
UCB with clustering and improved exploration

Abstract

In this paper, we present a novel algorithm for the stochastic multi-armed bandit problem. Our proposed method, referred to as ClusUCB, partitions the arms into clusters and then follows the UCB-Improved strategy with aggressive exploration factors to eliminate sub-optimal arms as well as clusters. Through a theoretical analysis, we establish that ClusUCB achieves a better gap-dependent regret upper bound than UCB-Improved[4] and MOSS[2] algorithms. Further, numerical experiments on test-cases with small gaps between optimal and sub-optimal mean rewards show that ClusUCB results in lower cumulative regret than several popular UCB variants as well as MOSS and Thompson sampling.

1 Introduction

In this paper, we consider the stochastic multi-armed bandit problem, a classical problem in sequential decision making. In this setting, a learning algorithm is provided with a set of decisions (or arms) with reward distributions unknown to the algorithm. The learning proceeds in an iterative fashion, where in each round, the algorithm chooses an arm and receives a stochastic reward that is drawn from a stationary distribution specific to the arm selected. Given the goal of maximizing the cumulative reward, the learning algorithm faces the exploration-exploitation dilemma, i.e., in each round should the algorithm select the arm which has the highest observed mean reward so far (*exploitation*), or should the algorithm choose a new arm to gain more knowledge of the true mean reward of the arms and thereby avert a sub-optimal greedy decision (*exploration*).

Formally, let $r_i, i = 1, \dots, K$ denote the mean rewards of the K arms and $r^* = \max_i r_i$ the optimal mean reward. The objective in the stochastic bandit problem is to mini-

mize the cumulative regret, which is defined as follows:

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

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

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

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

An early work involving a bandit setup is [19], where the author deals the problem of choosing between two treatments to administer on patients who come in sequentially. Following the seminal work of Robbins [18], bandit algorithms have been extensively studied in a variety of applications. From a theoretical standpoint, an asymptotic lower bound for the regret was established in [14]. In particular, it was shown there that for any consistent allocation strategy, we have $\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[R_T]}{\log T} \geq \sum_{\{i: r_i < r^*\}} \frac{(r^* - r_i)}{D(p_i || p^*)}$, where $D(p_i || p^*)$ is the Kullback-Leibler divergence between the reward densities p_i and p^* , corresponding to arms with mean r_i and r^* , respectively.

There have been several algorithms with strong regret guarantees. The foremost among them is UCB1 [5], which has a regret upper bound of $O\left(\frac{K \log T}{\Delta}\right)$, where $\Delta = \min_{i: \Delta_i > 0} \Delta_i$. This result is asymptotically order-optimal for the class of distributions considered. However, the worst case gap independent regret bound of UCB1 can be as bad as $O\left(\sqrt{TK \log T}\right)$. In [2], the authors propose the MOSS algorithm and establish that the worst case regret of MOSS is $O\left(\sqrt{TK}\right)$ which improves upon UCB1 by a factor of order $\sqrt{\log T}$. However, the gap-dependent regret of MOSS is $O\left(\frac{K^2 \log(T\Delta^2/K)}{\Delta}\right)$ and in certain regimes, this can be worse than even UCB1 (see [2],[15]).

The UCB-Improved algorithm, proposed in [4], is a round-based algorithm¹ variant of UCB1 that has a gap-dependent regret bound of $O\left(\frac{K \log T \Delta^2}{\Delta}\right)$, which is better than that of UCB1. On the other hand, the worst case regret of UCB-Improved is $O\left(\sqrt{TK \log K}\right)$.

Our Work

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

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

The motivation for our work stems from the remark in Section 2.4.3 of [8], where the authors conjecture that one should be able to obtain a bandit algorithm with a gap-dependent regret bound that is better than MOSS [2] and UCB-Improved [4], in particular, with a regret bound of the order $O\left(\frac{K \log(\frac{T}{H})}{\Delta}\right)$, where $H = \sum_{i: \Delta_i > 0} \frac{1}{\Delta_i^2}$. While ClusUCB does not achieve the conjectured regret bound, the theoretical analysis establishes that the gap-dependent regret of ClusUCB is always better than that of UCB-Improved and better than that of MOSS when $\sqrt{\frac{e}{T}} \leq \Delta \leq 1$ (see Table 1). Moreover, the gap-independent bound of ClusUCB is of the same order as UCB-Improved, i.e., $O(\sqrt{KT \log K})$. However, ClusUCB is not able to match the gap-independent bound of $O(\sqrt{KT})$ for MOSS.

On four synthetic setups with small gaps, we observe empirically that ClusUCB outperforms UCB-Improved[4] and MOSS[2] as well as other popular stochastic bandit algorithms such as DMED[13], UCB-V[3], Median Elimination[10], Thompson Sampling[1] and KL-UCB[12].

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

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

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

The rest of the paper is organized as follows: In Section 2, we present the ClusUCB algorithm. In Section 3, we present the associated regret bounds and prove the main theorem on the regret upper bound for ClusUCB in Section 4. In Section 5, we present the numerical experiments and provide concluding remarks in Section 6. Further proofs of corollaries and propositions presented in Section 4 are provided in the appendices.

2 Clustered UCB

Notation. We denote the set of arms by A , with the individual arms labeled $i, i = 1, \dots, K$. We denote an arbitrary round of ClusUCB by m . We denote an arbitrary cluster by s_k , the subset of arms within the cluster s_k by A_{s_k} and the set of clusters by S with $|S| = p \leq K$. Here p is a pre-specified limit for the number of clusters. For simplicity, we assume that the optimal arm is unique and denote it by $*$, with s^* denoting the corresponding cluster. The best arm in a cluster s_k is denoted by $a_{\max_{s_k}}$. We denote the sample mean of the rewards seen so far for arm i by \hat{r}_i and for the best arm within a cluster s_k by $\hat{r}_{a_{\max_{s_k}}}$. We assume that all arms' rewards are bounded in $[0, 1]$.

The algorithm

As mentioned in a recent work [16], UCB-Improved has two shortcomings:

- (i) A significant number of pulls are spent in early exploration, since each round m of UCB-Improved involves pulling every arm an identical $n_m = \left\lceil \frac{2 \log(T \epsilon_m^2)}{\epsilon_m^2} \right\rceil$ number of times. The quantity ϵ_m is initialized to 1 and halved after every round.
- (ii) In UCB-Improved, arms are eliminated conservatively, i.e., only after $\epsilon_m < \frac{\Delta_i}{2}$, the sub-optimal arm i is discarded

Algorithm 1 ClusUCB

Input: Number of clusters p , time horizon T , exploration parameters ρ_a, ρ_s and ψ .

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

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

for $m = 0, 1, \dots, \left\lfloor \frac{1}{2} \log_2 \frac{7T}{K} \right\rfloor$ **do**

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

Arm Elimination

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

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

Cluster Elimination

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

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

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

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

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

end for

with high probability. This is disadvantageous when K is large and the gaps are identical ($r_1 = r_2 = \dots = r_{K-1} < r^*$) and small.

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

As described in the pseudocode in Algorithm 1, ClusUCB begins with a initial clustering of arms that is performed by random uniform allocation. The set of clusters S thus obtained satisfies $|S| = p$, with individual clusters having a size that is bounded above by $\ell = \left\lceil \frac{K}{p} \right\rceil$. Each round of ClusUCB involves both individual arm as well as clus-

ter elimination conditions. These elimination conditions are inspired by UCB-Improved. Notice that, unlike UCB-Improved, there is no longer a single point of reference based on which we are eliminating arms. Instead now we have as many reference points to eliminate arms as number of clusters formed.

The exploration regulatory factor ψ governing the arm and cluster elimination conditions in ClusUCB is more aggressive than that in UCB-Improved. With appropriate choice of ψ and ρ_a and ρ_s we can achieve aggressive elimination even when the gaps Δ_i are small and K is large.

In [16], the authors recommend incorporating a factor of d_i inside the log-term of the UCB values, i.e., $\max\{\hat{r}_i + \sqrt{\frac{d_i \log T \epsilon_m^2}{2n_m}}\}$. The authors there examine the following choices for d_i : $\frac{T}{t_i}$, $\frac{\sqrt{T}}{t_i}$ and $\frac{\log T}{t_i}$, where t_i is the number of times an arm i has been sampled. Unlike [16], we employ cluster as well as arm elimination and establish from a theoretical analysis that the choice $\psi = \frac{T}{\log(KT)}$ helps in achieving a better gap-dependent regret upper bound for ClusUCB as compared to UCB-Improved and MOSS (see Corollary 1 in the next section).

3 Main results

We now state the main result that upper bounds the expected regret of ClusUCB.

Theorem 1 (Regret bound). *The regret R_T of ClusUCB satisfies*

$$\begin{aligned} \mathbb{E}[R_T] \leq & \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \left\{ \frac{C_1(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \Delta_i \right. \\ & + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \left. \right\} + \sum_{\substack{i \in A, \\ \Delta_i > b}} \left\{ 2\Delta_i + \frac{C_1(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right. \\ & + \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} + \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} \left. \right\} \\ & + \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \sum_{\substack{i \in A_{s^*}, \\ 0 < \Delta_i \leq b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{b^{4\rho_a-1}} \\ & + \sum_{\substack{i \in A \setminus A_{s^*}, \\ \Delta_i > b}} \frac{2C_2(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{\substack{i \in A \setminus A_{s^*}, \\ 0 < \Delta_i \leq b}} \frac{2C_2(\rho_s)T^{1-\rho_s}}{b^{4\rho_s-1}} \\ & + \max_{i: \Delta_i \leq b} \Delta_i T, \end{aligned}$$

where $b \geq \sqrt{\frac{K}{7T}}$, $C_1(x) = \frac{2^{1+4x}x^{2x}}{\psi^x}$, $C_2(x) = \frac{2^{2x+\frac{3}{2}}x^{2x}}{\psi^x}$ and A_{s^*} is the subset of arms in cluster s^* containing optimal arm a^* .

Proof. See Section 4. \square

We now specialize the result in the theorem above by substituting specific values for the exploration constants ρ_s , ρ_a and ψ .

Corollary 1 (Gap-dependent bound). *With $\psi = \frac{T}{\log(K)}$, $\rho_a = \frac{1}{2}$ and $\rho_s = \frac{1}{2}$, we have the following gap-dependent bound for the regret of ClusUCB:*

$$\begin{aligned} \mathbb{E}[R_T] \leq & \sum_{\substack{i \in A_{s^*}: \\ \Delta_i > b}} \left\{ \frac{6.8\sqrt{\log(KT)}}{\Delta_i} + \Delta_i + \right. \\ & \left. \frac{32 \log(T \frac{\Delta_i^2}{\sqrt{\log(K)}})}{\Delta_i} \right\} + \sum_{i \in A: \Delta_i > b} \left\{ \frac{4\sqrt{\log(K)}}{\Delta_i} + 2\Delta_i \right. \\ & \left. + \frac{64 \log(T \frac{\Delta_i^2}{\sqrt{\log(K)}})}{\Delta_i} \right\} + \sum_{\substack{i \in A_{s^*}: \\ 0 < \Delta_i \leq b}} \frac{2.8\sqrt{\log(K)}}{\Delta_i} \\ & + \sum_{\substack{i \in A \setminus A_{s^*}: \\ \Delta_i > b}} \frac{5.6\sqrt{\log(K)}}{\Delta_i} + \sum_{\substack{i \in A \setminus A \cup A_{s^*}: \\ 0 < \Delta_i \leq b}} \frac{5.6\sqrt{\log(K)}}{\Delta_i} \\ & + \max_{i \in A: \Delta_i \leq b} \Delta_i T, \quad \text{for all } b \geq \sqrt{\frac{K}{7T}}. \end{aligned}$$

Proof. See Appendix C. \square

The most significant term in the bound above is $\sum_{i \in A: \Delta_i \geq b} \frac{64 \log(T \frac{\Delta_i^2}{\sqrt{\log(K)}})}{\Delta_i}$ and hence, the regret upper bound for ClusUCB is of the order $O\left(\frac{K \log\left(\frac{T \Delta^2}{\sqrt{\log(K)}}\right)}{\Delta}\right)$. As shown in Table 1, the gap-dependent bound of ClusUCB is always better than UCB1 and UCB-Improved.

Since Corollary 1 holds for all $\Delta \geq \sqrt{\frac{K}{7T}}$, it can be clearly seen that for all $\sqrt{\frac{K}{7T}} \leq \Delta \leq 1$ and $K \geq 2$, the gap-dependent bound is better than that of MOSS and UCB-Improved. Also, since UCB-Improved holds for $\sqrt{\frac{e}{T}} \leq \Delta \leq 1$, in ClusUCB if we take γ such that $\frac{K}{\gamma} \approx e$ then $\Delta \geq \sqrt{\frac{K}{\gamma T}} \geq \sqrt{\frac{e}{T}}$ and we can easily mimic the same guarantee as UCB-Improved. In the theoretical results as well as in the experiments we have take $\gamma = 7$.

Corollary 2 (Gap-independent bound). *Considering the same gap of $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}}$ for all $i: i \neq *$ and with $\psi = \frac{T}{\log K}$, $p = \lceil \frac{K}{\log K} \rceil$, $\rho_a = \frac{1}{2}$ and $\rho_s = \frac{1}{2}$, we have the following gap-independent bound for the regret of ClusUCB:*

$$\mathbb{E}[R_T] \leq 70.8 \frac{\sqrt{T} \log K}{\sqrt{K}} + \frac{32 \sqrt{T} \log K \log(\log K)}{\sqrt{K}}$$

$$\begin{aligned} & + 4\sqrt{KT} + 128\sqrt{KT \log K} \\ & + \frac{64\sqrt{KT} \log(\log K)}{\sqrt{\log K}} + 2.8\sqrt{\frac{T \log K}{e}} \\ & + 5.6\sqrt{\frac{T}{e}} (\log K)^{\frac{3}{2}} + 5.6 \frac{K}{K + \log K} \sqrt{KT} \end{aligned}$$

Proof. See Appendix D. \square

From the above result, we observe that the order of the regret upper bound of ClusUCB is $O(\sqrt{KT \log K})$ and this matches the order of UCB-Improved. However, this is not as low as the order $O(\sqrt{KT})$ of MOSS.

Analysis of elimination error

Let \tilde{R}_T denote the contribution to the expected regret in the case when the optimal arm $*$ gets eliminated during one of the rounds of ClusUCB. This can happen if a sub-optimal arm eliminates $*$ or if a sub-optimal cluster eliminates the cluster s^* that contains $*$ - these correspond to cases b2 and b3 in the proof of Theorem 1 (see Section 4). We shall denote variants of ClusUCB that include arm elimination condition only and cluster elimination condition only as ClusUCB-AE and ClusUCB-CE, respectively, while ClusUCB corresponds to Algorithm 1, which uses both arm and cluster elimination conditions.

For ClusUCB-AE, the quantity \tilde{R}_T can be extracted from the proofs (in particular, case b2 in Appendix A) and simplified using the values $\rho_a = \frac{1}{2}$ and $\psi = \frac{T}{\log K}$, to obtain $\tilde{R}_T = 2\sqrt{KT \log K}$.

A similar exercise for ClusUCB-CE (see Case b2 in Appendix B) with $\rho_s = \frac{1}{2}$ and $\psi = \frac{T}{\log K}$ yields $\tilde{R}_T = 4\sqrt{KT \log K}$.

Finally, for ClusUCB, the relevant terms from Theorem 1 that corresponds to \tilde{R}_T can be simplified with $\rho_a = \frac{1}{2}$, $\rho_s = \frac{1}{2}$, $p = \lceil \frac{K}{\log K} \rceil$ and $\psi = \frac{T}{\log K}$ (as in Corollary 2 to obtain $\tilde{R}_T = \frac{5.3\sqrt{T} \log K^{\frac{3}{2}}}{\sqrt{K}} + \frac{5.3\sqrt{T} \log K}{\sqrt{K}} + 10.6 \frac{K}{K + \log K} \sqrt{KT \log K} + 10.6 \frac{K}{K + \log K} \sqrt{KT}$. Hence, in comparison to ClusUCB-AE and ClusUCB-CE which has an elimination regret bound of $O(\sqrt{KT \log K})$, the elimination error contribution to regret is lower in ClusUCB which has a bound of $O(\frac{K}{K + \log K} \sqrt{KT \log K})$.

Thus, we observe that clustering in conjunction with improved exploration via ρ_a, ρ_s, p and ψ helps in reducing the factor associated with \sqrt{KT} for the gap-independent error regret bound for ClusUCB. A table containing the regret error bound is shown in Appendix E.

4 Proof of Theorem 1

Proof. Let $A' = \{i \in A, \Delta_i > b\}$, $A'' = \{i \in A, \Delta_i > 0\}$, $A'_{s_k} = \{i \in A_{s_k}, \Delta_i > b\}$ and $A''_{s_k} = \{i \in A_{s_k}, \Delta_i > 0\}$. C_g is the cluster set containing max payoff arm from each cluster in g -th round. The arm having the highest payoff in a cluster s_k is denote by $a_{\max_{s_k}}$. Let for each sub-optimal arm $i \in A$, $m_i = \min\{m | \sqrt{\rho_a \epsilon_m} < \frac{\Delta_i}{2}\}$ and let for each cluster $s_k \in S$, $g_{s_k} = \min\{g | \sqrt{\rho_s \epsilon_g} < \frac{\Delta_{a_{\max_{s_k}}}}{2}\}$. Let $\check{A} = \{i \in A' | i \in s_k, \forall s_k \in S\}$.

The analysis proceeds by considering the contribution to the regret in each of the following cases:

Case a: Some sub-optimal arm i is not eliminated in round $\max(m_i, g_{s_k})$ or before, with the optimal arm $*$ $\in C_{\max(m_i, g_{s_k})}$.

We consider an arbitrary sub-optimal arm i and analyze the contribution to the regret when i is not eliminated in the following exhaustive sub-cases:

Case a1: In round $\max(m_i, g_{s_k})$, $i \in s^*$.

Similar to case (a) of [4], observe that when the following two conditions hold, arm i gets eliminated:

$$\hat{r}_i \leq r_i + c_{m_i} \text{ and } \hat{r}^* \geq r^* - c_{m_i}, \quad (1)$$

where $c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i}^2)}{2n_{m_i}}}$. The arm i gets eliminated because

$$\begin{aligned} \hat{r}_i + c_{m_i} &\leq r_i + 2c_{m_i} < r_i + \Delta_i - 2c_{m_i} = r^* - 2c_{m_i} \\ &\leq \hat{r}^* - c_{m_i}. \end{aligned}$$

In the above, we have used the fact that

$$c_{m_i} = \sqrt{\rho_a \epsilon_{m_i+1}} < \frac{\Delta_i}{4}, \text{ since } n_{m_i} = \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}} \text{ and } \rho_a \in (0, 1].$$

From the foregoing, we have to bound the events complementary to that in (1) for an arm i to not get eliminated. Considering Chernoff-Hoeffding bound this is done as follows:

$$\begin{aligned} \mathbb{P}(\hat{r}^* \leq r^* - c_{m_i}) &\leq \exp(-2c_{m_i}^2 n_{m_i}) \\ &\leq \exp(-2 * \frac{\rho_a \log(\psi T \epsilon_{m_i}^2)}{2n_{m_i}} * n_{m_i}) \\ &\leq \frac{1}{(\psi T \epsilon_{m_i}^2)^{\rho_a}} \end{aligned}$$

Along similar lines, we have $\mathbb{P}(\hat{r}_i \geq r_i + c_{m_i}) \leq \frac{1}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}$. Thus, the probability that a sub-optimal arm i is not eliminated in any round on or before m_i is bounded above by $\left(\frac{2}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}\right)$. Summing up over all arms in A'_{s^*} in conjunction with a simple bound of $T\Delta_i$ for each arm,

we obtain

$$\begin{aligned} \sum_{i \in A'_{s^*}} \left(\frac{2T\Delta_i}{(\psi T \epsilon_{m_i}^2)^{\rho_a}} \right) &\leq \sum_{i \in A'_{s^*}} \left(\frac{2T\Delta_i}{(\psi T \frac{\Delta_i^4}{16\rho_a^2})^{\rho_a}} \right) \\ &= \sum_{i \in A'_{s^*}} \left(\frac{C_1(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} \right), \text{ where } C_1(x) = \frac{2^{1+4x}x^{2x}}{\psi^x} \end{aligned}$$

Case a2: In round $\max(m_i, g_{s_k})$, $i \in s_k$ for some $s_k \neq s^*$.

Following a parallel argument like in Case a1, we have to bound the following two events of arm $a_{\max_{s_k}}$ not getting eliminated on or before g_{s_k} -th round,

$$\hat{r}_{a_{\max_{s_k}}} \geq r_{a_{\max_{s_k}}} + c_{g_{s_k}} \text{ and } \hat{r}^* \leq r^* - c_{g_{s_k}}$$

We can prove using Chernoff-Hoeffding bounds and considering independence of events mentioned above, that for $c_{g_{s_k}} = \sqrt{\frac{\rho_s \log(\psi T \epsilon_{g_{s_k}}^2)}{2n_{g_{s_k}}}}$ and $n_{g_{s_k}} = \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}}$ the probability of the above two events is bounded by $\left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}\right)$. Now, for any round g_{s_k} , all the elements of $C_{\max(m_i, g_{s_k})}$ are the respective maximum payoff arms of their cluster s_k , $\forall s_k \in S$, and since all the surviving arms are pulled equally in each round and since clusters are fixed so we can bound the maximum probability that a sub-optimal arm $i \in A'$ and $i \in s_k$ such that $a_{\max_{s_k}} \in C_{g_{s_k}}$ is not eliminated on or before the g_{s_k} -th round by the same probability as above.

Summing up over all p clusters and bounding the regret for each arm $i \in A'_{s_k}$ trivially by $T\Delta_i$,

$$\begin{aligned} \sum_{k=1}^p \sum_{i \in A'_{s_k}} \left(\frac{2T\Delta_i}{(\psi T \frac{\Delta_i^4}{16\rho_s^2})^{\rho_s}} \right) &= \sum_{i \in A'} \left(\frac{2T\Delta_i}{(\psi T \frac{\Delta_i^4}{16\rho_s^2})^{\rho_s}} \right) \\ &\leq \sum_{i \in A'} \left(\frac{2^{1+4\rho_s} \rho_s^{2\rho_s} T^{1-\rho_s}}{\psi^{\rho_s} \Delta_i^{4\rho_s-1}} \right) = \sum_{i \in A'} \frac{C_1(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \end{aligned}$$

Summing the bounds in Cases a1 – a2 and observing that the bounds in the aforementioned cases hold for any round $C_{\max\{m_i, g_{s_k}\}}$, we obtain the following contribution to the expected regret from case a:

$$\sum_{i \in A_{s^*}} \frac{C_1(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \sum_{i \in A'} \left(\frac{C_1(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right)$$

Case b: For each arm i , either i is eliminated in round $\max(m_i, g_{s_k})$ or before or there is no optimal arm $*$ in $C_{\max(m_i, g_{s_k})}$.

Case b1: $*$ $\in C_{\max(m_i, g_{s_k})}$ for each arm $i \in A'$ and cluster $s_k \in \check{A}$.

The condition in the case description above implies the following:

(i) each sub-optimal arm $i \in A'$ is eliminated on or before $\max(m_i, g_{s_k})$ and hence pulled not more than n_{m_i} number of times.

(ii) each sub-optimal cluster $s_k \in \check{A}$ is eliminated on or before $\max(m_i, g_{s_k})$ and hence pulled not more than $n_{g_{s_k}}$ number of times.

Hence, the maximum regret suffered due to pulling of a sub-optimal arm or a sub-optimal cluster is no more than the following:

$$\begin{aligned} & \sum_{i \in A'} \Delta_i \left\lceil \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}} \right\rceil + \sum_{k=1}^p \sum_{i \in A'_{s_k}} \Delta_i \left\lceil \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}} \right\rceil \\ & \leq \sum_{i \in A'} \Delta_i \left(1 + \frac{32 \rho_a \log \left(\psi T \left(\frac{\Delta_i}{2 \sqrt{\rho_a}} \right)^4 \right)}{\Delta_i^2} \right) \\ & \quad + \sum_{i \in A'} \Delta_i \left(1 + \frac{32 \rho_s \log \left(\psi T \left(\frac{\Delta_i}{2 \sqrt{\rho_s}} \right)^4 \right)}{\Delta_i^2} \right) \\ & \leq \sum_{i \in A'} \left[2 \Delta_i + \frac{32(\rho_a \log(\psi T \frac{\Delta_i^4}{16 \rho_a^2}) + \rho_s \log(\psi T \frac{\Delta_i^4}{16 \rho_s^2}))}{\Delta_i} \right] \end{aligned}$$

In the above, the first inequality follows since $\sqrt{\rho_a \epsilon_{m_i}} < \frac{\Delta_i}{2}$ and $\sqrt{\rho_s \epsilon_{n_{g_{s_k}}}} < \frac{\Delta_{a_{\max s_k}}}{2}$.

Case b2: $*$ is eliminated by some sub-optimal arm in s^*

Optimal arm a^* can get eliminated by some sub-optimal arm i only if arm elimination condition holds, i.e.,

$$\hat{r}_i - c_{m_i} > \hat{r}^* + c_{m_i},$$

where, as mentioned before, $c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i}^2)}{2 n_{m_i}}}$. From analysis in Case a1, notice that, if (1) holds in conjunction with the above, arm i gets eliminated. Also, recall from Case a1 that the events complementary to (1) have low-probability and can be upper bounded by $\frac{2}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}$. Moreover, a sub-optimal arm that eliminates $*$ has to survive until round m_* . In other words, all arms $j \in s^*$ such that $m_j < m_*$ are eliminated on or before m_* (this corresponds to case b1). Let, the arms surviving till m_* round be denoted by A'_{s^*} . 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 A''_{s^*} . Also maximal regret per step after eliminating $*$ is the maximal Δ_j among the remaining arms in A''_{s^*} with $m_j \geq m_*$. Let $m_b = \min\{m | \sqrt{\rho_a \epsilon_m} < \frac{\Delta_b}{2}\}$. Let $C_2(x) = \frac{2^{2x + \frac{3}{2}} x^{2x}}{\psi^x}$. Hence, the maximal regret after eliminating the arm $*$ is

upper bounded by,

$$\begin{aligned} & \max_{j \in A'_{s^*}} \sum_{m_*=0}^{m_j} \sum_{\substack{i \in A''_{s^*}: \\ m_i \geq m_*}} \left(\frac{2}{(\psi T \epsilon_{m_*}^2)^{\rho_a}} \right) T \max_{\substack{j \in A''_{s^*}: \\ m_j \geq m_*}} \Delta_j \\ & \leq \sum_{m_*=0}^{\max_{j \in A'_{s^*}} m_j} \sum_{i \in A''_{s^*}: m_i \geq m_*} \left(\frac{2}{(\psi T \epsilon_{m_*}^2)^{\rho_a}} \right) T \cdot 2 \sqrt{\rho_a \epsilon_{m_*}} \\ & \leq \sum_{m_*=0}^{\max_{j \in A'_{s^*}} m_j} \sum_{i \in A''_{s^*}: m_i \geq m_*} 4 \left(\frac{T^{1-\rho_a}}{\psi^{\rho_a} \epsilon_{m_*}^{2\rho_a - \frac{1}{2}}} \right) \\ & \leq \sum_{i \in A''_{s^*}: m_i \geq m_*} \sum_{m_*=0}^{\min\{m_i, m_b\}} \left(\frac{4 T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(2\rho_a - \frac{1}{2})m_*}} \right) \\ & \leq \sum_{i \in A'_{s^*}} \frac{4 T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(2\rho_a - \frac{1}{2})m_*}} + \sum_{i \in A''_{s^*} \setminus A'_{s^*}} \frac{4 T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(2\rho_a - \frac{1}{2})m_b}} \\ & \leq \sum_{i \in A'_{s^*}} \frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{2\rho_a + \frac{3}{2}}}{\psi^{\rho_a} \Delta_i^{4\rho_a - 1}} + \sum_{i \in A''_{s^*} \setminus A'_{s^*}} \frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{2\rho_a + \frac{3}{2}}}{\psi^{\rho_a} b^{4\rho_a - 1}} \\ & = \sum_{i \in A'_{s^*}} \frac{C_2(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a - 1}} + \sum_{i \in A''_{s^*} \setminus A'_{s^*}} \frac{C_2(\rho_a) T^{1-\rho_a}}{b^{4\rho_a - 1}}. \end{aligned}$$

Case b3: s^* is eliminated by some sub-optimal cluster.

Let $C'_g = \{a_{\max s_k} \in A' | \forall s_k \in S\}$ and $C''_g = \{a_{\max s_k} \in A'' | \forall s_k \in S\}$. A sub-optimal cluster s_k will eliminate s^* in round g_* only if the cluster elimination condition of Algorithm 1 holds, which is the following when $*$ $\in C_{g_*}$:

$$\hat{r}_{a_{\max s_k}} - c_{g_*} > \hat{r}^* + c_{g_*}. \quad (2)$$

Notice that when $*$ $\notin C_{g_*}$, since $r_{a_{\max s_k}} > r^*$, the inequality in (2) has to hold for cluster s_k to eliminate s^* . As in case b2, the probability that a given sub-optimal cluster s_k eliminates s^* is upper bounded by $\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}$ and all sub-optimal clusters with $g_{s_j} < g_*$ are eliminated before round g_* .

This leaves any arm $a_{\max s_b}$ such that $g_{s_b} \geq g_*$ to still survive and eliminate arm $*$ in round g_* . Let, such arms that survive $*$ belong to C''_g . Hence, following the same way as case b2, the maximal regret after eliminating $*$ is,

$$\begin{aligned} & \max_{a_{\max s_j} \in C'_g} \sum_{g_*=0}^{g_{s_j}} \sum_{\substack{a_{\max s_k} \in C''_g: \\ g_{s_k} \geq g_*}} \left(\frac{2}{(\psi T \epsilon_{g_*}^2)^{\rho_s}} \right) T \max_{\substack{a_{\max s_j} \in C''_g: \\ g_{s_j} \geq g_*}} \Delta_{a_{\max s_j}} \end{aligned}$$

Using $A' \supset C'_g$ and $A'' \supset C''_g$, we can bound the regret contribution from this case in a similar manner as Case b2

as follows:

$$\begin{aligned} & \sum_{i \in A' \setminus A'_{s^*}} \frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s + \frac{3}{2}}}{\psi^{\rho_s} \Delta_i^{4\rho_s-1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s + \frac{3}{2}}}{\psi^{\rho_s} b^{4\rho_s-1}} \\ &= \sum_{i \in A' \setminus A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s-1}} \end{aligned}$$

Case b4: $*$ is not in $C_{\max(m_i, g_{s_k})}$, but belongs to $B_{\max(m_i, g_{s_k})}$.

In this case the optimal arm $*$ $\in s^*$ is not eliminated, also s^* is not eliminated. So, for all sub-optimal arms i in A'_{s^*} which gets eliminated on or before $\max\{m_i, g_{s_k}\}$ will get pulled no less than $\left\lceil \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}} \right\rceil$ number of times, which leads to the following bound the contribution to the expected regret, as in Case b1:

$$\sum_{i \in A'_{s^*}} \left\{ \Delta_i + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \right\}$$

For arms $a_i \notin s^*$, the contribution to the regret cannot be greater than that in Case b3. So the regret is bounded by,

$$\sum_{i \in A' \setminus A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s-1}}$$

The main claim follows by summing the contributions to the expected regret from each of the cases above. \square

Proposition 1. The regret R_T for ClusUCB-AE satisfies

$$\begin{aligned} \mathbb{E}[R_T] \leq & \sum_{\substack{i \in A \\ \Delta_i > b}} \left\{ \frac{C_1(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \Delta_i + \frac{32\rho_a \log(\frac{\psi T \Delta_i^4}{16\rho_a^2})}{\Delta_i} \right. \\ & \left. + \frac{C_2(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} \right\} + \sum_{\substack{i \in A \\ 0 < \Delta_i \leq b}} \frac{C_2(\rho_a) T^{1-\rho_a}}{b^{4\rho_a-1}} + \max_{\substack{i \in A \\ \Delta_i \leq b}} \Delta_i T, \end{aligned}$$

for all $b \geq \sqrt{\frac{K}{T}}$. In the above, C_1, C_2 are as defined in Theorem 1.

Proof. See Appendix A. \square

Proposition 2. The regret R_T for ClusUCB-CE satisfies,

$$\mathbb{E}[R_T] \leq \sum_{\substack{i \in A \\ \Delta_i > b}} \left\{ \frac{C_1(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \frac{64\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} \right\}$$

$$\begin{aligned} & + 2\Delta_i \Big\} + \sum_{\substack{i \in A \setminus A_{s^*} \\ \Delta_i > b}} \frac{2C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \\ & + \sum_{\substack{i \in A \setminus A_{s^*} \\ 0 < \Delta_i \leq b}} \frac{2C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s-1}} + \max_{i \in A: \Delta_i \leq b} \Delta_i T, \end{aligned}$$

for all $b \geq \sqrt{\frac{K}{T}}$, with C_1 and C_2 as defined in Theorem 1.

Proof. See Appendix B. \square

5 Simulation experiments

For the purpose of performance comparison using cumulative regret as the metric, we implement the following algorithms: KL-UCB[12], DMED[13], MOSS[2], UCB1[5], UCB-Improved[4], Median Elimination[10], Thompson Sampling(TS)[1] and UCB-V[3]². The parameters of ClusUCB algorithm for all the experiments are set as follows: $\psi = \log T$, $\rho_s = 0.5$ and $\rho_a = 0.25$. When K is large and p is small it is advantageous to run $\rho_a < \rho_s$ (see Corollary 2) because this will aggressively eliminate arms within each cluster while full cluster elimination will be more conservative.

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

The second experiment is conducted over a testbed of 100 arms involving Gaussian reward distribution with expected rewards of the arms $r_{i \neq *:1-33} = 0.01$, $r_{i \neq *:34-99} = 0.06$ and $r_{i=100}^* = 0.1$ with variance set at $\sigma^2 = 0.3, \forall i \in A$. The horizon T is set for a large duration of 2×10^6 and the number of clusters $p = 20$. The regret is averaged over 100 independent runs and is shown in Figure 1b. In this case, in addition to EClusUCB, we also show the performance of ClusUCB algorithm (with $p = 10$). From the results in Figure 1b, we observe that EClusUCB with $p = 10$ outperforms ClusUCB with $p = 10$ as well as MOSS, UCB1,

²The implementation for KL-UCB and DMED were taken from [9]

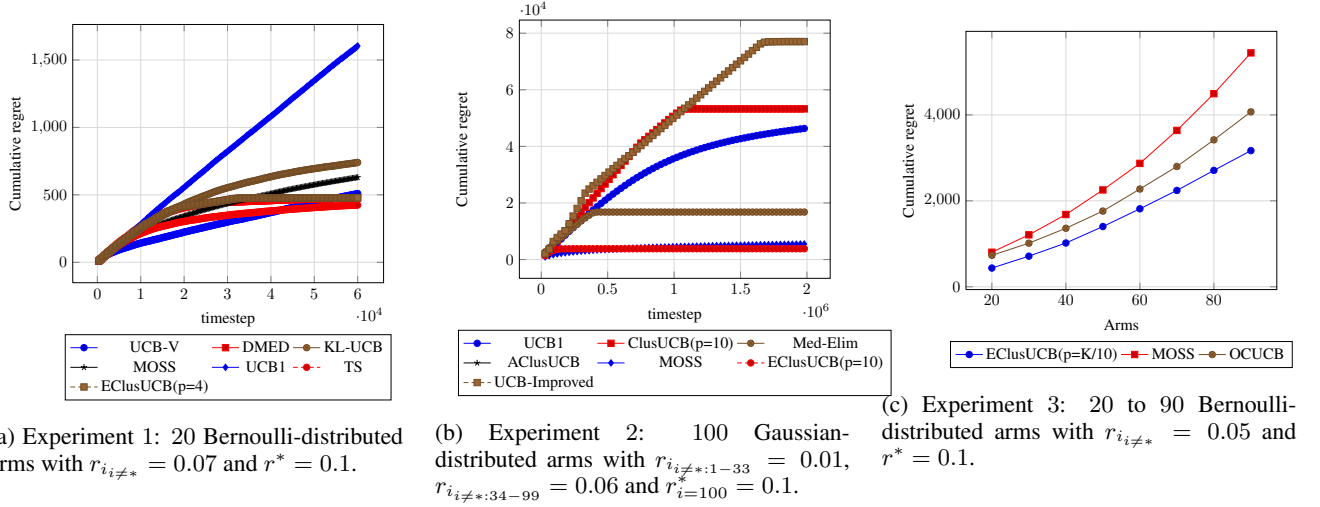
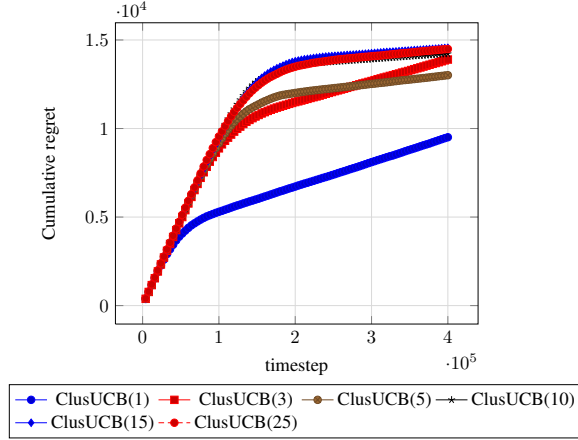


Figure 1: Cumulative regret for various bandit algorithms on three stochastic K-armed bandit environments.



UCB-Improved and Median-Elimination($\epsilon = 0.03, \delta = 0.1$). But as shown in Theorem 1, ClusUCB is better than UCB-Improved even though both are round based methods. Also the performance of UCB-Improved is poor in comparison to other algorithms, which is probably because of pulls wasted in initial exploration whereas ClusUCB with the choice of ψ, ρ_a and ρ_s performs much better.

The third experiment is conducted over a testbed of 20 – 90 (interval of 10) arms with Bernoulli reward distribution, where the expected rewards of the arms are $r_{i \neq *} = 0.05$ and $r^* = 0.1$. The horizon T is set to $10^5 + K^3$ and the number of arms are increased from $K = 20$ to 90. The proposed algorithm EClusUCB is run with $p = K/10$. The regret is averaged over 100 independent runs and is shown in Figure 1c. We report the performance of MOSS, OCUCB and EClusUCB only over this setup. From the results in Figure 1c, it is evident that the growth of regret for EClusUCB is lower than that of MOSS and OCUCB.

The fourth experiment is performed over a testbed having 50 Gaussian-distributed arms with $r_{i \neq *} = 0.8, \forall i \in A$, $r^* = 0.9$ and $\sigma^2 = 1.0$. In Figure 2a, we report the results with $T = 400000$ averaged over 100 independent runs for ClusUCB with $p = \{1, 3, 5, 10, 15, 25\}$. Also, in this experiment we take $\psi = K^2T$, $\rho_a = 0.25$ and $\rho_s = 0.5$ as stated in Corollary 2. The high variance leads to a greater number of errors committed by ClusUCB-AE that is ClusUCB($p = 1$) but as proved in Proposition 1 the cumulative regret is lesser than ClusUCB. But because of the increased errors committed in predicting the optimal arm and because of the large horizon, we eventually see that ClusUCB($p=5, 10, 25, 25$) outperforms ClusUCB-AE while ClusUCB($p = 3$) regret is worse than ClusUCB-AE. The error percentage in the 6 cases (in the order as shown in legend of Fig 2a) are 14, 12, 5, 3, 3 and 3. The range of p is shown to be between $\sqrt{\log K}$ to $\frac{K}{2}$ and as we approach $\frac{K}{2}$ we see that the error percentage stabilizes to 3%. More experiments are shown in Appendix K.

6 Conclusions and future work

From a theoretical viewpoint, we conclude that the gap-dependent regret bound of ClusUCB is lower than MOSS and UCB-Improved. From the numerical experiments on settings with small gaps between optimal and sub-optimal mean rewards, we observed that ClusUCB outperforms several popular bandit algorithms. While we exhibited better regret bounds for ClusUCB, it would be interesting future research to improve the theoretical analysis of ClusUCB to achieve the gap-independent regret bound of MOSS and possibly also the gap-dependent bound conjectured in Section 2.4.3 of [7].

References

- [1] Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit problem. *arXiv preprint arXiv:1111.1797*, 2011.
- [2] Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, pages 217–226, 2009.
- [3] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–1902, 2009.
- [4] Peter Auer and Ronald Ortner. Ucb revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Periodica Mathematica Hungarica*, 61(1-2):55–65, 2010.
- [5] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [6] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- [7] Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *arXiv preprint arXiv:1204.5721*, 2012.
- [8] Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *arXiv preprint arXiv:1209.1727*, 2012.
- [9] Olivier Cappe, Aurelien Garivier, and Emilie Kaufmann. pymabandits, 2012. <http://mloss.org/software/view/415/>.
- [10] Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *The Journal of Machine Learning Research*, 7:1079–1105, 2006.
- [11] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics Springer, Berlin, 2001.
- [12] Aurélien Garivier and Olivier Cappé. The kl-ucb algorithm for bounded stochastic bandits and beyond. *arXiv preprint arXiv:1102.2490*, 2011.
- [13] Junya Honda and Akimichi Takemura. An asymptotically optimal bandit algorithm for bounded support models. In *COLT*, pages 67–79. Citeseer, 2010.
- [14] Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- [15] Tor Lattimore. Optimally confident ucb: Improved regret for finite-armed bandits. *arXiv preprint arXiv:1507.07880*, 2015.
- [16] Yun-Ching Liu and Yoshimasa Tsuruoka. Modification of improved upper confidence bounds for regulating exploration in monte-carlo tree search. *Theoretical Computer Science*, 2016.
- [17] Shie Mannor and John N Tsitsiklis. The sample complexity of exploration in the multi-armed bandit problem. *Journal of Machine Learning Research*, 5(Jun):623–648, 2004.
- [18] Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*, pages 169–177. Springer, 1952.
- [19] William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, pages 285–294, 1933.
- [20] David Tolpin and Solomon Eyal Shimony. Mcts based on simple regret. In *AAAI*, 2012.

Appendix

A Proof of Proposition 1

Proof. Let $p = 1$ whereby we put all the arms in set A into one cluster, that is we have one UCB-Improved running throughout. So, for each sub-optimal arm i , $m_i = \min \{m | \sqrt{\rho_a \epsilon_m} < \frac{\Delta_i}{2}\}$ be the first round when $\sqrt{\rho_a \epsilon_m} < \frac{\Delta_i}{2}$. Also in this proof, since the clusters are fixed, so throughout the rounds each m_i is tied to a single arm. We also take $\rho_a \in (0, 1]$ as a constant in this proof whereby in Corollary 1 and 2 we use the different definitions. The theoretical analysis remains same as we have always bounded the values of $\rho_a \in (0, 1]$. Let $A' = \{i \in A : \Delta_i > b\}$ and $A'' = \{i \in A : \Delta_i > 0\}$.

Case a: Some sub-optimal arm i is not eliminated in round m_i or before and the optimal arm $*$ $\in B_{m_i}$

Following the steps of Theorem 1 Case a1, an arbitrary sub-optimal arm $i \in A'$ can get eliminated only when the event,

$$\hat{r}_i \leq r_i + c_{m_i} \text{ and } \hat{r}^* \geq r^* - c_{m_i} \quad (3)$$

takes place. So to bound the regret we need to bound the probability of the complementary event of these two conditions.

Putting the value of $n_{m_i} = \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}}$ in c_{m_i} , $c_{m_i} = \sqrt{\frac{\rho_a \epsilon_{m_i} \log(\psi T \epsilon_{m_i}^2)}{2 * 2 \log(\psi T \epsilon_{m_i}^2)}} = \frac{\sqrt{\rho_a \epsilon_{m_i}}}{2} = \sqrt{\rho_a \epsilon_{m_i+1}} < \frac{\Delta_i}{4}$, as $\rho_a \in (0, 1]$.

Again, for $i \in A'$,

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

Applying Chernoff-Hoeffding bound and considering independence of complementary of the two events in 3,

$$\begin{aligned} \mathbb{P}\{\hat{r}^* \leq r^* - c_{m_i}\} &\leq \exp(-2c_{m_i}^2 n_{m_i}) \\ &\leq \exp(-2 * \frac{\rho_a \log(\psi T \epsilon_{m_i}^2)}{2n_{m_i}} * n_{m_i}) \\ &\leq \frac{1}{(\psi T \epsilon_{m_i}^2)^{\rho_a}} \end{aligned}$$

Similarly, $\mathbb{P}\{\hat{r}_i \geq r_i + c_{m_i}\} \leq \frac{1}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}$

Summing, the two up, the probability that a sub-optimal arm i is not eliminated on or before m_i -th round is $\left(\frac{2}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}\right)$.

Summing up over all arms in A and bounding the regret for each arm $i \in A'$ trivially by $T \Delta_i$, we obtain

$$\begin{aligned} \sum_{i \in A'} \left(\frac{2T \Delta_i}{(\psi T \epsilon_{m_i}^2)^{\rho_a}} \right) &\leq \sum_{i \in A'} \left(\frac{2T \Delta_i}{(\psi T \frac{\Delta_i^4}{16 \rho_a^2})^{\rho_a}} \right) \\ &\leq \sum_{i \in A'} \left(\frac{2^{1+4\rho_a} T^{1-\rho_a} \rho_a^{2\rho_a} \Delta_i}{\psi \rho_a \Delta_i^{4\rho_a}} \right) \\ &\leq \sum_{i \in A'} \left(\frac{2^{1+4\rho_a} \rho_a^{2\rho_a} T^{1-\rho_a}}{\psi \rho_a \Delta_i^{4\rho_a-1}} \right) \end{aligned}$$

$$= \sum_{i \in A'} \left(\frac{C_1(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} \right), \text{ where } C_1(x) = \frac{2^{1+4x} x^{2x}}{\psi^x}$$

Case b: *Either an arm i is eliminated in round m_i or before or else there is no optimal arm $*$ $\in B_{m_i}$*

Case b1: *$*$ $\in B_{m_i}$ and each $i \in A'$ 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,

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

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

$$\begin{aligned} & \Delta_i \left\lceil \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}} \right\rceil \\ & \leq \Delta_i \left\lceil \frac{2 \log(\psi T (\frac{\Delta_i}{2\sqrt{\rho_a}})^4)}{(\frac{\Delta_i}{2\sqrt{\rho_a}})^2} \right\rceil, \text{ since } \sqrt{\rho_a \epsilon_{m_i}} < \frac{\Delta_i}{2} \\ & \leq \Delta_i \left(1 + \frac{32\rho_a \log(\psi T (\frac{\Delta_i}{2\sqrt{\rho_a}})^4)}{\Delta_i^2} \right) \\ & \leq \Delta_i \left(1 + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i^2} \right) \end{aligned}$$

Summing over all arms in A' the total regret is given by,

$$\sum_{i \in A'} \Delta_i \left(1 + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i^2} \right)$$

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 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 A' . Also maximal regret per step after eliminating $*$ is the maximal Δ_j among the remaining arms j with $m_j \geq m_*$. Let $m_b = \min\{m | \sqrt{\rho_a \epsilon_m} < \frac{\Delta_b}{2}\}$. Hence, the maximal regret after eliminating the arm $*$ is upper bounded by,

$$\begin{aligned} & \sum_{m_*=0}^{\max_{j \in A'} m_j} \sum_{i \in A'' : m_i > m_*} \left(\frac{2}{(\psi T \epsilon_{m_*}^2)^{\rho_a}} \right) \cdot T \max_{j \in A'' : m_j \geq m_*} \Delta_j \\ & \leq \sum_{m_*=0}^{\max_{j \in A'} m_j} \sum_{i \in A'' : m_i > m_*} \left(\frac{2}{(\psi T \epsilon_{m_*}^2)^{\rho_a}} \right) \cdot T \cdot 2\sqrt{\rho_a \epsilon_{m_*}} \\ & \leq \sum_{m_*=0}^{\max_{j \in A'} m_j} \sum_{i \in A'' : m_i > m_*} 4 \left(\frac{T^{1-\rho_a}}{\psi^{\rho_a} \epsilon_{m_*}^{2\rho_a - \frac{1}{2}}} \right) \end{aligned}$$

$$\begin{aligned}
 &\leq \sum_{i \in A'' : m_i > m_*} \sum_{m_*=0}^{\min\{m_i, m_b\}} \left(\frac{4T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(2\rho_a - \frac{1}{2})m_*}} \right) \\
 &\leq \sum_{i \in A'} \left(\frac{4T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(2\rho_a - \frac{1}{2})m_*}} \right) + \sum_{i \in A'' \setminus A'} \left(\frac{4T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(2\rho_a - \frac{1}{2})m_b}} \right) \\
 &\leq \sum_{i \in A'} \left(\frac{4\rho_a^{2\rho_a} T^{1-\rho_a} * 2^{2\rho_a - \frac{1}{2}}}{\psi^{\rho_a} \Delta_i^{4\rho_a - 1}} \right) + \sum_{i \in A'' \setminus A'} \left(\frac{4\rho_a^{2\rho_a} T^{1-\rho_a} * 2^{2\rho_a - \frac{1}{2}}}{\psi^{\rho_a} b^{4\rho_a - 1}} \right) \\
 &\leq \sum_{i \in A'} \left(\frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{2\rho_a + \frac{3}{2}}}{\psi^{\rho_a} \Delta_i^{4\rho_a - 1}} \right) + \sum_{i \in A'' \setminus A'} \left(\frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{2\rho_a + \frac{3}{2}}}{\psi^{\rho_a} b^{4\rho_a - 1}} \right) \\
 &= \sum_{i \in A'} \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a - 1}} \right) + \sum_{i \in A'' \setminus A'} \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{b^{4\rho_a - 1}} \right), \text{ where } C_2(x) = \frac{2^{2x + \frac{3}{2}} x^{2x}}{\psi^x}
 \end{aligned}$$

Summing up **Case a** and **Case b**, the total regret till round m is given by,

$$\begin{aligned}
 R_T \leq & \sum_{i \in A : \Delta_i > b} \left\{ \left(\frac{C_1(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a - 1}} \right) + \left(\Delta_i + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \right) \right. \\
 & \left. + \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a - 1}} \right) \right\} + \sum_{i \in A : 0 < \Delta_i \leq b} \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{\psi^{\rho_a} b^{4\rho_a - 1}} \right) + \max_{i \in A : \Delta_i \leq b} \Delta_i T
 \end{aligned}$$

□

Corollary 3. For $\rho_a = 1$ in the result of proposition 1 for ClusUCB-AE, we get a regret bound of

$$\sum_{i \in A : \Delta_i > b} \left(\Delta_i + \frac{44}{\psi(\Delta_i)^3} + \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i} \right) + \sum_{i \in A : 0 < \Delta_i \leq b} \frac{12}{\psi b^3}$$

Proof. In the result of Proposition 1 if we take $\rho_a = 1$ then the regret bound becomes $\sum_{i \in A : \Delta_i > b} \left(\Delta_i + \frac{44}{\psi(\Delta_i)^3} + \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i} \right) + \sum_{i \in A : 0 < \Delta_i \leq b} \frac{12}{\psi b^3}$. From the result we can see that for small Δ_i and large K , the terms like $\sum_{i \in A : \Delta_i > b} \left(\frac{44}{\psi(\Delta_i)^3} \right) + \sum_{i \in A : 0 < \Delta_i \leq b} \frac{12}{\psi b^3}$ can become the dominant term in the regret rather than $\sum_{i \in A : \Delta_i > b} \frac{32 \log(\psi T \Delta_i^4)}{\Delta_i}$. Intuitively, this actually suggests that the algorithm is trying to eliminate arms with too low exploration and so the probability of elimination is low and error(risk) is high. For this essentially we introduce ρ_a, ρ_s and ψ and by carefully defining their values enables us to eliminate arms and clusters aggressively and thereby reduce those two terms. □

B Proof of Proposition 2

An illustrative diagram explaining Cluster Elimination is given in **Figure 2**. A slight modification to the algorithm allows us to do cluster elimination without any arm elimination. By taking $p > 1$, removing the arm elimination condition, stopping when we are just left with one cluster and pulling the $\max\{\hat{r}_i\}$, where $i \in B_m$ we can achieve this. We also take $\rho_s \in (0, 1]$ as a constant in this proof whereby in Corollary 1 and 2 we use the different definitions. The theoretical analysis remains same as we have always bounded the values of $\rho_s \in (0, 1]$.

Proof. Let $C_{g_{s_k}} = \{\hat{r}_{a_{\max_{s_k}}} | \forall s_k \in S\}$, that is let $C_{g_{s_k}}$ be the set of all arms which has the maximum estimated payoff arms from their respective clusters in the g_{s_k} -th round. Let, for each sub-optimal cluster arm $a_{\max_{s_k}} \in C_{g_{s_k}}$, $g_{s_k} =$

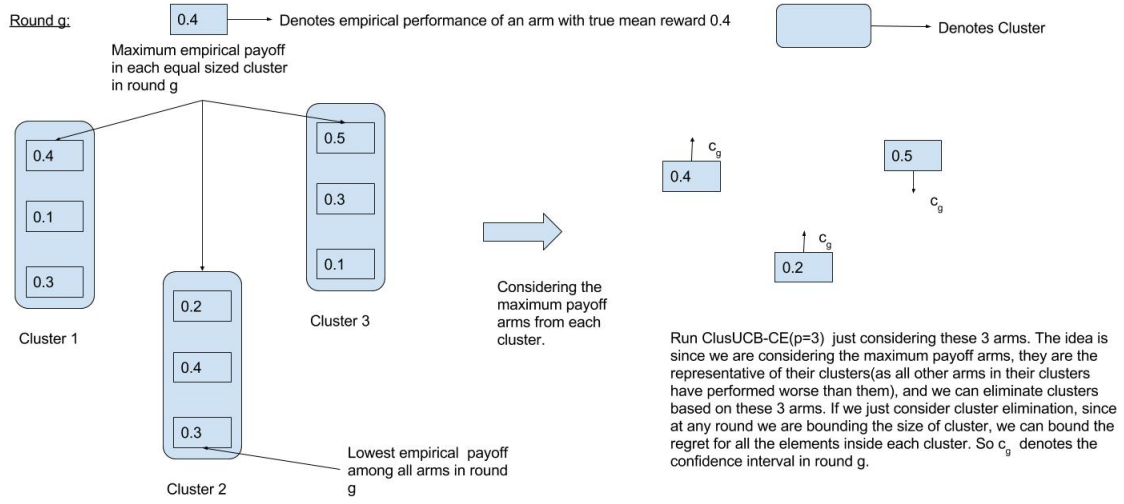


Figure 2: Cluster Elimination

$\min \{g | \sqrt{\rho_s \epsilon_g} < \frac{\Delta_{a_{\max_{s_k}}}}{2}\}$. So, g_{s_k} be the first round when $\sqrt{\rho_s \epsilon_{g_{s_k}}} < \frac{\Delta_{a_{\max_{s_k}}}}{2}$ where $a_{\max_{s_k}} \in C_{g_{s_k}}$ is the maximum payoff arm in cluster s_k . Here, $a_{\max_{s_k}}$ is called cluster arm. Here, A_{s_k} denotes the arm set in the cluster s_k . Let $A'_{s_k} = \{i \in A_{s_k} : \Delta_i > b\}$, $A''_{s_k} = \{i \in A_{s_k} : \Delta_i > 0\}$, $A' = \{i \in A : \Delta_i > b\}$ and $A'' = \{i \in A : \Delta_i > 0\}$.

Case a: Some sub-optimal cluster arm $a_{\max_{s_k}}$ is not eliminated in round g_{s_k} or before with $*$ $\in C_{g_{s_k}}$

Following the steps of Theorem 1 Case a2, an arbitrary sub-optimal arm $i \in A'$ can get eliminated only when the event,

$$\hat{r}_{a_{\max_{s_k}}} \leq r_{a_{\max_{s_k}}} + c_{m_i} \text{ and } \hat{r}^* \geq r^* - c_{m_i} \quad (4)$$

takes place. So to bound the regret we need to bound the probability of the complementary event of these two conditions.

Putting the value of $n_{g_{s_k}} = \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}}$ in $c_{g_{s_k}}$ we get,

$$\begin{aligned} c_{g_{s_k}} &= \sqrt{\frac{\rho_s * \epsilon_{g_{s_k}} \log(\psi T \epsilon_{g_{s_k}}^2)}{2 * 2 \log(\psi T \epsilon_{g_{s_k}}^2)}} \\ &= \sqrt{\frac{\rho_s \epsilon_{g_{s_k}}}{2}} \\ &= \sqrt{\rho_s \epsilon_{g_{s_k}+1}} < \frac{\sqrt{\rho_s} \Delta_{a_{\max_{s_k}}}}{4} < \frac{\Delta_{a_{\max_{s_k}}}}{4} \end{aligned}$$

Again, for $a_{\max_{s_k}}, * \in C_{g_{s_k}}$,

$$\begin{aligned} \hat{r}_{a_{\max_{s_k}}} + c_{g_{s_k}} &\leq r_{a_{\max_{s_k}}} + 2c_{g_{s_k}} = \hat{r}_{a_{\max_{s_k}}} + 4c_{g_{s_k}} - 2c_{g_{s_k}} \\ &< r_{a_{\max_{s_k}}} + \Delta_{a_{\max_{s_k}}} - 2c_{g_{s_k}} \\ &= r^* - 2c_{g_{s_k}} \\ &\leq \hat{r}^* - c_{g_{s_k}} \end{aligned}$$

Applying Chernoff-Hoeffding bound and considering independence of complementary of the two events in 4,

$$\begin{aligned}
 \mathbb{P}\left\{\hat{r}^* \leq r^* - c_{g_{s_k}}\right\} &\leq \exp(-2c_{g_{s_k}}^2 n_{g_{s_k}}) \\
 &\leq \exp(-2 * \frac{\rho_s \log(\psi T \epsilon_{g_{s_k}}^2)}{2n_{g_{s_k}}} * n_{g_{s_k}}) \\
 &\leq \frac{1}{(\psi T \epsilon_{g_k}^2)^{\rho_s}}
 \end{aligned}$$

Similarly, $\mathbb{P}\left\{\hat{r}_{a_{max_{s_k}}} \geq r_{a_{max_{s_k}}} + c_{g_{s_k}}\right\} \leq \frac{1}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}$

Summing, the two up, the probability that a sub-optimal cluster arm $a_{max_{s_k}} \in C_{g_{s_k}}$ is not eliminated in g_{s_k} -th round is $\left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}\right)$. Now, for each round g_{s_k} , all the elements of $C_{g_{s_k}}$ are the respective max payoff arms of their cluster s_k , that is all the other arms in their respective clusters have performed worse than them. Hence, since $A'_{s_k} \supset C_{g_{s_k}}$, we are pulling all the surviving arms equally in each round and since clusters are fixed so we can bound the maximum probability that a sub-optimal arm $j \in A'_{s_k}$ and $j \in s_k$ such that $a_{max_{s_k}} \in C_{g_{s_k}}$ is not eliminated on or before the g_{s_k} -th round by the same probability of

$$\left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}\right)$$

Summing up over all arms in s_k and bounding the regret trivially by $T\Delta_i$,

$$\sum_{i \in A'_{s_k}} \left(\frac{2T\Delta_i}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}\right)$$

Summing up over all p clusters and bounding the regret for each arm $i \in A'_{s_k}$ trivially by $T\Delta_i$,

$$\begin{aligned}
 \sum_{k=1}^p \sum_{i \in A'_{s_k}} \left(\frac{2T\Delta_i}{(\psi T \frac{\Delta_i^4}{16\rho_s^2})^{\rho_s}}\right) &= \sum_{i \in A'} \left(\frac{2T\Delta_i}{(\psi T \frac{\Delta_i^4}{16\rho_s^2})^{\rho_s}}\right) \\
 &\leq \sum_{i \in A'} \left(\frac{2^{1+4\rho_s} T^{1-\rho_s} \rho_s^{2\rho_s} \Delta_i}{\psi^{\rho_s} \Delta_i^{4\rho_s}}\right) \\
 &\leq \sum_{i \in A'} \left(\frac{2^{1+4\rho_s} \rho_s^{2\rho_s} T^{1-\rho_s}}{\psi^{\rho_s} \Delta_i^{4\rho_s-1}}\right) \\
 &= \sum_{i \in A'} \left(\frac{C_1(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}}\right), \text{ where } C_1(x) = \frac{2^{1+4x} x^{2x}}{\psi^x}
 \end{aligned}$$

Case b: For each arm i , either i is eliminated in round g_{s_k} or before or there is no optimal arm $*$ in $C_{g_{s_k}}$

Case b1: $*$ $\in C_{g_{s_k}}$ for each arm $i \in A'$ and cluster s_k eliminated on or before g_{s_k}

Again, in the g_{s_k} -th round, the maximum total elements in the cluster s_k can be no more than $\ell = \left\lceil \frac{K}{p} \right\rceil$.

Also, since we are eliminating a sub-optimal cluster arm $a_{max_{s_k}} \in C_{g_{s_k}}$ on or before round g_{s_k} , it is pulled (along with all the other arms in that cluster) no longer than,

$$n_{g_{s_k}} = \left\lceil \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}} \right\rceil$$

So, the total contribution of $a_{\max_{s_k}}$ along with all the other arms in the cluster till round g_{s_k} is given by,

$$\begin{aligned}
 & \sum_{i \in A_{s_k}} \Delta_i \left\lceil \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}} \right\rceil \\
 & \leq \sum_{i \in A'_{s_k}} \Delta_i \left\lceil \frac{2 \log(\psi T (\frac{\Delta_i}{2\sqrt{\rho_s}})^4)}{(\frac{\Delta_i}{2\sqrt{\rho_s}})^2} \right\rceil, \text{ since } \sqrt{\rho_s \epsilon_{g_{s_k}}} < \frac{\Delta_{a_{\max_{s_k}}}}{2} < \frac{\Delta_i}{2}, \text{ as } r_{a_{\max_{s_k}}} > r_i, \forall i \in s_k \\
 & \leq \sum_{i \in A'_{s_k}} \Delta_i \left(1 + \frac{32 * \rho_s * \log(\psi T (\frac{\Delta_i}{2\sqrt{\rho_s}})^4)}{\Delta_i^2} \right) \\
 & \leq \sum_{i \in A'_{s_k}} \Delta_i \left(1 + \frac{32 \rho_s \log(\psi T \frac{\Delta_i^4}{16 \rho_s^2})}{\Delta_i^2} \right)
 \end{aligned}$$

Summing over all p clusters the total regret is given by,

$$\begin{aligned}
 & \sum_{k=1}^p \sum_{i \in A'_{s_k}} \Delta_i \left(1 + \frac{32 \rho_s \log(\psi T \frac{\Delta_i^4}{16 \rho_s^2})}{\Delta_i^2} \right) \\
 & \leq \sum_{i \in A'} \Delta_i \left(1 + \frac{32 \rho_s \log(\psi T \frac{\Delta_i^4}{16 \rho_s^2})}{\Delta_i^2} \right)
 \end{aligned}$$

Case b2: s^* is eliminated by some sub-optimal cluster.

Let $C'_g = \{a_{\max_{s_k}} \in A' | \forall s_k \in S\}$ and $C''_g = \{a_{\max_{s_k}} \in A'' | \forall s_k \in S\}$. Firstly, if conditions of case b1 holds then the optimal arm $*$ $\in C_{g_{s_k}}$ will not be eliminated in round $g_{s_k} = g_*$ or it will lead to the contradiction that $r_{a_{\max_{s_k}}} > r^*$ where $a_{\max_{s_k}}, * \in C_{g_{s_k}}$. In any round g_* , if the optimal arm $*$ gets eliminated then for any round from 1 to g_{s_j} all arms $a_{\max_{s_j}} \in C_{g_{s_k}}, \forall s_j \neq s^*$ such that $g_{s_j} < g_*$ were eliminated according to assumption in Case a. Let, the arms surviving till g_* round be denoted by C'_g . This leaves any arm a_{s_b} such that $g_{s_b} \geq g_*$ to still survive and eliminate arm $*$ in round g_* . Let, such arms that survive $*$ belong to C''_g . Also maximal regret per step after eliminating $*$ is the maximal $\Delta_{a_{\max_{s_j}}}$ among the remaining arms $a_{\max_{s_j}} \in C_{g_{s_j}}$ with $g_{s_j} \geq g_*$. Hence, the maximal regret after eliminating the arm $*$ is upper bounded by,

$$\begin{aligned}
 & \max_{\substack{a_{\max_{s_j}} \in C'_g \\ g_{s_j} \geq g_*}} \sum_{g_{s_k} \geq g_*} \left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}} \right) \cdot T \max_{\substack{a_{\max_{s_j}} \in C''_g \\ g_{s_j} \geq g_*}} \Delta_{a_{\max_{s_j}}}
 \end{aligned}$$

But, we know that for any round g , elements of C_g are the best performers in their respective clusters. So, taking that into account and $A' \supset C'_g$ and $A'' \supset C''_g$ the regret can be bounded by,

$$\begin{aligned}
 & \max_{\substack{j \in A' \setminus A'_{s^*} \\ g_{s_j} \geq g_*}} \sum_{g_{s_k} \geq g_*} \left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}} \right) \cdot T \max_{\substack{j \in A'' \setminus A'_{s^*} \\ g_{s_j} \geq g_*}} \Delta_j \\
 & \leq \sum_{g_{s_k} \geq g_*} \sum_{i \in A'' \setminus A'_{s^*} : g_{s_k} > g_*} \left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}} \right) \cdot T \cdot 2\sqrt{\rho_s \epsilon_{g_{s_j}}}
 \end{aligned}$$

$$\begin{aligned}
 &\leq \sum_{g_*=0}^{\max_{j \in A' \setminus A'_{s^*}} g_{sj}} \sum_{i \in A'' \setminus A'_{s^*} : g_{s_k} > g_*} \left(\frac{4T^{1-\rho_s}}{\psi^{\rho_s} \epsilon_{g_{s_k}}^{2\rho_s - \frac{1}{2}}} \right) \\
 &\leq \sum_{i \in A'' \setminus A'_{s^*} : g_{s_k} > g_*} \sum_{g_*=0}^{\min\{g_{s_k}, g_{s_b}\}} \left(\frac{4T^{1-\rho_s}}{\psi^{\rho_s} 2^{(2\rho_s - \frac{1}{2})g_*}} \right) \\
 &\leq \sum_{i \in A' \setminus A'_{s^*}} \left(\frac{4T^{1-\rho_s}}{\psi^{\rho_s} 2^{(2\rho_s - \frac{1}{2})g_*}} \right) + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \left(\frac{4T^{1-\rho_s}}{\psi^{\rho_s} 2^{(2\rho_s - \frac{1}{2})g_{s_b}}} \right) \\
 &\leq \sum_{i \in A' \setminus A'_{s^*}} \left(\frac{4\rho_s^{2\rho_s} T^{1-\rho_s} * 2^{2\rho_s - \frac{1}{2}}}{\psi^{\rho_s} \Delta_i^{4\rho_s - 1}} \right) + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \left(\frac{4\rho_s^{2\rho_s} T^{1-\rho_s} * 2^{2\rho_s - \frac{1}{2}}}{\psi^{\rho_s} b^{4\rho_s - 1}} \right) \\
 &\leq \sum_{i \in A' \setminus A'_{s^*}} \left(\frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s + \frac{3}{2}}}{\psi^{\rho_s} \Delta_i^{4\rho_s - 1}} \right) + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \left(\frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s + \frac{3}{2}}}{\psi^{\rho_s} b^{4\rho_s - 1}} \right) \\
 &= \sum_{i \in A' \setminus A'_{s^*}} \left(\frac{C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s - 1}} \right) + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \left(\frac{C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s - 1}} \right), \text{ where } C_2(x) = \frac{2^{2x + \frac{3}{2}} x^{2x}}{\psi^x}
 \end{aligned}$$

Case b3: $*$ is not in $C_{g_{s_k}}$, but belongs to $B_{g_{s_k}}$

In this case the optimal arm $* \in s^*$ is not eliminated, also s^* is not eliminated. So, for all sub-optimal arms i in A' which gets eliminated on or before g_{s_k} will get pulled no less than $\left\lceil \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}} \right\rceil$ number of times, which leads to the following bound the contribution to the expected regret, as in Case b1:

$$\sum_{i \in A'} \left\{ \Delta_i + \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} \right\}$$

For arms $a_i \notin s^*$, the contribution to the regret cannot be greater than that in Case b2. So the regret is bounded by,

$$\sum_{i \in A' \setminus A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s - 1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s - 1}}$$

Summing up **Case a** and **Case b**, the total regret till round g is given by,

$$\begin{aligned}
 R_T &\leq \sum_{i \in A : \Delta_i > b} \left\{ \underbrace{\left(\frac{C_1(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s - 1}} \right)}_{\text{case a}} + \underbrace{\left(2\Delta_i + \frac{64\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} \right)}_{\text{case b1+b3}} \right\} \\
 &\quad + \underbrace{\sum_{i \in A \setminus A_{s^*} : \Delta_i > b} \left(\frac{2C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s - 1}} \right)}_{\text{case b2}} + \underbrace{\sum_{i \in A \setminus A_{s^*} : 0 < \Delta_i \leq b} \left(\frac{2C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s - 1}} \right)}_{\text{case b2}} + \max_{i \in A : \Delta_i \leq b} \Delta_i T
 \end{aligned}$$

□

C Proof of Corollary 1

Proof. Here we take $\psi = \frac{T}{\log(K)}$, $\rho_a = \frac{1}{2}$ and $\rho_s = \frac{1}{2}$. Taking into account Theorem 1 below for all $b \geq \sqrt{\frac{e}{T}}$

$$\begin{aligned}
 \mathbb{E}[R_T] \leq & \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \left\{ \frac{C_1(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \Delta_i + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \right\} + \sum_{\substack{i \in A, \\ \Delta_i > b}} \left\{ 2\Delta_i + \frac{C_1(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right. \\
 & + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} + \left. \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} \right\} + \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \sum_{\substack{i \in A_{s^*}, \\ 0 < \Delta_i \leq b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{b^{4\rho_a-1}} \\
 & + \sum_{i \in A \setminus A_{s^*}: \Delta_i > b} \frac{2C_2(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{i \in A \setminus A_{s^*}: 0 < \Delta_i \leq b} \frac{2C_2(\rho_s)T^{1-\rho_s}}{b^{4\rho_s-1}} + \max_{i: \Delta_i \leq b} \Delta_i T
 \end{aligned}$$

and putting the parameter values in the above Theorem 1 result,

$$\sum_{i \in A_{s^*}: \Delta_i > b} \left(\frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{1+4\rho_a}}{\psi^{\rho_a} \Delta_i^{4\rho_a-1}} \right) = \sum_{i \in A_{s^*}: \Delta_i > b} \left(\frac{T^{1-\frac{1}{2}} \frac{1}{2} 2^{*\frac{1}{2}} 2^{1+4*\frac{1}{2}}}{\left(\frac{T}{\log(K)}\right)^{\frac{1}{2}} \Delta_i^{4*\frac{1}{2}-1}} \right) = \sum_{i \in A_{s^*}: \Delta_i > b} \frac{4\sqrt{\log(K)}}{\Delta_i}$$

Similarly for the term,

$$\sum_{i \in A: \Delta_i > b} \left(\frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{1+4\rho_s}}{\psi^{\rho_s} \Delta_i^{4\rho_s-1}} \right) = \sum_{i \in A: \Delta_i > b} \frac{4\sqrt{\log(K)}}{\Delta_i}$$

For the term involving arm pulls,

$$\sum_{i \in A: \Delta_i > b} \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} = \sum_{i \in A: \Delta_i > b} \frac{16 \log(T^2 \frac{\Delta_i^4}{4 \log(K)})}{\Delta_i} \approx \sum_{i \in A: \Delta_i > b} \frac{32 \log(T \frac{\Delta_i^2}{\sqrt{\log(K)}})}{\Delta_i}$$

Similarly the term,

$$\sum_{i \in A: \Delta_i > b} \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \approx \sum_{i \in A: \Delta_i > b} \frac{32 \log(T \frac{\Delta_i^2}{\sqrt{\log(K)}})}{\Delta_i}$$

Lastly we can bound the error terms as,

$$\sum_{i \in A_{s^*}: 0 < \Delta_i \leq b} \left(\frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{2\rho_a + \frac{3}{2}}}{\psi^{\rho_a} \Delta_i^{4\rho_a-1}} \right) = \sum_{i \in A_{s^*}: 0 < \Delta_i \leq b} \frac{2.8\sqrt{\log(K)}}{\Delta_i}$$

Similarly for the term,

$$\sum_{i \in A \setminus A_{s^*}: 0 < \Delta_i \leq b} \left(\frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s + \frac{3}{2}}}{(\psi^{\rho_s}) \Delta_i^{4\rho_s-1}} \right) = \sum_{i \in A \setminus A_{s^*}: 0 < \Delta_i \leq b} \frac{2.8\sqrt{\log(K)}}{\Delta_i}$$

So, the total gap dependent regret bound for using both arm and cluster elimination comes of as

$$\sum_{i \in A_{s^*}: \Delta_i > b} \left\{ \frac{4\sqrt{\log(K)}}{\Delta_i} + \Delta_i + \frac{32 \log(T \frac{\Delta_i^2}{\sqrt{\log(K)}})}{\Delta_i} \right\} + \sum_{i \in A: \Delta_i > b} \left\{ \frac{4\sqrt{\log(K)}}{\Delta_i} + 2\Delta_i + \frac{64 \log(T \frac{\Delta_i^2}{\sqrt{\log(K)}})}{\Delta_i} \right\}$$

$$\begin{aligned}
 & + \sum_{i \in A_{s^*} : \Delta_i > b} \frac{2.8\sqrt{\log(K)}}{\Delta_i} + \sum_{i \in A_{s^*} : 0 < \Delta_i \leq b} \frac{2.8\sqrt{\log(K)}}{\Delta_i} \\
 & + \sum_{i \in A \setminus A_{s^*} : \Delta_i > b} \frac{5.6\sqrt{\log(K)}}{\Delta_i} + \sum_{i \in A \setminus A \cup A_{s^*} : 0 < \Delta_i \leq b} \frac{5.6\sqrt{\log(K)}}{\Delta_i} + \max_{i \in A : \Delta_i \leq b} \Delta_i T
 \end{aligned}$$

□

D Proof of Corollary 2

Proof. As stated in [4], we can have a bound on regret of the order of $\sqrt{KT \log K}$ in non-stochastic setting. This is shown in Exp4([6]) algorithm. Similarly, by choosing $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}}$ for all $i : i \neq * \in A$, in the bound of UCB1([5]) we get,

$$\sum_{i: r_i < r^*} \text{const} \frac{\log T}{\Delta_i} = \frac{\sqrt{KT} \log T}{\sqrt{\log K}}$$

So, this bound is worse than the non-stochastic setting and is clearly improvable and an upper bound regret of \sqrt{KT} is possible as shown in [2] for MOSS which is consistent with the lower bound as proposed by Mannor and Tsitsiklis([17]).

Hence, we take $b \approx \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{e}{T}}$ (the minimum value for b), $\psi = \frac{T}{\log K}$, $\rho_a = \frac{1}{2}$ and $\rho_s = \frac{1}{2}$.

Taking into account Theorem 1 below,

$$\begin{aligned}
 \mathbb{E}[R_T] & \leq \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \left\{ \frac{C_1(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \Delta_i + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \right\} + \sum_{\substack{i \in A, \\ \Delta_i > b}} \left\{ 2\Delta_i + \frac{C_1(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right. \\
 & + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} + \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} \left. \right\} + \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \sum_{\substack{i \in A_{s^*}, \\ 0 < \Delta_i \leq b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{b^{4\rho_a-1}} \\
 & + \sum_{i \in A \setminus A_{s^*} : \Delta_i > b} \frac{2C_2(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{i \in A \setminus A_{s^*} : 0 < \Delta_i \leq b} \frac{2C_2(\rho_s)T^{1-\rho_s}}{b^{4\rho_s-1}} + \max_{i: \Delta_i \leq b} \Delta_i T
 \end{aligned}$$

and putting the parameter values in the above Theorem 1 result,

$$\sum_{i \in A_{s^*} : \Delta_i > b} \left(\frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{1+4\rho_a}}{\psi^{\rho_a} \Delta_i^{4\rho_a-1}} \right) = \left(K \frac{T^{1-\frac{1}{2}} \frac{1}{2}^{\frac{1}{2}} 2^{1+4\frac{1}{2}}}{p(\frac{T}{\log K})^{\frac{1}{2}} \Delta_i^{\frac{4}{2}-1}} \right) = 4 \frac{\sqrt{KT}}{p}$$

Similarly, for the term,

$$\sum_{i \in A : \Delta_i > b} \left(\frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{1+4\rho_s}}{\psi^{\rho_s} \Delta_i^{4\rho_s-1}} \right) = 4\sqrt{KT}$$

For the term regarding number of pulls,

$$\sum_{i \in A : \Delta_i > b} \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} = \frac{32K\sqrt{T}^{\frac{1}{2}} \log(T^2 \frac{K^4(\log K)^2}{T^2 \log K})}{\sqrt{K \log K}}$$

$$\begin{aligned}
 &\leq \frac{32\sqrt{KT} \log(K^2(\sqrt{\log K}))}{\sqrt{\log K}} \\
 &= 64\sqrt{KT \log K} + \frac{32\sqrt{KT} \log(\sqrt{\log K})}{\sqrt{\log K}}
 \end{aligned}$$

Similarly for the term,

$$\sum_{i \in A: \Delta_i > b} \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} = 64\sqrt{KT \log K} + \frac{32\sqrt{KT} \log(\sqrt{\log K})}{\sqrt{\log K}}$$

Lastly we can bound the error terms as,

$$\sum_{i \in A_{s^*}: 0 \leq \Delta_i \leq b} \left(\frac{T^{1-\rho_a} \rho_a^{2\rho_a} 2^{2\rho_a + \frac{3}{2}}}{\psi^{\rho_a} \Delta_i^{4\rho_a - 1}} \right) = \frac{K}{p} \left(\frac{T^{1-\frac{1}{2}} \frac{1}{2}^{2\frac{1}{2}} 2^{2\frac{1}{2} + \frac{3}{2}}}{(\frac{T}{\log K})^{\frac{1}{2}} (\Delta_i)^{4\frac{1}{2} - 1}} \right) < \frac{5.3\sqrt{KT \log K}}{p}$$

Similarly for the term,

$$\sum_{i \in A \setminus A_{s^*}: \Delta_i > b} \left(\frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s + \frac{3}{2}}}{(\psi^{\rho_s}) \Delta_i^{4\rho_s - 1}} \right) < 5.6(K - \frac{K}{p}) \sqrt{\frac{T}{K \log K}}$$

Also, for all $b \geq \sqrt{\frac{K}{7T}}$,

$$\sum_{i \in A \setminus A_{s^*}: 0 < \Delta_i \leq b} \left(\frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s + \frac{3}{2}}}{(\psi^{\rho_s}) b^{4\rho_s - 1}} \right) < 5.6(K - \frac{K}{p}) \sqrt{\frac{T \log K}{K}}$$

After putting the value of $p = \left\lceil \frac{K}{\log K} \right\rceil$, we get,

$$\begin{aligned}
 \mathbb{E}[R_T] \leq & 4\frac{\sqrt{T} \log K}{\sqrt{K}} + 64\frac{\sqrt{T} \log K}{\sqrt{K}} + \frac{32\sqrt{T \log K} \log(\log K)}{\sqrt{K}} + 4\sqrt{KT} + 128\sqrt{KT \log K} \\
 & + \frac{64\sqrt{KT} \log(\log K)}{\sqrt{\log K}} + \frac{5.3\sqrt{T} \log K^{\frac{3}{2}}}{\sqrt{K}} + \frac{5.3\sqrt{T} \log K}{\sqrt{K}} + 10.6\frac{K}{K + \log K} \sqrt{KT \log K} + 10.6\frac{K}{K + \log K} \sqrt{KT}
 \end{aligned}$$

So, the total bound for using both arm and cluster elimination cannot be worse than,

$$\begin{aligned}
 \mathbb{E}[R_T] \leq & 74\frac{\sqrt{T} \log K}{\sqrt{K}} + \frac{32\sqrt{T \log K} \log(\log K)}{\sqrt{K}} + 4\sqrt{KT} + 128\sqrt{KT \log K} \\
 & + \frac{64\sqrt{KT} \log(\log K)}{\sqrt{\log K}} + \frac{5.3\sqrt{T} \log K^{\frac{3}{2}}}{\sqrt{K}} + 10.6\frac{K}{K + \log K} \sqrt{KT \log K} + 10.6\frac{K}{K + \log K} \sqrt{KT}
 \end{aligned}$$

□

E Why Clustering?

In this section we want to specify the apparent use of clustering. The error bounds are shown in Table 2.

Table 2: Error Bound

Elim Type	Error Bound	Remarks
Only Arm Elimination (ClusUCB-AE)	$\underbrace{\sum_{i \in A: \Delta_i > b} \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} \right)}_{\text{Case b2, Proposition 1}} + \underbrace{\sum_{i \in A: 0 < \Delta_i \leq b} \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{b^{4\rho_a-1}} \right)}_{\text{Case b2, Proposition 1}}$	With $\rho_a = \frac{1}{2}$, and $\psi = \frac{T}{\log K}$ this gives $2\sqrt{KT} + 2\sqrt{KT \log K}$. Hence, this has an order of $O(\sqrt{KT \log K})$.
Only Cluster Elimination (ClusUCB-CE)	$\underbrace{\sum_{i \in A \setminus A_{s^*}: \Delta_i > b} \left(\frac{2C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right)}_{\text{Case b2+b3, Proposition 2}} + \underbrace{\sum_{i \in A \setminus A_{s^*}: 0 \leq \Delta_i \leq b} \left(\frac{2C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s-1}} \right)}_{\text{Case b2+b3, Proposition 2}}$	With $\rho_s = \frac{1}{2}$ and $\psi = K^2 T$ this gives $4\sqrt{KT} + 4\sqrt{KT \log K}$. This is same as the bound using only arm elimination and has an order of $O(\sqrt{KT \log K})$.
Arm & Cluster Elimination (ClusUCB)	$\underbrace{\sum_{i \in A_{s^*}: \Delta_i > b} \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} \right) + \sum_{i \in A_{s^*}: 0 \leq \Delta_i \leq b} \left(\frac{C_2(\rho_a) T^{1-\rho_a}}{b^{4\rho_a-1}} \right)}_{\text{Case b2, Arm Elim, Theorem 1}} + \underbrace{\sum_{i \in A \setminus A_{s^*}: \Delta_i > b} \left(\frac{2C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right) + \sum_{i \in A \setminus A_{s^*}: 0 \leq \Delta_i \leq b} \left(\frac{2C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s-1}} \right)}_{\text{Case b3+b4, Clus Elim, Theorem 1}}$	With $\rho_a = \frac{1}{2}$, $\rho_s = \frac{1}{2}$, $p = \lceil \frac{K}{\log K} \rceil$ and $\psi = \frac{T}{\log K}$ this gives $\frac{5.3\sqrt{T \log K}^{\frac{3}{2}}}{\sqrt{K}} + \frac{5.3\sqrt{T \log K}}{\sqrt{K}} + 10.6 \frac{K}{K+\log K} \sqrt{KT \log K} + 10.6 \frac{K}{K+\log K} \sqrt{KT}$. So we can reduce the error bound to $O(\frac{K}{\log K} \sqrt{KT \log K})$.

Algorithm 2 EClusUCB

Input: Number of clusters p , time horizon T , exploration parameters ρ_a, ρ_s and ψ .

Initialization: Set $m := 0, B_0 := A, S_0 = S, \epsilon_0 := 1, M = \lfloor \frac{1}{2} \log_2 \frac{7T}{K} \rfloor, n_0 = \left\lceil \frac{2 \log(\psi T \epsilon_0^2)}{\epsilon_0} \right\rceil$ and $N_0 = K * n_0$.

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

Pull each arm once

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

Pull arm i in B_m such that $\max_{i \in B_m} \left\{ \hat{r}_i + \sqrt{\frac{\rho_s \log(\psi T \epsilon_m^2)}{2n_i}} \right\}$

$t := t + 1$

Arm Elimination

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

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

Cluster Elimination

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

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

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

Reset Parameters

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

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

$$n_m := \left\lceil \frac{2 \log(\psi T \epsilon_m^2)}{\epsilon_m} \right\rceil$$

$$N_m := t + |B_m| * n_m$$

$$m := m + 1$$

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

end if

end for

While looking at the error term for the 3 cases we see that using just Cluster elimination the error term is more than using just arm elimination while we can achieve a balance between the two by using both arm and cluster elimination simultaneously. From Table 2, we can see that the error term for using both arm and cluster elimination can become low depending on how we choose p since $|A_{s^*}| \leq \lceil \frac{A}{p} \rceil$ and in Corollary 2 we proved that taking $p = \lceil \frac{K}{\log K} \rceil$ reduces the elimination error bound.

F Efficient Clustered UCB

In ClusUCB we modified UCB-Improved to reduce early exploration and do faster elimination of sub-optimal arms and clusters with the help of clustering. With careful choosing of exploration regulatory factor ψ and elimination parameters ρ_a and ρ_s we were able to bring down the cumulative regret compared to UCB-Improved. One of the principal dis-advantage that still remain is that ClusUCB is a round based algorithm which samples all remaining arms equally in each round (still less compared to UCB-Improved). In this section, we introduce a further modification, as shown in Algorithm 2 called Efficient Clustered UCB hence referred to as EClusUCB where we introduce the idea of optimistic greedy sampling. A similar idea was also used in [16] which they used to modify the UCB-Improved algorithm. We further modify the idea by introducing clustering and arm elimination parameters. Also we use different exploration regulatory factor and we come

up with a cumulative regret bound whereas [16] only gives simple regret bound. In optimistic greedy sampling we only sample the arm with the highest upper confidence bound in each timestep. Also in EClusUCB we check arm elimination condition and cluster elimination condition in every timestep and update parameters when a round is complete. Drawing from the idea of UCB-Improved, we divide each round into $|B_m|n_m$ timesteps so that each surviving arms can be allocated atmost n_m pulls. In the next section we analyze the regret of EClusUCB.

G Proof of Theorem 2

Theorem 2 (Regret bound). *The regret R_T of EClusUCB satisfies*

$$\begin{aligned} \mathbb{E}[R_T] \leq & \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \left\{ \frac{C_1(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \Delta_i + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \right\} + \sum_{\substack{i \in A, \\ \Delta_i > b}} \left\{ 2\Delta_i + \frac{C_1(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right. \\ & + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} + \left. \frac{32\rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2})}{\Delta_i} \right\} + \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \sum_{\substack{i \in A_{s^*}, \\ 0 < \Delta_i \leq b}} \frac{C_2(\rho_a)T^{1-\rho_a}}{b^{4\rho_a-1}} \\ & + \sum_{\substack{i \in A \setminus A_{s^*}, \\ \Delta_i > b}} \frac{2C_2(\rho_s)T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{\substack{i \in A \setminus A_{s^*}, \\ 0 < \Delta_i \leq b}} \frac{2C_2(\rho_s)T^{1-\rho_s}}{b^{4\rho_s-1}} \\ & + \max_{i: \Delta_i \leq b} \Delta_i T, \end{aligned}$$

where $b \geq \sqrt{\frac{e}{T}}$, $C_1(x) = \frac{2^{1+4x}x^{2x}}{\psi^x}$, $C_2(x) = \frac{2^{2x+\frac{3}{2}}x^{2x}}{\psi^x}$, $\rho_a \leq \rho_s$ and A_{s^*} is the subset of arms in cluster s^* containing optimal arm a^* .

We can see from the result of this theorem that the regret upper bound for EClusUCB is same as ClusUCB. Also we can specialize the result of this Theorem 2 using Corollary 1 and 2 to achieve a similar bound as for Theorem 1.

Proof. Let $A' = \{i \in A, \Delta_i > b\}$, $A'' = \{i \in A, \Delta_i > 0\}$, $A'_{s_k} = \{i \in A_{s_k}, \Delta_i > b\}$ and $A''_{s_k} = \{i \in A_{s_k}, \Delta_i > 0\}$. C_g is the cluster set containing max payoff arm from each cluster in g -th round. The arm having the highest payoff in a cluster s_k is denote by $a_{\max_{s_k}}$. Let for each sub-optimal arm $i \in A$, $m_i = \min \{m | \sqrt{\rho_a \epsilon_m} < \frac{\Delta_i}{2}\}$ and let for each cluster $s_k \in S$, $g_{s_k} = \min \{g | \sqrt{\rho_s \epsilon_g} < \frac{\Delta_{a_{\max_{s_k}}}}{2}\}$. Let $\check{A} = \{i \in A' | i \in s_k, \forall s_k \in S\}$. Also n_i denotes total number of times an arm i has been pulled from time $t = 1$ to T . In the m -th round, n_m denotes the number of pulls allocated to the surviving arms in B_m .

The analysis proceeds by considering the contribution to the regret in each of the following cases:

Case a: *Some sub-optimal arm i is not eliminated in round $\max(m_i, g_{s_k})$ or before, with the optimal arm $*$ $\in C_{\max(m_i, g_{s_k})}$.*

We consider an arbitrary sub-optimal arm i and analyze the contribution to the regret when i is not eliminated in the following exhaustive sub-cases:

Case a1: *In round $\max(m_i, g_{s_k})$, $i \in s^*$.*

This case is similar to Case a1, Theorem 1. Let $c_i = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i}^2)}{2n_i}}$. For an arm to get eliminated it must satisfy the following condition,

$$\hat{r}_i \leq r_i + c_i \text{ and } \hat{r}^* \geq r^* - c^*, \quad (5)$$

Since each round now consist of $|B_m|n_m$ timesteps and since at every timestep the arm elimination condition is being checked, following a parallel argument as in Case a1, Theorem 1 we can show that when $n_i = n_{m_i}$ then

$$c_i = c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i}^2)}{2n_{m_i}}} = \sqrt{\rho_a \epsilon_{m_i+1}} < \frac{\Delta_i}{4}, \text{ since } n_i = n_{m_i} = \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}} \text{ and } \rho_a \in (0, 1].$$

Subsequently following the steps of Case a1, Theorem 1 we can upper bound the probability of the complementary of the events in 5 and show that the probability that a sub-optimal arm i is not eliminated in any round on or before m_i is bounded above by $\left(\frac{2}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}\right)$.

Case a2: In round $\max(m_i, g_{s_k})$, $i \in s_k$ for some $s_k \neq s^*$.

Let $c_{a_{\max s_k}} = \sqrt{\frac{\rho_s \log(\psi T \epsilon_m^2)}{2n_{a_{\max s_k}}}}$ and $n_{a_{\max s_k}}$ is the number of times $a_{\max s_k}$ is pulled from time $t = 1$ to T . Following a parallel argument like in Case a1, we have to bound the following two events of arm $a_{\max s_k}$ not getting eliminated on or before g_{s_k} -th round,

$$\hat{r}_{a_{\max s_k}} \geq r_{a_{\max s_k}} + c_{a_{\max s_k}} \text{ and } \hat{r}^* \leq r^* - c^*$$

As we argued in previous case, since cluster elimination condition being verified in each timestep, the above two conditions are possible when $n_{a_{\max s_k}} = n_{g_{s_k}}$. We can prove using Chernoff-Hoeffding bounds and considering independence of

events mentioned above, that for $c_{a_{\max s_k}} = c_{g_{s_k}} = \sqrt{\frac{\rho_s \log(\psi T \epsilon_{g_{s_k}}^2)}{2n_{g_{s_k}}}}$ and $n_{g_{s_k}} = \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}}$ the probability of the above two events is bounded by $\left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}\right)$. Now, for any round g_{s_k} , all the elements of $C_{\max(m_i, g_{s_k})}$ are the respective maximum payoff arms of their cluster s_k , $\forall s_k \in S$, and since clusters are fixed so we can bound the maximum probability that a sub-optimal arm $i \in A'$ and $i \in s_k$ such that $a_{\max s_k} \in C_{g_{s_k}}$ is not eliminated on or before the g_{s_k} -th round by the same probability as above.

Summing up over all p clusters and bounding the regret for each arm $i \in A'_{s_k}$ trivially by $T\Delta_i$ and following the steps of Case a2, Theorem 1, we can show that the regret for not eliminating clusters on or before round g_{s_k} is upper bounded by,

$$\sum_{i \in A'} \frac{C_1(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}}$$

Summing the bounds in Cases a1 – a2 and observing that the bounds in the aforementioned cases hold for any round $C_{\max\{m_i, g_{s_k}\}}$, we obtain the following contribution to the expected regret from Case a:

$$\sum_{i \in A_{s^*}} \frac{C_1(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \sum_{i \in A'} \left(\frac{C_1(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} \right)$$

Case b: For each arm i , either i is eliminated in round $\max(m_i, g_{s_k})$ or before or there is no optimal arm $*$ in $C_{\max(m_i, g_{s_k})}$.

Case b1: $*$ $\in C_{\max(m_i, g_{s_k})}$ for each arm $i \in A'$ and cluster $s_k \in \check{A}$.

The condition in the case description above implies the following:

- (i) each sub-optimal arm $i \in A'$ is eliminated on or before $\max(m_i, g_{s_k})$ and hence pulled not more than n_i number of times. But according to the condition of Case a1, $n_i \leq n_{m_i}$ number of times.
- (ii) each sub-optimal cluster $s_k \in \check{A}$ is eliminated on or before $\max(m_i, g_{s_k})$ and hence pulled not more than $n_{a_{\max s_k}}$ number of times. But again according to condition Case a2, $n_{a_{\max s_k}} \leq n_{g_{s_k}}$ number of times.

Hence, following the same steps as in Case b1, Theorem 1, the maximum regret suffered due to pulling of a sub-optimal arm or a sub-optimal cluster is no more than the following:

$$\sum_{i \in A'} \left[2\Delta_i + \frac{32(\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2}) + \rho_s \log(\psi T \frac{\Delta_i^4}{16\rho_s^2}))}{\Delta_i} \right]$$

Case b2: $*$ is eliminated by some sub-optimal arm in s^*

Optimal arm a^* can get eliminated by some sub-optimal arm i only if arm elimination condition holds, i.e.,

$$\hat{r}_i - c_i > \hat{r}^* + c^*,$$

where, as mentioned before, $c_i = c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i}^2)}{2n_{m_i}}}$. From analysis in Case a1, notice that, if (5) holds in conjunction with the above, arm i gets eliminated. Also, recall from Case a1 that the events complementary to (5) have low-probability and can be upper bounded by $\frac{2}{(\psi T \epsilon_{m_*}^2)^{\rho_a}}$. Moreover, a sub-optimal arm that eliminates $*$ has to survive until round m_* . In other words, all arms $j \in s^*$ such that $m_j < m_*$ are eliminated on or before m_* (this corresponds to case b1). Let, the arms surviving till m_* round be denoted by A'_{s^*} . 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 A''_{s^*} . Also maximal regret per step after eliminating $*$ is the maximal Δ_j among the remaining arms in A''_{s^*} with $m_j \geq m_*$. Let $m_b = \min\{m | \sqrt{\rho_a \epsilon_m} < \frac{\Delta_b}{2}\}$. Let $C_2(x) = \frac{2^{2x+\frac{3}{2}} x^{2x}}{\psi^x}$. Hence, the maximal regret after eliminating the arm $*$ is upper bounded by,

$$\sum_{m_*=0}^{\max_{j \in A'_{s^*}} m_j} \sum_{\substack{i \in A''_{s^*}: \\ m_i \geq m_*}} \left(\frac{2}{(\psi T \epsilon_{m_*}^2)^{\rho_a}} \right) \cdot T \max_{\substack{j \in A''_{s^*}: \\ m_j \geq m_*}} \Delta_j$$

Therefore, following the similar steps as in Case b2, Theorem 1, we can show that the above regret is upper bounded by,

$$\sum_{i \in A'_{s^*}} \frac{C_2(\rho_a) T^{1-\rho_a}}{\Delta_i^{4\rho_a-1}} + \sum_{i \in A''_{s^*} \setminus A'_{s^*}} \frac{C_2(\rho_a) T^{1-\rho_a}}{b^{4\rho_a-1}}$$

Case b3: s^* is eliminated by some sub-optimal cluster.

Let $C'_g = \{a_{\max_{s_k}} \in A' | \forall s_k \in S\}$ and $C''_g = \{a_{\max_{s_k}} \in A'' | \forall s_k \in S\}$. A sub-optimal cluster s_k will eliminate s^* in round g_* only if the cluster elimination condition of Algorithm 2 holds, which is the following when $*$ $\in C_{g_*}$:

$$\hat{r}_{a_{\max_{s_k}}} - c_{a_{\max_{s_k}}} > \hat{r}^* + c^*. \quad (6)$$

Notice that when $*$ $\notin C_{g_*}$, since $r_{a_{\max_{s_k}}} > r^*$, the inequality in (6) has to hold for cluster s_k to eliminate s^* . As in case b2, the probability that a given sub-optimal cluster s_k eliminates s^* is upper bounded by $\frac{2}{(\psi T \epsilon_{g_*}^2)^{\rho_s}}$ and all sub-optimal clusters with $g_{s_j} < g_*$ are eliminated before round g_* .

This leaves any arm $a_{\max_{s_b}}$ such that $g_{s_b} \geq g_*$ to still survive and eliminate arm $*$ in round g_* . Let, such arms that survive $*$ belong to C''_g . Hence, following the same way as case b2, the maximal regret after eliminating $*$ is,

$$\sum_{g_*=0}^{\max_{a_{\max_{s_j}} \in C'_g} g_{s_j}} \sum_{\substack{a_{\max_{s_k}} \in C''_g: \\ g_{s_k} \geq g_*}} \left(\frac{2}{(\psi T \epsilon_{g_*}^2)^{\rho_s}} \right) T \max_{\substack{a_{\max_{s_j}} \in C''_g: \\ g_{s_j} \geq g_*}} \Delta_{a_{\max_{s_j}}}$$

Using $A' \supset C'_g$ and $A'' \supset C''_g$, we can bound the regret contribution from this case in a similar manner as Case b2 as follows:

$$\begin{aligned} & \sum_{i \in A' \setminus A'_{s^*}} \frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s+\frac{3}{2}}}{\psi^{\rho_s} \Delta_i^{4\rho_s-1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{T^{1-\rho_s} \rho_s^{2\rho_s} 2^{2\rho_s+\frac{3}{2}}}{\psi^{\rho_s} b^{4\rho_s-1}} \\ &= \sum_{i \in A' \setminus A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s-1}} \end{aligned}$$

Case b4: $*$ is not in $C_{\max(m_i, g_{s_k})}$, but belongs to $B_{\max(m_i, g_{s_k})}$.

In this case the optimal arm $* \in s^*$ is not eliminated, also s^* is not eliminated. So, for all sub-optimal arms i in A'_{s^*} which gets eliminated on or before $\max\{m_i, g_{s_k}\}$ will get pulled no less than $n_i \leq \left\lceil \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}} \right\rceil$ number of times, which leads to the following bound the contribution to the expected regret, as in Case b1:

$$\sum_{i \in A'_{s^*}} \left\{ \Delta_i + \frac{32\rho_a \log(\psi T \frac{\Delta_i^4}{16\rho_a^2})}{\Delta_i} \right\}$$

For arms $a_i \notin s^*$, the contribution to the regret cannot be greater than that in Case b3. So the regret is bounded by,

$$\sum_{i \in A' \setminus A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{\Delta_i^{4\rho_s-1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{C_2(\rho_s) T^{1-\rho_s}}{b^{4\rho_s-1}}$$

The main claim follows by summing the contributions to the expected regret from each of the cases above. \square

H Proof of Theorem 3

Theorem 3. For every $0 < \eta < 1$ and $\gamma > 1$, there exists τ such that for all $T > \tau$ the simple regret of EClusUCB is upper bounded by,

$$SR_{EClusUCB} \leq 4K\gamma \sum_{i=1}^K \Delta_i \left\{ \exp\left(\frac{1}{2} - \rho_a - \frac{c_0\sqrt{e}}{4}\right) \left(\frac{\log(\psi T)}{T^{\frac{3}{2}}(\psi T^2)^{\rho_a}}\right) + \exp\left(\frac{1}{2} - \rho_s - \frac{c_0\sqrt{e}}{4}\right) \left(\frac{\log(\psi T)}{T^{\frac{3}{2}}(\psi T^2)^{\rho_s}}\right) \right\}$$

with probability at least $1 - \eta$, where $c_0 > 0$ is a constant.

Proof. We follow the same steps as in Theorem 2, [16]. First we will state the two facts used by this proof.

1. *Fact 1:* From Theorem 2 we know that the probability of elimination of a sub-optimal arm in the $\max(m_i, g_{s_k})$ round is $\left(\frac{2}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}\right)$ and of a sub-optimal cluster is $\left(\frac{2}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}\right)$.
2. *Fact 2:* From [20] we know that, for every $0 < \eta < 1$ and $\gamma > 1$, there exists τ such that for all $T > \tau$ the probability of a sub-optimal arm i being sampled in the m -th round is bounded by $P_m \leq 2\gamma \exp(-c_m \frac{\sqrt{T}}{2})$, where $c_m = \frac{c_0}{2^m}$.

We start with an upper bound on the number of plays $\delta_{\max(m_i, g_{s_k})}$ in the $\max(m_i, g_{s_k})$ -th round divided by the total number of plays T . We know from Fact 1 that the total number of arms surviving in the $\max(m_i, g_{s_k})$ -th arm is

$$|B_{\max(m_i, g_{s_k})}| = \left(\frac{2K}{(\psi T \epsilon_{m_i}^2)^{\rho_a}}\right) + \left(\frac{2K}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}}\right)$$

Again in EClusUCB, we know that the number of pulls allocated for each surviving arm i in the $\max(m_i, g_{s_k})$ -th round is $n_{\max(m_i, g_{s_k})} = \frac{2 \log(\psi T \epsilon_{\max(m_i, g_{s_k})}^2)}{\epsilon_{\max(m_i, g_{s_k})}}$. Therefore, the proportion of plays $\delta_{\max(m_i, g_{s_k})}$ in the $\max(m_i, g_{s_k})$ -th round can be written as,

$$\delta_{\max(m_i, g_{s_k})} = \frac{(|B_{\max(m_i, g_{s_k})}| \cdot n_{\max(m_i, g_{s_k})})}{T} \leq \left(\frac{1}{T} \cdot \frac{2K}{(\psi T \epsilon_{m_i}^2)^{\rho_a}} \cdot \frac{2 \log(\psi T \epsilon_{m_i}^2)}{\epsilon_{m_i}}\right) + \left(\frac{1}{T} \cdot \frac{2K}{(\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}} \cdot \frac{2 \log(\psi T \epsilon_{g_{s_k}}^2)}{\epsilon_{g_{s_k}}}\right)$$

$$\leq \left(\frac{4K \log(\psi T \epsilon_{m_i}^2)}{T \epsilon_{m_i} (\psi T \epsilon_{m_i}^2)^{\rho_a}} \right) + \left(\frac{4K \log(\psi T \epsilon_{g_{s_k}}^2)}{T \epsilon_{g_{s_k}} (\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}} \right)$$

Now, $\epsilon_{m_i} \geq \sqrt{\frac{e}{T}}$ and $\epsilon_{g_{s_k}} \geq \sqrt{\frac{e}{T}}$ for all rounds $m = 0, 1, 2, \dots, \lfloor \frac{1}{2} \log_2 \frac{T}{e} \rfloor$.

$$\begin{aligned} \delta_{\max(m_i, g_{s_k})} &\leq \left(\frac{4K \log(\psi T \epsilon_{m_i}^2)}{T \epsilon_{m_i} (\psi T \epsilon_{m_i}^2)^{\rho_a}} \right) + \left(\frac{4K \log(\psi T \epsilon_{g_{s_k}}^2)}{T \epsilon_{g_{s_k}} (\psi T \epsilon_{g_{s_k}}^2)^{\rho_s}} \right) \\ &\leq \left(\frac{4K \log(\psi T)}{T \epsilon_M (\psi T \epsilon_M^2)^{\rho_a}} \right) + \left(\frac{4K \log(\psi T)}{T \epsilon_M (\psi T \epsilon_M^2)^{\rho_s}} \right) \\ &\leq \left(\frac{4K e^{\frac{1}{2} - \rho_a} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_a}} \right) + \left(\frac{4K e^{\frac{1}{2} - \rho_s} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_s}} \right) \end{aligned}$$

Now, applying the bound from Fact 2, we can show that the probability of the sub-optimal arm i being pulled is bounded above by,

$$\begin{aligned} P_i &= \sum_{m=0}^M \delta_m \cdot P_m \leq \sum_{m=0}^M \left\{ \left(\frac{4K e^{\frac{1}{2} - \rho_a} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_a}} \right) + \left(\frac{4K e^{\frac{1}{2} - \rho_s} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_s}} \right) \right\} 2\gamma \exp\left(-\frac{c_m \sqrt{T}}{4}\right) \\ &\leq M \cdot \left\{ \left(\frac{4K e^{\frac{1}{2} - \rho_a} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_a}} \right) + \left(\frac{4K e^{\frac{1}{2} - \rho_s} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_s}} \right) \right\} 2\gamma \exp\left(-\frac{c_0 \sqrt{T}}{2M \cdot 4}\right) \\ &\leq \log_2 \frac{T}{e} \gamma \exp\left(\frac{c_0 \sqrt{e}}{4}\right) \left\{ \left(\frac{4K e^{\frac{1}{2} - \rho_a} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_a}} \right) + \left(\frac{4K e^{\frac{1}{2} - \rho_s} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_s}} \right) \right\}, \text{ for } M = \lfloor \frac{1}{2} \log_2 \frac{T}{e} \rfloor \end{aligned}$$

Hence, the simple regret of EClusUCB is upper bounded by,

$$\begin{aligned} SR_{EClusUCB} &= \sum_{i=1}^K \Delta_i \cdot P_i \leq \sum_{i=1}^K \Delta_i \cdot \log_2 \frac{T}{e} \gamma \exp\left(\frac{c_0 \sqrt{e}}{4}\right) \left\{ \left(\frac{4K e^{\frac{1}{2} - \rho_a} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_a}} \right) + \left(\frac{4K e^{\frac{1}{2} - \rho_s} \log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_s}} \right) \right\} \\ &\leq 4K \gamma \exp\left(\frac{1}{2} - \rho_a - \frac{c_0 \sqrt{e}}{4}\right) \sum_{i=1}^K \Delta_i \left(\frac{\log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_a}} \right) + 4K \gamma \exp\left(\frac{1}{2} - \rho_s - \frac{c_0 \sqrt{e}}{4}\right) \sum_{i=1}^K \Delta_i \left(\frac{\log(\psi T)}{T^{\frac{3}{2}} (\psi T^2)^{\rho_s}} \right) \end{aligned}$$

□

I Proof of Corollary 4

Corollary 4. For $\psi = \frac{T}{\log(K)}$, $\rho_a = \frac{1}{2}$ and $\rho_s = \frac{1}{2}$, the simple regret of EClusUCB is given by,

$$SR_{EClusUCB} \leq 8K \gamma \exp\left(-\frac{c_0 \sqrt{e}}{4}\right) \sum_{i=1}^K \Delta_i \left(\frac{2\sqrt{\log(K)} \log\left(\frac{T}{\sqrt{\log(K)}}\right)}{T^3} \right)$$

Proof. Putting $\psi = \frac{T}{\log(K)}$, $\rho_a = \frac{1}{2}$ and $\rho_s = \frac{1}{2}$ in the simple regret obtained in Theorem 3, we get

$$SR_{EClusUCB} \leq 8K \gamma \exp\left(-\frac{c_0 \sqrt{e}}{4}\right) \sum_{i=1}^K \Delta_i \left(\frac{\log\left(\frac{T^2}{\log(KT)}\right)}{T^{\frac{3}{2}} \left(\frac{T^3}{\log(KT)}\right)^{\frac{1}{2}}} \right)$$

Table 3: Simple regret upper bounds for different bandit algorithms

Algorithm	Upper bound
CCB	$O\left(K\gamma \sum_{i=1}^K \Delta_i \exp\left(2 - \frac{c_0\sqrt{e}}{4}\right) \log_2\left(\frac{T}{e}\right) \frac{\log T}{T^4}\right)$
EClusUCB	$O\left(K\gamma \exp\left(-\frac{c_0\sqrt{e}}{4}\right) \sum_{i=1}^K \Delta_i \left(\frac{\sqrt{\log(K)} \log\left(\frac{T}{\sqrt{\log(K)}}\right)}{T^3}\right)\right)$

$$\leq 8K\gamma \exp\left(-\frac{c_0\sqrt{e}}{4}\right) \sum_{i=1}^K \Delta_i \left(\frac{2\sqrt{\log(K)} \log\left(\frac{T}{\sqrt{\log(K)}}\right)}{T^3}\right)$$

□

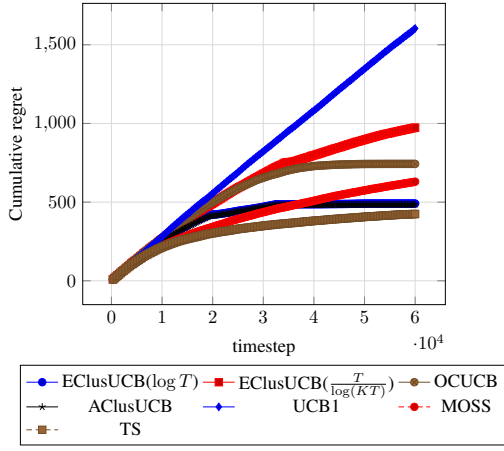
J Adaptive Clustered UCB

In the previous section we saw that EClusUCB deals with too much early exploration through optimistic greedy sampling. This reduces the cumulative regret, but still one of the principal disadvantage that EClusUCB suffers from is the lack of knowledge of the number of clusters p . One way to handle this is to estimate the number of clusters on the fly. In Algorithm 3, named Adaptive Clustered UCB, hence referred to as AClusUCB, we explore this idea. AClusUCB uses *hierarchical clustering* (see [11]) to find the number of clusters present. AClusUCB is similar to EClusUCB with 2 major differences. The first difference is the call to procedure CreateClusters at every timestep. CreateClusters sub-routine first creates a singleton cluster for each of the surviving arms in B_m and then clusters those singleton clusters $s_k, s_d \in S_m$ (say) if any arm $i \in s_k$ and $j \in s_d$ is such that $|\hat{r}_i - \hat{r}_j| \leq \epsilon_m$. We cluster based on ϵ_m because we have no prior knowledge of the gaps and we are estimating the gap by ϵ_m . Also, we are destroying the clusters after every timestep and re-constructing the clusters based on the condition specified. Since, the environment is stochastic, the initial clusters will have very poor purity (arms with ϵ_m -close expected means lying in a single cluster) whereas in the later rounds the purity becomes better which leads to the optimal arm * lying in a single cluster of its own which will eliminate all the other clusters based on the cluster elimination condition. The second difference is that, we limit the cluster size from start by $\ell_m = 2$ and then double it after every round. Since the environment is stochastic, if we do not limit the cluster size, then in the initial rounds it will result in huge chains of clusters because the initial estimates of $\hat{r}_i \forall i \in A$ will be poor. This condition just helps in stopping such large chains of clusters.

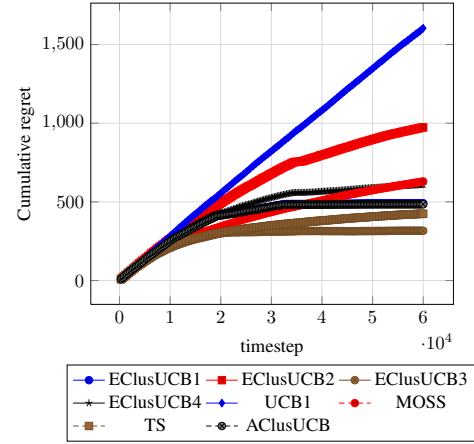
K More Experiments

In the first experiment, we check the performance of EClusUCB against 4 algorithms, UCB1, OCUCB[15], MOSS and TS. The testbed is similar to Experiment 1, with 20 arms involving Bernoulli reward distribution with expected rewards of the arms $r_{i \neq *}=0.07$ and $r^*=0.1$ and horizon T is set to 60000. The number of clusters p is set to 5 and the result is shown in Figure 3a. Here, we choose the $\psi = \{\frac{T}{\log(KT)}, \log T\}$. From the result we see that EClusUCB is sensitive to the exploratory regulatory factor ψ . The first choice of ψ is conservative and does too much exploration while the second choice is optimistic and explores less. Both ρ_a and ρ_s is set to $\frac{1}{2}$ in all the cases. We see that OCUCB will eventually perform better than MOSS which does not stabilize in 60000 timesteps but both of them is beaten by EClusUCB. The careful choosing of ρ_a, ρ_s and ψ enables this.

In the second Experiment the testbed is same as above and we check the performance of EClusUCB with various parameters against MOSS and TS and the result is shown in Figure 3b. In all the cases the number of clusters is set at 5. Here we see that a lesser exploration regulatory factor $\psi = \log T$ and $\rho_a = \rho_s = 0.25$ outperforms most of the other variants. Infact EClusUCB1, EClusUCB3 and EClusUCB4 are better than MOSS. In both the figures we see that EClusUCB(log T) performance is almost similar to AClusUCB.



(a) Experiment 5: Cumulative regret for EClusUCB and AClusUCB against various other UCB variants and Thompson Sampling



(b) Experiment 6: Cumulative regret for EClusUCB with different parameters and AClusUCB. EClusUCB1: $\psi = \log T, \rho_a = \rho_s = 0.5$. EClusUCB2: $\psi = \frac{T}{\log(KT)}, \rho_a = \rho_s = 0.5$. EClusUCB3: $\psi = \log T, \rho_a = \rho_s = 0.25$. EClusUCB4: $\psi = \frac{T}{\log(KT)}, \rho_a = \rho_s = 0.25$

Figure 3: Cumulative regret for various bandit algorithms on two stochastic K-armed bandit environments.

Algorithm 3 AClusUCB

Input: Time horizon T , exploration parameters ρ_a, ρ_s and ψ .

Initialization: Set $m := 0$, $B_0 := A$, $S_0 = S$, $\epsilon_0 := 1$, $M = \lfloor \frac{1}{2} \log_2 \frac{7T}{K} \rfloor$, $n_0 = \left\lceil \frac{2 \log(\psi T \epsilon_0^2)}{\epsilon_0} \right\rceil$, $\ell_0 := 2$ and

$N_0 = K * n_0$.

Pull each arm once

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

Pull arm i in B_m such that $\max_{i \in B_m} \left\{ \hat{r}_i + \sqrt{\frac{\rho_s \log(\psi T \epsilon_m^2)}{2n_i}} \right\}$

$t := t + 1$

Call CreateClusters()

Arm Elimination

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

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

Cluster Elimination

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

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

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

Reset Parameters

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

$$\ell_{m+1} = 2 * \ell_m$$

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

$$n_m := \left\lceil \frac{2 \log(\psi T \epsilon_m^2)}{\epsilon_m} \right\rceil$$

$$N_m := t + |B_m| * n_m$$

$$m := m + 1$$

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

end if

end for

procedure CREATECLUSTERS

Create singleton cluster $\{i\}$ for each arm $i \in B_m$ and call this partition as S_m .

For two cluster $s_k, s_d \in S_m$, join the clusters if any $|\hat{r}_i - \hat{r}_j| \leq \epsilon_m$ and $|s_k| + |s_d| \leq \ell_m$, where $i \in s_k$ and $j \in s_d$

end procedure
