

# Online Sequential Learning using Bandits

**Subhojyoti Mukherjee**

*IIT Madras, CSE Department*

SUBHOJYOTIMUKHERJEE22@GMAIL.COM

**Balaraman Ravindran (Guide)**

*IIT Madras, CSE Department*

RAVI@CSE.IITM.AC.IN

**Nandan Sudarsanam (Co-Guide)**

*IIT Madras, DoMS Department*

NANDAN@IITM.AC.IN

## 1. Introduction

In today's world a large number of problems in science and engineering, robotics and game playing, resource management, financial portfolio management, medical treatment design, ad placement, website optimization and packet routing can be modeled as sequential decision-making under uncertainty. Many of these real-world interesting sequential decision-making problems can be formulated as reinforcement learning (RL) problems ([Bertsekas and Tsitsiklis \(1996\)](#), [Sutton and Barto \(1998\)](#)). In an RL problem, an agent interacts with a dynamic, stochastic, and unknown environment, with the goal of finding an action-selection strategy or policy that optimizes some long-term performance measure. Every time when the agent interacts with the environment it receives a signal/reward from the environment based on which it modifies its policy. The agent learns to optimize the choice of actions over several time steps which is learned from the sequences of data that it receives from the environment. This is the crux of online sequential learning. This is in contrast to supervised learning methods that deal with labeled data which are independently and identically distributed (i.i.d.) samples from the domain and train some classifier on the entire training dataset to learn the pattern of this distribution to predict future samples (test dataset) with the assumption that it is sampled from the same domain, whereas the RL agent learns from the samples that are collected from the trajectories generated by its sequential interaction with the system. For an RL agent the trajectory consists of a series of sequential interactions whereby it transitions from one state to another following some dynamics intrinsic to the environment while collecting the reward till some stopping condition is reached. This is known as an episode. For a single-step interaction, i.e., when the episode terminates after a single transition, the problem is captured by the multi-armed bandit (MAB) model. Our work will focus on this idea of MAB model.

## 2. Motivation

The MAB model fits very well in various real-world scenarios that can be modeled as sequential decision-making problems. Some of which are mentioned as follows:-

1. *Online Shop Domain* ([Ghavamzadeh et al. \(2015\)](#)): In the online shop domain, a retailer aims to maximize profit by sequentially suggesting products to online shopping customers. In this scenario, at every timestep, the retailer displays an item to a customer from a pool of items

which has the highest probability of being selected by the customer. The episode ends when the customer selects or does not select a product (which will be considered as a loss to the retailer) and the process is again repeated till a pre-specified number of times with the retailer gathering valuable information regarding the customer from this behaviour and modifying its policy to display the next item.

2. *Medical Treatment Design* ([Thompson \(1933\)](#)): Here at every timestep, the agent chooses to administer one out of several treatments sequentially on a patient. Here, the episode ends when the patient responds well or does not respond well to the treatment whereby the agent modifies its policy for the next suggestion.
3. *Financial Portfolio Management*: In financial portfolio management MAB model can be used. Here, the agent is faced with the choice of selecting the most profitable stock option out of several stock options. The simplest strategy where we can employ a bandit model is this; at the start of every trading session the agent suggests a stock to purchase worth Re 1, while at the closing of the trading session it sells off the stock to witness its value after a day's trading. The profit recorded is treated as the reward revealed by the environment and the agent modifies its policy for the next day.

The thresholding bandit problem is a special case of combinatorial MAB problem where the learner has to suggest the best set of arms above a real valued threshold. This has several relevant industrial applications. The variants of TopM problem (identifying the best  $M$  arms from  $K$  given arms) can be readily used in the thresholding problem.

1. *Product Selection*: A company wants to introduce a new product in market and there is a clear separation of the test phase from the commercialization phase. In this case the company tries to minimize the loss it might incur in the commercialization phase by testing as much as possible in the test phase. So from the several variants of the product that are in the test phase the learning agent must suggest the product variant(s) that are above a particular threshold  $\tau$  at the end of the test phase that have the highest probability of minimizing loss in the commercialization phase. A similar problem has been discussed for single best product variant identification without threshold in [Bubeck et al. \(2011\)](#).
2. *Mobile Phone Channel Allocation*: Another similar problem as above concerns channel allocation for mobile phone communications ([Audibert et al. \(2009\)](#)). Here there is a clear separation between the allocation phase and communication phase whereby in the allocation phase a learning algorithm has to explore as many channels as possible to suggest the best possible set of channel(s) that are above a particular threshold  $\tau$ . The threshold depends on the subscription level of the customer. With higher subscription the customer is allowed better channel(s) with the  $\tau$  set high. Each evaluation of a channel is noisy and the learning algorithm must come up with the best possible suggestion within a very small number of attempts.
3. *Anomaly Detection and Classification*: Thresholding bandit can also be used for anomaly detection and classification where we define a cutoff level  $\tau$  and for any samples above this cutoff gets classified as an anomaly. For further reading we point the reader to section 3 of [Locatelli et al. \(2016\)](#).

In all the above examples the MAB model performs well mainly because all of them suffer from *exploration-exploitation dilemma*. This is characterized by action-selection choice faced by the agent where it must decide whether to stay with the action yielding highest reward till now or to explore newer actions which might be more profitable in the long run. MAB's are suited for such scenarios because

1. They are easy to implement.
2. The switch between exploration and exploitation is more well defined theoretically.
3. They perform well empirically.

### 3. Problem Definition

We formally define the MAB model in this section. We will mainly focus on the stochastic MAB model in our work. In this setting, the learning algorithm is provided with a set of decisions (or arms) with reward distributions unknown to the algorithm. These decisions can be the items to show the customers, or medical treatment, or the stock options. The learning proceeds in an iterative fashion, where in each timestep, the algorithm chooses an arm and receives a stochastic reward that is drawn from a stationary distribution specific to the arm selected. There are two types of goal we will be focusing our work on:-

1. Maximizing cumulative reward over the entire time horizon
2. Minimizing probability of error in suggesting the best arm(s) at any timestep

Given the goal of maximizing the cumulative reward, the learning algorithm faces the exploration-exploitation dilemma, i.e., in each round should the algorithm select the arm which has the highest observed mean reward so far (*exploitation*), or should the algorithm choose a new arm to gain more knowledge of the true mean reward of the arms and thereby avert a sub-optimal greedy decision (*exploration*).

Formally, let  $r_i, i = 1, \dots, K$  denote the mean rewards of the  $K$  arms and  $r^* = \max_i r_i$  the optimal mean reward. The objective in the stochastic bandit problem is to maximize cumulative reward by minimizing the cumulative regret, which is defined as follows:

$$R_T = r^*T - \sum_{i \in A} r_i N_i(T),$$

where  $T$  is the number of rounds,  $N_i(T) = \sum_{m=1}^T I(I_m = i)$  is the number of times the algorithm chose arm  $i$  up to round  $T$ . The expected regret of an algorithm after  $T$  rounds can be written as

$$\mathbb{E}[R_T] = \sum_{i=1}^K \mathbb{E}[N_i(T)] \Delta_i,$$

where  $\Delta_i = r^* - r_i$  denotes the gap between the means of the optimal arm and of the  $i$ -th arm.

In the pure exploration thresholding bandit setup the goal is different than minimizing the cumulative regret. Here the learning algorithm is provided with a threshold  $\tau$  and it has to output all such arms  $i$  whose mean of reward distribution  $r_i$  is above  $\tau$  after  $T$  rounds. This is a specific instance of

combinatorial pure exploration where the learning algorithm can explore as much as possible given a fixed horizon  $T$  and not be concerned with the usual exploration-exploitation dilemma. Let  $A$  be the set of all arms. Formally we can define a set  $S_\tau = \{i \in A : r_i \geq \tau\}$  and the complementary set  $S_\tau^C = \{i \in A : r_i < \tau\}$ . Also we define  $\hat{S}_\tau = \hat{S}_\tau(T) \subset A$  and its complementary set  $\hat{S}_\tau^C$  as the recommendation of the learning algorithm after  $T$  rounds. Given such sets exists, the performance of the learning agent is measured by how much accuracy it can discriminate between  $S_\tau$  and  $S_\tau^C$  after time horizon  $T$ . The loss  $\mathcal{L}$  is defined as:-

$$\mathcal{L}(T) = I(\{S_\tau \cap \hat{S}_\tau^C \neq \emptyset\} \cup \{\hat{S}_\tau \cap S_\tau^C \neq \emptyset\})$$

The goal of the learning agent is to minimize  $\mathcal{L}(T)$ . So, the expected loss after  $T$  rounds is

$$\mathbb{E}[\mathcal{L}(T)] = \mathbb{P}(\{S_\tau \cap \hat{S}_\tau^C \neq \emptyset\} \cup \{\hat{S}_\tau \cap S_\tau^C \neq \emptyset\})$$

which we can say is the probability of making mistake, that is whether the learning agent at the end of round  $T$  rejects arms from  $S_\tau$  or accepts arms from  $S_\tau^C$  in its final recommendation. Also, we are looking at an anytime algorithm, so the knowledge of  $T$  may not be known to the learner.

#### 4. Literature Review

An early work involving a bandit setup is [Thompson \(1933\)](#), where the author deals the problem of choosing between two treatments to administer on patients who come in sequentially. Following the seminal work of [Robbins \(1952\)](#), bandit algorithms have been extensively studied in a variety of applications. From a theoretical standpoint, an asymptotic lower bound for the regret was established in [Lai and Robbins \(1985\)](#). In particular, it was shown that for any consistent allocation strategy, we have  $\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[R_T]}{\log T} \geq \sum_{\{i: r_i < r^*\}} \frac{(r^* - r_i)}{D(p_i || p^*)}$ , where  $D(p_i || p^*)$  is the Kullback-Leibler divergence between the reward densities  $p_i$  and  $p^*$ , corresponding to arms with mean  $r_i$  and  $r^*$ , respectively.

There have been several algorithms with strong regret guarantees. The foremost among them is UCB1 by [Auer et al. \(2002\)](#), which has a regret upper bound of  $O\left(\frac{K \log T}{\Delta}\right)$ , where  $\Delta = \min_{i: \Delta_i > 0} \Delta_i$ . This result is asymptotically order-optimal for the class of distributions considered. However, the worst case gap independent regret bound of UCB1 can be as bad as  $O\left(\sqrt{TK \log T}\right)$ . In [Audibert and Bubeck \(2009\)](#), the authors propose the MOSS algorithm and establish that the worst case regret of MOSS is  $O\left(\sqrt{TK}\right)$  which improves upon UCB1 by a factor of order  $\sqrt{\log T}$ . However, the gap-dependent regret of MOSS is  $O\left(\frac{K^2 \log(T \Delta^2 / K)}{\Delta}\right)$  and in certain regimes, this can be worse than even UCB1 (see [Audibert and Bubeck \(2009\)](#), [Lattimore \(2015\)](#)). The UCB-Improved algorithm, proposed in [Auer and Ortner \(2010\)](#), is a round-based algorithm<sup>1</sup> variant of UCB1 that has a gap-dependent regret bound of  $O\left(\frac{K \log T \Delta^2}{\Delta}\right)$ , which is better than that of UCB1. On the other hand, the worst case regret of UCB-Improved is  $O\left(\sqrt{TK \log K}\right)$ .

---

1. An algorithm is *round-based* if it pulls all the arms equal number of times in each round and then proceeds to eliminate one or more arms that it identifies to be sub-optimal.

In the pure exploration setup, a significant amount of research has been done on finding the best arm(s) from a set of arms. The pure exploration setup has been explored in mainly two settings:-

1. Fixed Budget setting: In this setting the learning algorithm has to suggest the best arm(s) within a fixed number of attempts that is given as an input. The objective here is to maximize the probability of returning the best arm(s). One of the foremost papers to deal with single best arm identification is [Audibert et al. \(2009\)](#) where the authors come up with the algorithm UCB and Successive Reject(SR) with simple regret guarantees. The relationship between cumulative regret and simple regret is proved in [Bubeck et al. \(2011\)](#) where the authors prove that minimizing the simple regret necessarily results in maximizing the cumulative regret. In the combinatorial fixed budget setup [Gabillon et al. \(2011\)](#) come up with Gap-E and Gap-EV algorithm which suggests the best  $m$  (given as input) arms at the end of the budget with high probability. Similarly, [Bubeck et al. \(2013\)](#) comes up with the algorithm Successive Accept Reject(SAR) which is an extension of the SR algorithm. SAR is a round based algorithm whereby at the end of round an arm is either accepted or rejected based on certain conditions till the required top  $m$  arms are suggested at the end of the budget with high probability.
2. Fixed Confidence setting: In this setting the the learning algorithm has to suggest the best arm(s) with a fixed (given as input) confidence with as less number of attempts as possible. The single best arm identification has been handled in [Even-Dar et al. \(2006\)](#) where they come up with an algorithm called Successive Elimination (SE) which comes up with an arm that is  $\epsilon$  close to the optimal arm. In the combinatorial setup recently [Kalyanakrishnan et al. \(2012\)](#) have suggested the LUCB algorithm which on termination returns  $m$  arms which are atleast  $\epsilon$  close to the true top  $m$  arms with  $1 - \delta$  probability.

Apart from these two settings some unified approach has also been suggested in [Gabillon et al. \(2012\)](#) which proposes the algorithms UGapEb and UGapEc which can work in both the above two settings. A similar combinatorial setup was also explored in [Chen et al. \(2014\)](#) where the authors come up with more similarities and dissimilarities between these two settings in a more general setup. In their work, the learning algorithm, called Combinatorial Successive Accept Reject (CSAR) is similar to SAR with a more general setup. The thresholding bandit problem is a specific instance of the pure exploration setup of [Chen et al. \(2014\)](#). In the latest work in [Locatelli et al. \(2016\)](#) the algorithm Anytime Parameter-Free Thresholding (APT) algorithm comes up with a better anytime guarantee than CSAR for the thresholding bandit problem.

## 5. Objective and Scope

The main objectives of our work.

1. UCB-Improved has several shortfalls. It conducts too much early exploration and the arm elimination is very conservative. Our proposed algorithm tries to remedy this. We also explore the idea of employing clustering techniques in MAB which, till now has only been studied in the contextual bandit setup (see [Li et al. \(2010\)](#), [Bui et al. \(2012\)](#), [Gentile et al. \(2014\)](#), [Nguyen and Lauw \(2014\)](#)). We propose an algorithm that minimizes cumulative regret in the multi-armed stochastic bandit setup using clustering and improved exploration and performs at par with the best algorithms currently available. We also show that there is a distinct advantage of using clustering in action elimination type of algorithms in MABs.

2. We also propose an algorithm that minimizes probability of error in a combinatorial stochastic bandit setup whereby the learning algorithm, provided with a threshold  $\tau$ , proposes the best set of arms such that  $r_i \geq \tau$  at every timestep using action elimination method. This problem is similar to the pure exploration setup with a major difference being that the algorithm has to suggest as many arms whose  $r_i$  is above  $\tau$  which may be more than one arm. Our proposed methodology is in contrast with APT which uses an UCB1 type of strategy to suggest the best arms above the threshold whereas we employ an UCB-Improved type of strategy with better exploration parameters.
3. In both the above scenarios we intend to explore both the theoretical guarantees and empirical performance and verify against the best algorithms available in the field.

## 6. Proposed Methodology

We propose a variant of UCB algorithm, henceforth referred to as ClusUCB, that incorporates clustering and an improved exploration scheme. ClusUCB is a round-based algorithm that starts with a partition of the arms into small clusters, each having same number of arms. The clustering is done at the start with a pre-specified number of clusters. Each round of ClusUCB involves both (individual) arm elimination as well as cluster elimination.

The clustering of arms provides two benefits. First, it creates a context where UCB-Improved like algorithm can be run in parallel on smaller sets of arms with limited exploration, which could lead to fewer pulls of sub-optimal arms with the help of more aggressive elimination of sub optimal arms. Second, the cluster elimination leads to whole sets of sub-optimal arms being simultaneously eliminated when they are found to yield poor results. These two simultaneous criteria for arm elimination can be seen as borrowing the strengths of UCB-Improved as well as other popular round based approaches.

The motivation for our work stems from the remark in Section 2.4.3 of [Bubeck et al. \(2012\)](#), where the authors conjecture that one should be able to obtain a bandit algorithm with a gap-dependent regret bound that is better than MOSS ([Audibert and Bubeck \(2009\)](#)) and UCB-Improved ([Auer and Ortner \(2010\)](#)), in particular, with a regret bound of the order  $O\left(\frac{K \log(\frac{T}{H})}{\Delta}\right)$ , where

$H = \sum_{i: \Delta_i > 0} \frac{1}{\Delta_i^2}$ . While ClusUCB does not achieve the conjectured regret bound, the theoretical analysis establishes that the gap-dependent regret of ClusUCB is always better than that of UCB-Improved and better than that of MOSS when  $\sqrt{\frac{\epsilon}{T}} \leq \Delta \leq 1$  (see Table 1). Moreover, the gap-independent bound of ClusUCB is of the same order as UCB-Improved, i.e.,  $O(\sqrt{KT \log K})$ . However, ClusUCB is not able to match the gap-independent bound of  $O(\sqrt{KT})$  for MOSS.

As mentioned in a recent work [Liu and Tsuruoka \(2016\)](#), UCB-Improved has two shortcomings:

- A significant number of pulls are spent in early exploration, since each round  $m$  of UCB-Improved involves pulling every arm an identical  $n_m = \left\lceil \frac{2 \log(T \epsilon_m^2)}{\epsilon_m^2} \right\rceil$  number of times. The quantity  $\epsilon_m$  is initialized to 1 and halved after every round.

Table 1: Gap-dependent regret bounds for different bandit algorithms

Algorithm	Upper bound
UCB1	$O\left(\frac{K \log T}{\Delta}\right)$
UCB-Improved	$O\left(\frac{K \log(T\Delta^2)}{\Delta}\right)$
MOSS	$O\left(\frac{K^2 \log(T\Delta^2/K)}{\Delta}\right)$
ClusUCB	$O\left(\frac{K \log\left(\frac{T\Delta^2}{\sqrt{\log(KT)}}\right)}{\Delta}\right)$

- In UCB-Improved, arms are eliminated conservatively, i.e, only after  $\epsilon_m < \frac{\Delta_i}{2}$ , the sub-optimal arm  $i$  is discarded with high probability. This is disadvantageous when  $K$  is large and the gaps are identical ( $r_1 = r_2 = \dots = r_{K-1} < r^*$ ) and small.

To reduce early exploration, the number  $n_m$  of times each arm is pulled per round in ClusUCB is lower than that of UCB-Improved and also that of Median-Elimination, which used  $n_m = \frac{4}{\epsilon^2} \log\left(\frac{3}{\delta}\right)$ , where  $\epsilon, \delta$  are confidence parameters. To handle the second problem mentioned above, ClusUCB partitions the larger problem into several small sub-problems using clustering and then performs local exploration aggressively to eliminate sub-optimal arms within each clusters with high probability.

As described in the pseudocode in Algorithm 1, ClusUCB begins with a initial clustering of arms that is performed by random uniform allocation. The set of clusters  $S$  thus obtained satisfies  $|S| = p$ , with individual clusters having a size that is bounded above by  $\ell = \left\lceil \frac{K}{p} \right\rceil$ . Each round of ClusUCB involves both individual arm as well as cluster elimination conditions. These elimination conditions are inspired by UCB-Improved. Notice that, unlike UCB-Improved, there is no longer a single point of reference based on which we are eliminating arms. Instead now we have as many reference points to eliminate arms as number of clusters formed.

The exploration regulatory factor  $\psi$  governing the arm and cluster elimination conditions in ClusUCB is more aggressive than that in UCB-Improved. With appropriate choice of  $\psi$  and  $\rho_a$  and  $\rho_s$  we can achieve aggressive elimination even when the gaps  $\Delta_i$  are small and  $K$  is large.

In Liu and Tsuruoka (2016), the authors recommend incorporating a factor of  $d_i$  inside the log-term of the UCB values, i.e.,  $\max\{\hat{r}_i + \sqrt{\frac{d_i \log T \epsilon_m^2}{2n_m}}\}$ . The authors there examine the following choices for  $d_i$ :  $\frac{T}{t_i}$ ,  $\frac{\sqrt{T}}{t_i}$  and  $\frac{\log T}{t_i}$ , where  $t_i$  is the number of times an arm  $i$  has been sampled.

**Algorithm 1** ClusUCB

**Input:** Number of clusters  $p$ , time horizon  $T$ , exploration parameters  $\rho_a, \rho_s$  and  $\psi$ .

**Initialization:** Set  $B_0 := A, S_0 = S$  and  $\epsilon_0 := 1$ .

Create a partition  $S_0$  of the arms at random into  $p$  clusters of size up to  $\ell = \left\lceil \frac{K}{p} \right\rceil$  each.

**for**  $m = 0, 1, \dots, \left\lfloor \frac{1}{2} \log_2 \frac{T}{e} \right\rfloor$  **do**

Pull each arm in  $B_m$  so that the total number of times it has been pulled is  $n_m = \left\lceil \frac{2 \log(\psi T \epsilon_m^2)}{\epsilon_m} \right\rceil$ .

**Arm Elimination**

For each cluster  $s_k \in S_m$ , delete arm  $i \in s_k$  from  $B_m$  if

$$\hat{r}_i + \sqrt{\frac{\rho_a \log(\psi T \epsilon_m^2)}{2n_m}} < \max_{j \in s_k} \left\{ \hat{r}_j - \sqrt{\frac{\rho_a \log(\psi T \epsilon_m^2)}{2n_m}} \right\}$$

**Cluster Elimination**

Delete cluster  $s_k \in S_m$  and remove all arms  $i \in s_k$  from  $B_m$  if

$$\begin{aligned} & \max_{i \in s_k} \left\{ \hat{r}_i + \sqrt{\frac{\rho_s \log(\psi T \epsilon_m^2)}{2n_m}} \right\} \\ & < \max_{j \in B_m} \left\{ \hat{r}_j - \sqrt{\frac{\rho_s \log(\psi T \epsilon_m^2)}{2n_m}} \right\}. \end{aligned}$$

Set  $\epsilon_{m+1} := \frac{\epsilon_m}{2}$

Set  $B_{m+1} := B_m$

Stop if  $|B_m| = 1$  and pull  $i \in B_m$  till  $T$  is reached.

**end**

Unlike [Liu and Tsuruoka \(2016\)](#), we employ cluster as well as arm elimination and establish from a theoretical analysis that the choice  $\psi = \frac{T}{\log(KT)}$  helps in achieving a better gap-dependent regret upper bound for ClusUCB as compared to UCB-Improved and MOSS.

We also approach the threshold bandit setup with a similar type of algorithm. The algorithm Augmented UCB is presented below:-

In algorithm 2, hence referred to as AugUCB, we have three input parameters,  $\rho$  which is the arm elimination parameter,  $\psi$  which is the exploration regulatory factor and the threshold  $\tau$ . The salient features of the algorithm are listed below:-

- AugUCB combines the power of UCB-Improved ([Auer and Ortner \(2010\)](#)), APT ([Locatelli et al. \(2016\)](#)) and SAR ([Gabillon et al. \(2011\)](#)). The main approach is based on UCB-Improved with modifications suited for the thresholding bandit problem. The active set  $B_0$  is initialized with all the arms from  $A$ .



---

**Algorithm 2** AugmentedUCB
 

---

**Input:** Time horizon  $T$ , exploration parameters  $\rho$  and  $\psi$ , threshold  $\tau$ .

**Initialization:** Set  $B_0 := A$ ,  $M = \left\lfloor \frac{1}{2} \log_2 \frac{T}{e} \right\rfloor$ ,  $m := 0$ ,  $\epsilon_0 := 1$ ,  $\ell_0 = \left\lceil \frac{2 \log(\psi T \epsilon_0^2)}{\epsilon_0} \right\rceil$  and  $N_0 = K * \ell_0$ .

Pull each arm once

**for**  $t = K + 1, \dots, T$  **do**

Pull arm  $i$  in  $B_m$  such that  $\min_{i \in B_m} \left\{ |\hat{r}_i - \tau| - \sqrt{\frac{\rho \log(\psi T \epsilon_m^2)}{2n_i}} \right\}$ , where  $n_i$  is the number of times the arm  $i$  has been pulled.

$t := t + 1$

**Arm Elimination**

For each arm  $i \in B_m$ , remove arm  $i$  from  $B_m$  if

$$\hat{r}_i + \sqrt{\frac{\rho \log(\psi T \epsilon_m^2)}{2n_i}} < \tau - \sqrt{\frac{\rho \log(\psi T \epsilon_m^2)}{2n_i}}$$

For each arm  $i \in B_m$ , remove arm  $i$  from  $B_m$  if

$$\hat{r}_i - \sqrt{\frac{\rho \log(\psi T \epsilon_m^2)}{2n_i}} > \tau + \sqrt{\frac{\rho \log(\psi T \epsilon_m^2)}{2n_i}}$$

**if**  $t \geq N_m$  **and**  $m \leq M$  **then**

**Reset Parameters**

$$\epsilon_{m+1} := \frac{\epsilon_m}{2}$$

$$B_{m+1} := B_m$$

$$\ell_{m+1} = \left\lceil \frac{2 \log(\psi T \epsilon_{m+1}^2)}{\epsilon_{m+1}} \right\rceil$$

$$N_{m+1} := t + |B_{m+1}| \ell_{m+1}$$

$$m := m + 1$$

**end**

**end**

Output  $\hat{S}_\tau = \{i : \hat{r}_i \geq \tau\}$ .

---

- We divide the entire budget  $T$  into rounds/phases as like UCB-Improved, SAR and CSAR. The choice of  $M$  comes from UCB-Improved which necessarily entails that the  $\epsilon_m \geq \sqrt{\frac{\epsilon}{T}}$ . So,  $M$  is the total number of rounds and is the same as UCB-Improved. After the end of each such round  $m$  we eliminate arm(s) from active set  $B_m$  and update parameters.
- As suggested by [Liu and Tsuruoka \(2016\)](#) to make AugUCB an anytime algorithm and to overcome too much early exploration, we no longer pull all the arms equal number of times

in each round but pull the arm that minimizes

$$\left\{ |\hat{r}_i - \tau| - \sqrt{\frac{\rho \log(\psi T \epsilon_m^2)}{2n_i}} \right\}$$

in the active set  $B_m$ .

- $\min_{i \in B_m} \left\{ |\hat{r}_i - \tau| - \sqrt{\frac{\rho \log(\psi T \epsilon_m^2)}{2n_i}} \right\}$  condition actually makes it possible to pull the arms closer to the threshold  $\tau$ . This is a strategy used by APT.
- This also gets rid of the excessive initial exploration employed by UCB-Improved and yet with suitable choice of  $\rho$  and  $\psi$  we can fine tune the exploration.
- The arm elimination condition simply removes arm(s) if the algorithm is sufficiently sure that the mean of the arms are very high or very low about the threshold. This although is a tactic similar to SAR or CSAR, but here at any round, an arbitrary number of arms can be accepted or rejected thereby improving upon SAR and CSAR which accepts/rejects one arm in every round.
- At the end of the budget  $T$  the algorithm outputs all the arms whose estimated average payoff  $\hat{r}_i$  is above the threshold  $\tau$  thereby making this an anytime algorithm whereby we need not finish every round.
- The arm elimination condition(s) helps in re-allocating the remaining budget/pulls among the surviving arms. Those among the remaining arms are pulled which are closer to the threshold. Arms lying far to the either side of the threshold are eliminated from the active set  $B_m$ .

## 7. Work Done

1. On four synthetic setups with small gaps, we observe empirically that ClusUCB outperforms UCB-Improved (Auer and Ortner (2010)) and MOSS (Audibert and Bubeck (2009)) as well as other popular stochastic bandit algorithms such as DMED (Honda and Takemura (2010)), UCB-V (Audibert et al. (2009)), Median Elimination (Even-Dar et al. (2006)), Thompson Sampling (Agrawal and Goyal (2011)) and KL-UCB (Garivier and Cappé (2011)). The results can be seen in Appendix 9 section. We also prove strong theoretical guarantees for Clustered UCB, and the manuscript is currently under review in AISTATS 2017.
2. We have started work on AugUCB algorithm and are currently trying to prove it's simple regret bounds. Empirically the algorithm has performed well and the result is shown in Appendix 9.

## 8. Contribution

From a theoretical viewpoint, we conclude that the gap-dependent regret bound of ClusUCB is lower than that of MOSS and UCB-Improved. From the numerical experiments on settings with small gaps between optimal and sub-optimal mean rewards, we observed that ClusUCB outperforms several popular bandit algorithms. While we exhibited better regret bounds for ClusUCB, it

would be interesting future research to improve the theoretical analysis of ClusUCB to achieve the gap-independent regret bound of MOSS and possibly also the gap-dependent bound conjectured in Section 2.4.3 of [Bubeck and Cesa-Bianchi \(2012\)](#).

## References

- Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit problem. *arXiv preprint arXiv:1111.1797*, 2011.
- Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, pages 217–226, 2009.
- Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–1902, 2009.
- Peter Auer and Ronald Ortner. Ucb revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Periodica Mathematica Hungarica*, 61(1-2):55–65, 2010.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- Dimitri P Bertsekas and John N Tsitsiklis. Neuro-dynamic programming (optimization and neural computation series, 3). *Athena Scientific*, 7:15–23, 1996.
- Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *arXiv preprint arXiv:1204.5721*, 2012.
- Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412(19):1832–1852, 2011.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *arXiv preprint arXiv:1209.1727*, 2012.
- Sébastien Bubeck, Tengyao Wang, and Nitin Viswanathan. Multiple identifications in multi-armed bandits. In *ICML (1)*, pages 258–265, 2013.
- Loc Bui, Ramesh Johari, and Shie Mannor. Clustered bandits. *arXiv preprint arXiv:1206.4169*, 2012.
- Olivier Cappe, Aurelien Garivier, and Emilie Kaufmann. pymabandits, 2012. <http://mloss.org/software/view/415/>.
- Shouyuan Chen, Tian Lin, Irwin King, Michael R Lyu, and Wei Chen. Combinatorial pure exploration of multi-armed bandits. In *Advances in Neural Information Processing Systems*, pages 379–387, 2014.
- Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *The Journal of Machine Learning Research*, 7:1079–1105, 2006.

- Victor Gabillon, Mohammad Ghavamzadeh, Alessandro Lazaric, and Sébastien Bubeck. Multi-bandit best arm identification. In *Advances in Neural Information Processing Systems*, pages 2222–2230, 2011.
- Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In *Advances in Neural Information Processing Systems*, pages 3212–3220, 2012.
- Aurélien Garivier and Olivier Cappé. The kl-ucb algorithm for bounded stochastic bandits and beyond. *arXiv preprint arXiv:1102.2490*, 2011.
- Claudio Gentile, Shuai Li, and Giovanni Zappella. Online clustering of bandits. In *ICML*, pages 757–765, 2014.
- Mohammad Ghavamzadeh, Shie Mannor, Joelle Pineau, Aviv Tamar, et al. *Bayesian reinforcement learning: a survey*. World Scientific, 2015.
- Junya Honda and Akimichi Takemura. An asymptotically optimal bandit algorithm for bounded support models. In *COLT*, pages 67–79. Citeseer, 2010.
- Shivaram Kalyanakrishnan, Ambuj Tewari, Peter Auer, and Peter Stone. Pac subset selection in stochastic multi-armed bandits. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 655–662, 2012.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- Tor Lattimore. Optimally confident ucb: Improved regret for finite-armed bandits. *arXiv preprint arXiv:1507.07880*, 2015.
- Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670. ACM, 2010.
- Yun-Ching Liu and Yoshimasa Tsuruoka. Modification of improved upper confidence bounds for regulating exploration in monte-carlo tree search. *Theoretical Computer Science*, 2016.
- Andrea Locatelli, Maurilio Gutzeit, and Alexandra Carpentier. An optimal algorithm for the thresholding bandit problem. *arXiv preprint arXiv:1605.08671*, 2016.
- Trong T Nguyen and Hady W Lauw. Dynamic clustering of contextual multi-armed bandits. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pages 1959–1962. ACM, 2014.
- Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*, pages 169–177. Springer, 1952.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
- William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, pages 285–294, 1933.

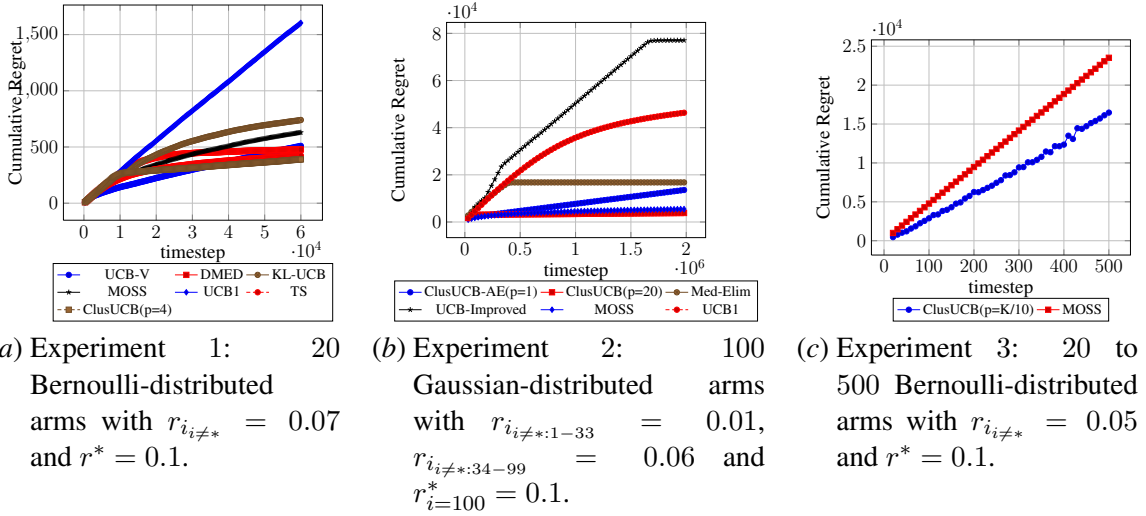


Figure 1: Cumulative regret for various bandit algorithms on three stochastic K-armed bandit environments.

## 9. Appendix

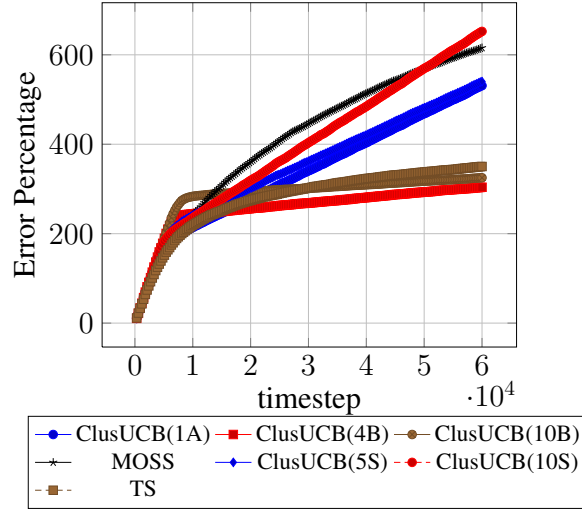
In this appendix, we show the different experiments carried on various test-beds. Appendix A contains experiments on ClusUCB while Appendix B contains experiments on AugUCB.

### Appendix A

For the purpose of performance comparison using cumulative regret as the metric, we implement the following algorithms: KL-UCB [Garivier and Cappé \(2011\)](#), DMED [Honda and Takemura \(2010\)](#), MOSS [Audibert and Bubeck \(2009\)](#), UCB1 [Auer et al. \(2002\)](#), UCB-Improved [Auer and Ortner \(2010\)](#), Median Elimination [Even-Dar et al. \(2006\)](#), Thompson Sampling (TS) [Agrawal and Goyal \(2011\)](#) and UCB-V [Audibert et al. \(2009\)](#)<sup>2</sup>. The parameters of ClusUCB algorithm for all the experiments are set as follows:  $\psi = \log T$ ,  $\rho_s = \frac{1}{2^{2m+1}}$  and  $\rho_a = \frac{1}{2^{4m+1}}$  for  $m = 0, 1, \dots, \lfloor \frac{1}{2} \log_2 \frac{T}{e} \rfloor$ . So in all the experiments  $\rho_a$  and  $\rho_s$  are initialized to 1 and then reduced after every round. By this definition of  $\rho_a, \rho_s$  we have made sure that their value always remain bounded  $\in (0, 1]$ . When  $K$  is large and  $p$  is small it is advantageous to run  $\rho_a < \rho_s$  because this will aggressively eliminate arms within each cluster while full cluster elimination will be more conservative.

The first experiment is conducted over a testbed of 20 arms for the test-cases involving Bernoulli reward distribution with expected rewards of the arms  $r_{i \neq *} = 0.07$  and  $r^* = 0.1$ . These type of cases are frequently encountered in web-advertising domain. The horizon  $T$  is set to 60000. After limited exploratory experimentation the number of clusters  $p$  for ClusUCB is set to 4. The regret is averaged over 100 independent runs and is shown in Figure 1(a). ClusUCB, MOSS, UCB1, UCB-V, KL-UCB, TS and DMED are run in this experimental setup and we observe that ClusUCB performs better than all the aforementioned algorithms. Because of the short horizon  $T$ , we do not implement UCB-Improved and Median Elimination on this test-case. We also observe that in this case the

2. The implementation for KL-UCB and DMED were taken from [Cappé et al. \(2012\)](#)



(a)

Figure 2: Experiment 4: Cumulative regret for ClusUCB variants: 1, 4, 5, 10 correspond to number  $p$  of clusters and A, S, B correspond to having only arm elimination, only cluster elimination and having both in Algorithm 1. For using just cluster elimination in ClusUCB, we stop when we left with only one cluster and play the max payoff arm of that cluster for the remaining time horizon.

cumulative regret of ClusUCB and TS are very similar to each other with ClusUCB being slightly better.

The second experiment is conducted over a testbed of 100 arms involving Gaussian reward distribution with expected rewards of the arms  $r_{i \neq *:1-33} = 0.01$ ,  $r_{i \neq *:34-99} = 0.06$  and  $r_{i=100}^* = 0.1$  with variance set at  $\sigma = 0.3, \forall i \in A$ . The horizon  $T$  is set for a large duration of  $2 \times 10^6$  and the number of clusters  $p = 20$ . The regret is averaged over 100 independent runs and is shown in Figure 1(b). In this case, in addition to ClusUCB, we also show the performance of no-clustering version of ClusUCB algorithm (i.e.,  $p = 1$ ). From the results in Figure 1(b), we observe that ClusUCB with  $p = 20$  outperforms ClusUCB with  $p = 1$  as well as MOSS, UCB1, UCB-Improved and Median-Elimination( $\epsilon = 0.03, \delta = 0.1$ ). We also observed that the ClusUCB variant that uses only the arm elimination condition in Algorithm 1 performs worse than the variant that employs both cluster and arm elimination conditions. Also the performance of UCB- Improved is poor in comparison to other algorithms, which is probably because of pulls wasted in initial exploration.

The third experiment is conducted over a testbed of 20–500 (interval of 10) arms with Bernoulli reward distribution, where the expected rewards of the arms are  $r_{i \neq *} = 0.05$  and  $r^* = 0.1$ . The horizon  $T$  is set to  $10^5 + K^2 \times 10^4$  and the number of arms are increased from  $K = 20$  to 500. The proposed algorithm ClusUCB is run with  $p = K/10$ . The regret is averaged over 500 independent runs and is shown in Figure 1(c). We report the performance of MOSS and ClusUCB only over this setup. From the results in Figure 1(c), it is evident that the growth of regret for ClusUCB is lower than that of MOSS.

The fourth experiment is performed over a testbed having 20 Bernoulli-distributed arms with  $r_{i \neq *} = 0.06, \forall i \in A$  and  $r^* = 0.1$ . In Figure 2, we report the results with  $T = 60000$  averaged over 100 independent runs for ClusUCB with  $p = \{1, 4, 10\}$ . Also, in this experiment we take

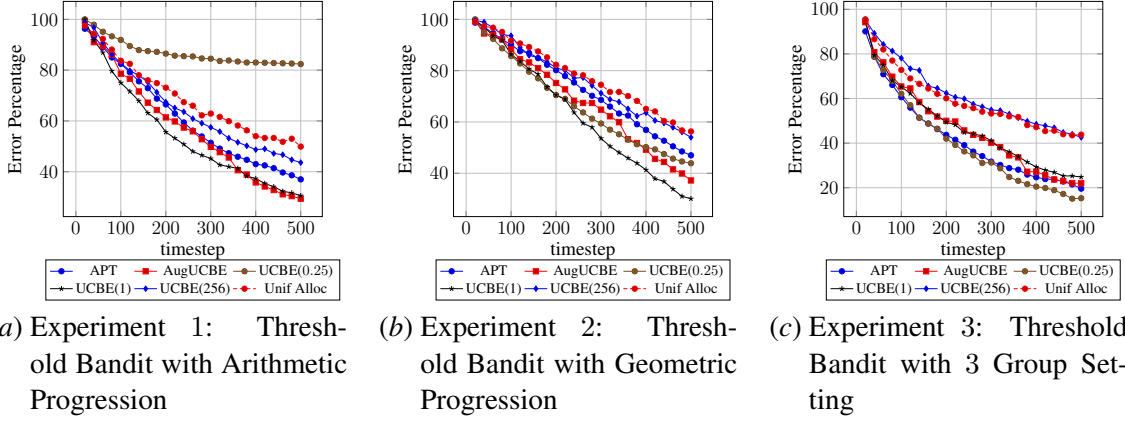


Figure 3: Experiments with thresholding bandit

$\psi = 1$  to make sure that the exploration is low. Under these conditions, exploration is low and the elimination parameters are decreased very fast. This leads to a greater number of errors committed by ClusUCB-AE that is ClusUCB( $p = 1$ ). Two other benchmark algorithms are considered here, MOSS which is one of our main competitors and TS which has performed near equivalently in experiment 1 (see Fig 1(a)). In this case we see that, since  $\rho_a$  is decreased very fast, the optimal arm  $a^*$  gets eliminated most of the time for no clustering  $p = 1$ . While a balance of  $p$ ,  $\rho_a$  and  $\rho_s$  gives a much better result. ClusUCB with  $p = 4$  and 10 perform better than MOSS, while  $p = 1$  with just arm elimination does not converge and  $p = 5, 10$  with just Cluster elimination and no arm elimination also does not converge, implying both cluster and arm elimination are necessary. This also concurs with the theoretical observations. Note also that ClusUCB( $p = 4$ ) outperforms TS in this setup.

## Appendix B

In this section we compare the empirical performance of AugUCB against APT, Uniform Allocation and UCBE algorithm. The threshold  $\tau$  is set at 0.5 for all experiments. Each algorithm is run independently a 1000 times for 500 timesteps and the output set of arms suggested by the algorithms at every timestep is recorded. The output is considered erroneous if the correct set of arms is not  $i = \{6, 7, 8, 9, 10\}$  (true for all the experiments). The error percentage over 1000 runs is plotted against 500 timesteps. For the uniform allocation algorithm, for each  $t = 1, 2, \dots, T$  the arms are sampled uniformly. For UCBE algorithm (Audibert et al. (2009)) which was built for single best arm identification, we modify it according to Locatelli et al. (2016) to suit the goal of finding arms above the threshold  $\tau$ . So the exploration parameter  $a$  in UCBE is set to  $a_i = 4^i \frac{T-K}{H}$  for  $i \in \{-1, 0, 4\}$  and  $H = \sum_{i=1}^K \frac{1}{\Delta_i^2}$  is defined as the problem complexity. Then for each timestep  $t = 1, 2, \dots, T$  we pull the arm that maximizes  $\{|\hat{r}_i - \tau| - \sqrt{\frac{a_i}{n_i}}\}$ , where  $n_i$  is the number of times the arm  $i$  is pulled till  $t - 1$  timestep. Also, APT is run with  $\epsilon = 0$ , which denotes the precision with which the algorithm suggests the best set of arms. So when  $\epsilon$  is set to 0 APT has to suggest the exact set of arms above the threshold. For AugUCB we take  $\psi = K^2 T$  and we initialize  $\rho = \frac{1}{2^m}$  for  $m = 0, 1, 2, \dots, \gamma$ .

The high value of  $\psi$  helps in improved exploration whereas we decrease  $\rho$  sufficiently after every round to facilitate arm elimination.

The first experiment is conducted on a testbed of 10 arms involving Bernoulli reward distribution with expected rewards of the arms  $r_{1:4} = 0.2 + (0 : 3) * 0.05$ ,  $r_5 = 0.45$ ,  $r_6 = 0.55$  and  $r_{7:10} = 0.65 + (0 : 3) * 0.05$ . The means are set as arithmetic progression. In this experiment we see that AugUCB performs better than all the other algorithms mentioned. Only UCBE(1) catches up with AugUCB and that is because it has access to the exact problem complexity. The result is shown in Figure 3(a).

The second experiment is conducted on a testbed of 10 arms with the means set as Geometric Progression. The testbed involves Bernoulli reward distribution with expected rewards of the arms as  $r_{1:4} = 0.4 - (0.2)^{1:4}$ ,  $r_5 = 0.45$ ,  $r_6 = 0.55$  and  $r_{7:10} = 0.6 + (0.2)^{5-(1:4)}$ . AugUCB, APT, Uniform Allocation and UCBE with the same settings as experiment 1 are run on this testbed. The result is shown in Figure 3(b). Here, we see that AugUCB beats APT with only UCBE(1) performing at par with AugUCB.

The third experiment is conducted on a testbed of 10 arms with the means divided into 3 groups. Again the testbed involves Bernoulli reward distribution with expected rewards of the arms as  $r_{1:3} = 0.1$ ,  $r_{4:7} = \{0.35, 0.45, 0.55, 0.65\}$  and  $r_{8:10} = 0.9$ . AugUCB, APT, Uniform Allocation and UCBE with the same settings as experiment 1 are run on this testbed. The result is shown in Figure 3(c). Here, also we see that AugUCB beats APT.