrate good performance of our approach across the board.

## 7.1 Service Quality Improvement Study

**Setup** This cohort of beneficiaries registered in the program between Feb 16, 2021 and March 15, 2021 as  $D_{test}$  and started receiving automated voice messages few days post enrolment as per their gestational age. Additionally, as per the current standard of care, any of these beneficiaries could initiate a service call by placing a "missed call". The 23003 beneficiaries are randomly distributed across 3 groups, each group adding to the CSOC as follows:

- Current-Standard-of-Care (CSOC) Group: The beneficiaries in this group follow the original standard of care, where there are no ARMMAN initiated service calls. The listenership behavior of beneficiaries in this group is used as a benchmark for the RR and RMAB groups.
- RMAB group: In this group, beneficiaries are selected for ARMMANinitiated service call per week via the Whittle Index policy described in Section 4. Even though all beneficiaries within a cluster are modeled by identical MDP parameters, their states may evolve independently, and so the Whittle indices are tracked for each beneficiary separately, leading to an RMAB with 7668 arms.
- Round Robin (RR) group: By default, NGOs including ARMMAN often conduct service calls using some systematic set order – the idea here is to have an easily executable policy, that services enough of a cross-section of

Table 2: Beneficiary distribution in the three groups and their start states during week 0 of the study.

Group	Engaging (E)	Non-Engaging (NE)	Total
RMAB	3571	4097	7668
RR	3647	4021	7668
CSOC	3661	4006	7667

beneficiaries and can be scaled up or down per week based on available resources. To recreate this setting, we generate service calls to beneficiaries based on the ascending order of their date of enrollment for this RR group, as recommended by ARMMAN. If this method succeeds compared to CSOC, then a simple manual strategy is enough; RMAB style optimization may not be needed.

Table 2 shows the absolute number of beneficiaries in states E or NE, where the state is computed using one week of engagement data between April 19 - April 26, 2021.

Beneficiaries across all three groups receive the same automated voice messages regarding pregnancy and post-birth care throughout the program, and no health related information is withheld from any beneficiary. The study only aims to evaluate the effectiveness of ARMMAN-initiated outbound service calls with respect to improving engagement with the program across the three groups. No interviews or research data or feedback was collected from the beneficiaries.

The study started on April 26, 2021, with m beneficiaries selected from the RMAB and RR group each  $(m \ll N)$  per week for ARMMAN-initiated service calls. ARMMAN staff performing service calls were blind to the experimental groups that the beneficiaries belonged to. Recall, the goal of the service calls is to encourage the beneficiaries to engage with the health information message program in the future. For this study, number of service calls m was on average 225 per week for each of RMAB and RR groups to reflect real-world constraints on service calls. The study was scheduled for a total of 7 weeks, during which 20% of the RMAB (and RR) group had received a service calls by ARMMAN. 4

**Results** We present our key results from the study in Figure 4. The results are computed at the end of 7 weeks from the start of the quality improvement study on April 26, 2021.

Figure 4 measures the impact of service calls by the RMAB and RR policies in comparison to the CSOC Group. Beneficiaries' engagement with the program typically starts to dwindle with time.

 $<sup>^4</sup>$  Each beneficiary group also received very similar beneficiary-initiated calls, but these were less than 10% of the ARMMAN-initiated calls in RMAB or RR groups over 7 weeks.



Fig. 4: Cumulative number of weekly engagement drops prevented (in comparison to the CSOC group) by RMAB far exceed those prevented by RR.

In Figure 4, we measure the impact of a service call policy as the cumulative drop in engagement prevented compared to the CSOC Group. We consider drop in engagement instead of the raw engagement numbers themselves, because of the slight difference in the numbers of beneficiaries in engaging (E) state at the start of the study. The drop in engagement under a policy  $\pi$  at time t can be measured as the change in engagement:

$$\Delta_{current}^{\pi}(t) \coloneqq \sum_{n \in N} (R_n(s_0) - R_n(s_t)) \tag{2}$$

where  $R_n(s_t)$  represents the reward for  $n^{th}$  beneficiary in state  $s_t$  at time step t and cumulative drop in engagement is:

$$\Delta_{cumulative}^{\pi}(t) \coloneqq \sum_{n \in N} \sum_{\zeta=0}^{\zeta=t} (R_n(s_0) - R_n(s_{\zeta}))$$
(3)

The cumulative drop in engagement prevented by a policy  $\pi$ , in comparison to the CSOC Group is thus simply:

$$\Delta_{cumulative}^{\pi}(t) - \Delta_{cumulative}^{CSOC}(t) \tag{4}$$

and is plotted on the y-axis of Figure 4.

Figure 4 shows that the RMAB policy prevents a total 622 instances of a drop in automated health message engagement, at the end of 7 weeks, as compared to CSOC. RR group, on the other hand, only prevents 101 engagement drops by the end of week 7. Given that there are a total of 1944 engagement drops in the CSOC group, we show in the first row of Table 3, that the RMAB group has 32.0% and 28.3% less cumulative engagement drops as compared to the CSOC and RR groups respectively by the end of the study.

a mical regression model. We use: $p < 0.05$ , $p < 0.1$								
	RMAB vs	RR vs	RMAB vs					
	CSOC	CSOC	RR					
% reduction in cumulative en- gagement drops	32.0%	5.2%	28.3%					
p-value	$0.044^{*}$	0.740	$0.098^{\dagger}$					
Coefficient $\beta$	-0.0819	-0.0137	-0.0068					

Table 3: Statistical significance for service call policy impact at week 7 is tested using a linear regression model. We use: p < 0.05; p < 0.1

**Statistical Analysis** To investigate the benefit from use of RMAB policy over policies in the RR and CSOC groups, we use regression analysis [1]. Specifically, we fit a linear regression model to predict number of cumulative engagement drops at week 7 while controlling for treatment assignment and covariates specified by beneficiary registration features. The model is given by:

$$Y_i = k + \beta T_i + \sum_{j=1}^J \gamma_j x_{ij} + \epsilon_i$$

where for the  $i_{th}$  beneficiary,  $Y_i$  is the outcome variable defined as number of cumulative engagement drops at week 7, k is the constant term,  $\beta$  is the treatment effect,  $T_i$  is the treatment indicator variable,  $x_i$  is a vector of length J representing the  $i_{th}$  beneficiary's registration features,  $\gamma_j$  represents the impact of the  $j_{th}$  feature on the outcome variable and  $\epsilon_i$  is the error term. For evaluating the effect of RMAB service calls as compared to CSOC group, we fit the regression model only for the subset of beneficiaries assigned to either of these two groups.  $T_i$  is set to 1 for beneficiaries belonging to the RMAB group and 0 for those in CSOC group. We repeat the same experiment to compare RR vs CSOC group and RMAB vs RR group.

The results are summarized in Table 3. We find that RMAB has a statistically significant treatment effect in reducing cumulative engagement drop (negative  $\beta, p < 0.05$ ) as compared to CSOC group. However, the treatment effect is not statistically significant when comparing RR with CSOC group (p = 0.740). Additionally, comparing RMAB group with RR, we find  $\beta$ , the RMAB treatment effect, to be significant (p < 0.1). This shows that RMAB policy has a statistically significant effect on reducing cumulative engagement drop as compared to both the RR policy and CSOC. RR fails to achieve statistical significance against CSOC. Together these results illustrate the importance of RMAB's optimization of service calls, and that without such optimization, service calls may not yield any benefits.

**RMAB Strategies** We analyse RMAB's strategic selection of beneficiaries in comparison to RR using Figure 5, where we group beneficiaries according to their whittle indices, equivalently their (cluster, state). Figure 5 plots the frequency distribution of beneficiaries (shown via corresponding clusters) who were



Fig. 5: Distributions of clusters picked for service calls by RMAB and RR are significantly different. RMAB is very strategic in picking only a few clusters with a promising probability of success, RR displays no such selection.

selected by RMAB and RR in the first two weeks. For example, the top plot in Figure 5a shows that RMAB selected 60 beneficiaries from cluster 29 (NE state). First, we observe that RMAB was clearly more selective, choosing beneficiaries from just four (Figure 5a) or seven (Figure 5b) clusters, rather than RR that chose from 20. Further, we assign each cluster a hue based on their probability of transitioning to engaging state from their current state given a service call. Figure 5 reveals that RMAB consistently prioritizes clusters with high probability of success (blue hues) while RR deploys no such selection; its distribution emulates the overall distribution of beneficiaries across clusters (mixed blue and red hues).

Furthermore, Figure 6a further highlights the situation in week 1, where RMAB spent 100% of its service calls on beneficiaries in the non-engaging state while RR spent the same on only 64%. Figure 6b shows that RMAB converts 31.2% of the beneficiaries shown in Figure 6a from non-engaging to engaging state by week 7, while RR does so for only 13.7%. This further illustrates the need for optimizing service calls for them to be effective, as done by RMAB.

## 7.2 Synthetic Results

We run additional simulations to test other service call policies beyond those included in the quality improvement study and confirm the superior performance of RMAB. Specifically, we compare to the following baselines: (1) RANDOM is a naive baseline that selects m arms at random. (2) MYOPIC is a greedy algorithm that pulls arms optimizing for the reward in the immediate next time step. WHITTLE is our algorithm.

We compute a normalized reward of an algorithm ALG as:  $\frac{100 \times (\overline{R}^{\text{ALG}} - \overline{R}^{\text{CSOC}})}{\overline{R}^{\text{WHITTLE}} - \overline{R}^{\text{CSOC}}}$ 

where  $\overline{R}$  is the total discounted reward. Simulation results are averaged over 30 independent trials and run over 40 weeks.



Fig. 6: (a) % of week 1 service calls on non-engaging beneficiaries (b) % of non-engaging beneficiaries of week 1 receiving service calls that converted to engaging by week 7



Fig. 7: Performance of MYOPIC can be arbitrarily bad and even worse than RAN-DOM, unlike the Whittle policy.

Figure 7 presents simulation of an adversarial example [25] consisting of x% of non-recoverable and 100-x% of self-correcting beneficiaries for different values of x. Self-correcting beneficiaries tend to miss automated voice messages sporadically, but revert to engaging ways without needing a service call. Non-recoverable beneficiaries are those who may drop out for good, if they stop engaging. We find that in such situations, MYOPIC proves brittle, as it performs even worse than RANDOM while WHITTLE performs well consistently. The actual quality improvement study cohort consists of 48.12% non-recoverable beneficiaries (defined by  $P_{01}^p < 0.2$ ) and the remaining comprised of self-correcting and other types of beneficiaries.

# 8 Lessons Learned

The widespread use of cell-phones, particularly in the global south, has enabled non-profits to launch massive programs delivering key health messages to a broad population of beneficiaries in a cost-effective manner. We present an RMAB based system to assist these non-profits in optimizing their limited service resources. To the best of our knowledge, ours is the first study to demonstrate the effectiveness of such RMAB-based resource optimization in real-world public health contexts. These encouraging results have initiated the transition of our RMAB software to ARMMAN for real-world deployment. We hope this work paves the way for use of RMABs in many other health service applications.

Some key lessons learned from this research, which complement some of the lessons outlined in [10, 37, 41] include the following. First, social-impact driven engagement and design iterations with the NGOs on the ground is crucial to understanding the right AI model for use and appropriate research challenges. As discussed in footnote 1, our initial effort used a one-shot prediction model, and only after some design iterations we arrived at the current RMAB model. Next, given the missing parameters in RMAB, we found that the assumptions made in literature for learning such parameters did not apply in our domain, exposing new research challenges in RMABs. In short, domain partnerships with NGOs to achieve real social impact automatically revealed requirements for use of novel application of an AI model (RMAB) and new research problems in this model.

Second, data and compute limitations of non-profits are a real-world constraint, and must be seen as genuine research challenges in AI for social impact, rather than limitations. In our domain, one key technical contribution in our RMAB system is deploying clustering methods on offline historical data to infer unknown RMAB parameters. Data is limited as not enough samples are available for any given beneficiary, who may stay in the program for a limited time. Non-profit partners also cannot bear the burden of massive compute requirements. Our clustering approach allows efficient offline mapping to Whittle indices, addressing both data and compute limits, enabling scale-up to service 10s if not 100s of thousands of beneficiaries. Third, in deploying AI systems for social impact, there are many technical challenges that may not need innovative solutions, but they are critical to deploying solutions at scale. Indeed, deploying any system in the real world is challenging, but even more so in domains where NGOs may be interacting with low-resource communities. We hope this work serves as a useful example of deploying an AI based system for social impact in partnership with non-profits in the real world and will pave the way for more such solutions with real-world impact.

Finally, there are also some important topics for future work in improving the RMAB system, which include handling fairness [26], changing the current RMAB model with two actions to incorporate multiple actions [18], and improving the RMAB model from interactions with beneficiaries [5].

## Acknowledgement

We would like to thank all our collaborators at Google Research India, Harvard, ARMMAN and Google.org for supporting this work and making the field study possible.

## References

1. ANGRIST, J. D., AND PISCHKE, J.-S. *Mostly harmless econometrics*. Princeton university press, 2008.

- 2. ARMMAN. mmitra. https://armman.org/mmitra/, 2020.
- 3. AVRACHENKOV, K., AND BORKAR, V. S. Whittle index based q-learning for restless bandits with average reward. arXiv preprint arXiv:2004.14427 (2020).
- AYER, T., ZHANG, C., BONIFONTE, A., SPAULDING, A. C., AND CHHATWAL, J. Prioritizing hepatitis c treatment in us prisons. *Operations Research* 67, 3 (2019), 853–873.
- BISWAS, A., AGGARWAL, G., VARAKANTHAM, P., AND TAMBE, M. Learn to intervene: An adaptive learning policy for restless bandits in application to preventive healthcare. In Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021 (2021), Z. Zhou, Ed., ijcai.org, pp. 4039–4046.
- CHEN, R., SANTO, K., WONG, G., SOHN, W., SPALLEK, H., CHOW, C., AND IRVING, M. Mobile apps for dental caries prevention: Systematic search and quality evaluation. *JMIR mHealth and uHealth 9* (01 2021).
- COROTTO, P. S., MCCAREY, M. M., ADAMS, S., KHAZANIE, P., AND WHELLAN, D. J. Heart failure patient adherence: epidemiology, cause, and treatment. *Heart failure clinics 9*, 1 (2013), 49–58.
- DAHIYA, K., AND BHATIA, S. Customer churn analysis in telecom industry. In 2015 4th International Conference on Reliability, Infocom Technologies and Optimization (ICRITO)(Trends and Future Directions) (2015), IEEE, pp. 1–6.
- EAGLE, N., MACY, M., AND CLAXTON, R. Network diversity and economic development. *Science* 328, 5981 (2010), 1029–1031.
- FLORIDI, L., COWLS, J., KING, T., AND TADDEO, M. How to design ai for social good: Seven essential factors. *Science and Engineering Ethics* 26 (06 2020).
- 11. GENTILE, C., LI, S., AND ZAPPELLA, G. Online clustering of bandits. In International Conference on Machine Learning (2014), PMLR, pp. 757–765.
- GILBERT, E. N. Capacity of a burst-noise channel. Bell system technical journal 39, 5 (1960), 1253–1265.
- GLAZEBROOK, K. D., RUIZ-HERNANDEZ, D., AND KIRKBRIDE, C. Some indexable families of restless bandit problems. *Advances in Applied Probability* 38, 3 (2006), 643–672.
- 14. HELPMUM. Preventing maternal and infant mortality in nigeria. https://helpmum.org/, 2021.
- 15. JOHNSON, J. . Momente: Connecting women to care, one text at a time. https://www.jnj.com/our-giving/ momentec-connecting-women-to-care-one-text-at-a-time, 2017.
- JUNG, Y. H., AND TEWARI, A. Regret bounds for thompson sampling in episodic restless bandit problems. Advances in Neural Information Processing Systems (2019).
- KAUR, J., KAUR, M., CHAKRAPANI, V., WEBSTER, J., SANTOS, J., AND KU-MAR, R. Effectiveness of information technology-enabled 'smart eating' health promotion intervention: A cluster randomized controlled trial. *PLOS ONE 15* (01 2020), e0225892.
- KILLIAN, J. A., BISWAS, A., SHAH, S., AND TAMBE, M. Q-learning lagrange policies for multi-action restless bandits. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining* (2021), pp. 871–881.
- KILLIAN, J. A., WILDER, B., SHARMA, A., CHOUDHARY, V., DILKINA, B., AND TAMBE, M. Learning to prescribe interventions for tuberculosis patients using digital adherence data. *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Jul 2019).

- 20. LESKOVEC, J., ADAMIC, L. A., AND HUBERMAN, B. A. The dynamics of viral marketing. *ACM Trans. Web* 1, 1 (may 2007), 5–es.
- LI, C., WU, Q., AND WANG, H. Unifying clustered and non-stationary bandits. In International Conference on Artificial Intelligence and Statistics (2021), PMLR, pp. 1063–1071.
- LI, S., CHEN, W., AND LEUNG, K.-S. Improved algorithm on online clustering of bandits. arXiv preprint arXiv:1902.09162 (2019).
- LIAO, P., GREENEWALD, K., KLASNJA, P., AND MURPHY, S. Personalized heartsteps: A reinforcement learning algorithm for optimizing physical activity. *Proceed*ings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies 4, 1 (2020), 1–22.
- MARTIN, L. R., WILLIAMS, S. L., HASKARD, K. B., AND DIMATTEO, M. R. The challenge of patient adherence. *Therapeutics and clinical risk management 1*, 3 (2005), 189.
- MATE, A., KILLIAN, J., XU, H., PERRAULT, A., AND TAMBE, M. Collapsing bandits and their application to public health intervention. *Advances in Neural Information Processing Systems* 34 (2020).
- MATE, A., PERRAULT, A., AND TAMBE, M. Risk-aware interventions in public health: Planning with restless multi-armed bandits. *Autonomous Agents and Multii-Agent Systems (AAMAS)* (2021).
- MINTZ, Y., ASWANI, A., KAMINSKY, P., FLOWERS, E., AND FUKUOKA, Y. Nonstationary bandits with habituation and recovery dynamics. *Operations Research* 68, 5 (2020), 1493–1516.
- 28. NISHTALA, S., KAMARTHI, H., THAKKAR, D., NARAYANAN, D., GRAMA, A., HEGDE, A., PADMANABHAN, R., MADHIWALLA, N., CHAUDHARY, S., RAVIN-DRAN, B., AND TAMBE, M. Missed calls, automated calls and health support: Using ai to improve maternal health outcomes by increasing program engagement, 2020.
- PFAMMATTER, A., SPRING, B., SALIGRAM, N., DAVÉ, R., GOWDA, A., BLAIS, L., ARORA, M., RANJANI, H., GANDA, O., HEDEKER, D., REDDY, S., AND RAMALINGAM, S. mhealth intervention to improve diabetes risk behaviors in india: A prospective, parallel group cohort study. *Journal of Medical Internet Research* 18 (08 2016), e207.
- PILOTE, L., TULSKY, J. P., ZOLOPA, A. R., HAHN, J. A., SCHECTER, G. F., AND MOSS, A. R. Tuberculosis Prophylaxis in the Homeless: A Trial to Improve Adherence to Referral. Archives of Internal Medicine 156, 2 (01 1996), 161–165.
- PUTERMAN, M. L. Markov Decision Processes: Discrete Stochastic Dynamic Programming. Wiley Series in Probability and Statistics. Wiley, 1994.
- QIAN, Y., ZHANG, C., KRISHNAMACHARI, B., AND TAMBE, M. Restless poachers: Handling exploration-exploitation tradeoffs in security domains. In Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems, Singapore, May 9-13, 2016 (2016), C. M. Jonker, S. Marsella, J. Thangarajah, and K. Tuyls, Eds., ACM, pp. 123–131.
- SHAABAN, E., HELMY, Y., KHEDR, A., AND NASR, M. A proposed churn prediction model. *International Journal of Engineering Research and Applications 2*, 4 (2012), 693–697.
- SOMBABU, B., MATE, A., MANJUNATH, D., AND MOHARIR, S. Whittle index for AoI-aware scheduling. In 2020 12th International Conference on Communication Systems & Networks (2020), IEEE.

- 18 Authors Suppressed Due to Excessive Length
- SON, Y.-J., KIM, H.-G., KIM, E.-H., CHOI, S., AND LEE, S.-K. Application of support vector machine for prediction of medication adherence in heart failure patients. *Healthcare informatics research* 16, 4 (2010), 253–259.
- THIRUMURTHY, H., AND LESTER, R. T. M-health for health behaviour change in resource-limited settings: applications to hiv care and beyond. *Bulletin of the World Health Organization 90* (2012), 390–392.
- 37. TOMAŠEV, N., CORNEBISE, J., HUTTER, F., MOHAMED, S., PICCIARIELLO, A., CONNELLY, B., BELGRAVE, D., EZER, D., HAERT, F., MUGISHA, F., ABILA, G., ARAI, H., ALMIRAAT, H., PROSKURNIA, J., SNYDER, K., OTAKE, M., OTHMAN, M., GLASMACHERS, T., WEVER, W., AND CLOPATH, C. Ai for social good: unlocking the opportunity for positive impact. *Nature Communications 11* (05 2020), 2468.
- TULDRÀ, A., FERRER, M. J., FUMAZ, C. R., BAYÉS, R., PAREDES, R., BURGER, D. M., AND CLOTET, B. Monitoring Adherence to HIV Therapy. Archives of Internal Medicine 159, 12 (06 1999), 1376–1377.
- 39. WEBER, R. R., AND WEISS, G. On an index policy for restless bandits. *Journal* of applied probability (1990), 637–648.
- WHITTLE, P. Restless bandits: Activity allocation in a changing world. Journal of applied probability (1988), 287–298.
- 41. WILDER, B., ONASCH-VERA, L., DIGUISEPPI, G., PETERING, R., HILL, C., YA-DAV, A., RICE, E., AND TAMBE, M. Clinical trial of an ai-augmented intervention for hiv prevention in youth experiencing homelessness. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2021), vol. 35, pp. 14948–14956.
- XIE, Y., LI, X., NGAI, E., AND YING, W. Customer churn prediction using improved balanced random forests. *Expert Systems with Applications 36*, 3, Part 1 (2009), 5445 – 5449.
- YANG, L., LIU, B., LIN, L., XIA, F., CHEN, K., AND YANG, Q. Exploring clustering of bandits for online recommendation system. In *Fourteenth ACM Conference* on Recommender Systems (2020), pp. 120–129.
- 44. ZHOU, M., MINTZ, Y., FUKUOKA, Y., GOLDBERG, K., FLOWERS, E., KAMIN-SKY, P., CASTILLEJO, A., AND ASWANI, A. Personalizing mobile fitness apps using reinforcement learning. In *CEUR workshop proceedings* (2018), vol. 2068, NIH Public Access.