

# Efficient-UCBV: An Almost Optimal Algorithm using Variance Estimates

Author names withheld

## Abstract

We propose a novel variant of the UCB algorithm (referred to as Efficient-UCB-Variance (EUCBV)) for minimizing cumulative regret in the stochastic multi-armed bandit (MAB) 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, Munos, and Szepesvári, 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  $O(\sqrt{KT})$  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.

## 1 Introduction

In this paper, we deal with the stochastic multi-armed bandit (MAB) 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. The mean of the reward distribution 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. With this formulation 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 or explore arms more thoroughly to ensure that a correct decision is being made. The objective in the stochastic bandit problem is 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 recent years the MAB 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.

### 1.1 Related Work

Bandit problems has been extensively studied in several earlier works such as Thompson (1933), Robbins (1952) and Lai and Robbins (1985). Lai and Robbins in Lai and Robbins (1985) established an asymptotic lower bound for the cumulative regret. Over the years stochastic MABs has seen several algorithms with strong regret guarantees. For further reference an interested reader can look into Bubeck and Cesa-Bianchi (2012). The upper confidence bound algorithms balance the exploration-exploitation dilemma by linking the uncertainty in estimate of an arm with the number of times an arm is pulled, and therefore ensuring sufficient exploration. One of the earliest among these algorithms is UCB1 (Auer, Cesa-Bianchi, and Fischer, 2002), which 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  $O(\sqrt{KT \log T})$ . 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 gap-dependent regret of MOSS is  $O\left(\frac{K^2 \log(T\Delta^2/K)}{\Delta}\right)$  and in certain regimes, this can be worse than even UCB1 (see Audibert and Bubeck (2009); Lattimore (2015)).

The UCB-Improved algorithm, proposed in Auer and Ortner (2010), is a round-based<sup>1</sup> variant of UCB1, that 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  $O(\sqrt{KT \log K})$ . 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, 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.

In contrast to the above work, the UCBV (Audibert, Munos, and Szepesvári, 2009) algorithm utilizes variance estimates to compute the confidence intervals for each arm. 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, Munos, and Szepesvári (2009) showed that UCBV outperforms UCB1 in several scenarios.

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, Cappé, and Garivier, 2012). The TS algorithm maintains a posterior reward distribution for each arm; at each round, the algorithm samples values from these distribution 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.

The final design principle we will state is the information theoretic approach of DMED (Honda and Takemura, 2010) and KLUCB (Garivier and Cappé, 2011) algorithms. The algorithm KLUCB uses Kullback-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

<sup>1</sup>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 sub-optimal.

Table 1: Regret upper bound of different algorithms

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

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.

## 1.2 Our Contributions

In this paper 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 performs significantly better than the existing algorithms in the stochastic MAB setting. EUCBV outperforms UCBV (Audibert, Munos, and Szepesvári, 2009) which also takes into account empirical variance but is less powerful than EUCBV because of the usage of exploration regulatory factor, arm elimination parameter and non-uniform arm selection (as opposed to UCB-Improved) by EUCBV. 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 enjoys a worst case gap-independent regret bound of  $O(\sqrt{KT})$  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 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-Imp, 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. Also 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). In Auer and Ortner (2010) the authors conjectured that an UCB-Improved like arm-elimination algorithm can greatly benefit by taking into consideration the variance of the arms. 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  $O(\sqrt{KT \log K})$ .

The rest of the paper is organized as follows. In section 2 we present the EUCBV algorithm. Our main theoretical results are stated in section 3, while the proofs are established in section 4. Section 5 contains results and discussions from our numerical experiments. We draw our conclusions in section 6 and section 7 is Appendix (supplementary material).

## 2 Algorithm: Efficient UCB Variance

**2.1 Notations:** We denote the set of arms by  $\mathcal{A}$ , with the individual arms labeled  $i$ , where  $i = 1, \dots, K$ . We denote an arbitrary round of EUCBV by  $m$ . For simplicity, we assume that the optimal arm is unique and denote it by  $*$ . 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. We assume the rewards of all arms are bounded in  $[0, 1]$ .

**2.2 The algorithm:** Earlier round-based arm elimination algorithms like Median Elimination (Even-Dar, Mannor, and Mansour, 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.

EUCBV algorithm which is mainly based on the arm elimination technique of the UCB-Improved algorithm

---

### Algorithm 1 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**

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**

$$\begin{aligned} \epsilon_{m+1} &:= \frac{\epsilon_m}{2} \\ B_{m+1} &:= B_m \\ n_{m+1} &:= \left\lceil \frac{\log(\psi T \epsilon_{m+1}^2)}{2\epsilon_{m+1}} \right\rceil \\ N_{m+1} &:= t + |B_{m+1}| n_{m+1} \\ m &:= m + 1 \end{aligned}$$

**end if**

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

**end for**

---

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, Munos, and Szepesvári (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, Mannor, and Mansour (2006)) where the arm elimination takes place only at the end of the respective exploration rounds.

## 3 Main Results

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

**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

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

for all  $b \geq \sqrt{\frac{e}{T}}$  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, Munos, and Szepesvári (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, Munos, and Szepesvári (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:* 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 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.

Next, 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.,  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{e}{T}}, \forall i \in \mathcal{A}$  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 7.7.

*Discussion:* In the non-stochastic scenario, Auer et al. (2002) showed that the bound on the cumulative regret for EXP-4 is  $O(\sqrt{KT \log K})$ . However, in the stochastic case, UCB1 proposed in Auer, Cesa-Bianchi, and Fischer (2002) incurred a regret of order of  $O(\sqrt{KT \log T})$  which is clearly improvable. From the above result we see that in the gap-independent bound of EUCBV the most significant term is  $O(\sqrt{KT})$  which matches the upper bound of MOSS and OCUCB, and is better than UCB-Improved, UCB1 and UCBV (see Table 1).

## 4 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  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i})}{2\epsilon_{m_i}}$  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  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i})}{2\epsilon_{m_i}}$  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 7.1, 7.2, 7.3, 7.4, 7.5 and 7.6 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

$$\hat{r}_i \leq r_i + c_i, \hat{r}^* \geq r^* - c^*, c_i \geq c^* \text{ and } z_i \geq n_{m_i}. \quad (1)$$

Indeed, in round  $m_i$  suppose (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 (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, Munos, and Szepesvári (2009) we can show that a bound on the complementary of the first condition is given by,

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

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 (1) is,

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

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,

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

Also, for eq. (3) we can show that for any  $\epsilon_{m_i} \in [\sqrt{\frac{e}{T}}, 1]$

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

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 (1) by (4) + (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,

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

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,

$$\sum_{i \in \mathcal{A}'} \Delta_i \left( 1 + \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i^2} \right) = \sum_{i \in \mathcal{A}'} \left( \Delta_i + \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i} \right)$$

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

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,

$$\begin{aligned}
&\sum_{m_*=0}^{\max_{j \in \mathcal{A}'} m_j} \sum_{i \in \mathcal{A}'' : m_i > m_*} \left( \frac{368K^4}{(T^{\frac{5}{4}} \sqrt{\epsilon_{m_*}})} \right) \cdot T \max_{j \in \mathcal{A}'' : m_j \geq m_*} \Delta_j \\
&\leq \sum_{m_*=0}^{\max_{j \in \mathcal{A}'} m_j} \sum_{i \in \mathcal{A}'' : m_i > m_*} \left( \frac{368K^4 \sqrt{4}}{(T^{\frac{5}{4}} \sqrt{\epsilon_{m_*}})} \right) \cdot T \cdot 4\sqrt{\epsilon_{m_*}} \\
&\stackrel{(a)}{\leq} \sum_{m_*=0}^{\max_{j \in \mathcal{A}'} m_j} \sum_{i \in \mathcal{A}'' : m_i > m_*} \left( \frac{C_2 K^4}{T^{\frac{1}{4}} \epsilon_{m_*}^{\frac{1}{2} - \frac{1}{2}}} \right) \\
&\leq \sum_{i \in \mathcal{A}'' : m_i > m_*} \sum_{m_*=0}^{\min \{m_i, m_b\}} \left( \frac{C_2 K^4}{T^{\frac{1}{4}}} \right) \\
&\leq \sum_{i \in \mathcal{A}'} \left( \frac{C_2 K^4}{T^{\frac{1}{4}}} \right) + \sum_{i \in \mathcal{A}'' \setminus \mathcal{A}'} \left( \frac{C_2 K^4}{T^{\frac{1}{4}}} \right).
\end{aligned}$$

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.  $C_0 = C_1 + C_2$ .

## 5 Experiments

In this section, we conduct extensive empirical evaluations of EUCEV against several other popular bandit algorithms. We use expected cumulative regret as the metric of comparison. We compare with the following algorithms: KLUCB+ (Garivier and Cappé, 2011), DMED (Honda and Takemura, 2010), MOSS (Audibert and Bubeck, 2009), UCB1 (Auer, Cesa-Bianchi, and Fischer, 2002), UCB-Improved (Auer and Ortner, 2010), Median Elimination (Even-Dar, Mannor, and Mansour, 2006), Thompson Sampling (TS) (Agrawal and Goyal, 2011), OCUCB (Lattimore, 2015), Bayes-UCB (BU) (Kaufmann, Cappé, and Garivier, 2012) and UCB-V (Audibert, Munos, and Szepesvári, 2009)<sup>2</sup>. The parameters of EUCEV 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)).

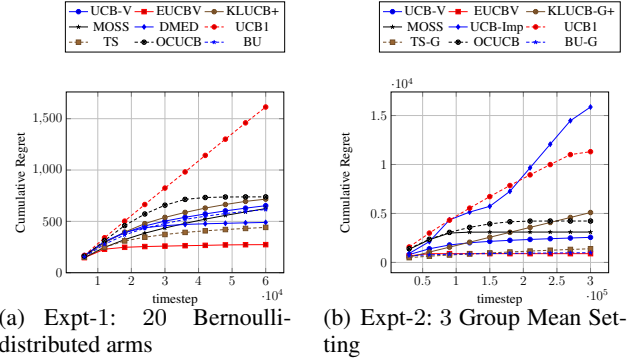


Figure 1: A comparison of the cumulative regret incurred by the various bandit algorithms.

**Experiment-1 (Bernoulli with uniform gaps):** This experiment is conducted to observe the performance of EUCEV 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 1(a). EUCEV, MOSS, OCUCB, UCB1, UCB-V, KLUCB+, TS, BU and DMED are run in this experimental setup. Not only do we observe that EUCEV 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 EUCEV 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

<sup>2</sup>The implementation for KLUCB, Bayes-UCB and DMED were taken from Cappé, Garivier, and Kaufmann (2012)

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 1(b). From the results in Figure 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 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 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 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  $r_{1:10} = 0.045$ ,  $r_{34:99} = 0.04$ ,  $r_{100}^* = 0.05$  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 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  $\sigma_{100}^2 = 0.25$ . 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 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.

## 6 Conclusion and Future Works

In this paper, 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  $O(\sqrt{KT})$ . 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.



## References

- Agrawal, S., and Goyal, N. 2011. Analysis of thompson sampling for the multi-armed bandit problem. *arXiv preprint arXiv:1111.1797*.
- Audibert, J.-Y., and Bubeck, S. 2009. Minimax policies for adversarial and stochastic bandits. In *COLT*, 217–226.
- Audibert, J.-Y., and Bubeck, S. 2010. Best arm identification in multi-armed bandits. In *COLT-23th Conference on Learning Theory-2010*, 13–p.
- Audibert, J.-Y.; Munos, R.; and Szepesvári, C. 2009. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science* 410(19):1876–1902.
- Auer, P., and Ortner, R. 2010. Ucb revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Periodica Mathematica Hungarica* 61(1-2):55–65.
- Auer, P.; Cesa-Bianchi, N.; Freund, Y.; and Schapire, R. E. 2002. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing* 32(1):48–77.
- Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine learning* 47(2-3):235–256.
- Bubeck, S., and Cesa-Bianchi, N. 2012. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *arXiv preprint arXiv:1204.5721*.
- Bubeck, S.; Munos, R.; and Stoltz, G. 2011. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science* 412(19):1832–1852.
- Cappe, O.; Garivier, A.; and Kaufmann, E. 2012. pymabandits. <http://mloss.org/software/view/415/>.
- Even-Dar, E.; Mannor, S.; and Mansour, Y. 2006. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *The Journal of Machine Learning Research* 7:1079–1105.
- Garivier, A., and Cappé, O. 2011. The kl-ucb algorithm for bounded stochastic bandits and beyond. *arXiv preprint arXiv:1102.2490*.
- Honda, J., and Takemura, A. 2010. An asymptotically optimal bandit algorithm for bounded support models. In *COLT*, 67–79. Citeseer.
- Kaufmann, E.; Cappé, O.; and Garivier, A. 2012. On bayesian upper confidence bounds for bandit problems. In *AISTATS*, 592–600.
- Lai, T. L., and Robbins, H. 1985. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics* 6(1):4–22.
- Lattimore, T. 2015. Optimally confident ucb: Improved regret for finite-armed bandits. *arXiv preprint arXiv:1507.07880*.
- Liu, Y.-C., and Tsuruoka, Y. 2016. Modification of improved upper confidence bounds for regulating exploration in monte-carlo tree search. *Theoretical Computer Science*.
- Robbins, H. 1952. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*. Springer. 169–177.
- Thompson, W. R. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 285–294.



## 7 Appendix

### 7.1 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 2** 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. ■

### 7.2 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 3** 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. ■

### 7.3 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,

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

**Proof 4** We start by recalling from equation (2) that,

$$\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \mathbb{P}(\hat{r}_i > r_i + \bar{c}_i) + \mathbb{P}(\hat{v}_i \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}) \quad (6)$$

where

$$\begin{aligned} c_i &= \sqrt{\frac{\rho(\hat{v}_i+2) \log(\psi T \epsilon_{m_i})}{4z_i}} \text{ and} \\ \bar{c}_i &= \sqrt{\frac{\rho(\sigma_i^2 + \sqrt{\epsilon_{m_i}} + 2) \log(\psi T \epsilon_{m_i})}{4z_i}}. \end{aligned}$$

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,

$$\bar{c}_i \leq \sqrt{\frac{\rho\epsilon_{m_i}(\sigma_i^2 + \sqrt{\epsilon_{m_i}} + 2)}{2}} \leq \sqrt{\epsilon_{m_i}}. \quad (7)$$

The first term in the LHS of (6) can be bounded using the Bernstein inequality as below:

$$\begin{aligned} \mathbb{P}(\hat{r}_i > r_i + \bar{c}_i) &\leq \exp\left(-\frac{(\bar{c}_i)^2 z_i}{2\sigma_i^2 + \frac{2}{3}\bar{c}_i}\right) \\ &\stackrel{(a)}{\leq} \exp\left(-\rho\left(\frac{3\sigma_i^2 + 3\sqrt{\epsilon_{m_i}} + 6}{6\sigma_i^2 + 2\sqrt{\epsilon_{m_i}}}\right) \log(\psi T \epsilon_{m_i})\right) \\ &\stackrel{(b)}{\leq} \exp(-\rho \log(\psi T \epsilon_{m_i})) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \quad (8) \end{aligned}$$

where, (a) is obtained by substituting equation 7 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 (6) can be simplified as follows:

$$\begin{aligned} &\mathbb{P}\left\{\hat{v}_i \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}\right\} \\ &\leq \mathbb{P}\left\{\frac{1}{n_i} \sum_{t=1}^{n_i} (X_{i,t} - r_i)^2 - (\hat{r}_i - r_i)^2 \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}\right\} \\ &\leq \mathbb{P}\left\{\frac{\sum_{t=1}^{n_i} (X_{i,t} - r_i)^2}{n_i} \geq \sigma_i^2 + \sqrt{\epsilon_{m_i}}\right\} \\ &\stackrel{(a)}{\leq} \mathbb{P}\left\{\frac{\sum_{t=1}^{n_i} (X_{i,t} - r_i)^2}{n_i} \geq \sigma_i^2 + \bar{c}_i\right\} \\ &\stackrel{(b)}{\leq} \exp\left(-\rho\left(\frac{3\sigma_i^2 + 3\sqrt{\epsilon_{m_i}} + 6}{6\sigma_i^2 + 2\sqrt{\epsilon_{m_i}}}\right) \log(\psi T \epsilon_{m_i})\right) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \quad (9) \end{aligned}$$

where inequality (a) is obtained using (7), while (b) follows from the Bernstein inequality.

Thus, using (8) and (9) in (6) we obtain  $\mathbb{P}(\hat{r}_i > r_i + c_i) \leq \frac{2}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}}$ . ■

#### 7.4 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 5** 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,

$$\begin{aligned} \mathbb{P}\{c^* > c_i\} &\leq \sum_{m=0}^{m_i} \sum_{z^*=1}^{n_m} \sum_{z_i=1}^{n_m} \left( \mathbb{P}\{\hat{r}^* < r^* - c^*\} + \mathbb{P}\{\hat{r}_i > r_i + c_i\} \right) \\ &\stackrel{(a)}{\leq} \sum_{m=0}^{m_i} |B_m| n_m \left( \mathbb{P}\{\hat{r}^* < r^* - c^*\} + \mathbb{P}\{\hat{r}_i > r_i + c_i\} \right) \\ &\stackrel{(b)}{\leq} \sum_{m=0}^{m_i} \frac{4K}{(\psi T \epsilon_{m_i})^{\frac{3\rho}{2}}} \frac{\log(\psi T \epsilon_m^2)}{2\epsilon_m} \times \\ &\quad \left( \mathbb{P}\{\hat{r}^* < r^* - c^*\} + \mathbb{P}\{\hat{r}_i > r_i + c_i\} \right) \\ &\stackrel{(c)}{\leq} \sum_{m=0}^{m_i} \frac{4K}{(\psi T \epsilon_m)^{\frac{3\rho}{2}}} \frac{\log(T)}{\epsilon_m} \left[ \frac{4}{(\psi T \epsilon_m)^{\frac{3\rho}{2}}} + \frac{4}{(\psi T \epsilon_m)^{\frac{3\rho}{2}}} \right] \\ &\leq \sum_{m=0}^{m_i} \frac{32K \log T}{(\psi T \epsilon_m)^{3\rho} \epsilon_m} \leq \frac{32K \log T}{(\psi T)^{3\rho}} \sum_{m=0}^{m_i} \frac{1}{\epsilon_m^{3\rho+1}} \\ &\stackrel{(d)}{\leq} \sum_{m=0}^{m_i} \frac{32K \log T}{(\psi T)^{3\rho}} \left( \sum_{m=0}^{m_i} \frac{1}{\epsilon_m} \right)^{3\rho+1} \\ &\stackrel{(e)}{\leq} \frac{32K \log T}{\left(\frac{T^2}{K^2}\right)^{\frac{3}{2}}} \left[ \left( 1 + \frac{2(2^{\frac{1}{2} \log_2 \frac{T}{\epsilon}} - 1)}{2 - 1} \right)^{\frac{5}{2}} \right] \end{aligned}$$

$$\leq \frac{182K^4 T^{\frac{5}{4}} \log T}{T^3} \stackrel{(f)}{\leq} \frac{182K^4}{T^{\frac{5}{4}}} \stackrel{(g)}{\leq} \frac{182K^4}{T^{\frac{5}{4}} \sqrt{\epsilon_{m_i}}}$$

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) and then we substitute the value of  $n_m$ , (c) happens by substituting the value of  $\psi$  and considering  $\epsilon_m \in [\sqrt{\frac{\epsilon}{T}}, 1]$ , (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  $\epsilon_{m_i} \in [\sqrt{\frac{\epsilon}{T}}, 1]$ . ■

#### 7.5 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 6** Following a similar argument as in Lemma 4, we can show that in the  $m_i$ -th round,

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

#### 7.6 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 7** 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,

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

■

## 7.7 Proof of Corollary 1

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

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

**Proof 8** *From Bubeck, Munos, and Stoltz (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  $C = \lfloor \frac{T}{e} \rfloor$  and choose  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{e}{T}}$  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  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{e}{T}}$  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  $\sigma_i^2 \in [0, \frac{1}{4}], \forall i \in \mathcal{A}$ . ■