

# **A study on online sequential learning using Bandits**

*A THESIS*

*submitted by*

**SUBHOJYOTI MUKHERJEE**

*for the award of the degree*

*of*

**MASTER OF SCIENCE**

(by Research)



**DEPARTMENT OF COMPUTER SCIENCE AND  
ENGINEERING  
INDIAN INSTITUTE OF TECHNOLOGY MADRAS**

**December 2017**

# THESIS CERTIFICATE

This is to certify that the thesis titled **A study on online sequential learning using Bandits**, submitted by **Subhojyoti Mukherjee**, to the Indian Institute of Technology, Madras, for the award of the degree of **Master of Science (Research)**, is a bona fide record of the research work done by him under our supervision. The contents of this thesis, in full or in parts, have not been submitted to any other Institute or University for the award of any degree or diploma.

**Dr. Balaraman Ravindran**  
Research Guide  
Associate Professor  
Dept. of Computer Science  
IIT-Madras, 600 036

**Dr. Nandan Sudarsanam**  
Research Co-Guide  
Assistant Professor  
Dept. of Management Studies  
IIT-Madras, 600 036

Place: Chennai

Date: 22nd December 2017

# **ACKNOWLEDGEMENTS**

Thanks to all those who made  $\text{T}_{\text{E}}\text{X}$  and  $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$  what it is today.

# ABSTRACT

**KEYWORDS:** Reinforcement Learning, Bandits, UCB.

The thesis studies the following topics in the area of Reinforcement Learning: Multi-armed Bandits, Multi-armed bandits in stationary distribution with the goal of cumulative regret minimization, Thresholding bandits in pure exploration setting, and analysis of bandit theory in piece-wise stationary distributions. The common underlying theme is the study of bandit theory and its application in various types of environments. We start with a general introduction to Multi-armed bandits, its connection to the wider reinforcement learning theory and then we discuss the various types of bandits available in the literature. The subsequent chapter deals with the classic multi-armed bandit problem in stationary distribution, one of the first setting studied by the bandit community and which successively gave rise to several new directions in bandit theory. We propose a novel algorithm in this setting and compare both theoretically and empirically its performance against the available algorithms in this setting. In the next chapter, we move onto a very specific type of bandit setup called the thresholding bandit problem and discuss extensively on its usage, available state-of-the-art algorithms on this setting and our solution to the problem. We give theoretical guarantees on the expected loss of our algorithm and also analyze its performance against state-of-the-art algorithms in numerical simulations in multiple synthetic environments. The final chapter deals with the notion of piece-wise stationary distribution and how available bandit algorithms can be modified to perform well in this setting. We propose a set of algorithms for this setting with the goal of minimizing cumulative regret which uses various techniques ranging from changepoint detection mechanism to aggregation of experts.

# Contents

<b>ACKNOWLEDGEMENTS</b>	<b>i</b>
<b>ABSTRACT</b>	<b>ii</b>
<b>LIST OF TABLES</b>	<b>vi</b>
<b>LIST OF FIGURES</b>	<b>vii</b>
<b>ABBREVIATIONS</b>	<b>viii</b>
<b>NOTATION</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Motivation . . . . .	2
1.3 Types of Information Feedback . . . . .	3
1.3.1 Full information feedback . . . . .	4
1.3.2 Partial information feedback . . . . .	4
1.3.3 Bandit feedback . . . . .	4
1.4 Different types of Bandits . . . . .	4
1.4.1 Types of Bandits based on Environment . . . . .	5
1.4.2 Types of Bandits based on goal . . . . .	6
1.4.3 Collaborative Bandits . . . . .	6
<b>2 Stochastic Multi-armed Bandits</b>	<b>7</b>
2.1 Introduction . . . . .	7
2.2 Notations and assumptions . . . . .	8
2.3 Problem Definition . . . . .	8
2.4 Motivation . . . . .	10
2.5 Related Work in SMAB . . . . .	10

2.5.1	Lower bound in SMAB . . . . .	10
2.5.2	The Upper Confidence Bound approach . . . . .	10
2.5.3	Bayesian Approach . . . . .	14
2.5.4	Information Theoretic approach . . . . .	14
2.5.5	Discussion on the various confidence intervals . . . . .	15
2.6	Conclusion . . . . .	16
<b>3</b>	<b>Efficient UCB Variance: An almost optimal algorithm in SMAB setting</b>	<b>17</b>
3.1	Introduction . . . . .	17
3.2	Our Contributions . . . . .	18
3.3	Algorithm: Efficient UCB Variance . . . . .	19
3.4	Main Results . . . . .	21
3.5	Proofs . . . . .	23
3.6	Experiments . . . . .	28
3.7	Conclusion and Future Works . . . . .	32
<b>4</b>	<b>Thresholding Bandits</b>	<b>33</b>
4.1	Introduction . . . . .	33
4.2	Notations . . . . .	34
4.3	Problem Definition . . . . .	34
4.4	Motivation . . . . .	35
4.5	Related Work in Pure Exploration . . . . .	35
4.5.1	Fixed Budget setting . . . . .	36
4.5.2	Fixed Confidence setting . . . . .	37
4.5.3	Unified Setting . . . . .	37
4.6	Related Work in Thresholding Bandits . . . . .	37
4.7	Conclusion . . . . .	38
<b>5</b>	<b>Augmented UCB for TBP</b>	<b>39</b>
5.1	Introduction . . . . .	39
5.2	Our Contribution . . . . .	39
5.3	Augmented-UCB Algorithm . . . . .	41
5.4	Theoretical Results . . . . .	43

5.5	Numerical Experiments . . . . .	47
5.6	Conclusion and Future Works . . . . .	51
<b>A</b>	<b>Appendix on Concentration Inequalities</b>	<b>52</b>
A.0.1	Sample space . . . . .	54
A.0.2	Events . . . . .	54
A.0.3	Sigma-algebra . . . . .	54
A.0.4	Measure . . . . .	54
A.0.5	The triplet . . . . .	54
A.0.6	Filtration . . . . .	54
A.1	Martingale . . . . .	54
A.1.1	Super-martingale . . . . .	54
A.1.2	Sub-martingale . . . . .	54
A.2	Convergence theorems . . . . .	54
A.2.1	Monotone convergence theorem . . . . .	54
A.2.2	Dominated convergence theorem . . . . .	54
A.2.3	Fatou's Lemma . . . . .	54
A.3	Sub-Gaussian distribution . . . . .	54
A.4	Concentration Inequalities . . . . .	55
A.4.1	Markov's inequality . . . . .	55
A.4.2	Chernoff-Hoeffding Bound . . . . .	56
A.4.3	Empirical Bernstein inequality . . . . .	56
<b>B</b>	<b>Appendix for EUCBV</b>	<b>58</b>
B.0.4	Proof of Lemma 1 . . . . .	58
B.0.5	Proof of Lemma 2 . . . . .	58
B.0.6	Proof of Lemma 3 . . . . .	59
B.0.7	Proof of Lemma 4 . . . . .	61
B.0.8	Proof of Lemma 5 . . . . .	62
B.0.9	Proof of Lemma 6 . . . . .	63
B.0.10	Proof of Corollary 1 . . . . .	64
	<b>LIST OF PAPERS BASED ON THESIS</b>	<b>69</b>

## List of Tables

2.1	Confidence interval of different algorithms . . . . .	15
3.1	Regret upper bound of different algorithms . . . . .	19
5.1	AugUCB vs. State of the art . . . . .	40



## List of Figures

3.1	A comparison of the cumulative regret incurred by the various bandit algorithms. . . . .	29
3.2	Further Experiments with EUCEV . . . . .	30
5.1	Performances of the various TBP algorithms in terms of error percentage vs. time-step, for six different experimental scenarios. . . . .	50

# **ABBREVIATIONS**

<b>IITM</b>	Indian Institute of Technology, Madras
<b>RTFM</b>	Read the Fine Manual

## NOTATION

$r$	Radius, $m$
$\alpha$	Angle of thesis in degrees
$\beta$	Flight path in degrees

# Chapter 1

## Introduction

### 1.1 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.

## 1.2 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.

### 1.3 Types of Information Feedback

In an online sequential setting, the feedback that the learner receives from the environment can be characterized into three broad categories, full information feedback, partial information feedback and bandit feedback.

To illustrate the different types of feedback we will take help of the following example. Let a learner be given a set of actions  $i \in \mathcal{A}$  such that  $|\mathcal{A}| = K$ . Let, the environment be such that each action has a probability distribution  $D_i$  attached to it which is fixed throughout the time horizon  $T$ . The learning proceeds as follows:-

---

**Algorithm 1** An online sequential game

---

**Input:** Time horizon  $T$ ,  $K$  number of arms with unknown parameters of reward distribution

**for** each timestep  $t = 1, 2, \dots, T$  **do**

    The environment chooses a reward  $r_{i,t}, \forall i \in \mathcal{A}$ .

    The learner chooses  $m$  actions such that  $m < K$ , where  $\mathcal{A}$  is the set of arms and  $|\mathcal{A}| = K$ .

    The learner observes the reward  $R_{m,t} = F(r_{i,t})$ .

**end for**

---

### 1.3.1 Full information feedback

In full information feedback, when a learner selects  $m$  actions then the environment reveals the rewards of all the actions  $i \in \mathcal{A}$ . Hence, in this form of feedback the learner observes  $R_{m,t} = \{r_{i,t}, \forall i \in \mathcal{A}\}$ .

### 1.3.2 Partial information feedback

In partial information feedback, when a learner selects  $m$  actions then the environment reveals the rewards of only those  $m$  actions for  $m \in \mathcal{A}$ . Hence, in this form of feedback the learner observes  $R_{m,t} = \{r_{m,t}, \forall m \in \mathcal{A}\}$ . This is also sometimes called the semi-bandit feedback.

### 1.3.3 Bandit feedback

In bandit feedback, when a learner selects  $m$  actions then the environment reveals a cumulative reward of those  $m$  actions for  $m \in \mathcal{A}$ . Hence, in this form of feedback the learner observes  $R_{m,t} = \sum_{q=1}^m r_{q,t}$ . Note, that when  $m = 1$ , then the learner observes the reward of only that action that it has chosen out of  $K$  actions.

## 1.4 Different types of Bandits

In this section we discuss on the various types of bandits that are available in literature.

### 1.4.1 Types of Bandits based on Environment

#### Stochastic Bandits

In stochastic bandits, the distribution associated with each of the arms remains fixed throughout the time horizon  $T$ . Some of the notable papers associated with this type of setup are Robbins (1952), Lai and Robbins (1985), Agrawal (1995), Auer *et al.* (2002a), Auer and Ortner (2010), Audibert and Bubeck (2009), Lattimore (2015), etc. Chapter 2 and Chapter 3 is based on this setup where we discuss extensively on the latest state-of-the-art algorithms.

#### Non-stochastic Bandits

In non-stochastic setting the distribution associated with each arm varies over the duration of the play. Two notable examples of this are:-

- **Adversarial bandits:** In adversarial bandits, an adversary decides the payoff for each arm before the learner selects an arm. This adversary may or may not be oblivious to the learning algorithm employed by the learner. In each of these cases a different guarantee on the performance of the learner can be arrived at. Some of the important papers in this setting are Auer *et al.* (2002b), Auer *et al.* (2002b), Auer (2002), Kocák *et al.* (2014).
- **Piece-wise stationary:** Another setup under this setting can be the piece-wise stationary setting. In this setting, the distribution associated with each arm is not fixed throughout the time horizon and changes either arbitrarily at particular changepoints, or changes at a fixed period. The distribution associated with each arm then remains fixed till the next changepoint is encountered. Chapter (to be added) is based on this.

#### Contextual Bandits

The idea of clustering has been extensively studied in the contextual bandit setup, an extension of the MAB where side information or features are attached to each arm. The clustering is done over the features representing the arms to capture the complexity of the problem better when a large-number of arms are involved. Typical examples of this setting are in web-advertising domain, news article selection, etc. Some notable papers



available for this setting are Auer (2002), Langford and Zhang (2008), Li *et al.* (2010), Beygelzimer *et al.* (2011), Slivkins (2014), etc.

## **1.4.2 Types of Bandits based on goal**

**Cumulative regret minimization**

**Simple regret minimization**

**External Regret minimization**

## **1.4.3 Collaborative Bandits**

Distributed bandits are specific setup of MAB where a network of bandits collaborate with each other to identify the optimal arm(s) (see Awerbuch and Kleinberg (2008); Liu and Zhao (2010); Szörényi *et al.* (2013); Hillel *et al.* (2013)). In our setting we can assign each of the  $p$  clusters to individual bandits and at the end of each round they can share information synchronously to identify the optimal arm. This naturally results in a speedup of operation and helps in identifying the best arm faster. The clustering in this case is typically done over the feature space Bui *et al.* (2012), Cesa-Bianchi *et al.* (2013), Gentile *et al.* (2014).

# Chapter 2

## Stochastic Multi-armed Bandits

### 2.1 Introduction

In this chapter, we deal with the stochastic multi-armed bandit (SMAB) setting. In its classical form, stochastic MABs represent a sequential learning problem where a learner is exposed to a finite set of actions (or arms) and needs to choose one of the actions at each timestep. After choosing (or pulling) an arm the learner receives a reward, which is conceptualized as an independent random draw from stationary distribution associated with the selected arm. Also, note that in SMAB, the distribution associated with each arm is fixed throughout the entire duration of the horizon denoted by  $T$ .

---

**Algorithm 2** SMAB formulation

---

**Input:** Time horizon  $T$ ,  $K$  number of arms with unknown parameters of reward distribution

**for** each timestep  $t = 1, 2, \dots, T$  **do**

    The learner chooses an arm  $i \in \mathcal{A}$ , where  $\mathcal{A}$  is the set of arms and  $|\mathcal{A}| = K$ .

    The learner observes the reward  $X_{i,t} \sim^{i.i.d} D_i$  where,  $D_i$  is the distribution associated with the arm  $i$ .

**end for**

---

The rest of the chapter is organized as follows. We specify all the notations and assumptions in section 2.2. Then we define the problem statement for the SMAB setting in section 2.3. In the next section 2.4 we discuss the motivations behind the SMAB setting. In section 2.5 we discuss extensively on the various state-of-the-art algorithms available for the SMAB setting. Finally, we draw our conclusions in section 2.6.

## 2.2 Notations and assumptions

**Assumption 1** *In the considered SMAB setting we assume the optimal arm to be unique and it is denoted by  $*$ .*

**Assumption 2** *We assume the rewards of all arms are bounded in  $[0, 1]$ .*

**Notations:** The mean of the reward distribution  $D_i$  associated with an arm  $i$  is denoted by  $r_i$  whereas the mean of the reward distribution of the optimal arm  $*$  is denoted by  $r^*$  such that  $r_i < r^*, \forall i \in \mathcal{A}$ , where  $\mathcal{A}$  is the set of arms such that  $|\mathcal{A}| = K$ . We denote the individual arms labeled  $i$ , where  $i = 1, \dots, K$ . We denote the sample mean of the rewards for an arm  $i$  at time instant  $t$  by  $\hat{r}_i(t) = \frac{1}{z_i(t)} \sum_{\ell=1}^{z_i(t)} X_{i,\ell}$ , where  $X_{i,\ell}$  is the reward sample received when arm  $i$  is pulled for the  $\ell$ -th time, and  $z_i(t)$  is the number of times arm  $i$  has been pulled until timestep  $t$ . We denote the true variance of an arm by  $\sigma_i^2$  while  $\hat{v}_i(t)$  is the estimated variance, i.e.,  $\hat{v}_i(t) = \frac{1}{z_i(t)} \sum_{\ell=1}^{z_i(t)} (X_{i,\ell} - \hat{r}_i)^2$ . Whenever there is no ambiguity about the underlying time index  $t$ , for simplicity we neglect  $t$  from the notations and simply use  $\hat{r}_i, \hat{v}_i$ , and  $z_i$  to denote the respective quantities. Also,  $\Delta$  denotes the minimum gap such that  $\Delta = \min_{i \in \mathcal{A}} \{\Delta_i\}$ .

## 2.3 Problem Definition

With the formulation of SMAB stated in algorithm 2, the learner seeks to identify the optimal arm as quickly as possible to maximize its rewards. In the pursuit of this, the learner faces the task of balancing exploitation and exploration. In other words, should the learner pull the arm which currently has the best-known estimates (exploit) or explores arms more thoroughly to ensure that a correct decision is being made. This is termed as the *exploration-exploitation dilemma*, one of the fundamental challenges of reinforcement learning as discussed in chapter 1.

The objective of the learner in the SMAB setting is to maximize his rewards or in other words, to minimize the cumulative regret, which is defined as follows:

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

where  $T$  is the number of timesteps, and  $z_i(T)$  is the number of times the algorithm has chosen arm  $i$  up to timestep  $T$ . The expected regret of an algorithm after  $T$  timesteps can be written as,

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

where  $\Delta_i = r^* - r_i$  is the gap between the means of the optimal arm and the  $i$ -th arm. In the theoretical analysis of each algorithm, we try to obtain bounds on this cumulative regret. These bounds can be both asymptotic or for a finite horizon. Again, these regret bounds can be either gap-dependent or gap-independent bounds.

1. **Asymptotic regret bounds:** These type of regret bounds are valid for a large horizon  $T$  tending to infinity. In other words, if the guarantees of these bounds to be held true then an infinite number of samples needs to be collected.
2. **Finite horizon regret bounds:** These type of regret bounds are valid for a finite horizon when a limited number of samples are allowed to be collected. Note, that the knowledge of horizon may or may not be known to the learner.
3. **Gap-Dependent regret bounds:** In gap-dependent or problem dependent regret bounds the regret is obtained as a measure of the gap  $\Delta_i = r^* - r_i$  for an arm  $i \in \mathcal{A}$  along with the time horizon and number of arms. It is so called because the regret bound depends explicitly on the means of the arms considered for that environment along with the stated assumptions on the distribution.
4. **Gap-Independent regret bounds:** In gap-independent regret bound the regret does not contain the gaps and is stated explicitly in terms of the number of arms and the horizon. This is because the regret depends only on the distributional assumption, but not on the means of the arms considered. In fact, gap-independent regret bounds point to something more general and informative. These type of bounds actually give us the maximum possible regret such that no matter what is the policy, there will be an environment on which the policy achieves almost the same regret as the gap-independent regret upper bound. This leads to the notion of minimax regret.
5. **Minimax regret bounds:** For a finite horizon  $T$ ,  $K$  number of arms, for all set of possible policies  $\pi_{T,K}$  over  $T$  and  $K$  and all possible environment class  $\mathcal{E}$  the minimax regret is given by,

$$R_T(\mathcal{E}) = \inf_{\pi \in \pi_{T,K}} \sup_{E \in \mathcal{E}} R_T(\pi, E).$$

Hence, this value is independent of any specific choice of a policy  $\pi$  but only depends on  $T$ ,  $K$  and  $\mathcal{E}$  where the dependence on  $K$  is hidden in  $\mathcal{E}$ .

## 2.4 Motivation

There has been a significant amount of research in the area of stochastic MABs. One of the earliest work can be traced to Thompson (1933), which deals with the problem of choosing between two treatments to administer on patients who come in sequentially. In Thompson (1935) this work was extended to include more general cases of finitely many treatments. In recent years the SMAB setting has garnered extensive popularity because of its simple learning model and its practical applications in a wide-range of industries, including, but not limited to, mobile channel allocations, online advertising and computer simulation games. Some of these problems have been already discussed in chapter 1, section 1.2 and an interested reader can refer to it.

## 2.5 Related Work in SMAB

### 2.5.1 Lower bound in SMAB

SMAB problems have been extensively studied in several earlier works such as Thompson (1933), Thompson (1935), Robbins (1952) and Lai and Robbins (1985). Lai and Robbins in Lai and Robbins (1985) established an asymptotic lower bound for the cumulative regret. It showed that for any consistent allocation strategy, we can have

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[R_T]}{\log T} \geq \sum_{\{i: r_i < r^*\}} \frac{(r^* - r_i)}{KL(Q_i || Q^*)}$$

where  $KL(Q_i || Q^*)$  is the Kullback-Leibler divergence between the reward densities  $Q_i$  and  $Q^*$ , corresponding to arms with mean  $r_i$  and  $r^*$ , respectively.

### 2.5.2 The Upper Confidence Bound approach

Over the years SMABs have seen several algorithms with strong regret guarantees. For further reference, an interested reader can look into Bubeck and Cesa-Bianchi (2012). In the next few subsections, we will explicitly focus on the upper confidence bound algorithms which is a type of non-Bayesian algorithm widely used in SMAB setting.

The upper confidence bound or UCB algorithms balance the exploration-exploitation dilemma by linking the uncertainty in the estimate of an arm with the number of times an arm is pulled and therefore ensuring sufficient exploration.

### UCB1 Algorithm

One of the earliest among these algorithms is UCB1 algorithm proposed first in Agrawal (1995) and subsequently analyzed in Auer *et al.* (2002a). The UCB1 algorithm (as stated in Auer *et al.* (2002a)) is mentioned in algorithm 3.

---

#### Algorithm 3 UCB1

---

- 1: **Input:**  $K$  number of arms with unknown parameters of reward distribution
  - 2: Pull each arm once
  - 3: **for**  $t = K + 1, \dots, T$  **do**
  - 4:     Pull the arm such that  $\arg \max_{i \in A} \left\{ \hat{r}_i + \sqrt{\frac{2 \log(t)}{n_i}} \right\}$
  - 5:      $t := t + 1$
  - 6: **end for**
- 

The intuition behind this algorithm is simple and it follows from the ideas of concentration inequalities in probability measure theory. The term  $\sqrt{\frac{2 \log(t)}{n_i}}$  is called the confidence interval of the arm  $i$  and it signifies a measure of uncertainty over the arm  $i$  based on the history of observed rewards for that arm. Therefore, lesser the confidence interval, higher is our confidence that the estimated mean  $\hat{r}_i$  is lying close to the expected mean  $r_i$  of the arm  $i$ . Also, note that the confidence interval decreases at the rate of  $O\left(\frac{1}{\sqrt{n_i}}\right)$  which signifies the rate of convergence of  $\hat{r}_i$  to  $r_i$  and depends on the number of time the arm has been pulled.

UCB1 has a gap-dependent 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. But, the worst case gap-independent regret bound of UCB1 is found to be  $O\left(\sqrt{KT \log T}\right)$ .

### UCB-Improved Algorithm

The UCB-Improved stated in algorithm 4, proposed in Auer and Ortner (2010), is a round-based variant of UCB1. An algorithm is *round-based* if it pulls all the arms equal number of times in each round and then eliminates one or more arms that it deems to be

---

**Algorithm 4** UCB-Improved

---

- 1: **Input:** Time horizon  $T$ ,  $K$  number of arms with unknown parameters of reward distribution
  - 2: **Initialization:** Set  $B_0 := \mathcal{A}$  and  $\epsilon_0 := 1$ .
  - 3: **for**  $m = 0, 1, \dots, \lfloor \frac{1}{2} \log_2 \frac{T}{e} \rfloor$  **do**
  - 4:     Pull each arm in  $B_m$ ,  $n_m = \left\lceil \frac{2 \log(T\epsilon_m^2)}{\epsilon_m} \right\rceil$  number of times.
  - 5:     ***Arm Elimination by Mean Estimation***
  - 6:     For each  $i \in B_m$ , delete arm  $i$  from  $B_m$  if,
$$\hat{r}_i + \sqrt{\frac{\log(T\epsilon_m^2)}{2n_m}} < \max_{j \in B_m} \left\{ \hat{r}_j - \sqrt{\frac{\log(T\epsilon_m^2)}{2n_m}} \right\}$$
  - 7:     Set  $\epsilon_{m+1} := \frac{\epsilon_m}{2}$ , Set  $B_{m+1} := B_m$
  - 8:     Stop if  $|B_m| = 1$  and pull  $i \in B_m$  till  $n$  is reached.
  - 9: **end for**
- 

sub-optimal. Note, that in this algorithm the confidence interval term is  $\sqrt{\frac{\log(T\epsilon_m^2)}{2n_m}}$  which is constant in the  $m$ -th round as  $n_m$  is fixed for that round and all arms are being pulled an equal number of times in each round. This is unlike UCB1 algorithm where the confidence interval term depends on  $n_i$  which is a random variable. Also, note that in UCB-Improved the knowledge of horizon is required before-hand to calculate the confidence intervals whereas no such input is required for UCB1.

UCB-Improved incurs 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 gap-independent regret bound of UCB-Improved is  $O(\sqrt{KT \log K})$ . Empirically, UCB-Improved is outperformed by UCB1 in almost all environments. This stems from the fact that UCB-Improved is pulling all arms equal number of times in each round and hence spends a significant number of pulls in initial exploration as opposed to UCB1 thereby incurring higher regret.

## MOSS Algorithm

In the later work of Audibert and Bubeck (2009), the authors propose the MOSS algorithm and showed that the worst case gap-independent regret bound of MOSS is  $O(\sqrt{KT})$  which improves upon UCB1 by a factor of order  $\sqrt{\log T}$ . However, the

---

**Algorithm 5** MOSS

---

- 1: **Input:** Time horizon  $T$ ,  $K$  number of arms with unknown parameters of reward distribution
  - 2: Pull each arm once
  - 3: **for**  $t = K + 1, \dots, T$  **do**
  - 4:     Pull the arm such that  $\arg \max_{i \in \mathcal{A}} \left\{ \hat{r}_i + \sqrt{\frac{\max\{0, \log(\frac{T}{Kn_i})\}}{n_i}} \right\}$
  - 5:      $t := t + 1$
  - 6: **end for**
- 

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)).

Recently in Lattimore (2015), the authors showed that the algorithm OCUCB achieves order-optimal gap-dependent regret bound of  $O\left(\sum_{i=2}^K \frac{\log(T/H_i)}{\Delta_i}\right)$  where  $H_i = \sum_{j=1}^K \min\left\{\frac{1}{\Delta_i^2}, \frac{1}{\Delta_j^2}\right\}$ , and a gap-independent regret bound of  $O(\sqrt{KT})$ . This is the best known gap-dependent and gap-independent regret bounds in the stochastic MAB framework. However, unlike our proposed EUCBV algorithm (in chapter 3), OCUCB does not take into account the variance of the arms; as a result, empirically we find that our algorithm outperforms OCUCB in all the environments considered.

### UCB-Variance algorithm

---

**Algorithm 6** UCBV

---

- 1: **Input:**  $K$  number of arms with unknown parameters of reward distribution
  - 2: Pull each arm once
  - 3: **for**  $t = K + 1, \dots, T$  **do**
  - 4:     Pull the arm such that  $\max_{i \in \mathcal{A}} \left\{ \hat{r}_i + \sqrt{\frac{2\hat{v}_i \log(t)}{s_i}} + \frac{3 \log(t)}{2} \right\}$
  - 5:      $t := t + 1$
  - 6: **end for**
- 

In contrast to the above work, the UCBV (Audibert *et al.*, 2009) algorithm utilizes variance estimates to compute the confidence intervals for each arm. In UCBV the confidence interval term is given by  $\sqrt{\frac{2\hat{v}_i \log(t)}{s_i}} + \frac{3 \log(t)}{2}$  where  $\hat{v}_i$  denotes the empirical variance of the arm  $i$ . Hence, the confidence interval makes sure that the arms whose variances are high are pulled more often to get a better estimates of their  $\hat{r}_i$ .

UCBV has a gap-dependent regret bound of  $O\left(\frac{K\sigma_{\max}^2 \log T}{\Delta}\right)$ , where  $\sigma_{\max}^2$  denotes



the maximum variance among all the arms  $i \in \mathcal{A}$ . Its gap-independent regret bound can be inferred to be same as that of UCB1 i.e  $O(\sqrt{KT \log T})$ . Empirically, Audibert *et al.* (2009) showed that UCBV outperforms UCB1 in several scenarios.

### 2.5.3 Bayesian Approach

---

#### Algorithm 7 Bernoulli Thompson Sampling

---

**Input:**  $K$  number of arms with unknown parameters of reward distribution

**Initialization:** For each arm  $i := 1$  to  $K$  set  $S_i = 0$  and  $F_i = 0$

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

**for**  $i = 1, \dots, K$  **do**

        Sample  $\theta_i(t)$  from the  $Beta(S_i + 1, F_i + 1)$  distribution.

**end for**

    Play the arm  $i(t) := \arg \max_i \theta_i(t)$  and observe reward  $X_{i,t}$ .

**if**  $X_{i,t} = 1$  **then**  $S_i(t) = S_i(t) + 1$

**else**  $F_i(t) = F_i(t) + 1$

**end if**

**end for**

---

Another notable design principle which has recently gained a lot of popularity is the Thompson Sampling (TS) algorithm ((Thompson, 1933), (Agrawal and Goyal, 2011)) and Bayes-UCB (BU) algorithm (Kaufmann *et al.*, 2012). This TS is stated in algorithm 7. The TS algorithm maintains a posterior reward distribution for each arm; at each round, the algorithm samples values from these distributions and the arm corresponding to the highest sample value is chosen. Although TS is found to perform extremely well when the reward distributions are Bernoulli, it is established that with Gaussian priors the worst-case regret can be as bad as  $\Omega(\sqrt{KT \log T})$  (Lattimore, 2015). The BU algorithm is an extension of the TS algorithm that takes quartile deviations into consideration while choosing arms.

### 2.5.4 Information Theoretic approach

The final design principle we state is the information theoretic approach of DMED (Honda and Takemura, 2010) and KLUCB (Garivier and Cappé, 2011), (Cappé *et al.*,

2013) algorithms. The algorithm KLUCB uses Kullbeck-Leibler divergence to compute the upper confidence bound for the arms. KLUCB is stable for a short horizon and is known to reach the Lai and Robbins (1985) lower bound in the special case of Bernoulli distribution. However, Garivier and Cappé (2011) showed that KLUCB, MOSS and UCB1 algorithms are empirically outperformed by UCBV in the exponential distribution as they do not take the variance of the arms into consideration.

### 2.5.5 Discussion on the various confidence intervals

A comparative analysis of the confidence interval of the UCB algorithms is discussed in table 2.1.

Table 2.1: Confidence interval of different algorithms

Algorithm	Confidence interval	Horizon as input	Remarks
UCB1	$\sqrt{\frac{2 \log(t)}{n_i}}$	No	Loose confidence interval leading to high regret upper bounds.
UCBV	$\sqrt{\frac{2\hat{v}_i \log(t)}{s_i}} + \frac{3 \log(t)}{2}$	No	Confidence interval uses variance estimation.
UCB-Imp	$\sqrt{\frac{\log(T\epsilon_m^2)}{2n_m}}$	Yes	Same confidence interval for all arms in a particular round.
MOSS	$\sqrt{\frac{\max\{0, \log(\frac{T}{Kn_i})\}}{n_i}}$	Yes	Confidence interval is based on dividing the horizon uniformly for $K$ arms.
OCUCB	$\sqrt{\frac{2 \log(\frac{2T}{t})}{n_i}}$	Yes	Tightest confidence interval leading to order-optimal regret bounds.

## 2.6 Conclusion

In this chapter, we looked at the stochastic multi-armed bandit (SMAB) setting and discussed how it is important in the general reinforcement learning setup. We also looked at the various state-of-the-art algorithms in the literature for the SMAB setting and discussed the advantages and disadvantages of them. The regret bounds that have been proven for the said algorithms have also been discussed at length and their confidence intervals have also been compared against each other. In the next chapter, we provide our solution to this SMAB setting which achieves an almost order-optimal regret bound.

## Chapter 3

# Efficient UCB Variance: An almost optimal algorithm in SMAB setting

### 3.1 Introduction

In this chapter, we look at a novel variant of the UCB algorithm (referred to as Efficient-UCB-Variance (EUCBV)) for minimizing cumulative regret in the stochastic multi-armed bandit (SMAB) setting. EUCBV incorporates the arm elimination strategy proposed in UCB-Improved (Auer and Ortner, 2010) while taking into account the variance estimates to compute the arms' confidence bounds, similar to UCBV (Audibert *et al.*, 2009). Through a theoretical analysis we establish that EUCBV incurs a *gap-dependent* regret bound of  $O\left(\frac{K\sigma_{\max}^2 \log(T\Delta^2/K)}{\Delta}\right)$  after  $T$  trials, where  $\Delta$  is the minimal gap between optimal and sub-optimal arms; the above bound is an improvement over that of existing state-of-the-art UCB algorithms (such as UCB1, UCB-Improved, UCBV, MOSS). Further, EUCBV incurs a *gap-independent* regret bound of  $O\left(\sqrt{KT}\right)$  which is an improvement over that of UCB1, UCBV and UCB-Improved, while being comparable with that of MOSS and OCUCB. Through an extensive numerical study, we show that EUCBV significantly outperforms the popular UCB variants (like MOSS, OCUCB, etc.) as well as Thompson sampling and Bayes-UCB algorithms.

The rest of the chapter is organized as follows. We elaborate our contributions in section 3.2 and in section 3.3 we present the EUCBV algorithm. Our main theoretical results are stated in section 3.4, while the proofs are established in section 3.5. Section 3.6 contains results and discussions from our numerical experiments. We draw our conclusions in section 3.7 and Appendix B contains the proofs of the lemmas that have been used for proving the main result.

## 3.2 Our Contributions

We propose the Efficient-UCB-Variance (henceforth referred to as EUCBV) algorithm for the stochastic MAB setting. EUCBV combines the approach of UCB-Improved, CCB (Liu and Tsuruoka, 2016) and UCBV algorithms. EUCBV, by virtue of taking into account the empirical variance of the arms, exploration parameters and non-uniform arm selection (as opposed to UCB-Improved), performs significantly better than the existing algorithms in the stochastic MAB setting. EUCBV outperforms UCBV (Audibert *et al.*, 2009) which also takes into account empirical variance but is less powerful than EUCBV because of the usage of exploration regulatory factor by UCBV. Also, we carefully design the confidence interval term with the variance estimates along with the pulls allocated to each arm to balance the risk of eliminating the optimal arm against excessive optimism. Theoretically we refine the analysis of Auer and Ortner (2010) and prove that for  $T \geq K^{2.4}$  our algorithm is order optimal and achieves a worst case gap-independent regret bound of  $O\left(\sqrt{KT}\right)$  which is same as that of MOSS and OCUCB but better than that of UCBV, UCB1 and UCB-Improved. Also, the gap-dependent regret bound of EUCBV is better than UCB1, UCB-Improved and MOSS but is poorer than OCUCB. However, EUCBV's gap-dependent bound matches OCUCB in the worst case scenario when all the gaps are equal. Through our theoretical analysis we establish the exact values of the exploration parameters for the best performance of EUCBV. Our proof technique is highly generic and can be easily extended to other MAB settings. An illustrative table containing the bounds is provided in Table 3.1.

Empirically, we show that EUCBV, owing to its estimating the variance of the arms, exploration parameters and non-uniform arm pull, performs significantly better than MOSS, OCUCB, UCB-Improved, UCB1, UCBV, TS, BU, DMED, KLUCB and Median Elimination algorithms. Note that except UCBV, TS, KLUCB and BU (the last three with Gaussian priors) all the aforementioned algorithms do not take into account the empirical variance estimates of the arms. Also, for the optimal performance of TS, KLUCB and BU one has to have the prior knowledge of the type of distribution, but EUCBV requires no such prior knowledge. EUCBV is the first arm-elimination algorithm that takes into account the variance estimates of the arm for minimizing cumulative regret and thereby answers an open question raised by Auer and Ortner (2010),

Table 3.1: Regret upper bound of different algorithms

Algorithm	Gap-Dependent	Gap-Independent
EUCBV	$O\left(\frac{K\sigma_{\max}^2 \log(\frac{T\Delta^2}{K})}{\Delta}\right)$	$O(\sqrt{KT})$
UCB1	$O\left(\frac{K \log T}{\Delta}\right)$	$O(\sqrt{KT \log T})$
UCBV	$O\left(\frac{K\sigma_{\max}^2 \log T}{\Delta}\right)$	$O(\sqrt{KT \log T})$
UCB-Imp	$O\left(\frac{K \log(T\Delta^2)}{\Delta}\right)$	$O(\sqrt{KT \log K})$
MOSS	$O\left(\frac{K^2 \log(T\Delta^2/K)}{\Delta}\right)$	$O(\sqrt{KT})$
OCUCB	$O\left(\frac{K \log(T/H_i)}{\Delta}\right)$	$O(\sqrt{KT})$

where the authors conjectured that an UCB-Improved like arm-elimination algorithm can greatly benefit by taking into consideration the variance of the arms. Also, it is the first algorithm that follows the same proof technique of UCB-Improved and achieves a gap-independent regret bound of  $O(\sqrt{KT})$  thereby, closing the gap of UCB-Improved which achieved a gap-independent regret bound of  $O(\sqrt{KT \log K})$ .

### 3.3 Algorithm: Efficient UCB Variance

**The algorithm:** Earlier round-based arm elimination algorithms like Median Elimination (Even-Dar *et al.*, 2006) and UCB-Improved mainly suffered from two basic problems:

- (i) *Initial exploration:* Both of these algorithms pull each arm equal number of times in each round, and hence waste a significant number of pulls in initial explorations.
- (ii) *Conservative arm-elimination:* In UCB-Improved, arms are eliminated conservatively, i.e, only after  $\epsilon_m < \frac{\Delta_i}{2}$ , where the quantity  $\epsilon_m$  is initialized to 1 and halved after every round. In the worst case scenario when  $K$  is large, and the gaps are uniform ( $r_1 = r_2 = \dots = r_{K-1} < r^*$ ) and small this results in very high regret.

---

**Algorithm 8** EUCBV

---

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

**Initialization:** Set  $m := 0$ ,  $B_0 := \mathcal{A}$ ,  $\epsilon_0 := 1$ ,  $M = \lfloor \frac{1}{2} \log_2 \frac{T}{e} \rfloor$ ,  $n_0 = \lceil \frac{\log(\psi T \epsilon_0^2)}{2\epsilon_0} \rceil$  and  $N_0 = K n_0$ .

Pull each arm once

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

    Pull arm  $i \in \arg \max_{j \in B_m} \left\{ \hat{r}_j + \sqrt{\frac{\rho(\hat{v}_j + 2) \log(\psi T \epsilon_m)}{4z_j}} \right\}$ , where  $z_j$  is the number of times arm  $j$  has been pulled.

**Arm Elimination by Mean Estimation**

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

$$\hat{r}_i + \sqrt{\frac{\rho(\hat{v}_i + 2) \log(\psi T \epsilon_m)}{4z_i}} < \max_{j \in B_m} \left\{ \hat{r}_j - \sqrt{\frac{\rho(\hat{v}_j + 2) \log(\psi T \epsilon_m)}{4z_j}} \right\}$$

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

**Reset Parameters**

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

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

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

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

$$m := m + 1$$

**end if**

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

**end for**

---

The EUCBV algorithm, which is mainly based on the arm elimination technique of the UCB-Improved algorithm, remedies these by employing exploration regulatory factor  $\psi$  and arm elimination parameter  $\rho$  for aggressive elimination of sub-optimal arms. Along with these, similar to CCB (Liu and Tsuruoka, 2016) algorithm, EUCBV uses optimistic greedy sampling whereby at every timestep it only pulls the arm with the highest upper confidence bound rather than pulling all the arms equal number of times in each round. Also, unlike the UCB-Improved, UCB1, MOSS and OCUCB algorithms (which are based on mean estimation) EUCBV employs mean and variance estimates (as in Audibert *et al.* (2009)) for arm elimination. Further, we allow for arm-elimination at every time-step, which is in contrast to the earlier work (e.g., Auer and Ortner (2010); Even-Dar *et al.* (2006)) where the arm elimination takes place only at the end of the respective exploration rounds.

### 3.4 Main Results

The main result of this chapter is presented in the following theorem, where we establish a regret upper bound for the proposed EUCBV algorithm.

#### Gap-Dependent bound of EUCBV

**Theorem 1 (Gap-Dependent Bound)** *For  $T \geq K^{2.4}$ ,  $\rho = \frac{1}{2}$  and  $\psi = \frac{T}{K^2}$ , the regret  $R_T$  for EUCBV satisfies*

$$\begin{aligned} \mathbb{E}[R_T] \leq & \sum_{i \in \mathcal{A}: \Delta_i > b} \left\{ \frac{C_0 K^4}{T^{\frac{1}{4}}} + \left( \Delta_i + \frac{320 \sigma_i^2 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right) \right\} \\ & + \sum_{i \in \mathcal{A}: 0 < \Delta_i \leq b} \frac{C_2 K^4}{T^{\frac{1}{4}}} + \max_{i \in \mathcal{A}: 0 < \Delta_i \leq b} \Delta_i T. \end{aligned}$$

for all  $b \geq \sqrt{\frac{\epsilon}{T}}$  and  $C_0, C_2$  are integer constants.

**Proof 1 (Outline)** *The proof is along the lines of the technique in Auer and Ortner (2010). It comprises of three modules. In the first module we prove the necessary conditions for arm elimination within a specified number of rounds. However, here we require some additional technical results (see Lemma 1 and Lemma 2) to bound the length of the confidence intervals. Further, note that our algorithm combines the variance-estimate based approach of Audibert et al. (2009) with the arm-elimination technique of Auer and Ortner (2010) (see Lemma 3). Also, while Auer and Ortner (2010) uses Chernoff-Hoeffding bound to derive their regret bound whereas in our work we use Bernstein inequality (as in Audibert et al. (2009)) to obtain the bound. To bound the probability of the non-uniform arm selection before it gets eliminated we use Lemma 4 and Lemma 5. In the second module we bound the number of pulls required if an arm is eliminated on or before a particular number of rounds. Note that the number of pulls allocated in a round  $m$  for each arm is  $n_m := \left\lceil \frac{\log(\psi T \epsilon_m^2)}{2\epsilon_m} \right\rceil$  which is much lower than the number of pulls of each arm required by UCB-Improved or Median-Elimination. We introduce the variance term in the most significant term in the bound by Lemma 6. Finally, the third module deals with case of bounding the regret, given that a sub-optimal arm eliminates the optimal arm. ■*



**Discussion 1** From the above result we see that the most significant term in the gap-dependent bound is of the order  $O\left(\frac{K\sigma_{\max}^2 \log(T\Delta^2/K)}{\Delta}\right)$  which is better than the existing results for UCB1, UCBV, MOSS and UCB-Improved (see Table 3.1). Also as like UCBV, this term scales with the variance. Audibert and Bubeck (2010) have defined the term  $H_1 = \sum_{i=1}^K \frac{1}{\Delta_i^2}$ , which is referred to as the hardness of a problem; Bubeck and Cesa-Bianchi (2012) have conjectured that the gap-dependent regret upper bound can match  $O\left(\frac{K \log(T/H_1)}{\Delta}\right)$ . However, in Lattimore (2015) it is proved that the gap-dependent regret bound cannot be lower than  $O\left(\sum_{i=2}^K \frac{\log(T/H_i)}{\Delta_i}\right)$ , where  $H_i = \sum_{j=1}^K \min\left\{\frac{1}{\Delta_i^2}, \frac{1}{\Delta_j^2}\right\}$  (OCUCB proposed in Lattimore (2015) achieves this bound). Further, in Lattimore (2015) it is shown that only in the worst case scenario when all the gaps are equal (so that  $H_1 = H_i = \sum_{i=1}^K \frac{1}{\Delta^2}$ ) the above two bounds match. In the latter scenario, considering  $\sigma_{\max}^2 \leq \frac{1}{4}$  as all rewards are bounded in  $[0, 1]$ , we see that the gap-dependent bound of EUCBV simplifies to  $O\left(\frac{K \log(T/H_1)}{\Delta}\right)$ , thus matching the gap-dependent bound of OCUCB which is order optimal.

### Gap-Independent bound of EUCBV

In this section, we specialize the result of Theorem 1 in Corollary 1 to obtain the gap-independent worst case regret bound.

**Corollary 1 (Gap-Independent Bound)** When the gaps of all the sub-optimal arms are identical, i.e.,  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{\epsilon}{T}}, \forall i \in \mathcal{A}$  and  $C_3$  being an integer constant, the regret of EUCBV is upper bounded by the following gap-independent expression:

$$\mathbb{E}[R_T] \leq \frac{C_3 K^5}{T^{\frac{1}{4}}} + 320\sqrt{KT}.$$

The proof is given in Appendix B.0.10.

**Discussion 2** In the non-stochastic scenario, Auer et al. (2002b) showed that the bound on the cumulative regret for EXP-4 is  $O\left(\sqrt{KT \log K}\right)$ . However, in the stochastic case, UCB1 proposed in Auer et al. (2002a) incurred a regret of order of  $O\left(\sqrt{KT \log T}\right)$  which is clearly improvable. From the above result we see that in the gap-independent bound of EUCBV the most significant term is  $O\left(\sqrt{KT}\right)$  which matches the upper

bound of MOSS and OCUCB, and is better than UCB-Improved, UCB1 and UCBV (see Table 3.1).

### 3.5 Proofs

We first present a few technical lemmas that is required to prove the result in Theorem 1.

**Lemma 1** *If  $T \geq K^{2.4}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$  and  $m \leq \frac{1}{2} \log_2 \left( \frac{T}{e} \right)$ , then,*

$$\frac{\rho m \log(2)}{\log(\psi T) - 2m \log(2)} \leq \frac{3}{2}.$$

**Lemma 2** *If  $T \geq K^{2.4}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$ ,  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$  and  $c_i = \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}}$ , then,*

$$c_i < \frac{\Delta_i}{4}.$$

**Lemma 3** *If  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $c_i = \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i})}{2\epsilon_{m_i}}$  then we can show that,*

$$\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \frac{2}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}}.$$

**Lemma 4** *If  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$ ,  $c_i = \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$  then in the  $m_i$ -th round,*

$$\mathbb{P}\{c^* > c_i\} \leq \frac{182K^4}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}}.$$

**Lemma 5** *If  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$ ,  $c_i = \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$  then in the  $m_i$ -th round,*

$$\mathbb{P}\{z_i < n_{m_i}\} \leq \frac{182K^4}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}}.$$

**Lemma 6** For two integer constants  $c_1$  and  $c_2$ , if  $20c_1 \leq c_2$  then,

$$c_1 \frac{4\sigma_i^2 + 4}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right) \leq c_2 \frac{\sigma_i^2}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right).$$

The proofs of lemmas 1 - 6 can be found in Appendix B.0.4, B.0.5, B.0.6, B.0.7, B.0.8 and B.0.9 respectively.

## Proof of Theorem 1

**Proof 1** For each sub-optimal arm  $i \in \mathcal{A}$ , let  $m_i = \min \{m | \sqrt{4\epsilon_{m_i}} < \frac{\Delta_i}{4}\}$ . Also, let  $\mathcal{A}' = \{i \in \mathcal{A} : \Delta_i > b\}$  and  $\mathcal{A}'' = \{i \in \mathcal{A} : \Delta_i > 0\}$ . Note that as all rewards are bounded in  $[0, 1]$ , it implies that  $0 \leq \sigma_i^2 \leq \frac{1}{4}, \forall i \in \mathcal{A}$ . Now, as in Auer and Ortner (2010), we bound the regret under the following two cases:

- Case (a): some sub-optimal arm  $i$  is not eliminated in round  $m_i$  or before and the optimal arm  $* \in B_{m_i}$
- Case (b): an arm  $i \in B_{m_i}$  is eliminated in round  $m_i$  (or before), or there is no optimal arm  $* \in B_{m_i}$

The details of each case are contained in the following sub-sections.

**Case (a):** For simplicity, let  $c_i := \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}}$  denote the length of the confidence interval corresponding to arm  $i$  in round  $m_i$ . Thus, in round  $m_i$  (or before) whenever  $z_i \geq n_{m_i} \geq \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$ , by applying Lemma 2 we obtain  $c_i < \frac{\Delta_i}{4}$ . Now, the sufficient conditions for arm  $i$  to get eliminated by an optimal arm in round  $m_i$  is given by

$$\hat{r}_i \leq r_i + c_i, \hat{r}^* \geq r^* - c^*, c_i \geq c^* \text{ and } z_i \geq n_{m_i}. \quad (3.1)$$

Indeed, in round  $m_i$  suppose (3.1) holds, then we have

$$\begin{aligned} \hat{r}_i + c_i &\leq r_i + 2c_i = r_i + 4c_i - 2c_i \\ &< r_i + \Delta_i - 2c_i \leq r^* - 2c^* \leq \hat{r}^* - c^* \end{aligned}$$

so that a sub-optimal arm  $i \in \mathcal{A}'$  gets eliminated. Thus, the probability of the complementary event of these four conditions in (3.1) yields a bound on the probability that arm  $i$  is not eliminated in round  $m_i$ . Following the proof of Lemma 1 of Audibert et al. (2009) we can show that a bound on the complementary of the first condition is given by,

$$\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \mathbb{P}(\hat{r}_i > r_i + \bar{c}_i) + \mathbb{P}(\hat{v}_i \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}) \quad (3.2)$$

where

$$\bar{c}_i = \sqrt{\frac{\rho(\sigma_i^2 + \sqrt{\epsilon_{m_i}} + 2) \log(\psi T \epsilon_{m_i})}{4n_{m_i}}}.$$

From Lemma 3 we can show that  $\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \mathbb{P}(\hat{r}_i > r_i + \bar{c}_i) + \mathbb{P}(\hat{v}_i \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}) \leq \frac{2}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}}$ . Similarly,  $\mathbb{P}\{\hat{r}^* < r^* - c^*\} \leq \frac{2}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}}$ . Summing the above two contributions, the probability that a sub-optimal arm  $i$  is not eliminated on or before  $m_i$ -th round by the first two conditions in (3.1) is,

$$\left( \frac{4}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \right). \quad (3.3)$$

Again, from Lemma 4 and Lemma 5 we can bound the probability of the complementary of the event  $c_i \geq c^*$  and  $z_i \geq n_{m_i}$  by,

$$\frac{182K^4}{T^{\frac{5}{4}}\sqrt{\epsilon_{m_i}}} + \frac{182K^4}{T^{\frac{5}{4}}\sqrt{\epsilon_{m_i}}} \leq \frac{364K^4}{T^{\frac{5}{4}}\sqrt{\epsilon_{m_i}}}. \quad (3.4)$$

Also, for eq. (3.3) we can show that for any  $\epsilon_{m_i} \in [\sqrt{\frac{e}{T}}, 1]$

$$\begin{aligned} \left( \frac{4}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \right) &\stackrel{(a)}{\leq} \left( \frac{4}{\left(\frac{T^2}{K^2} \epsilon_{m_i}\right)^{\frac{3}{4}}} \right) \leq \left( \frac{4K^{\frac{3}{2}}}{(T^{\frac{3}{2}} \epsilon_{m_i}^{\frac{1}{4}} \sqrt{\epsilon_{m_i}})} \right) \\ &\stackrel{(b)}{\leq} \left( \frac{4K^{\frac{3}{2}}}{(T^{\frac{3}{2} - \frac{1}{8}} \sqrt{\epsilon_{m_i}})} \right) \leq \frac{4K^4}{T^{\frac{5}{4}}\sqrt{\epsilon_{m_i}}}. \end{aligned} \quad (3.5)$$

Here, in (a) we substitute the values of  $\psi$  and  $\rho$  and (b) follows from the identity  $\epsilon_{m_i}^{\frac{1}{4}} \geq (\frac{e}{T})^{\frac{1}{8}}$  as  $\epsilon_{m_i} \geq \sqrt{\frac{e}{T}}$ .

Summing up over all arms in  $\mathcal{A}'$  and bounding the regret for all the four arm elimination conditions in (3.1) by (3.4) + (3.5) for each arm  $i \in \mathcal{A}'$  trivially by  $T\Delta_i$ , we obtain

$$\begin{aligned} & \sum_{i \in \mathcal{A}'} \left( \frac{4K^4 T \Delta_i}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}} \right) + \sum_{i \in \mathcal{A}'} \left( \frac{364K^4 T \Delta_i}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}} \right) \\ & \stackrel{(a)}{\leq} \sum_{i \in \mathcal{A}'} \left( \frac{368K^4 T \Delta_i}{T^{\frac{5}{4}} \left( \frac{\Delta_i^2}{4.16} \right)^{\frac{1}{2}}} \right) \stackrel{(b)}{\leq} \sum_{i \in \mathcal{A}'} \left( \frac{C_1 K^4}{(T)^{\frac{1}{4}}} \right). \end{aligned}$$

Here, (a) happens because  $\sqrt{4\epsilon_{m_i}} < \frac{\Delta_i}{4}$ , and in (b),  $C_1$  denotes a constant integer value.

**Case (b):** Here, there are two sub-cases to be considered.

**Case (b1) ( $* \in B_{m_i}$  and each  $i \in \mathcal{A}'$  is eliminated on or before  $m_i$ ):** Since we are eliminating a sub-optimal arm  $i$  on or before round  $m_i$ , it is pulled no longer than,

$$z_i < \left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil$$

So, the total contribution of  $i$  till round  $m_i$  is given by,

$$\begin{aligned} & \Delta_i \left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil \stackrel{(a)}{\leq} \Delta_i \left\lceil \frac{\log(\psi T (\frac{\Delta_i}{16 \times 256})^4)}{2(\frac{\Delta_i}{4\sqrt{4}})^2} \right\rceil \\ & \leq \Delta_i \left( 1 + \frac{32 \log(\psi T (\frac{\Delta_i}{16384})^4)}{\Delta_i^2} \right) \leq \Delta_i \left( 1 + \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i^2} \right). \end{aligned}$$

Here, (a) happens because  $\sqrt{4\epsilon_{m_i}} < \frac{\Delta_i}{4}$ . Summing over all arms in  $\mathcal{A}'$  the total

regret is given by,

$$\begin{aligned}
& \sum_{i \in \mathcal{A}'} \Delta_i \left( 1 + \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i^2} \right) = \sum_{i \in \mathcal{A}'} \left( \Delta_i + \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i} \right) \\
& \stackrel{(a)}{\leq} \sum_{i \in \mathcal{A}'} \left( \Delta_i + \frac{64 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right) \\
& \stackrel{(b)}{\leq} \sum_{i \in \mathcal{A}'} \left( \Delta_i + \frac{16(4\sigma_i^2 + 4) \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right) \\
& \stackrel{(c)}{\leq} \sum_{i \in \mathcal{A}'} \left( \Delta_i + \frac{320\sigma_i^2 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right).
\end{aligned}$$

We obtain (a) by substituting the value of  $\psi$ , (b) from  $0 \leq \sigma_i^2 \leq \frac{1}{4}, \forall i \in \mathcal{A}$  and (c) from Lemma 6.

**Case (b2) (Optimal arm  $*$  is eliminated by a sub-optimal arm):** Firstly, if conditions of Case a holds then the optimal arm  $*$  will not be eliminated in round  $m = m_*$  or it will lead to the contradiction that  $r_i > r^*$ . In any round  $m_*$ , if the optimal arm  $*$  gets eliminated then for any round from 1 to  $m_j$  all arms  $j$  such that  $m_j < m_*$  were eliminated according to assumption in Case a. Let the arms surviving till  $m_*$  round be denoted by  $\mathcal{A}'$ . This leaves any arm  $a_b$  such that  $m_b \geq m_*$  to still survive and eliminate arm  $*$  in round  $m_*$ . Let such arms that survive  $*$  belong to  $\mathcal{A}''$ . Also maximal regret per step after eliminating  $*$  is the maximal  $\Delta_j$  among the remaining arms  $j$  with  $m_j \geq m_*$ . Let  $m_b = \min \{m | \sqrt{4\epsilon_m} < \frac{\Delta_b}{4}\}$ . Hence, the maximal regret after eliminating the arm  $*$  is upper bounded by,

$$\begin{aligned}
& \sum_{m_*=0}^{\max_{j \in \mathcal{A}'} m_j} \sum_{i \in \mathcal{A}'' : m_i > m_*} \left( \frac{368K^4}{(T^{\frac{5}{4}} \sqrt{\epsilon_{m_*}})} \right) \cdot T \max_{j \in \mathcal{A}'' : m_j \geq m_*} \Delta_j \\
& \leq \sum_{m_*=0}^{\max_{j \in \mathcal{A}'} m_j} \sum_{i \in \mathcal{A}'' : m_i > m_*} \left( \frac{368K^4 \sqrt{4}}{(T^{\frac{5}{4}} \sqrt{\epsilon_{m_*}})} \right) \cdot T \cdot 4 \sqrt{\epsilon_{m_*}}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} \sum_{m_*=0}^{\max_{j \in \mathcal{A}'} m_j} \sum_{i \in \mathcal{A}'' : m_i > m_*} \left( \frac{C_2 K^4}{T^{\frac{1}{4}} \epsilon_{m_*}^{\frac{1}{2} - \frac{1}{2}}} \right) \\
&\leq \sum_{i \in \mathcal{A}'' : m_i > m_*} \sum_{m_*=0}^{\min\{m_i, m_b\}} \left( \frac{C_2 K^4}{T^{\frac{1}{4}}} \right) \\
&\leq \sum_{i \in \mathcal{A}'} \left( \frac{C_2 K^4}{T^{\frac{1}{4}}} \right) + \sum_{i \in \mathcal{A}'' \setminus \mathcal{A}'} \left( \frac{C_2 K^4}{T^{\frac{1}{4}}} \right).
\end{aligned}$$

Here at (a),  $C_2$  denotes an integer constant.

Finally, summing up the regrets in **Case a** and **Case b**, the total regret is given by

$$\begin{aligned}
\mathbb{E}[R_T] &\leq \sum_{i \in \mathcal{A} : \Delta_i > b} \left\{ \frac{C_0 K^4}{T^{\frac{1}{4}}} + \left( \Delta_i + \frac{320 \sigma_i^2 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right) \right\} \\
&\quad + \sum_{i \in \mathcal{A} : 0 < \Delta_i \leq b} \frac{C_2 K^4}{T^{\frac{1}{4}}} + \max_{i \in \mathcal{A} : 0 < \Delta_i \leq b} \Delta_i T
\end{aligned}$$

where  $C_0, C_1, C_2$  are integer constants s.t.  $C_0 = C_1 + C_2$ .

### 3.6 Experiments

In this section, we conduct extensive empirical evaluations of EUCBV against several other popular MAB algorithms. We use expected cumulative regret as the metric of comparison. The comparison is conducted against the following algorithms: KLUCB+ (Garivier and Cappé, 2011), DMED (Honda and Takemura, 2010), MOSS (Audibert and Bubeck, 2009), UCB1 (Auer *et al.*, 2002a), UCB-Improved (Auer and Ortner, 2010), Median Elimination (Even-Dar *et al.*, 2006), Thompson Sampling (TS) (Agrawal and Goyal, 2011), OCUCB (Lattimore, 2015), Bayes-UCB (BU) (Kaufmann *et al.*, 2012) and UCB-V (Audibert *et al.*, 2009)<sup>1</sup>. The parameters of EUCBV algorithm for all the experiments are set as follows:  $\psi = \frac{T}{K^2}$  and  $\rho = 0.5$  (as in Corollary 1). Note that KLUCB+ empirically outperforms KLUCB (as shown in Garivier and Cappé (2011)).

---

<sup>1</sup>The implementation for KLUCB, Bayes-UCB and DMED were taken from Cappé *et al.* (2012)

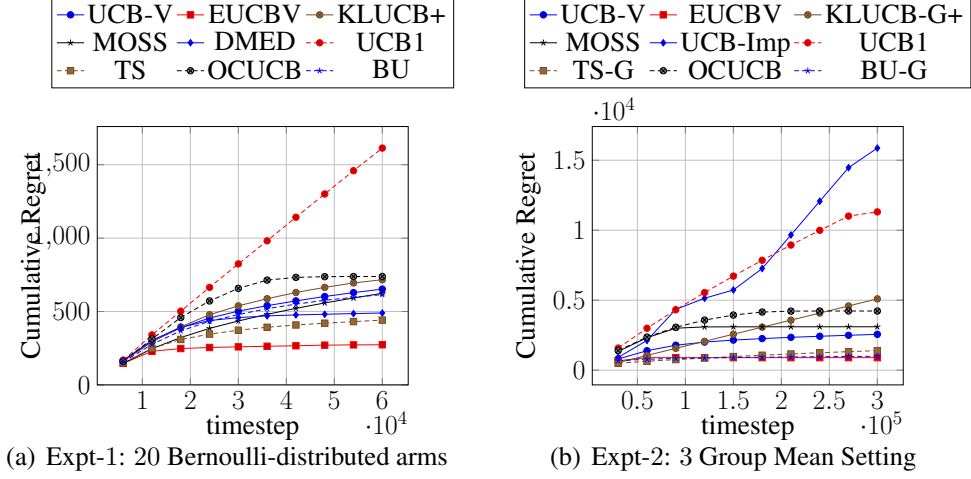


Figure 3.1: A comparison of the cumulative regret incurred by the various bandit algorithms.

**Experiment-1 (Bernoulli with uniform gaps):** This experiment is conducted to observe the performance of EUCBV over a short horizon. The horizon  $T$  is set to 60000. The testbed comprises of 20 Bernoulli distributed arms with expected rewards of the arms as  $r_{1:19} = 0.07$  and  $r_{20}^* = 0.1$  and these type of cases are frequently encountered in web-advertising domain (see Garivier and Cappé (2011)). The regret is averaged over 100 independent runs and is shown in Figure 3.1(a). EUCBV, MOSS, OCUCB, UCB1, UCB-V, KLUCB+, TS, BU and DMED are run in this experimental setup. Not only do we observe that EUCBV performs better than all the non-variance based algorithms such as MOSS, OCUCB, UCB-Improved and UCB1, but it also outperforms UCBV because of the choice of the exploration parameters. Because of the small gaps and short horizon  $T$ , we do not compare with UCB-Improved and Median Elimination for this test-case.

**Experiment-2 (Gaussian 3 Group Mean Setting):** This experiment is conducted to observe the performance of EUCBV over a large horizon in Gaussian distribution testbed. This setting comprises of a large horizon of  $T = 3 \times 10^5$  timesteps and a large set of arms. This testbed comprises of 100 arms involving Gaussian reward distributions with expected rewards of the arms in 3 groups,  $r_{1:66} = 0.07$ ,  $r_{67:99} = 0.01$  and  $r_{100}^* = 0.09$  with variance set as  $\sigma_{1:66}^2 = 0.01$ ,  $\sigma_{67:99}^2 = 0.25$  and  $\sigma_{100}^2 = 0.25$ . The regret is averaged over 100 independent runs and is shown in Figure 3.1(b). From the results in Figure 3.1(b), we observe that since the gaps are small and the variances of the optimal arm and the arms farthest from the optimal arm are the highest, EUCBV,



which allocates pulls proportional to the variances of the arms, outperforms all the non-variance based algorithms MOSS, OCUCB, UCB1, UCB-Improved and Median-Elimination ( $\epsilon = 0.1, \delta = 0.1$ ). The performance of Median-Elimination is extremely weak in comparison with the other algorithms and its plot is not shown in Figure 3.1(b). We omit its plot in order to more clearly show the difference between EUCBV, MOSS and OCUCB. Also note that the order of magnitude in the y-axis (cumulative regret) of Figure 3.1(b) is  $10^4$ . KLUCB-Gauss+ (denoted by KLUCB-G+), TS-G and BU-G are initialized with Gaussian priors. Both KLUCB-G+ and UCBV which is a variance-aware algorithm perform much worse than TS-G and EUCBV. The performance of DMED is similar to KLUCB-G+ in this setup and its plot is omitted.

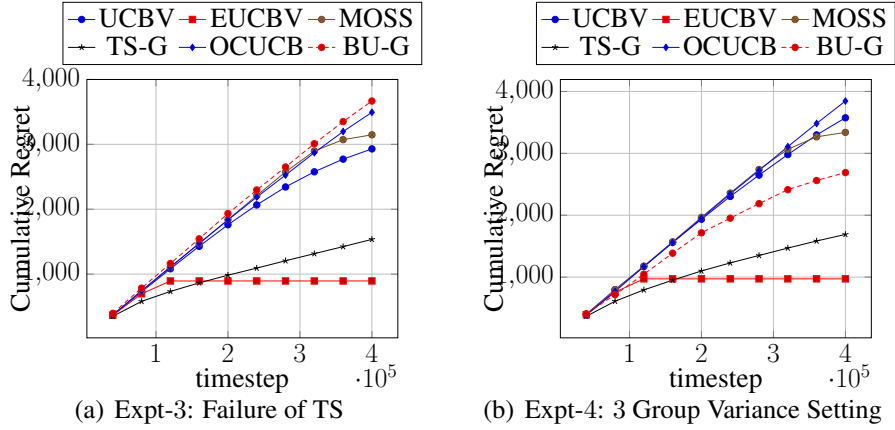


Figure 3.2: Further Experiments with EUCBV

**Experiment-3 (Failure of TS):** This experiment is conducted to demonstrate that in certain environments when the horizon is large, gaps are small and the variance of the optimal arm is high, the Bayesian algorithms (like TS) do not perform well but EUCBV performs exceptionally well. This experiment is conducted on 100 Gaussian distributed arms such that expected rewards of the arms  $r_{1:10} = 0.045$ ,  $r_{11:99} = 0.04$ ,  $r_{100}^* = 0.05$  and the variance is set as  $\sigma_{1:10}^2 = 0.01$ ,  $\sigma_{100}^2 = 0.25$  and  $T = 4 \times 10^5$ . The variance of the arms  $i = 11 : 99$  are chosen uniform randomly between  $[0.2, 0.24]$ . TS and BU with Gaussian priors fail because here the chosen variance values are such that only variance-aware algorithms with appropriate exploration factors will perform well or otherwise it will get bogged down in costly exploration. The algorithms that are not variance-aware will spend a significant amount of pulls trying to find the optimal arm. The result is shown in Figure 3.2(a). Predictably EUCBV, which allocates pulls proportional to the variance of the arms, outperforms its closest competitors TS-G,

BU-G, UCBV, MOSS and OCUCB. The plots for KLUCB-G+, DMED, UCB1, UCB-Improved and Median Elimination are omitted from the figure as their performance is extremely weak in comparison with other algorithms. We omit their plots to clearly show how EUCBV outperforms its nearest competitors. Note that EUCBV by virtue of its aggressive exploration parameters outperforms UCBV in all the experiments even though UCBV is a variance-based algorithm. The performance of TS-G is also weak and this is in line with the observation in Lattimore (2015) that the worst case regret of TS when Gaussian prior is used is  $\Omega(\sqrt{KT \log T})$ .

**Experiment-4 (Gaussian 3 Group Variance setting):** This experiment is conducted to show that when the gaps are uniform and variance of the arms are the only discriminative factor then the EUCBV performs extremely well over a very large horizon and over a large number of arms. This testbed comprises of 100 arms with Gaussian reward distributions, where the expected rewards of the arms are  $r_{1:99} = 0.09$  and  $r_{100}^* = 0.1$ . The variances of the arms are divided into 3 groups. The group 1 consist of arms  $i = 1 : 49$  where the variances are chosen uniform randomly between  $[0.0, 0.05]$ , group 2 consist of arms  $i = 50 : 99$  where the variances are chosen uniform randomly between  $[0.19, 0.24]$  and for the optimal arm  $i = 100$  (group 3) the variance is set as  $\sigma_*^2 = 0.25$ . We report the cumulative regret averaged over 100 independent runs. The horizon is set at  $T = 4 \times 10^5$  timesteps. We report the performance of MOSS, BU-G, UCBV, TS-G and OCUCB who are the closest competitors of EUCBV over this uniform gap setup. From the results in Figure 3.2(b), it is evident that the growth of regret for EUCBV is much lower than that of TS-G, MOSS, BU-G, OCUCB and UCBV. Because of the poor performance of KLUCB-G+ in the last two experiments we do not implement it in this setup. Also, note that for optimal performance BU-G, TS-G and KLUCB-G+ require the knowledge of the type of distribution to set their priors. Also, in all the experiments with Gaussian distributions EUCBV significantly outperforms all the Bayesian algorithms initialized with Gaussian priors.

### 3.7 Conclusion and Future Works

In this chapter, we studied the EUCBV algorithm which takes into account the empirical variance of the arms and employs aggressive exploration parameters in conjunction with non-uniform arm selection (as opposed to UCB-Improved) to eliminate sub-optimal arms. Our theoretical analysis conclusively established that EUCBV exhibits an order-optimal gap-independent regret bound of  $O\left(\sqrt{KT}\right)$ . Empirically, we show that EUCBV performs superbly across diverse experimental settings and outperforms most of the bandit algorithms in a stochastic MAB setup. Our experiments show that EUCBV is extremely stable for larger horizons and performs consistently well across different types of distributions. One avenue for future work is to remove the constraint of  $T \geq K^{2.4}$  required for EUCBV to reach the order optimal regret bound. Another future direction is to come up with an anytime version of EUCBV. An anytime algorithm does not need the horizon  $T$  as an input parameter.

# Chapter 4

## Thresholding Bandits

### 4.1 Introduction

In the previous chapters 2 and 3 we studied the stochastic multi-armed bandit (SMAB) setting with the goal of minimizing cumulative regret. In this chapter we will study another setting called Pure-exploration multi-armed bandits. An interested reader can read through the previous chapters or can continue from here. Though we re-use the ideas from SMABs, the goal of pure exploration setup is distinctly different from that of cumulative regret minimization of SMABs and the required algorithms to understand this setup are mentioned in this chapter itself. Pure-exploration MAB problems are unlike their traditional (exploration vs. exploitation) counterparts, the SMABs, where the objective is to minimize the cumulative regret. The cumulative regret is the total loss incurred by the learner for not playing the optimal arm throughout the time horizon  $T$ . In pure-exploration problems a learning algorithm, until time  $T$ , can invest entirely on exploring the arms without being concerned about the loss incurred while exploring; the objective is to minimize the probability that the arm recommended at time  $T$  is not the best arm. In this chapter, we further consider a combinatorial version of the pure-exploration MAB, called the thresholding bandit problem (TBP). Here, the learning algorithm is provided with a threshold  $\tau$ , and the objective, after exploring for  $T$  rounds, is to output all arms  $i$  whose  $r_i$  is above  $\tau$ . It is important to emphasize that the *thresholding* bandit problem is different from the *threshold* bandit setup studied in Abernethy *et al.* (2016), where the learner receives an unit reward whenever the value of an observation is above a threshold.

The rest of the chapter is organized as follows. We specify all the notations and assumptions in section 4.2. Then we define the problem statement for the TBP setting in section 4.3. In the next section 4.4 we discuss the motivations behind the TBP setting. In section 4.5 we discuss extensively on the various state-of-the-art algorithms available

for the pure exploration setting and then in section 4.6 we discuss the latest works done in the TBP setting. Finally, we draw our conclusions in section 4.7.

## 4.2 Notations

To benefit the reader, we again recall the notations we stated in chapter 2 and also a few additional notations.  $\mathcal{A}$  denotes the set of arms, and  $|\mathcal{A}| = K$  is the number of arms in  $\mathcal{A}$ . For arm  $i \in \mathcal{A}$ , we use  $r_i$  to denote the true mean of the distribution from which the rewards are sampled, while  $\hat{r}_i(t)$  denotes the estimated mean at time  $t$ . Formally, using  $n_i(t)$  to denote the number of times arm  $i$  has been pulled until time  $t$ , we have  $\hat{r}_i(t) = \frac{1}{n_i(t)} \sum_{z=1}^{n_i(t)} X_{i,z}$ , where  $X_{i,z}$  is the reward sample received when arm  $i$  is pulled for the  $z$ -th time. Similarly, we use  $\sigma_i^2$  to denote the true variance of the reward distribution corresponding to arm  $i$ , while  $\hat{v}_i(t)$  is the estimated variance, i.e.,  $\hat{v}_i(t) = \frac{1}{n_i(t)} \sum_{z=1}^{n_i(t)} (X_{i,z} - \hat{r}_i)^2$ . Whenever there is no ambiguity about the underlying time index  $t$ , for simplicity we neglect  $t$  from the notations and simply use  $\hat{r}_i$ ,  $\hat{v}_i$ , and  $n_i$ , to denote the respective quantities. Let  $\Delta_i = |\tau - r_i|$  denote the distance of the true mean from the threshold  $\tau$ . Also, the rewards are assumed to take values in  $[0, 1]$ .

## 4.3 Problem Definition

Formally, the problem we consider is the following. First, we define the set  $S_\tau = \{i \in \mathcal{A} : r_i \geq \tau\}$ . Note that,  $S_\tau$  is the set of all arms whose reward mean is greater than  $\tau$ . Let  $S_\tau^c$  denote the complement of  $S_\tau$ , i.e.,  $S_\tau^c = \{i \in \mathcal{A} : r_i < \tau\}$ . Next, let  $\hat{S}_\tau = \hat{S}_\tau(T) \subseteq \mathcal{A}$  denote the recommendation of a learning algorithm (under consideration) after  $T$  time units of exploration, while  $\hat{S}_\tau^c$  denotes its complement.

The performance of the learning agent is measured by the accuracy with which it can classify the arms into  $S_\tau$  and  $S_\tau^c$  after time horizon  $T$ . Equivalently, using  $\mathbb{I}(E)$  to denote the indicator of an event  $E$ , the *loss*  $\mathcal{L}(T)$  is defined as

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

Finally, the goal of the learning agent is to minimize the expected loss:

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

Note that the expected loss is simply the *probability of mis-classification* (i.e., error), that occurs either if a good arm is rejected or a bad arm is accepted as a good one.

## 4.4 Motivation

The above TBP formulation has several applications, for instance, from areas ranging from anomaly detection and classification (see Locatelli *et al.* (2016)) to industrial application. Particularly in industrial applications a learners objective is to choose (i.e., keep in operation) all machines whose productivity is above a threshold. The TBP also finds applications in mobile communications (see Audibert and Bubeck (2010)) where the users are to be allocated only those channels whose quality is above an acceptable threshold.

## 4.5 Related Work in Pure Exploration

Significant amount of literature is available on the stochastic MAB setting with respect to minimizing the cumulative regret. Chapter 2 and 3 deals with that. In this work we are particularly interested in *pure-exploration MABs*, where the focus is primarily on simple regret rather than the cumulative regret. 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. The pure exploration problem has been explored mainly under the following two settings:

### 4.5.1 Fixed Budget setting

Here the learning algorithm has to suggest the best arm(s) within a fixed time-horizon  $T$ , that is usually given as an input. The objective is to maximize the probability of returning the best arm(s). This is the scenario we consider in this chapter. Some of the important algorithms used in pure exploration setting are discussed in the next part.

#### UCB-Exploration Algorithm

---

##### Algorithm 9 UCBE

---

- 1: **Input:** The budget  $T$
  - 2: Pull each arm once
  - 3: **for**  $t = K + 1, \dots, T$  **do**
  - 4:     Pull the arm such that  $\arg \max_{i \in \mathcal{A}} \left\{ \hat{r}_i + \sqrt{\frac{a}{n_i}} \right\}$ , where  $a = \frac{25(T - K)}{36H_1}$  and
 
$$H_1 = \sum_{i=1}^K \frac{1}{\Delta_i^2}.$$
  - 5:      $t := t + 1$
  - 6: **end for**
- 

In Audibert and Bubeck (2010) the authors propose the UCBE and the Successive Reject (SR) algorithm, and prove simple-regret guarantees for the problem of identifying the single best arm. In the combinatorial fixed budget setup Gabillon *et al.* (2011) propose the GapE and GapE-V algorithms that suggest, with high probability, the best  $m$  arms at the end of the time budget.

#### Successive Reject Algorithm

Similarly, Bubeck *et al.* (2013) introduce the Successive Accept Reject (SAR) algorithm, which is an extension of the SR algorithm; SAR is a round based algorithm whereby at the end of each round an arm is either accepted or rejected (based on certain confidence conditions) until the top  $m$  arms are suggested at the end of the budget with high probability. A similar combinatorial setup was explored in Chen *et al.* (2014) where the authors propose the Combinatorial Successive Accept Reject (CSAR) algorithm, which is similar in concept to SAR but with a more general setup.

---

**Algorithm 10** Successive Reject(SR)

---

- 1: **Input:** The budget  $T$
  - 2: **Initialization:**  $n_0 = 0$
  - 3: **Definition:**  $\log^- K = \frac{1}{2} + \sum_{i=2}^K \frac{1}{i}$ ,  $n_k = \frac{1}{\log^- K} \frac{T - K}{K + 1 - m}$
  - 4: **for** For each phase  $m = 1, \dots, K - 1$  **do**
  - 5:     For each  $i \in B_m$ , select arm  $i$  for  $n_k - n_{k-1}$  rounds.
  - 6:     Let  $B_{m+1} = B_m \setminus \arg \min_{i \in B_m} \hat{r}_i$  (remove one element from  $B_m$ , if there is a tie, select randomly the arm to dismiss among the worst arms).
  - 7:      $m := m + 1$
  - 8: **end for**
  - 9: Output the single remaining  $i \in B_m$ .
- 

### 4.5.2 Fixed Confidence setting

In this setting the learning algorithm has to suggest the best arm(s) with a fixed confidence (given as input) with as fewer number of attempts as possible. The single best arm identification has been studied in Even-Dar *et al.* (2006), while for the combinatorial setup Kalyanakrishnan *et al.* (2012) have proposed the LUCB algorithm which, on termination, returns  $m$  arms which are at least  $\epsilon$  close to the true top- $m$  arms with probability at least  $1 - \delta$ . For a detail survey of this setup we refer the reader to Jamieson and Nowak (2014).

### 4.5.3 Unified Setting

Apart from these two settings some unified approaches 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. The thresholding bandit problem is a specific instance of the pure-exploration setup of Chen *et al.* (2014).

## 4.6 Related Work in Thresholding Bandits

In the latest work of Locatelli *et al.* (2016) Anytime Parameter-Free Thresholding (APT) algorithm comes up with an improved anytime guarantee than CSAR for the thresholding bandit problem.



---

**Algorithm 11** APT

---

**Input:** Time horizon  $T$ , threshold  $\tau$ , tolerance factor  $\epsilon \geq 0$

Pull each arm once

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

    Pull arm  $j \in \arg \min_{i \in A} \{ (|\hat{r}_i - \tau| + \epsilon) \sqrt{n_i} \}$  and observe the reward for arm  $j$ .

**end for**

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

---

## 4.7 Conclusion

In this chapter, we looked at the pure exploration MAB and thresholding bandit (TBP) setting which is a special case of combinatorial pure exploration MAB. We then looked at the various state-of-the-art algorithms in the literature for the pure-exploration setting and discussed the advantages and disadvantages of them. Then we looked at the latest algorithm for the TBP setting. The expected loss that has been proven for the said algorithms have also been discussed at length and their exploration parameters have also been compared against each other. In the next chapter, we provide our solution to this TBP setting which uses variance estimation to find the set of arms above the threshold.

# Chapter 5

## Augmented UCB for TBP

### 5.1 Introduction

In this chapter we look at the Augmented-UCB (AugUCB) algorithm for a fixed-budget version of the thresholding bandit problem (TBP), where the objective is to identify a set of arms whose expected mean is above a threshold. A key feature of AugUCB is that it uses both mean and variance estimates to eliminate arms that have been sufficiently explored; to the best of our knowledge this is the first algorithm to employ such an approach for the considered TBP. Theoretically, we obtain an upper bound on the loss (probability of mis-classification) incurred by AugUCB. Although UCBEV in literature provides a better guarantee, it is important to emphasize that UCBEV has access to problem complexity (whose computation requires arms' mean and variances), and hence is not realistic in practice; this is in contrast to AugUCB whose implementation does not require any such complexity inputs. We conduct extensive simulation experiments to validate the performance of AugUCB. Through our simulation work, we establish that AugUCB, owing to its utilization of variance estimates, performs significantly better than the state-of-the-art APT, CSAR and other non variance-based algorithms.

### 5.2 Our Contribution

We propose the Augmented UCB (AugUCB) algorithm for the fixed-budget setting of a specific combinatorial, pure-exploration, stochastic MAB called the thresholding bandit problem. AugUCB essentially combines the approach of UCB-Improved, CCB (Liu and Tsuruoka, 2016) and APT algorithms. Our algorithm takes into account the empirical variances of the arms along with mean estimates; to the best of our knowledge this is the first variance-based algorithm for the considered TBP. Thus, we also address an open problem discussed in Auer and Ortner (2010) of designing an algorithm that

can eliminate arms based on variance estimates. In this regard, note that both CSAR and APT are not variance-based algorithms.

Our theoretical contribution comprises proving an upper bound on the expected loss incurred by AugUCB (Theorem 2). In Table 5.1 we compare the upper bound on the losses incurred by the various algorithms, including AugUCB. The terms  $H_1$ ,  $H_2$ ,  $H_{CSAR,2}$ ,  $H_{\sigma,1}$  and  $H_{\sigma,2}$  represent various problem complexities, and are as defined in Section 5.4. From Section 5.4 we note that, for all  $K \geq 8$ , we have

$$\log(K \log K) H_{\sigma,2} > \log(2K) H_{\sigma,2} \geq H_{\sigma,1}.$$

Thus, it follows that the upper bound for UCBEV is better than that for AugUCB. However, implementation of UCBEV algorithm requires  $H_{\sigma,1}$  as input, whose computation is not realistic in practice. In contrast, our AugUCB algorithm requires no such complexity factor as input.

Proceeding with the comparisons, we emphasize that the upper bound for AugUCB is, in fact, not comparable with that of APT and CSAR; this is because the complexity term  $H_{\sigma,2}$  is not explicitly comparable with either  $H_1$  or  $H_{CSAR,2}$ . However, through extensive simulation experiments we find that AugUCB significantly outperforms both APT, CSAR and other non variance-based algorithms. AugUCB also outperforms UCBEV under explorations where non-optimal values of  $H_{\sigma,1}$  are used. In particular, we consider experimental scenarios comprising large number of arms, with the variances of arms in  $S_\tau$  being large. AugUCB, being variance based, exhibits superior

Table 5.1: AugUCB vs. State of the art

Algorithm	Upper Bound on Expected Loss
AugUCB	$\exp \left( -\frac{T}{4096 \log(K \log K) H_{\sigma,2}} + \log(2KT) \right)$
UCBEV	$\exp \left( -\frac{1}{512} \frac{T-2K}{H_{\sigma,1}} + \log(6KT) \right)$
APT	$\exp \left( -\frac{T}{64H_1} + 2 \log((\log(T) + 1)K) \right)$
CSAR	$\exp \left( -\frac{T-K}{72 \log(K) H_{CSAR,2}} + 2 \log(K) \right)$

performance under these settings.

The remainder of the paper is organized as follows. In section 5.3 we present our AugUCB algorithm. Section 5.4 contains our main theorem on expected loss, while section 5.5 contains simulation experiments. We finally draw our conclusions in section 5.6.

### 5.3 Augmented-UCB Algorithm

**The Algorithm:** The Augmented-UCB (AugUCB) algorithm is presented in Algorithm 12. AugUCB is essentially based on the arm elimination method of the UCB-Improved Auer and Ortner (2010), but adapted to the thresholding bandit setting proposed in Locatelli *et al.* (2016). However, unlike the UCB improved (which is based on mean estimation) our algorithm employs *variance estimates* (as in Audibert *et al.* (2009)) for arm elimination; to the best of our knowledge this is the first variance-aware algorithm for the thresholding bandit problem. Further, we allow for arm-elimination at each time-step, which is in contrast to the earlier work (e.g., Auer and Ortner (2010); Chen *et al.* (2014)) where the arm elimination task is deferred to the end of the respective exploration rounds. The details are presented below.

The active set  $B_0$  is initialized with all the arms from  $\mathcal{A}$ . We divide the entire budget  $T$  into rounds/phases like in UCB-Improved, CCB, SAR and CSAR. At every time-step AugUCB checks for arm elimination conditions, while updating parameters at the end of each round. As suggested by Liu and Tsuruoka (2016) to make AugUCB to overcome too much early exploration, we no longer pull all the arms equal number of times in each round. Instead, we choose an arm in the active set  $B_m$  that minimizes  $(|\hat{r}_i - \tau| - 2s_i)$  where

$$s_i = \sqrt{\frac{\rho\psi_m(\hat{v}_i + 1) \log(T\epsilon_m)}{4n_i}}$$

with  $\rho$  being the arm elimination parameter and  $\psi_m$  being the exploration regulatory factor. The above condition ensures that an arm closer to the threshold  $\tau$  is pulled; parameter  $\rho$  can be used to fine tune the elimination interval. The choice of exploration factor,  $\psi_m$ , comes directly from Audibert and Bubeck (2010) and Bubeck *et al.* (2011)

---

**Algorithm 12** AugUCB

---

**Input:** Time budget  $T$ ; parameter  $\rho$ ; threshold  $\tau$

**Initialization:**  $B_0 = \mathcal{A}$ ;  $m = 0$ ;  $\epsilon_0 = 1$ ;

$$M = \left\lfloor \frac{1}{2} \log_2 \frac{T}{e} \right\rfloor; \quad \psi_0 = \frac{T \epsilon_0}{128 \left( \log(\frac{3}{16} K \log K) \right)^2};$$
$$\ell_0 = \left\lceil \frac{2 \psi_0 \log(T \epsilon_0)}{\epsilon_0} \right\rceil; \quad N_0 = K \ell_0$$

Pull each arm once

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

    Pull arm  $j \in \arg \min_{i \in B_m} \left\{ |\hat{r}_i - \tau| - 2s_i \right\}$

$t \leftarrow t + 1$

**for**  $i \in B_m$  **do**

**if**  $(\hat{r}_i + s_i < \tau - s_i)$  or  $(\hat{r}_i - s_i > \tau + s_i)$  **then**

$B_m \leftarrow B_m \setminus \{i\}$  (Arm deletion)

**end if**

**end for**

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

**Reset Parameters**

$\epsilon_{m+1} \leftarrow \frac{\epsilon_m}{2}$

$B_{m+1} \leftarrow B_m$

$\psi_{m+1} \leftarrow \frac{T \epsilon_{m+1}}{128 \left( \log(\frac{3}{16} K \log K) \right)^2}$

$\ell_{m+1} \leftarrow \left\lceil \frac{2 \psi_{m+1} \log(T \epsilon_{m+1})}{\epsilon_{m+1}} \right\rceil$

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

$m \leftarrow m + 1$

**end if**

**end for**

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

---

where it is stated that in pure exploration setup, the exploring factor must be linear in  $T$  (so that an exponentially small probability of error is achieved) rather than being logarithmic in  $T$  (which is more suited for minimizing cumulative regret).

## 5.4 Theoretical Results

### Problem Complexity

Let us begin by recalling the following definitions of the *problem complexity* as introduced in Locatelli *et al.* (2016):

$$H_1 = \sum_{i=1}^K \frac{1}{\Delta_i^2} \text{ and } H_{CSAR,2} = \min_{i \in \mathcal{A}} \frac{i}{\Delta_{(i)}^2}$$

where  $(\Delta_{(i)} : i \in \mathcal{A})$  is obtained by arranging  $(\Delta_i : i \in \mathcal{A})$  in an increasing order. Also, from Chen *et al.* (2014) we have

$$H_{CSAR,2} = \max_{i \in \mathcal{A}} \frac{i}{\Delta_{(i)}^2}.$$

$H_{CSAR,2}$  is the complexity term appearing in the bound for the CSAR algorithm. The relation between the above complexity terms are as follows (see Locatelli *et al.* (2016)):

$$H_1 \leq \log(2K)H_2 \text{ and } H_1 \leq \log(K)H_{CSAR,2}.$$

As ours is a variance-aware algorithm, we require  $H_1^\sigma$  (as defined in Gabillon *et al.* (2011)) that incorporates reward variances into its expression as given below:

$$H_{\sigma,1} = \sum_{i=1}^K \frac{\sigma_i + \sqrt{\sigma_i^2 + (16/3)\Delta_i}}{\Delta_i^2}.$$

Finally, analogous to  $H_{CSAR,2}$ , in this paper we introduce the complexity term  $H_{\sigma,2}$ , which is given by

$$H_{\sigma,2} = \max_{i \in \mathcal{A}} \frac{i}{\tilde{\Delta}_{(i)}^2}$$

where  $\tilde{\Delta}_i^2 = \frac{\Delta_i^2}{\sigma_i + \sqrt{\sigma_i^2 + (16/3)\Delta_i}}$ , and  $(\tilde{\Delta}_{(i)})$  is an increasing ordering of  $(\tilde{\Delta}_i)$ . Following the results in Audibert and Bubeck (2010), we can show that

$$H_{\sigma,2} \leq H_{\sigma,1} \leq \overline{\log}(K)H_{\sigma,2} \leq \log(2K)H_{\sigma,2}.$$

## Proof of expected loss of AugUCB

Our main result is summarized in the following theorem where we prove an upper bound on the expected loss.

**Theorem 2** *For  $K \geq 4$  and  $\rho = 1/3$ , the expected loss of the AugUCB algorithm is given by,*

$$\mathbb{E}[\mathcal{L}(T)] \leq 2KT \exp\left(-\frac{T}{4096 \log(K \log K) H_{\sigma,2}}\right).$$

**Proof 2 Discussion 3** *The proof comprises of two modules. In the first module we investigate the necessary conditions for arm elimination within a specified number of rounds, which is motivated by the technique in Auer and Ortner (2010). Bounds on the arm-elimination probability is then obtained; however, since we use variance estimates, we invoke the Bernstein inequality (as in Audibert et al. (2009)) rather than the Chernoff-Hoeffding bounds (which is appropriate for the UCB-Improved (Auer and Ortner, 2010)). In the second module, as in Locatelli et al. (2016), we first define a favourable event that will yield an upper bound on the expected loss. Using union bound, we then incorporate the result from module-1 (on the arm elimination probability), and finally derive the result through a series of simplifications.*

*The details of the proof as stated in discussion 3 are as follows.*

**Arm Elimination:** Recall the notations used in the algorithm, Also, for each arm  $i \in \mathcal{A}$ , define  $m_i = \min \{m | \sqrt{\rho \epsilon_m} < \frac{\Delta_i}{2}\}$ . In the  $m_i$ -th round, whenever  $n_i = \ell_{m_i} \geq \frac{2\psi_{m_i} \log(T\epsilon_{m_i})}{\epsilon_{m_i}}$ , we obtain (as  $\hat{v}_i \in [0, 1]$ )

$$s_i \leq \sqrt{\frac{\rho(\hat{v}_i + 1)\epsilon_{m_i}}{8}} \leq \frac{\sqrt{\rho\epsilon_{m_i}}}{2} < \frac{\Delta_i}{4}. \quad (5.1)$$

First, let us consider a bad arm  $i \in \mathcal{A}$  (i.e.,  $r_i < \tau$ ). We note that, in the  $m_i$ -th round whenever  $\hat{r}_i \leq r_i + 2s_i$ , then arm  $i$  is eliminated as a bad arm. This is easy to verify as

follows: using (5.1) we obtain,

$$\hat{r}_i \leq r_i + 2s_i < r_i + \Delta_i - 2s_i = \tau - 2s_i$$

which is precisely one of the elimination conditions in Algorithm 12. Thus, the probability that a bad arm is not eliminated correctly in the  $m_i$ -th round (or before) is given by

$$\mathbb{P}(\hat{r}_i > r_i + 2s_i) \leq \mathbb{P}(\hat{r}_i > r_i + 2\bar{s}_i) + \mathbb{P}(\hat{v}_i \geq \sigma_i^2 + \sqrt{\rho\epsilon_{m_i}}) \quad (5.2)$$

where

$$\bar{s}_i = \sqrt{\frac{\rho\psi_{m_i}(\sigma_i^2 + \sqrt{\rho\epsilon_{m_i}} + 1) \log(T\epsilon_{m_i})}{4n_i}}$$

Note that, substituting  $n_i = \ell_{m_i} \geq \frac{2\psi_{m_i} \log(T\epsilon_{m_i})}{\epsilon_{m_i}}$ ,  $\bar{s}_i$  can be simplified to obtain,

$$2\bar{s}_i \leq \frac{\sqrt{\rho\epsilon_{m_i}(\sigma_i^2 + \sqrt{\rho\epsilon_{m_i}} + 1)}}{2} \leq \sqrt{\rho\epsilon_{m_i}}. \quad (5.3)$$

The first term in the LHS of (3.2) can be bounded using the Bernstein inequality as below:

$$\begin{aligned} & \mathbb{P}(\hat{r}_i > r_i + 2\bar{s}_i) \\ & \leq \exp\left(-\frac{(2\bar{s}_i)^2 n_i}{2\sigma_i^2 + \frac{4}{3}\bar{s}_i}\right) \\ & \leq \exp\left(-\frac{\rho\psi_{m_i}(\sigma_i^2 + \sqrt{\rho\epsilon_{m_i}} + 1) \log(T\epsilon_{m_i})}{2\sigma_i^2 + \frac{2}{3}\sqrt{\rho\epsilon_{m_i}}}\right) \\ & \stackrel{(a)}{\leq} \exp\left(-\frac{3\rho T\epsilon_{m_i}}{256a^2} \left(\frac{\sigma_i^2 + \sqrt{\rho\epsilon_{m_i}} + 1}{3\sigma_i^2 + \sqrt{\rho\epsilon_{m_i}}}\right) \log(T\epsilon_{m_i})\right) \\ & := \exp(-Z_i) \end{aligned} \quad (5.4)$$

where, for simplicity, we have used  $\alpha_i$  to denoted the exponent in the inequality (a).

Also, note that (a) is obtained by using  $\psi_{m_i} = \frac{T\epsilon_{m_i}}{128a^2}$ , where  $a = (\log(\frac{3}{16}K \log K))$ .



The second term in the LHS of (5.2) can be simplified as follows:

$$\begin{aligned}
& \mathbb{P}\left\{\hat{v}_i \geq \sigma_i^2 + \sqrt{\rho\epsilon_{m_i}}\right\} \\
& \leq \mathbb{P}\left\{\frac{1}{n_i} \sum_{t=1}^{n_i} (X_{i,t} - r_i)^2 - (\hat{r}_i - r_i)^2 \geq \sigma_i^2 + \sqrt{\rho\epsilon_{m_i}}\right\} \\
& \leq \mathbb{P}\left\{\frac{\sum_{t=1}^{n_i} (X_{i,t} - r_i)^2}{n_i} \geq \sigma_i^2 + \sqrt{\rho\epsilon_{m_i}}\right\} \\
& \stackrel{(a)}{\leq} \mathbb{P}\left\{\frac{\sum_{t=1}^{n_i} (X_{i,t} - r_i)^2}{n_i} \geq \sigma_i^2 + 2\bar{s}_i\right\} \\
& \stackrel{(b)}{\leq} \exp\left(-\frac{3\rho\psi_{m_i}}{2} \left(\frac{\sigma_i^2 + \sqrt{\rho\epsilon_{m_i}} + 1}{3\sigma_i^2 + \sqrt{\rho\epsilon_{m_i}}}\right) \log(T\epsilon_{m_i})\right) \\
& = \exp(-Z_i)
\end{aligned} \tag{5.5}$$

where inequality (a) is obtained using (5.3), while (b) follows from the Bernstein inequality.

Thus, using (5.4) and (5.5) in (5.2) we obtain  $\mathbb{P}(\hat{r}_i > r_i + 2s_i) \leq 2\exp(-Z_i)$ . Proceeding similarly, for a good arm  $i \in \mathcal{A}$ , the probability that it is not correctly eliminated in the  $m_i$ -th round (or before) is also bounded by  $\mathbb{P}(\hat{r}_i < r_i - 2s_i) \leq 2\exp(-Z_i)$ . In general, for any  $i \in \mathcal{A}$  we have

$$\mathbb{P}(|\hat{r}_i - r_i| > 2s_i) \leq 4\exp(-Z_i). \tag{5.6}$$

**Favourable Event:** Following the notation in Locatelli et al. (2016) we define the event

$$\xi = \left\{ \forall i \in \mathcal{A}, \forall t = 1, 2, \dots, T : |\hat{r}_i - r_i| \leq 2s_i \right\}.$$

Note that, on  $\xi$  each arm  $i \in \mathcal{A}$  is eliminated correctly in the  $m_i$ -th round (or before). Thus, it follows that  $\mathbb{E}[\mathcal{L}(T)] \leq P(\xi^c)$ . Since  $\xi^c$  can be expressed as an union of the events  $(|\hat{r}_i - r_i| > 2s_i)$  for all  $i \in \mathcal{A}$  and all  $t = 1, 2, \dots, T$ , using union bound we can write

$$\mathbb{E}[\mathcal{L}(T)] \leq \sum_{i \in \mathcal{A}} \sum_{t=1}^T \mathbb{P}(|\hat{r}_i - r_i| > 2s_i)$$

$$\begin{aligned}
&\leq \sum_{i \in \mathcal{A}} \sum_{t=1}^T 4 \exp(-Z_i) \\
&\leq 4T \sum_{i \in \mathcal{A}} \exp \left( -\frac{3\rho T \epsilon_{m_i}}{256a^2} \left( \frac{\sigma_i^2 + \sqrt{\rho \epsilon_{m_i}} + 1}{3\sigma_i^2 + \sqrt{\rho \epsilon_{m_i}}} \right) \log(T \epsilon_{m_i}) \right) \\
&\stackrel{(a)}{\leq} 4T \sum_{i \in \mathcal{A}} \exp \left( -\frac{3T \Delta_i^2}{4096a^2} \left( \frac{4\sigma_i^2 + \Delta_i + 4}{12\sigma_i^2 + \Delta_i} \right) \log\left(\frac{3}{16} T \Delta_i^2\right) \right) \\
&\stackrel{(b)}{\leq} 4T \sum_{i \in \mathcal{A}} \exp \left( -\frac{12T \Delta_i^2}{(12\sigma_i + 12\Delta_i)} \frac{\log(\frac{3}{16} K \log K)}{4096a^2} \right) \\
&\stackrel{(c)}{\leq} 4T \sum_{i \in \mathcal{A}} \exp \left( -\frac{T \Delta_i^2 \log(\frac{3}{16} K \log K)}{4096(\sigma_i + \sqrt{\sigma_i^2 + (16/3)\Delta_i})a^2} \right) \\
&\stackrel{(d)}{\leq} 4T \sum_{i \in \mathcal{A}} \exp \left( -\frac{T \log(\frac{3}{16} K \log K)}{4096 \tilde{\Delta}_i^{-2} a^2} \right) \\
&\stackrel{(e)}{\leq} 4T \sum_{i \in \mathcal{A}} \exp \left( -\frac{T \log(\frac{3}{16} K \log K)}{4096 \max_j(j \tilde{\Delta}_{(j)}^{-2}) (\log(\frac{3}{16} K \log K))^2} \right) \\
&\stackrel{(f)}{\leq} 4KT \exp \left( -\frac{T}{4096 \log(K \log K) H_{\sigma,2}} \right).
\end{aligned}$$

The justification for the above simplifications are as follows:

- (a) is obtained by noting that in round  $m_i$  we have  $\frac{\Delta_i}{4} \leq \sqrt{\rho \epsilon_{m_i}} < \frac{\Delta_i}{2}$ .
- For (b), we note that the function  $x \mapsto x \exp(-Cx^2)$ , where  $x \in [0, 1]$ , is decreasing on  $[1/\sqrt{2C}, 1]$  for any  $C > 0$  (see Bubeck et al. (2011); Auer and Ortner (2010)). Thus, using  $C = \lfloor T/e \rfloor$  and  $\min_{j \in \mathcal{A}} \Delta_j = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{e}{T}}$ , we obtain (b).
- To obtain (c) we have used the inequality  $\Delta_i \leq \sqrt{\sigma_i^2 + (16/3)\Delta_i}$  (which holds because  $\Delta_i \in [0, 1]$ ).
- (d) is obtained simply by substituting  $\tilde{\Delta}_i = \frac{\Delta_i^2}{\sigma_i + \sqrt{\sigma_i^2 + (16/3)\Delta_i}}$  and  $a = \log(\frac{3}{16} K \log K)$ .
- Finally, to obtain (e) and (f), note that  $\tilde{\Delta}_i^{-2} \leq i \tilde{\Delta}_i^{-2} \leq \max_{j \in \mathcal{A}} j \Delta_{(j)}^{-2} = H_{\sigma,2}$ .

## 5.5 Numerical Experiments

In this section, we empirically compare the performance of AugUCB against APT, UCBE, UCBEV, CSAR and the uniform-allocation (UA) algorithms. A brief note about these algorithms are as follows:

- APT: This algorithm is from Locatelli *et al.* (2016); we set  $\epsilon = 0.05$ , which is the margin-of-error within which APT suggests the set of good arms.

- AugUCB: This is the Augmented-UCB algorithm proposed in this paper; as in Theorem 2 we set  $\rho = \frac{1}{3}$ .

- UCBE: This is a modification of the algorithm in Audibert *et al.* (2009) (as it was originally proposed for the best arm identification problem); here, we set  $a = \frac{T-K}{H_1}$ , and at each time-step an arm  $i \in \arg \min \left\{ |\hat{r}_i - \tau| - \sqrt{\frac{a}{n_i}} \right\}$  is pulled.

- UCBEV: This is a modification of the algorithm in Gabillon *et al.* (2011) (proposed for the TopM problem); its implementation is identical to UCBE, but with  $a = \frac{T-2K}{H_{\sigma,1}}$ . As mentioned earlier, note that UCBEV's implementation would not be possible in real scenarios, as it requires computing the problem complexity  $H_{\sigma,1}$ . However, for theoretical reasons we show the best performance achievable by UCBEV. In experiment 6 we perform further explorations of UCBEV with alternate settings of  $a$ .

- CSAR: Modification of the successive-reject algorithm in Chen *et al.* (2014); here, we reject the arm farthest from  $\tau$  after each round.

- UA: The naive strategy where at each time-step an arm is uniformly sampled from  $\mathcal{A}$  (the set of all arms); however, UA is known to be optimal if all arms are equally difficult to classify.

Motivated by the settings considered in Locatelli *et al.* (2016), we design six different experimental scenarios that are obtained by varying the arm means and variances. Across all experiments consists of  $K = 100$  arms (indexed  $i = 1, 2, \dots, 100$ ) of which  $S_\tau = \{6, 7, \dots, 10\}$ , where we have fixed  $\tau = 0.5$ . In all the experiments, each algorithm is run independently for 10000 time-steps. At every time-step, the output set,  $\hat{S}_\tau$ , suggested by each algorithm is recorded; the output is counted as an error if  $\hat{S}_\tau \neq S_\tau$ . In Figure 1, for each experiment, we have reported the percentage of error incurred by the different algorithms as a function of time; Error percentage is obtained by repeating each experiment independently for 500 iterations, and then respectively computing the fraction of errors. The details of the considered experiments are as follows.

**Experiment-1:** The reward distributions are Gaussian with means  $r_{1:4} = 0.2 + (0 : 3) \cdot 0.05$ ,  $r_5 = 0.45$ ,  $r_6 = 0.55$ ,  $r_{7:10} = 0.65 + (0 : 3) \cdot 0.05$  and  $r_{11:100} = 0.4$ . Thus, the

means of the first 10 arms follow an arithmetic progression. The remaining arms have identical means; this setting is chosen because now a significant budget is required in exploring these arms, thus increasing the problem complexity.

The corresponding variances are  $\sigma_{1:5}^2 = 0.5$  and  $\sigma_{6:10}^2 = 0.6$ , while  $\sigma_{11:100}^2$  is chosen independently and uniform in the interval  $[0.38, 0.42]$ ; note that, the variances of the arms in  $S_\tau$  are higher than those of the other arms. The corresponding results are shown in Figure 5.1(a), from where we see that UCBEV, which has access to the problem complexity while being variance-aware, outperforms all other algorithm (including UCBE which also has access to the problem complexity but does not take into account the variances of the arms). Interestingly, the performance of our AugUCB (without requiring any complexity input) is comparable with UCBEV, while it outperforms UCBE, APT and the other non variance-aware algorithms that we have considered.

**Experiment-2:** We again consider Gaussian reward distributions. However, here the means of the first 10 arms constitute a geometric progression. Formally, the reward means are  $r_{1:4} = 0.4 - (0.2)^{1:4}$ ,  $r_5 = 0.45$ ,  $r_6 = 0.55$ ,  $r_{7:10} = 0.6 + (0.2)^{5-(1:4)}$  and  $r_{11:100} = 0.4$ ; the arm variances are as in experiment-1. The corresponding results are shown in Figure 5.1(b). We again observe AugUCB outperforming the other algorithms, except UCBEV.

**Experiment-3:** Here, the first 10 arms are partitioned into three groups, with all arms in a group being assigned the same mean; the reward distributions are again Gaussian. Specifically, the reward means are  $r_{1:3} = 0.1$ ,  $r_{4:7} = \{0.35, 0.45, 0.55, 0.65\}$  and  $r_{8:10} = 0.9$ ; as before,  $r_{11:100} = 0.4$  and all the variances are as in Experiment-1. The results for this scenario are presented in Figure 5.1(c). The observations are inline with those made in the previous experiments.

**Experiment-4:** The setting is similar to that considered in Experiment-3, but with the first 10 arms partitioned into two groups; the respective means are  $r_{1:5} = 0.45$ ,  $r_{6:10} = 0.55$ . The corresponding results are shown in Figure 5.1(d), from where the good performance of AugUCB is again validated.

**Experiment-5:** This is again the two group setting involving Gaussian reward distributions. The reward means are as in Experiment-4, while the variances are  $\sigma_{1:5}^2 = 0.3$

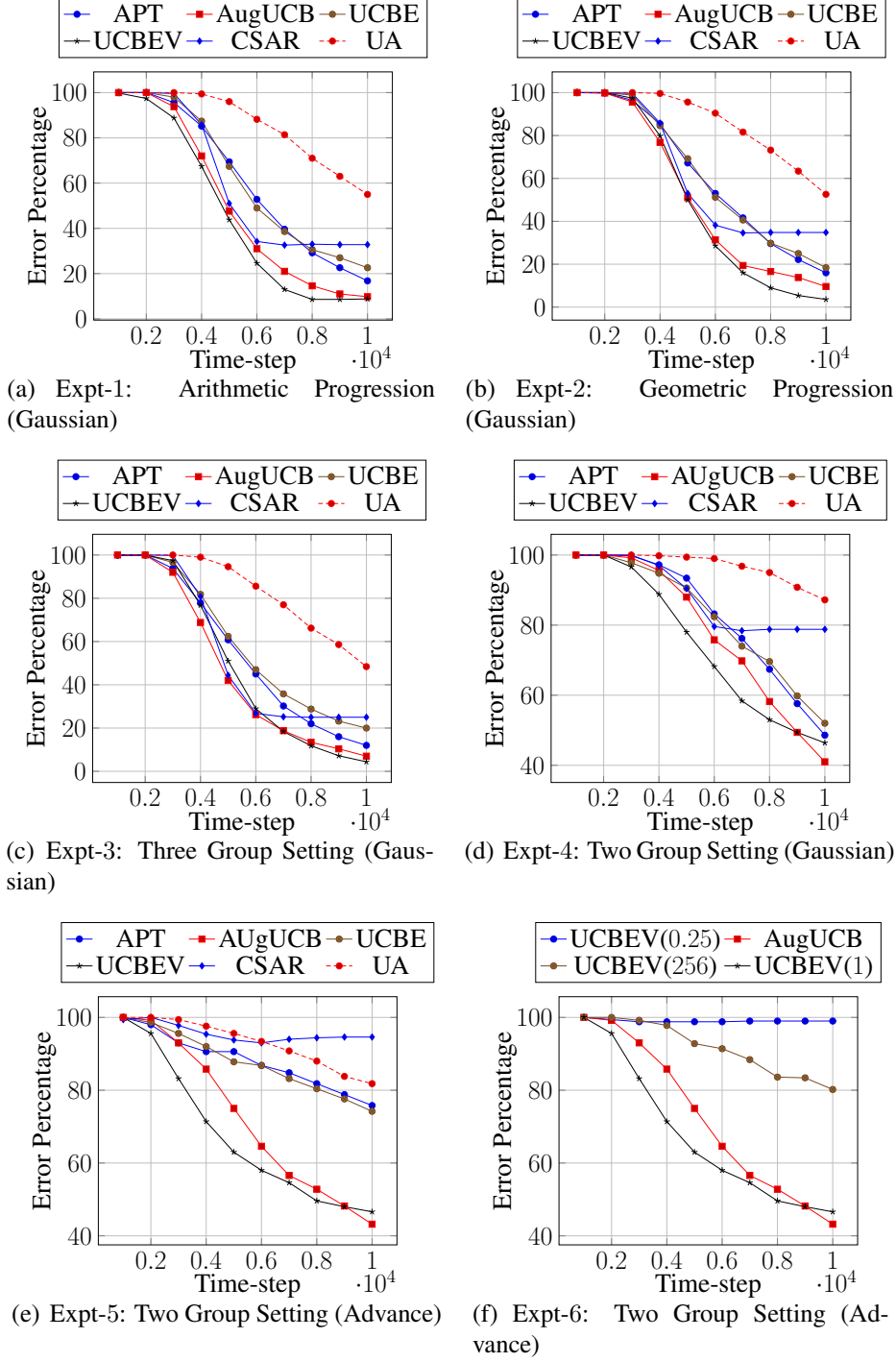


Figure 5.1: Performances of the various TBP algorithms in terms of error percentage vs. time-step, for six different experimental scenarios.

and  $\sigma_{6:10}^2 = 0.8$ ;  $\sigma_{11:100}^2$  are independently and uniformly chosen in the interval  $[0.2, 0.3]$ . The corresponding results are shown in Figure 5.1(e). We refer to this setup as *Advanced* because here the chosen variance values are such that only variance-aware algorithms will perform well. Hence, we see that UCBEV performs very well in comparison with the other algorithms. However, it is interesting to note that the performance of

AugUCB catches-up with UCBEV as the time-step increases, while significantly outperforming the other non-variance aware algorithms.

**Experiment-6:** We use the same setting as in Experiment-5, but conduct more exploration of UCBEV with different values of the exploration parameter  $a$ . The corresponding results are shown in Figure 5.1(f). As studied in Locatelli *et al.* (2016), we implement UCBEV with  $a_i = 4^i \frac{T-2K}{H_{\sigma,1}}$  for  $i = -1, 0, 4$ . Here,  $a_0$  corresponds to UCBEV(1) (in Figure 5.1(f)) which is UCBEV run with the optimal choice of  $H_{\sigma,1}$ . For other choices of  $a_i$  we see that UCBEV( $a_i$ ) is significantly outperformed by AugUCB.

Finally, note that in all the above experiments, the CSAR algorithm, although performs well initially, quickly exhausts its budget and saturates at a higher error percentage. This is because it pulls all arms equally in each round, with the round lengths being non-adaptive.

## 5.6 Conclusion and Future Works

We proposed the AugUCB algorithm for a fixed-budget, pure-exploration TBP. Our algorithm employs both mean and variance estimates for arm elimination. This, to our knowledge is the first variance-based algorithm for the specific TBP that we have considered. We first prove an upper bound on the expected loss incurred by AugUCB. We then conduct simulation experiments to validate the performance of AugUCB. In comparison with APT, CSAR and other non variance-based algorithms, we find that the performance of AugUCB is significantly better. Further, the performance of AugUCB is comparable with UCBEV (which is also variance-based), although the latter exhibits a slightly better performance. However, UCBEV is not implementable in practice as it requires computing problem complexity,  $H_{\sigma,1}$ , while AugUCB (requiring no such inputs) can be easily deployed in real-life scenarios. It would be an interesting future work to design an anytime version of the AugUCB algorithm.

## **Appendix A**

### **Appendix on Concentration Inequalities**





## **A.0.1 Sample space**

## **A.0.2 Events**

## **A.0.3 Sigma-algebra**

### **Borel-sigma algebra**

## **A.0.4 Measure**

### **The probability measure**

## **A.0.5 The triplet**

## **A.0.6 Filtration**

# **A.1 Martingale**

## **A.1.1 Super-martingale**

## **A.1.2 Sub-martingale**

# **A.2 Convergence theorems**

## **A.2.1 Monotone convergence theorem**

## **A.2.2 Dominated convergence theorem**

## **A.2.3 Fatou's Lemma**

# **A.3 Sub-Gaussian distribution**

Let a random variable  $X \in \mathbb{R}$  with variance as  $\sigma^2$ . Then  $X$  is said to be  $\sigma$ -sub-gaussian for  $\sigma \geq 0$  such that  $\mathbb{E}[X] = 0$  and its moment generating function satisfies for all  $\lambda \in \mathbb{R}$

the following condition,

$$\mathbb{E}[\exp \lambda X] \leq \exp \left( -\frac{\lambda^2 \sigma^2}{2} \right)$$

Also, note that sub-gaussian distribution is a class of distribution rather than a distribution itself.

**Remark 1** A random variable  $X \in [0, 1]$  is said to be  $\frac{1}{2}$ -sub-gaussian with its moment generating function satisfying the condition,

$$\mathbb{E}[\exp \lambda X] \leq \exp \left( -\frac{\lambda^2}{8} \right), \forall \lambda \in \mathbb{R}$$

## A.4 Concentration Inequalities

In this section we state some of the concentration inequalities used in the proofs in several chapters of the thesis. Concentration inequality deals with the control of the tail of the average of independent random variables from their expected mean.

Let,  $X_1, X_2, \dots, X_n$  be a sequence of independent random variables defined on a probability space  $(\omega, \mathcal{F}, \mathbb{P})$ , is bounded in  $[a_i, b_i], \forall i = 1, 2, \dots, n$ . Let  $S_n$  denote the sum of the random variables such that  $S_n = X_1 + X_2 + \dots + X_n$ ,  $\hat{r} = \frac{S_n}{n}$  and  $E[S_n] = r$ . Let  $\mathcal{F}_n$  be an increasing sequence of  $\sigma$ -fields of  $\mathcal{F}$  such that for each  $n$ ,  $\sigma(X_1, \dots, X_n) \subset \mathcal{F}_t$  and for  $q > t$ ,  $X_q$  is independent of  $\mathcal{F}_n$ .

### A.4.1 Markov's inequality

Markov's inequality states that, for any  $\epsilon > 0$ ,

$$\mathbb{P}[S_n > \epsilon] \leq \frac{\mathbb{E}[S_n]}{\epsilon}$$

### A.4.2 Chernoff-Hoeffding Bound

Chernoff-Hoeffding gives us the following inequality regarding the sums of independent random variables  $S_n$  and their deviation from their expectation  $\mathbb{E}[S_n] = r$ , for any  $\epsilon > 0$ ,

$$\begin{aligned}\mathbb{P}\{S_n - n\mathbb{E}[S_n] \geq \epsilon\} &\leq \exp\left(-\frac{2\epsilon^2}{n \sum_{i=1}^n (a_i - b_i)}\right) \\ &= \mathbb{P}\{S_n - n\mathbb{E}[S_n] \leq -\epsilon\} \leq \exp\left(-\frac{2\epsilon^2}{n \sum_{i=1}^n (a_i - b_i)}\right)\end{aligned}$$

Considering all the random variables bounded in  $[0, 1]$ , the above two inequalities can be reduced to,

$$\begin{aligned}\mathbb{P}\left\{\left|\frac{S_n}{n} - \mathbb{E}[S_n]\right| \geq \epsilon\right\} &\leq 2 \exp(-2\epsilon^2 n) \\ &= \mathbb{P}\{|\hat{r} - r| \geq \epsilon\} \leq 2 \exp(-2\epsilon^2 n)\end{aligned}$$

### A.4.3 Empirical Bernstein inequality

Similar to Chernoff-Hoeffding bound, empirical Bernstein inequality gives us the following inequality regarding the sums of independent random variables  $S_n$  and their deviation from their expectation  $\mathbb{E}[S_n] = r$ , for any  $\epsilon > 0$ ,

$$\begin{aligned}\mathbb{P}\{S_n - n\mathbb{E}[S_n] \geq \epsilon\} &\leq \exp\left(-\frac{2\epsilon^2}{\left(2\sigma^2 + \frac{2b_{\max}\epsilon}{3}\right) n \sum_{i=1}^n (a_i - b_i)}\right), \\ \mathbb{P}\{S_n - n\mathbb{E}[S_n] \leq -\epsilon\} &\leq \exp\left(-\frac{2\epsilon^2}{\left(2\sigma^2 + \frac{2b_{\max}\epsilon}{3}\right) n \sum_{i=1}^n (a_i - b_i)}\right)\end{aligned}$$

Considering all the random variables bounded in  $[0, 1]$ , the above two inequalities can be reduced to,

$$\begin{aligned}
\mathbb{P}\left\{\left|\frac{S_n}{n} - \mathbb{E}[S_n]\right| \geq \epsilon\right\} &\leq 2 \exp\left(-\frac{2\epsilon^2 n}{(2\sigma^2 + \frac{2\epsilon}{3})}\right) \\
&= \mathbb{P}\{|\hat{r} - r| \geq \epsilon\} \leq 2 \exp\left(-\frac{2\epsilon^2 n}{(2\sigma^2 + \frac{2\epsilon}{3})}\right)
\end{aligned}$$

## Appendix B

### Appendix for EUCEB

#### B.0.4 Proof of Lemma 1

**Lemma 1** *If  $T \geq K^{2.4}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$  and  $m \leq \frac{1}{2} \log_2 \left( \frac{T}{e} \right)$ , then,*

$$\frac{\rho m \log(2)}{\log(\psi T) - 2m \log(2)} \leq \frac{3}{2}.$$

**Proof 3** *The proof is based on contradiction. Suppose*

$$\frac{\rho m \log(2)}{\log(\psi T) - 2m \log(2)} > \frac{3}{2}.$$

*Then, with  $\psi = \frac{T}{K^2}$  and  $\rho = \frac{1}{2}$ , we obtain*

$$\begin{aligned} 6 \log(K) &> 6 \log(T) - 7m \log(2) \\ &\stackrel{(a)}{\geq} 6 \log(T) - \frac{7}{2} \log_2 \left( \frac{T}{e} \right) \log(2) \\ &= 2.5 \log(T) + 3.5 \log_2(e) \log(2) \\ &\stackrel{(b)}{=} 2.5 \log(T) + 3.5 \end{aligned}$$

*where (a) is obtained using  $m \leq \frac{1}{2} \log_2 \left( \frac{T}{e} \right)$ , while (b) follows from the identity  $\log_2(e) \log(2) = 1$ . Finally, for  $T \geq K^{2.4}$  we obtain,  $6 \log(K) > 6 \log(K) + 3.5$ , which is a contradiction. ■*

#### B.0.5 Proof of Lemma 2

**Lemma 2** *If  $T \geq K^{2.4}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$ ,  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$  and  $c_i = \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}}$ , then,*

$$c_i < \frac{\Delta_i}{4}$$

**Proof 4** In the  $m_i$ -th round since  $z_i \geq n_{m_i}$ , by substituting  $z_i$  with  $n_{m_i}$  we can show that,

$$\begin{aligned}
c_i &\leq \sqrt{\frac{\rho(\hat{v}_i + 2)\epsilon_{m_i} \log(\psi T \epsilon_{m_i})}{2 \log(\psi T \epsilon_{m_i}^2)}} \stackrel{(a)}{\leq} \sqrt{\frac{2\rho\epsilon_{m_i} \log(\frac{\psi T \epsilon_{m_i}^2}{\epsilon_{m_i}})}{\log(\psi T \epsilon_{m_i}^2)}} \\
&= \sqrt{\frac{2\rho\epsilon_{m_i} \log(\psi T \epsilon_{m_i}^2) - 2\rho\epsilon_{m_i} \log(\epsilon_{m_i})}{\log(\psi T \epsilon_{m_i}^2)}} \\
&\leq \sqrt{2\rho\epsilon_{m_i} - \frac{2\rho\epsilon_{m_i} \log(\frac{1}{2^{m_i}})}{\log(\psi T \frac{1}{2^{2m_i}})}} \\
&\leq \sqrt{2\rho\epsilon_{m_i} + \frac{2\rho\epsilon_{m_i} \log(2^{m_i})}{\log(\psi T) - \log(2^{2m_i})}} \\
&\leq \sqrt{2\rho\epsilon_{m_i} + \frac{2\rho\epsilon_{m_i} m_i \log(2)}{\log(\psi T) - 2m_i \log(2)}} \\
&\stackrel{(b)}{\leq} \sqrt{2\rho\epsilon_{m_i} + 2 \cdot \frac{3}{2} \epsilon_{m_i}} < \sqrt{4\epsilon_{m_i}} < \frac{\Delta_i}{4}.
\end{aligned}$$

In the above simplification, (a) is due to  $\hat{v}_i \in [0, 1]$ , while (b) is obtained using Lemma 1.

■

### B.0.6 Proof of Lemma 3

**Lemma 3** If  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $c_i = \sqrt{\frac{\rho(\hat{v}_i + 2) \log(\psi T \epsilon_{m_i})}{4z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i})}{2\epsilon_{m_i}}$  then we can show that,

$$\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \frac{2}{(\psi T \epsilon_{m_i})^{\frac{3p}{2}}}.$$

**Proof 5** We start by recalling from equation (3.2) that,

$$\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \mathbb{P}(\hat{r}_i > r_i + \bar{c}_i) + \mathbb{P}(\hat{v}_i \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}) \quad (\text{B.1})$$

where

$$c_i = \sqrt{\frac{\rho(\hat{v}_i + 2) \log(\psi T \epsilon_{m_i})}{4z_i}} \text{ and}$$

$$\bar{c}_i = \sqrt{\frac{\rho(\sigma_i^2 + \sqrt{\epsilon_{m_i}} + 2) \log(\psi T \epsilon_{m_i})}{4z_i}}.$$

Note that, substituting  $z_i \geq n_{m_i} \geq \frac{\log(\psi T \epsilon_{m_i})}{2\epsilon_{m_i}}$ ,  $\bar{c}_i$  can be simplified to obtain,

$$\bar{c}_i \leq \sqrt{\frac{\rho \epsilon_{m_i} (\sigma_i^2 + \sqrt{\epsilon_{m_i}} + 2)}{2}} \leq \sqrt{\epsilon_{m_i}}. \quad (\text{B.2})$$

The first term in the LHS of (B.1) can be bounded using the Bernstein inequality as below:

$$\begin{aligned} \mathbb{P}(\hat{r}_i > r_i + \bar{c}_i) &\leq \exp\left(-\frac{(\bar{c}_i)^2 z_i}{2\sigma_i^2 + \frac{2}{3}\bar{c}_i}\right) \\ &\stackrel{(a)}{\leq} \exp\left(-\rho\left(\frac{3\sigma_i^2 + 3\sqrt{\epsilon_{m_i}} + 6}{6\sigma_i^2 + 2\sqrt{\epsilon_{m_i}}}\right) \log(\psi T \epsilon_{m_i})\right) \\ &\stackrel{(b)}{\leq} \exp(-\rho \log(\psi T \epsilon_{m_i})) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \end{aligned} \quad (\text{B.3})$$

where, (a) is obtained by substituting equation B.2 and (b) occurs because for all  $\sigma_i^2 \in [0, \frac{1}{4}]$ ,  $\left(\frac{3\sigma_i^2 + 3\sqrt{\epsilon_{m_i}} + 6}{6\sigma_i^2 + 2\sqrt{\epsilon_{m_i}}}\right) \geq \frac{3}{2}$ .

The second term in the LHS of (B.1) can be simplified as follows:

$$\begin{aligned} &\mathbb{P}\left\{\hat{v}_i \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}\right\} \\ &\leq \mathbb{P}\left\{\frac{1}{n_i} \sum_{t=1}^{n_i} (X_{i,t} - r_i)^2 - (\hat{r}_i - r_i)^2 \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}\right\} \\ &\leq \mathbb{P}\left\{\frac{\sum_{t=1}^{n_i} (X_{i,t} - r_i)^2}{n_i} \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}\right\} \\ &\stackrel{(a)}{\leq} \mathbb{P}\left\{\frac{\sum_{t=1}^{n_i} (X_{i,t} - r_i)^2}{n_i} \geq \sigma_i^2 + \bar{c}_i\right\} \\ &\stackrel{(b)}{\leq} \exp\left(-\rho\left(\frac{3\sigma_i^2 + 3\sqrt{\epsilon_{m_i}} + 6}{6\sigma_i^2 + 2\sqrt{\epsilon_{m_i}}}\right) \log(\psi T \epsilon_{m_i})\right) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \end{aligned} \quad (\text{B.4})$$

where inequality (a) is obtained using (B.2), while (b) follows from the Bernstein in-

equality.

Thus, using (B.3) and (B.4) in (B.1) we obtain  $\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \frac{2}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}}$ .  $\blacksquare$

### B.0.7 Proof of Lemma 4

**Lemma 4** If  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$ ,  $c_i = \sqrt{\frac{\rho(\hat{v}_i + 2) \log(\psi T \epsilon_{m_i})}{4z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$  then in the  $m_i$ -th round,

$$\mathbb{P}\{c^* > c_i\} \leq \frac{182K^4}{T^{\frac{5}{4}}\sqrt{\epsilon_{m_i}}}.$$

**Proof 6** From the definition of  $c_i$  we know that  $c_i \propto \frac{1}{z_i}$  as  $\psi$  and  $T$  are constants. Therefore in the  $m_i$ -th round,

$$\begin{aligned} \mathbb{P}\{c^* > c_i\} &\leq \mathbb{P}\{z^* < z_i\} \\ &\leq \sum_{m=0}^{m_i} \sum_{z^*=1}^{n_m} \sum_{z_i=1}^{n_m} \left( \mathbb{P}\{\hat{r}^* < r^* - c^*\} + \mathbb{P}\{\hat{r}_i > r_i + c_i\} \right) \end{aligned}$$

Now, applying Bernstein inequality and following the same way as in Lemma 3 we can show that,

$$\begin{aligned} \mathbb{P}\{\hat{r}^* < r^* - c^*\} &\leq \exp\left(-\frac{(c^*)^2}{2\sigma_*^2 + \frac{2c^*}{3}} z^*\right) \leq \frac{4}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \\ \mathbb{P}\{\hat{r}_i > r_i + c_i\} &\leq \exp\left(-\frac{(c_i)^2}{2\sigma_i^2 + \frac{2c_i}{3}} z_i\right) \leq \frac{4}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \end{aligned}$$

Hence, summing everything up,

$$\begin{aligned} \mathbb{P}\{c^* > c_i\} &\leq \sum_{m=0}^{m_i} \sum_{z^*=1}^{n_m} \sum_{z_i=1}^{n_m} \left( \mathbb{P}\{\hat{r}^* < r^* - c^*\} + \mathbb{P}\{\hat{r}_i > r_i + c_i\} \right) \\ &\stackrel{(a)}{\leq} \sum_{m=0}^{m_i} |B_m| n_m \left( \mathbb{P}\{\hat{r}^* < r^* - c^*\} + \mathbb{P}\{\hat{r}_i > r_i + c_i\} \right) \\ &\stackrel{(b)}{\leq} \sum_{m=0}^{m_i} \frac{4K}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \frac{\log(\psi T \epsilon_m^2)}{2\epsilon_m} \times \end{aligned}$$



$$\begin{aligned}
& \left( \mathbb{P}\{\hat{r}^* < r^* - c^*\} + \mathbb{P}\{\hat{r}_i > r_i + c_i\} \right) \\
& \stackrel{(c)}{\leq} \sum_{m=0}^{m_i} \frac{4K}{(\psi T \epsilon_m)^{\frac{3\rho}{2}}} \frac{\log(T)}{\epsilon_m} \left[ \frac{4}{(\psi T \epsilon_m)^{\frac{3\rho}{2}}} + \frac{4}{(\psi T \epsilon_m)^{\frac{3\rho}{2}}} \right] \\
& \leq \sum_{m=0}^{m_i} \frac{32K \log T}{(\psi T \epsilon_m)^{3\rho} \epsilon_m} \leq \frac{32K \log T}{(\psi T)^{3\rho}} \sum_{m=0}^{m_i} \frac{1}{\epsilon_m^{3\rho+1}} \\
& \stackrel{(d)}{\leq} \sum_{m=0}^{m_i} \frac{32K \log T}{(\psi T)^{3\rho}} \left( \sum_{m=0}^{m_i} \frac{1}{\epsilon_m} \right)^{3\rho+1} \\
& \stackrel{(e)}{\leq} \frac{32K \log T}{\left(\frac{T^2}{K^2}\right)^{\frac{3}{2}}} \left[ \left( 1 + \frac{2(2^{\frac{1}{2} \log_2 \frac{T}{e}} - 1)}{2 - 1} \right)^{\frac{5}{2}} \right] \\
& \leq \frac{182K^4 T^{\frac{5}{4}} \log T}{T^3} \stackrel{(f)}{\leq} \frac{182K^4}{T^{\frac{5}{4}}} \stackrel{(g)}{\leq} \frac{182K^4}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}}
\end{aligned}$$

where, (a) comes from the total pulls allocated for all  $i \in B_m$  till the  $m$ -th round, in (b) the arm count  $|B_m|$  can be bounded by using equation (3.3) and then we substitute the value of  $n_m$ , (c) happens by substituting the value of  $\psi$  and considering  $\epsilon_m \in [\sqrt{\frac{e}{T}}, 1]$ , (d) follows as  $\frac{1}{\epsilon_m} \geq 1, \forall m$ , in (e) we use the standard geometric progression formula and then we substitute the values of  $\rho$  and  $\psi$ , (f) follows from the inequality  $\log T \leq \sqrt{T}$  and (g) is valid for any  $\epsilon_{m_i} \in [\sqrt{\frac{e}{T}}, 1]$ . ■

## B.0.8 Proof of Lemma 5

**Lemma 5** If  $m_i = \min\{m | \sqrt{4\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$ ,  $c_i = \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$  then in the  $m_i$ -th round,

$$\mathbb{P}\{z_i < n_{m_i}\} \leq \frac{182K^4}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}}.$$

**Proof 7** Following a similar argument as in Lemma 4, we can show that in the  $m_i$ -th round,

$$\mathbb{P}\{z_i < n_{m_i}\}$$

$$\begin{aligned}
&\leq \sum_{m=0}^{m_i} \sum_{z_i=1}^{n_m} \sum_{z^*=1}^{n_m} \left( \mathbb{P}\{\hat{r}^* > r^* - c^*\} + \mathbb{P}\{\hat{r}_i < r_i + c_i\} \right) \\
&\leq \frac{32K \log T}{(\psi T)^{3\rho}} \sum_{m=0}^{m_i} \frac{1}{\epsilon_m^{3\rho+1}} \leq \frac{182K^4}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}}.
\end{aligned}$$

■

### B.0.9 Proof of Lemma 6

**Lemma 6** *For two integer constants  $c_1$  and  $c_2$ , if  $20c_1 \leq c_2$  then,*

$$c_1 \frac{4\sigma_i^2 + 4}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right) \leq c_2 \frac{\sigma_i^2}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right).$$

**Proof 8** *We again prove this by contradiction. Suppose,*

$$c_1 \frac{4\sigma_i^2 + 4}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right) > c_2 \frac{\sigma_i^2}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right).$$

*Further reducing the above two terms we can show that,*

$$\begin{aligned}
4c_1\sigma_i^2 + 4c_1 &> c_2\sigma_i^2 \\
\Rightarrow 4c_1 \cdot \frac{1}{4} + 4c_1 &\stackrel{(a)}{>} \frac{c_2}{4} \\
\Rightarrow 20c_1 &> c_2.
\end{aligned}$$

*Here, (a) occurs because  $0 \leq \sigma_i^2 \leq \frac{1}{4}, \forall i \in \mathcal{A}$ . But, we already know that  $20c_1 \leq c_2$ . Hence,*

$$c_1 \frac{4\sigma_i^2 + 4}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right) \leq c_2 \frac{\sigma_i^2}{\Delta_i} \log \left( \frac{T\Delta_i^2}{K} \right).$$

■

### B.0.10 Proof of Corollary 1

**Corollary 1 (Gap-Independent Bound)** When the gaps of all the sub-optimal arms are identical, i.e.,  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{\epsilon}{T}}, \forall i \in \mathcal{A}$  and  $C_3$  being an integer constant, the regret of EUCEB is upper bounded by the following gap-independent expression:

$$\mathbb{E}[R_T] \leq \frac{C_3 K^5}{T^{\frac{1}{4}}} + 320 \sqrt{KT}.$$

**Proof 9** From Bubeck et al. (2011) we know that the function  $x \in [0, 1] \mapsto x \exp(-Cx^2)$  is decreasing on  $[\frac{1}{\sqrt{2C}}, 1]$  for any  $C > 0$ . Thus, we take  $C = \lfloor \frac{T}{\epsilon} \rfloor$  and choose  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{\epsilon}{T}}$  for all  $i$ .

First, let us recall the result in Theorem 1 below:

$$\begin{aligned} \mathbb{E}[R_T] &\leq \sum_{i \in \mathcal{A}: \Delta_i > b} \left\{ \frac{C_0 K^4}{T^{\frac{1}{4}}} + \left( \Delta_i + \frac{320 \sigma_i^2 \log(\frac{T \Delta_i^2}{K})}{\Delta_i} \right) \right\} \\ &\quad + \sum_{i \in \mathcal{A}: 0 < \Delta_i \leq b} \frac{C_2 K^4}{T^{\frac{1}{4}}} + \max_{i \in \mathcal{A}: 0 < \Delta_i \leq b} \Delta_i T. \end{aligned}$$

Now, with  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{\epsilon}{T}}$  we obtain,

$$\begin{aligned} \sum_{i \in \mathcal{A}: \Delta_i > b} \frac{320 \sigma_i^2 \log(\frac{T \Delta_i^2}{K})}{\Delta_i} &\leq \frac{320 \sigma_{\max}^2 K \sqrt{T} \log(T \frac{K(\log K)}{TK})}{\sqrt{K \log K}} \\ &\leq \frac{320 \sigma_{\max}^2 \sqrt{KT} \log(\log K)}{\sqrt{\log K}} \stackrel{(a)}{\leq} 320 \sigma_{\max}^2 \sqrt{KT} \end{aligned}$$

where (a) follows from the identity  $\frac{\log(\log K)}{\sqrt{\log K}} \leq 1$  for  $K \geq 2$ .

Thus, the total worst case gap-independent bound is given by

$$\begin{aligned} \mathbb{E}[R_T] &\stackrel{(a)}{\leq} \frac{C_3 K^5}{T^{\frac{1}{4}}} + 320 \sigma_{\max}^2 \sqrt{KT} \\ &\stackrel{(b)}{\leq} \frac{C_3 K^5}{T^{\frac{1}{4}}} + 320 \sqrt{KT} \end{aligned}$$

where, in (a),  $C_3$  is an integer constant such that  $C_3 = C_0 + C_2$  and (b) occurs because  $\sigma_i^2 \in [0, \frac{1}{4}], \forall i \in \mathcal{A}$ .



## Bibliography

1. **Abernethy, J. D., K. Amin, and R. Zhu**, Threshold bandits, with and without censored feedback. *In Advances In Neural Information Processing Systems*. 2016.
2. **Agrawal, R.** (1995). Sample mean based index policies by  $o(\log n)$  regret for the multi-armed bandit problem. *Advances in Applied Probability*, **27**(4), 1054–1078.
3. **Agrawal, S. and N. Goyal** (2011). Analysis of thompson sampling for the multi-armed bandit problem. *arXiv preprint arXiv:1111.1797*.
4. **Audibert, J.-Y. and S. Bubeck**, Minimax policies for adversarial and stochastic bandits. *In COLT*. 2009.
5. **Audibert, J.-Y. and S. Bubeck**, Best arm identification in multi-armed bandits. *In COLT-23th Conference on Learning Theory-2010*. 2010.
6. **Audibert, J.-Y., R. Munos, and C. Szepesvári** (2009). Exploration–exploitation trade-off using variance estimates in multi-armed bandits. *Theoretical Computer Science*, **410**(19), 1876–1902.
7. **Auer, P.** (2002). Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, **3**(Nov), 397–422.
8. **Auer, P., N. Cesa-Bianchi, and P. Fischer** (2002a). Finite-time analysis of the multi-armed bandit problem. *Machine learning*, **47**(2-3), 235–256.
9. **Auer, P., N. Cesa-Bianchi, Y. Freund, and R. E. Schapire** (2002b). The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, **32**(1), 48–77.
10. **Auer, P. and R. Ortner** (2010). Ucb revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Periodica Mathematica Hungarica*, **61**(1-2), 55–65.
11. **Awerbuch, B. and R. Kleinberg** (2008). Competitive collaborative learning. *Journal of Computer and System Sciences*, **74**(8), 1271–1288.
12. **Bertsekas, D. P. and J. N. Tsitsiklis** (1996). Neuro-dynamic programming (optimization and neural computation series, 3). *Athena Scientific*, **7**, 15–23.
13. **Beygelzimer, A., J. Langford, L. Li, L. Reyzin, and R. E. Schapire**, Contextual bandit algorithms with supervised learning guarantees. *In AISTATS*. 2011.
14. **Bubeck, S. and N. Cesa-Bianchi** (2012). Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *arXiv preprint arXiv:1204.5721*.
15. **Bubeck, S., R. Munos, and G. Stoltz** (2011). Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, **412**(19), 1832–1852.

16. **Bubeck, S., T. Wang, and N. Viswanathan**, Multiple identifications in multi-armed bandits. *In ICML (1)*. 2013.
17. **Bui, L., R. Johari, and S. Mannor** (2012). Clustered bandits. *arXiv preprint arXiv:1206.4169*.
18. **Cappe, O., A. Garivier, and E. Kaufmann** (2012). pymabandits. <http://mloss.org/software/view/415/>.
19. **Cappé, O., A. Garivier, O.-A. Maillard, R. Munos, G. Stoltz, et al.** (2013). Kullback–leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, **41**(3), 1516–1541.
20. **Cesa-Bianchi, N., C. Gentile, and G. Zappella**, A gang of bandits. *In Advances in Neural Information Processing Systems*. 2013.
21. **Chen, S., T. Lin, I. King, M. R. Lyu, and W. Chen**, Combinatorial pure exploration of multi-armed bandits. *In Advances in Neural Information Processing Systems*. 2014.
22. **Even-Dar, E., S. Mannor, and Y. Mansour** (2006). Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *The Journal of Machine Learning Research*, **7**, 1079–1105.
23. **Gabillon, V., M. Ghavamzadeh, and A. Lazaric**, Best arm identification: A unified approach to fixed budget and fixed confidence. *In Advances in Neural Information Processing Systems*. 2012.
24. **Gabillon, V., M. Ghavamzadeh, A. Lazaric, and S. Bubeck**, Multi-bandit best arm identification. *In Advances in Neural Information Processing Systems*. 2011.
25. **Garivier, A. and O. Cappé** (2011). The kl-ucb algorithm for bounded stochastic bandits and beyond. *arXiv preprint arXiv:1102.2490*.
26. **Gentile, C., S. Li, and G. Zappella**, Online clustering of bandits. *In ICML*. 2014.
27. **Ghavamzadeh, M., S. Mannor, J. Pineau, A. Tamar, et al.**, *Bayesian reinforcement learning: a survey*. World Scientific, 2015.
28. **Hillel, E., Z. S. Karnin, T. Koren, R. Lempel, and O. Somekh**, Distributed exploration in multi-armed bandits. *In Advances in Neural Information Processing Systems*. 2013.
29. **Honda, J. and A. Takemura**, An asymptotically optimal bandit algorithm for bounded support models. *In COLT*. Citeseer, 2010.
30. **Jamieson, K. and R. Nowak**, Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. *In Information Sciences and Systems (CISS), 2014 48th Annual Conference on*. IEEE, 2014.
31. **Kalyanakrishnan, S., A. Tewari, P. Auer, and P. Stone**, Pac subset selection in stochastic multi-armed bandits. *In Proceedings of the 29th International Conference on Machine Learning (ICML-12)*. 2012.

32. **Kaufmann, E., O. Cappé, and A. Garivier**, On bayesian upper confidence bounds for bandit problems. *In AISTATS*. 2012.
33. **Kocák, T., G. Neu, M. Valko, and R. Munos**, Efficient learning by implicit exploration in bandit problems with side observations. *In Advances in Neural Information Processing Systems*. 2014.
34. **Lai, T. L. and H. Robbins** (1985). Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, **6**(1), 4–22.
35. **Langford, J. and T. Zhang**, The epoch-greedy algorithm for multi-armed bandits with side information. *In Advances in neural information processing systems*. 2008.
36. **Lattimore, T.** (2015). Optimally confident ucb: Improved regret for finite-armed bandits. *arXiv preprint arXiv:1507.07880*.
37. **Li, L., W. Chu, J. Langford, and R. E. Schapire**, A contextual-bandit approach to personalized news article recommendation. *In Proceedings of the 19th international conference on World wide web*. ACM, 2010.
38. **Liu, K. and Q. Zhao** (2010). Distributed learning in multi-armed bandit with multiple players. *IEEE Transactions on Signal Processing*, **58**(11), 5667–5681.
39. **Liu, Y.-C. and Y. Tsuruoka** (2016). Modification of improved upper confidence bounds for regulating exploration in monte-carlo tree search. *Theoretical Computer Science*.
40. **Locatelli, A., M. Gutzeit, and A. Carpentier** (2016). An optimal algorithm for the thresholding bandit problem. *arXiv preprint arXiv:1605.08671*.
41. **Robbins, H.**, Some aspects of the sequential design of experiments. *In Herbert Robbins Selected Papers*. Springer, 1952, 169–177.
42. **Slivkins, A.** (2014). Contextual bandits with similarity information. *Journal of Machine Learning Research*, **15**(1), 2533–2568.
43. **Sutton, R. S. and A. G. Barto**, *Reinforcement learning: An introduction*. MIT press, 1998.
44. **Szörényi, B., R. Busa-Fekete, I. Hegedüs, R. Ormándi, M. Jelasity, and B. Kégl**, Gossip-based distributed stochastic bandit algorithms. *In ICML (3)*. 2013.
45. **Thompson, W. R.** (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 285–294.
46. **Thompson, W. R.** (1935). On the theory of apportionment. *American Journal of Mathematics*, **57**(2), 450–456.

## **LIST OF PAPERS BASED ON THESIS**

1. Authors.... Title... *Journal*, Volume, Page, (year).