
MaxGap Bandit: Adaptive Algorithms for Approximate Ranking

Sumeet Katariya *
UW-Madison and Amazon
sumeetsk@gmail.com

Ardhendu Tripathy *
UW-Madison
astripathy@wisc.edu

Robert Nowak
UW-Madison
rdnowak@wisc.edu

Abstract

This paper studies the problem of adaptively sampling from K distributions (arms) in order to identify the largest gap between any two adjacent means. We call this the MaxGap-bandit problem. This problem arises naturally in approximate ranking, noisy sorting, outlier detection, and top-arm identification in bandits. The key novelty of the MaxGap bandit problem is that it aims to adaptively determine the natural partitioning of the distributions into a subset with larger means and a subset with smaller means, where the split is determined by the largest gap rather than a pre-specified rank or threshold. Estimating an arm’s gap requires sampling its neighboring arms in addition to itself, and this dependence results in a novel hardness parameter that characterizes the sample complexity of the problem. We propose elimination and UCB-style algorithms and show that they are minimax optimal. Our experiments show that the UCB-style algorithms require 6-8x fewer samples than non-adaptive sampling to achieve the same error.

1 Introduction

Consider an algorithm that can draw i.i.d. samples from K unknown distributions. The goal is to partially rank the distributions according to their (unknown) means. This model encompasses many problems including best-arms identification (BAI) in multi-armed bandits, noisy sorting and ranking, and outlier detection. Partial ranking is often preferred to complete ranking because correctly ordering distributions with nearly equal means is an expensive task (in terms of number of required samples). Moreover, in many applications it is arguably unnecessary to resolve the order of such close distributions. This observation motivates algorithms that aim to recover a partial ordering into groups/clusters of distributions with similar means. This entails identifying large “gaps” in the ordered sequence of means. The focus of this paper is the fundamental problem of finding the *largest gap* by sampling adaptively. Identification of the largest gap separates the distributions into two groups, and thus recursive application would allow one to identify any number of groupings in a partial order.

As illustration, consider a subset of images from the Chicago streetview dataset [17] shown in Fig. 1. In this study, people were asked to judge how safe each scene looks [18], and a larger mean indicates a safer looking scene. While each person has a different sense of how safe an image looks, when aggregated there are clear trends in the safety scores (denoted by $\mu_{(i)}$) of the images. Fig. 1 schematically shows the distribution of scores given by people as a bell curve below each image. Assuming the sample means are close to their true means, one can nominally classify them as ‘safe’, ‘maybe unsafe’ and ‘unsafe’ as indicated in Fig. 1. Here we have implicitly used the large gaps $\mu_{(2)} - \mu_{(3)}$ and $\mu_{(4)} - \mu_{(5)}$ to mark the boundaries. Note that finding the safest image (BAI) is hard as we need a lot of human responses to decide the larger mean between the two rightmost distributions; it is also arguably unnecessary. A common way to address this problem is to specify a tolerance ϵ [7],

* Authors contributed equally and are listed alphabetically.

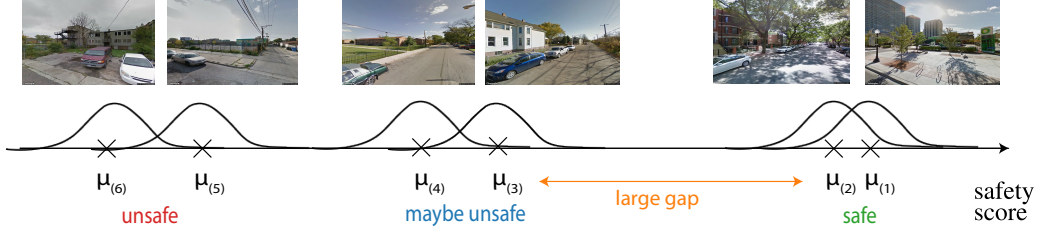


Figure 1: Six representative images from Chicago streetview dataset and their safety (Borda) scores.

and stop sampling if the means are less than ϵ apart; however determining this can require $\Omega(1/\epsilon^2)$ samples. Distinguishing the top 2 distributions from the rest is easy and can be efficiently done using top- m arm identification [15], however this requires the experimenter to prescribe the location $m = 2$ where a large gap exists which is unknown. *Automatically identifying natural splits in the set of distributions is the aim of the new theory and algorithms we propose. We call this problem of adaptive sampling to find the largest gap the MaxGap-bandit problem.*

1.1 Notation and Problem Statement

We will use multi-armed bandit terminology and notation throughout the paper. The K distributions will be called *arms* and drawing a sample from a distribution will be referred to as *sampling the arm*. Let $\mu_i \in \mathbb{R}$ denote the mean of the i -th arm, $i \in \{1, 2, \dots, K\} =: [K]$. We add a parenthesis around the subscript j to indicate the j -th largest mean, i.e., $\mu_{(K)} \leq \mu_{(K-1)} \leq \dots \leq \mu_{(1)}$. For the i -th arm, we define its gap Δ_i to be the maximum of its left and right gaps, i.e.,

$$\Delta_i = \max\{\mu_{(\ell)} - \mu_{(\ell+1)}, \mu_{(\ell-1)} - \mu_{(\ell)}\} \quad \text{where } \mu_i = \mu_{(\ell)}. \quad (1)$$

We define $\mu_{(0)} = -\infty$ and $\mu_{(K+1)} = \infty$ to account for the fact that extreme arms have only one gap. The goal of the MaxGap-bandit problem is to (adaptively) sample the arms and return two clusters

$$C_1 = \{(1), (2), \dots, (m)\} \quad \text{and} \quad C_2 = \{(m+1), \dots, (K)\},$$

where m is the rank of the arm with the largest gap between *adjacent* means, i.e.,

$$m = \arg \max_{j \in [K-1]} \mu_{(j)} - \mu_{(j+1)}. \quad (2)$$

The mean values are unknown as is the ordering of the arms according to their means. A solution to the MaxGap-bandit problem is an algorithm which given a probability of error $\delta > 0$, samples the arms and upon stopping partitions $[K]$ into two clusters \hat{C}_1 and \hat{C}_2 such that

$$\mathbb{P}(\hat{C}_1 \neq C_1) \leq \delta. \quad (3)$$

This setting is known as the fixed-confidence setting [10], and the goal is to achieve the probably correct clustering using as few samples as possible. In the sequel, we assume that m is uniquely defined and let $\Delta_{\max} = \Delta_{i^*}$ where $\mu_{i^*} = \mu_{(m)}$.

1.2 Comparison to a Naive Algorithm: Sort then search for MaxGap

The MaxGap-bandit problem is not equivalent to BAI on $\binom{K}{2}$ gaps since the MaxGap-bandit problem requires identifying the largest gap between *adjacent* arm means (BAI on $\binom{K}{2}$ gaps would always identify $\mu_{(1)} - \mu_{(K)}$ as the largest gap). This suggests a naive two-step algorithm: we first sample the arms enough number of times so as to identify all pairs of adjacent arms (i.e., we sort the arms according to their means), and then run a BAI bandit algorithm [13] on the $(K-1)$ gaps between adjacent arms to identify the largest gap (an unbiased sample of the gap can be obtained by taking the difference of the samples of the two arms forming the gap).

We analyze the sample complexity of this naive algorithm in Appendix A, and discuss the results here for an example configuration. Consider the arrangement of means shown in Fig. 2 where there is one



Figure 2: Configuration with one large gap

large gap Δ_{\max} and all the other gaps are equal to $\Delta_{\min} \ll \Delta_{\max}$. The naive algorithm has a sample complexity $\Omega(K/\Delta_{\min}^2)$ (the first sorting step requires these many samples) which can be very large. Is this sorting of the arm means necessary? For instance, we do not need to sort K real numbers in order to cluster them according to the largest gap¹. The algorithms we propose in this paper solve the MaxGap-bandit problem without necessarily sorting the arm means. For the configuration in Fig. 2 they require $\tilde{O}(K/\Delta_{\max}^2)$ samples, giving a saving of approximately $(\Delta_{\max}/\Delta_{\min})^2$ samples.

The analysis of our algorithms suggests a novel hardness parameter for the MaxGap-bandit problem that we discuss next. We let $\Delta_{i,j} := \mu_j - \mu_i$ for all $i, j \in [K]$. We show in Section 5 that the number of samples taken from distribution i due to its right gap is inversely proportional to the square of

$$\gamma_i^r := \max_{j: \Delta_{i,j} > 0} \min \{ \Delta_{i,j}, \Delta_{\max} - \Delta_{i,j} \}. \quad (4)$$

For the left gap of i we define γ_i^l analogously. The total number of samples drawn from distribution i is inversely proportional to the square of $\gamma_i := \min\{\gamma_i^r, \gamma_i^l\}$. The intuition for Eq. (4) is that distribution i can be eliminated quickly if there is another distribution j that has a moderately large gap from i (so that this gap can be quickly detected), but not too large (so that the gap is easy to distinguish from Δ_{\max}), and (4) chooses the best j . We discuss (4) in detail in Section 5, where we show that our algorithms use $\tilde{O}(\sum_{i \in [K] \setminus \{(m), (m+1)\}} \gamma_i^{-2} \log(K/\delta\gamma_i))$ samples to find the largest gap with probability at least $1 - \delta$. This sample complexity is minimax optimal.

1.3 Summary of Main Results and Paper Organization

In addition to motivating and formulating the MaxGap-bandit problem, we make the following contributions. First, we design elimination and UCB-style algorithms as solutions to the MaxGap-bandit problem that do not require sorting the arm means (Section 3). These algorithms require computing upper bounds on the gaps Δ_i , which can be formulated as a mixed integer optimization problem. We design a computationally efficient dynamic programming subroutine to solve this optimization problem and this is our second contribution (Section 4). Third, we analyze the sample complexity of our proposed algorithms, and discover a novel problem-hardness parameter (Section 5). This parameter arises because of the arm interactions in the MaxGap-bandit problem where, in order to reduce uncertainty in the value of an arm's gap, we not only need to sample the said arm but also its neighboring arms. Fourth, we show that this sample complexity is minimax optimal (Section 6). Finally, we evaluate the empirical performance of our algorithms on simulated and real datasets and observe that they require 6-8x fewer samples than non-adaptive sampling to achieve the same error (Section 7).

2 Related Work

One line of related research is best-arm identification (BAI) in multi-armed bandits. A typical goal in this setting is to identify the top- m arms with largest means, where m is a *prespecified* number [15, 16, 1, 3, 9, 4, 14, 7, 20]. As explained in Section 1, our motivation behind formulating the MaxGap-bandit problem is to have an adaptive algorithm which finds the “natural” set of top arms as delineated by the largest gap in consecutive mean values. Our work can also be used to automatically detect “outlier” arms [23].

The MaxGap-bandit problem is different from the standard multi-armed bandit because of the local dependence of an arm's gap on other arms. Other best-arm settings where an arm's reward can inform the quality of other arms include linear bandits [22] and combinatorial bandits [5, 11]. In these problems, the decision space is known to the learner, i.e., the vectors corresponding to the arms in linear bandits and the subsets of arms over which the objective function is to be optimized in combinatorial bandits is known to the learner. However in our problem, we do not know the sorted order of the arm means, i.e., the set of all valid gaps is unknown *a priori*. Our problem does not reduce to these settings.

¹First find the smallest and largest numbers, say a and b respectively. Divide the interval $[a, b]$ into $K + 1$ equal-width bins and map each number to its corresponding bin, while maintaining the smallest and largest number in each bin. Since at least one bin is empty by the pigeonhole principle, the largest gap is between two numbers belonging to different bins. Calculate all gaps between bins and cluster based on the largest of those.

Another related problem is noisy sorting and ranking. Here the typical goal is to sort a list using noisy pairwise comparisons. Our framework encompasses noisy ranking based on Borda scores [1]. The Borda score of an item is the probability that it is ranked higher in a pairwise comparison with another item chosen uniformly at random. In our setting, the Borda score is the mean of each distribution. Much of the theoretical computer science literature on this topic assumes a bounded noise model for comparisons (i.e., comparisons are probably correct with a positive margin) [8, 6, 2, 21]. This is unrealistic in many real-world applications since near equals or outright ties are not uncommon. The largest gap problem we study can be used to (partially) order items into two natural groups, one with large means and one with small means. Previous related work considered a similar problem with prescribed (non-adaptive) quantile groupings [18].

3 MaxGap Bandit Algorithms

We propose elimination [7] and UCB [13] style algorithms for the MaxGap-bandit problem. These algorithms operate on the arm *gaps* instead of the arm *means*. The subroutine to construct confidence intervals on the gaps (denoted by $\mathcal{U}_{\Delta_a}(t)$) using confidence intervals on the arm means (denoted by $[l_a(t), r_a(t)]$) is described in Algorithm 4 in Section 4, and this subroutine is used by all three algorithms described in this section.

3.1 Elimination Algorithm: MaxGapElim

At each time step, MaxGapElim (Algorithm 1) samples all arms in an active set consisting of arms a whose gap upper bound \mathcal{U}_{Δ_a} is larger than the global lower bound $L\Delta$ on the maximum gap, and stops when there are only two arms in the active set.

Algorithm 1 MaxGapElim

```

1: Initialize active set  $A = [K]$ 
2: for  $t = 1, 2, \dots$  do // rounds
3:    $\forall a \in A$ , sample arm  $a$ , compute  $[l_a(t), r_a(t)]$  using (5). // arm confidence intervals
4:    $\forall a \in A$ , compute  $\mathcal{U}_{\Delta_a}(t)$  using Algorithm 4. // upper bound on arm max gap
5:   Compute  $L\Delta(t)$  using (9). // lower bound on max gap
6:    $\forall a \in A$ , if  $\mathcal{U}_{\Delta_a}(t) \leq L\Delta(t)$ ,  $A = A \setminus a$ . // Elimination
7:   If  $|A| = 2$ , stop. Return clusters using max gap in the empirical means. // Stopping condition

```

3.2 UCB algorithms: MaxGapUCB and MaxGapTop2UCB

MaxGapUCB (Algorithm 2) is motivated from the principle of “optimism in the face of uncertainty”. It samples *all* arms with the highest gap upper bound. Note that there are at least two arms with the highest gap upper bound because any gap is shared by at least two arms (one on the right and one on the left). The stopping condition is akin to the stopping condition in Jamieson et al. [13].

Algorithm 2 MaxGapUCB

```

1: Initialize  $\mathcal{U} = [K]$ .
2: for  $t = 1, 2, \dots$  do
3:    $\forall a \in \mathcal{U}$ , sample  $a$  and update  $[l_a(t), r_a(t)]$  using (5).
4:    $\forall a \in [K]$ , compute  $\mathcal{U}_{\Delta_a}(t)$  using Algorithm 4.
5:   Let  $M_1(t) = \max_{j \in [K]} \mathcal{U}_{\Delta_j}(t)$ . Set  $\mathcal{U} = \{a : \mathcal{U}_{\Delta_a}(t) = M_1(t)\}$ . // highest gap-UCB arms
6:   If  $\exists i, j$  such that  $T_i(t) + T_j(t) \geq c \sum_{a \notin \{i, j\}} T_a(t)$ , stop. // stopping condition

```

Alternatively, we can use an LUCB [16]-type algorithm that samples arms which have the two highest gap upper bounds, and stops when the second-largest gap upper bound is smaller than the global lower bound $L\Delta(t)$. We refer to this algorithm as MaxGapTop2UCB (Algorithm 3).

Algorithm 3 MaxGapTop2UCB

- 1: Initialize $\mathcal{U}_1 \cup \mathcal{U}_2 = [K]$.
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: $\forall a \in \mathcal{U}_1 \cup \mathcal{U}_2$, sample a and update $[l_a(t), r_a(t)]$ using (5).
 - 4: $\forall a \in [K]$, compute $\text{U}\Delta_a(t)$ using Algorithm 4.
 - 5: Let $M_1(t) = \max_{j \in [K]} \text{U}\Delta_j(t)$. Set $\mathcal{U}_1 = \{a : \text{U}\Delta_a(t) = M_1(t)\}$. // highest gap-UCB arms
 - 6: Let $M_2(t) = \max_{j \in [K] \setminus \mathcal{U}_1} \text{U}\Delta_j(t)$. Set $\mathcal{U}_2 = \{a : \text{U}\Delta_a(t) = M_2(t)\}$. // 2nd highest gap-UCB
 - 7: Compute $\text{L}\Delta(t)$ using (9). If $M_2(t) < \text{L}\Delta(t)$, stop.
-

Algorithm 4 Procedure to find $\text{U}\Delta_a(t)$

- 1: Set $P_a^r = \{i : l_i(t) \in [l_a(t), r_a(t)]\}$.
 - 2: $\text{U}\Delta_a^r(t) = \max_{i \in P_a^r} \{G_a^r(l_i(t), t)\}$, where $G_a^r(x, t)$ is given by (7). // eqn. (8)
 - 3: Set $P_a^l = \{i : r_i(t) \in [l_a(t), r_a(t)]\}$.
 - 4: $\text{U}\Delta_a^l(t) = \max_{i \in P_a^l} \{G_a^l(r_i(t), t)\}$, where $G_a^l(x, t)$ is given by (19). // eqn. (20)
 - 5: **return** $\text{U}\Delta_a(t) \leftarrow \max\{\text{U}\Delta_a^r(t), \text{U}\Delta_a^l(t)\}$
-

4 Confidence Bounds for Gaps

In this section we explain how to construct confidence bounds for the arm gaps (denoted by $\text{U}\Delta_a$ and $\text{L}\Delta$) using confidence bounds for the arm means (denoted by $[l_a, r_a]$). These bounds are key ingredients for the algorithms described in Section 3.

Given i.i.d. samples from arm a , an empirical mean $\hat{\mu}_a$ and confidence interval on the arm mean can be constructed using standard methods. Let $T_a(t)$ denote the number of samples from arm a after t time steps of the algorithm. Throughout our analysis and experimentation we use confidence intervals on the mean of the form

$$l_a(t) = \hat{\mu}_a(t) - c_{T_a(t)} \text{ and } r_a(t) = \hat{\mu}_a(t) + c_{T_a(t)}, \text{ where } c_s = \sqrt{\frac{\log(4Ks^2/\delta)}{s}}. \quad (5)$$

The confidence intervals are chosen so that [12]

$$\mathbb{P}(\forall t \in \mathbb{N}, \forall a \in [K], \mu_a \in [l_a(t), r_a(t)]) \geq 1 - \delta. \quad (6)$$

Conceptually, the confidence intervals on the arm means can be used to construct upper confidence bounds on the mean gaps $\{\Delta_i\}_{i \in [K]}$ in the following manner. Consider all possible configurations of the arm means that satisfy the confidence interval constraints in (5). Each configuration fixes the gaps associated with any arm $a \in [K]$. Then the maximum gap value over all configurations is the upper confidence bound on arm a 's gap; we denote it as $\text{U}\Delta_a$. The above procedure can be formulated as a mixed integer linear program (see Appendix B.1). In the algorithms in Section 3, this optimization problem needs to be solved at every time t and for every arm $a \in [K]$ before querying a new sample, which can be practically infeasible. In Algorithm 4, we give an efficient $O(K^2)$ time dynamic programming algorithm to compute $\text{U}\Delta_a$. We next explain the main ideas used in this algorithm, and refer the reader to Appendix B.2 for the proofs.

Each arm a has a right and left gap, $\Delta_a^r := \mu_{(\ell-1)} - \mu_{(\ell)}$ and $\Delta_a^l := \mu_{(\ell)} - \mu_{(\ell+1)}$, where ℓ is the rank of a , i.e., $\mu_a = \mu_{(\ell)}$. We construct separate upper bounds $\text{U}\Delta_a^r(t)$ and $\text{U}\Delta_a^l(t)$ for these gaps and then define $\text{U}\Delta_a(t) = \max\{\text{U}\Delta_a^r(t), \text{U}\Delta_a^l(t)\}$. Here we provide an intuitive description for how the bounds are computed, focusing on $\text{U}\Delta_a^r(t)$ as an example. To start, suppose the true mean of arm a is known exactly, while the means of other arms are only known to lie within their confidence intervals. If there exist arms that cannot go to the left of arm a , one can see that the largest right gap for a is obtained by placing all arms that can go to the left of a at their leftmost positions, and all remaining arms at their rightmost positions, as shown in Fig. 3(a). If however all arms can go to the left of arm a , the configuration that gives the largest right gap for a is obtained by placing the arm with the largest upper bound at its right boundary, and all other arms at their left boundaries, as illustrated in Fig. 3(b). We define a function $G_a^r(x, t)$ that takes as input a known position x for the mean of arm a



Figure 3: Computing maximum right gap of blue arm when its true mean is known (at position indicated by blue x), while the other means are known only to lie within their confidence intervals. (a) If there exist arms that cannot go to the left of blue (red, green, purple), the largest right gap for blue is obtained by placing all arms that can go to the left of blue at their left boundaries and the remaining arms at their rightmost positions. (b) If all arms can go to the left of blue, the largest right gap for blue is obtained by placing the arm with the largest right confidence bound (purple) at its right boundary and all other arms at their left boundaries.

and the confidence intervals of all other arms at time t , and returns the maximum right gap for arm a using the above idea as follows.

$$G_a^r(x, t) = \begin{cases} \min_{j: l_j(t) > x} r_j(t) - x & \text{if } \{j : l_j(t) > x\} \neq \emptyset, \\ \max_{j \neq a} r_j(t) - x & \text{otherwise.} \end{cases} \quad (7)$$

However, the true mean of arm a is not known exactly but only that it lies within its confidence interval. The insight that helps here is that $G_a^r(x, t)$ must achieve its maximum when x is at one of the finite locations in $\{l_j(t) : l_a(t) \leq l_j(t) \leq r_a(t)\}$. We define $P_a^r := \{j : l_a(t) \leq l_j(t) \leq r_a(t)\}$ as the set of arms relevant for the right gap of a , and then the maximum possible right gap of a is

$$\text{UD}_a^r(t) = \max\{G_a^r(l_j(t), t) : j \in P_a^r\}. \quad (8)$$

An upper bound for the left gap UD_a^l can be similarly obtained. We explain this and give a proof of correctness in Appendix B.2.

The algorithms also use a single global lower bound on the maximum gap. To do so, we sort the items according to their empirical means, and find partitions of items that are clearly separated in terms of their confidence intervals. At time t , let $(i)_t$ denote the arm with the i^{th} -largest empirical mean, i.e., $\hat{\mu}_{(K)_t}(t) \leq \dots \leq \hat{\mu}_{(2)_t}(t) \leq \hat{\mu}_{(1)_t}(t)$ (this can be different from the true ranking which is denoted by (\cdot) without the subscript t). We *detect* a nonzero gap at arm k if $\max_{a \in \{(k+1)_t, \dots, (K)_t\}} r_a(t) < \min_{a \in \{(1)_t, \dots, (k)_t\}} l_a(t)$. Thus, a lower bound on the largest gap is

$$\text{LD}(t) = \max_{k \in [K-1]} \left(\min_{a \in \{(1)_t, \dots, (k)_t\}} l_a(t) - \max_{a \in \{(k+1)_t, \dots, (K)_t\}} r_a(t) \right). \quad (9)$$

5 Analysis

In this section, we first state the accuracy and sample complexity guarantees for MaxGapElim and MaxGapUCB, and then discuss our results. The proofs can be found in the Supplementary material.

Theorem 1. *With probability $1 - \delta$, MaxGapElim, MaxGapUCB and MaxGapTop2UCB cluster the arms according to the maximum gap, i.e., they satisfy (3).*

The number of times arm a is sampled by both the algorithms depends on $\gamma_a = \min\{\gamma_a^l, \gamma_a^r\}$ where

$$\gamma_a^r = \max_{j: 0 < \Delta_{a,j} < \Delta_{\max}} \min\{\Delta_{a,j}, (\Delta_{\max} - \Delta_{a,j})\} \quad (10)$$

$$\gamma_a^l = \max_{j: 0 < \Delta_{j,a} < \Delta_{\max}} \min\{\Delta_{j,a}, (\Delta_{\max} - \Delta_{j,a})\}. \quad (11)$$

The maxima is assumed to be ∞ in (10) and (11) if there is no j that satisfies the constraint to account for edge arms. The quantity γ_a acts as a measure of hardness for arm a ; Theorem 2 states that MaxGapElim and MaxGapUCB sample arm a at most $\tilde{O}(1/\gamma_a^2)$ number of times (up to log factors).

Theorem 2. *With probability $1 - \delta$, the sample complexity of MaxGapElim and MaxGapUCB is bounded by*

$$O\left(\sum_{a \in [K] \setminus \{(m), (m+1)\}} \frac{\log(K/\delta\gamma_a)}{\gamma_a^2}\right)$$

Next, we provide intuition for why the sample complexity depends on the parameters in (10) and (11). In particular, we show that $O((\gamma_a^r)^{-2})$ (resp. $O((\gamma_a^l)^{-2})$) is the number of samples of a required to rule out arm a 's right (resp. left) gap from being the largest gap.

Let us focus on the right gap for simplicity. To understand how (10) naturally arises, consider Fig. 4, which denotes the confidence intervals on the means at some time t . A lower bound on the gap $L\Delta(t)$ can be computed between the left and right confidence bounds of arms 10 and 11 respectively as shown. Consider the computation of the upper bound $U\Delta_7^r(t)$ on the right gap of arm $a = 7$. Arm 4 lies to the right of arm 7 with high probability (unlike the arms with dashed confidence intervals), so the upper bound $U\Delta_7^r(t) \leq r_4(t) - l_7(t)$. Considering only the right gap for simplicity, as soon as $U\Delta_7^r(t) < L\Delta(t)$, arm 7 can be eliminated as a candidate for the maximum gap. Thus, an arm a is removed from consideration as soon as we find a *helper* arm j (arm 4 in Fig. 4) that satisfies two properties: (1) the confidence interval of arm j is disjoint from that of arm a , and (2) the upper bound $U\Delta_a^r(t) = r_j(t) - l_a(t) < L\Delta(t)$.

The first of these conditions gives rise to the term $\Delta_{a,j}$ in (10), and the second condition gives rise to the term $(\Delta_{\max} - \Delta_{a,j})$. Since any arm j that satisfies these conditions can act as a helper for arm a , we take the maximum over all arms j to yield the smallest sample complexity for arm a .

This also shows that if all arms are either very close to a or at a distance approximately Δ_{\max} from a , then the upper bound $U\Delta_7^r(t) = r_4(t) - l_7(t) > L\Delta(t)$ and arm 7 cannot be eliminated. Thus arm a could have a small gap with respect to its adjacent arms, but if there is a large gap in the vicinity of arm a , it cannot be eliminated quickly. This illustrates that the maximum gap identification problem is not equivalent to best-arm identification (BAI) on gaps. Section 6 formalizes this intuition.

Key Differences compared to BAI Analysis: The analysis of MaxGapUCB is very different from the standard UCB analysis. On a high-level, in BAI, the number of samples of a sub-optimal arm i is bounded by observing that

$$\begin{aligned} \text{Arm } i \text{ is pulled} &\implies \mu_i + 2c_{T_i(t)} \geq \hat{\mu}_i + c_{T_i(t)} \geq \hat{\mu}_{(1)} + c_{T_{(1)}(t)} \geq \mu_{(1)} \\ &\implies 2c_{T_i(t)} \geq \mu_{(1)} - \mu_i = \Delta_i. \end{aligned} \quad (12)$$

The last inequality *directly* bounds the number of samples $T_i(t)$ of a sub-optimal arm i . In MaxGapUCB, the gap upper bound is obtained using the confidence intervals of two arms, and the fact that a sub-optimal gap (i, j) has the highest gap-UCB implies that

$$\begin{aligned} (\mu_j + 2c_{T_j(t)}) - (\mu_i - 2c_{T_i(t)}) &\geq (\hat{\mu}_j + c_{T_j(t)}) - (\hat{\mu}_i - 2c_{T_i(t)}) \geq \Delta_{\max} \\ &\implies 2(c_{T_j(t)} + c_{T_i(t)}) \geq \Delta_{\max} - \Delta_{ij}. \end{aligned}$$

Thus unlike the reasoning in (12), the number of samples from arm i is coupled to the number of samples from arm j , and $T_i(t) \rightarrow \infty$ if j is not sampled enough. We show in our analysis that this cannot happen in MaxGapUCB. Furthermore, any arm i is coupled with multiple other arms since the ordering of the arms is unknown, and may have to be sampled even if its own gap is small - a phenomenon absent in standard BAI analysis because of the independence of the arm means. In our proof, we account for all samples of an arm by defining states the arm can belong to (called levels), and arguing about the confidence intervals of the arms in unison.

6 Minimax Lower Bound

In this section, we demonstrate that the MaxGap problem is fundamentally different from best-arm identification (BAI) on gaps. We construct a problem instance and prove a lower bound on the number of samples needed by any probably correct algorithm. The lower bound matches the upper bounds in the previous section for this instance.

Lemma 1. Consider a model \mathcal{B} with $K = 4$ normal distributions $\mathcal{P}_i = \mathcal{N}(\mu_i, 1)$, where

$$\mu_4 = 0, \quad \mu_3 = \epsilon, \quad \mu_2 = \nu + 2\epsilon, \quad \mu_1 = 2\nu + 2\epsilon,$$

for some $\nu \gg \epsilon > 0$. Then any algorithm that is correct with probability at least $1 - \delta$ must collect $\Omega(1/\epsilon^2)$ samples of arm 4 in expectation.

Proof Outline: The proof uses a standard change of measure argument [10]. We construct another problem instance \mathcal{B}' which has a different maximum gap clustering compared to \mathcal{B} (see Fig. 5; the maxgap clustering in \mathcal{B} is $\{4, 3\} \cup \{2, 1\}$, while the maxgap clustering in \mathcal{B}' is $\{4, 3, 2\} \cup \{1\}$), and show that in order to distinguish between \mathcal{B} and \mathcal{B}' , any probably correct algorithm must collect at least $\Omega(1/\epsilon^2)$ samples of arm 4 in expectation (see Appendix E for details). From the definition of γ_a using (10),(11), it is easy to check that $\gamma_4 = \epsilon$. Therefore, for problem instance \mathcal{B} our algorithms find the maxgap clustering using at most $O(\log(\epsilon/\delta)/\epsilon^2)$ samples of arm 4 (c.f. Theorem 2). This essentially matches the lower bound above.

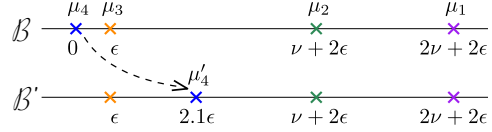


Figure 5: Changing the original bandit model \mathcal{B} to \mathcal{B}' . μ_4 is shifted to the right by 2.1ϵ . As a result, the maximum gap in \mathcal{B}' is between green and purple.

This example illustrates why the maximum gap identification problem is different from a simple BAI on gaps. Suppose an oracle told a BAI algorithm the ordering of the arm means. Using the ordering it can convert the 4-arm maximum gap problem \mathcal{B} to a BAI problem on 3 gaps, with distributions $\mathcal{P}_{4,3} = \mathcal{N}(\epsilon, 2)$, $\mathcal{P}_{3,2} = \mathcal{N}(\nu + \epsilon, 2)$, and $\mathcal{P}_{2,1} = \mathcal{N}(\nu, 2)$. The BAI algorithm can sample arms i and $i + 1$ to get a sample of the gap $(i + 1, i)$. We know from standard BAI analysis [13] that the gap $(4, 3)$ can be eliminated from being the largest by sampling it (and hence arm 4) $O(1/\nu^2)$ times, which can be arbitrarily lower than the $1/\epsilon^2$ lower bound in Lemma 1. Thus the ordering information given to the BAI algorithm is crucial for it to quickly identify the larger gaps. The problem we solve in this paper is identifying the maximum gap when the ordering information is *not* available.

7 Experiments

We conduct three experiments. First, we verify the validity of our sample complexity bounds in Section 7.1. We then study the performance of our adaptive algorithms on simulated data in Section 7.2, and on the Streetview dataset in Section 7.3. The code for all experiments is publicly available [19].

7.1 Sample Complexity

In Fig. 6(b) and Fig. 6(c), we plot the empirical stopping time against the theoretical sample complexity (Theorem 2) for different arm configurations. We choose the arm configuration in Fig. 6(a) containing $K = 15$ arms and three unique gaps - a small gap Δ_3 and two large gaps $\Delta_2 < \Delta_1 = \Delta_{\max} = 0.4$. The hardness parameter is changed by increasing Δ_2 (from 0.35 to 0.39) and bringing it closer to Δ_1 . The rewards are normally distributed with $\sigma = 0.05$. We see a linear relationship in Fig. 6(b) which suggests that the sample complexity expression in Theorem 2 is correct up to constants. In Fig. 6(c) we include random sampling and see that our adaptive algorithms require up to 5x fewer samples when run until completion. Fig. 6(c) also shows that our adaptive algorithms always outperform random sampling, and the gains increase with hardness. We used a lower bound based stopping condition for Random, Elimination, Top2UCB, and set $c = 5$ in the UCB stopping condition (value of c chosen empirically as in [13]).

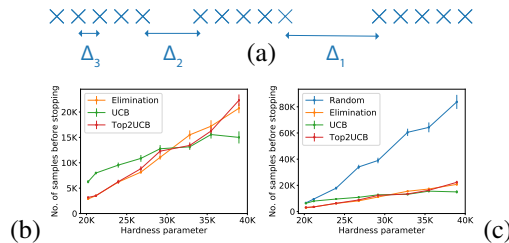


Figure 6: Stopping time experiments

7.2 Simulated Data

In the second experiment, we study the performance on a simulated set of means containing two large gaps. The mean distribution plotted in Fig. 7(a) has $K = 24$ arms ($\mathcal{N}(\cdot, 1)$), with two large mean

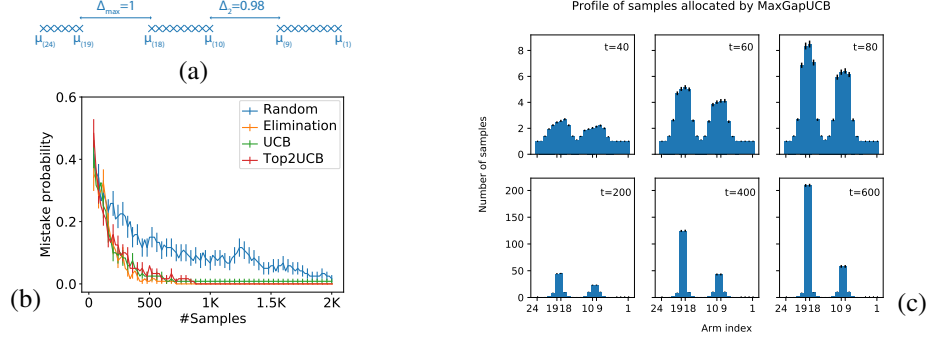


Figure 7: (a) Two large gaps. (b) Clustering error probability for means shown in Fig. 7(a). (c) The profile of samples allocated by MaxGapUCB to each arm in (a) at different time steps.

gaps $\Delta_{10,9} = 0.98, \Delta_{19,18} = 1.0$, and remaining small gaps ($\Delta_{i+1,i} = 0.2$ for $i \notin \{9, 18\}$). We expect to see a big advantage for adaptive sampling in this example because almost every sub-optimal arm has a *helper* arm (see Section 5) which can help eliminate it quickly, and adaptive algorithms can then focus on distinguishing the two large gaps. A non-adaptive algorithm on the other hand would continue sampling all arms. We plot the fraction of times $C_1 \neq \{1, \dots, 18\}$ in 120 runs in Fig. 7(b), and see that the active algorithms identify the largest gap in 8x fewer samples. To visualize the adaptive allocation of samples to the arms, we plot in Fig. 7(c) the number of samples queried for each arm at different time steps by MaxGapUCB. Initially, MaxGapUCB allocates samples uniformly over all the arms. After a few time steps, we see a bi-modal profile in the number of samples. Since all arms that achieve the largest UD are sampled, we see that several arms that are near the pairs (10, 9) and (19, 18) are also sampled frequently. As time progresses, only the pairs (10, 9) and (19, 18) get sampled, and eventually more samples are allocated to the larger gap (19, 18) among the two.

7.3 Streetview Dataset

For our third experiment we study performance on the Streetview dataset [17, 18] whose means are plotted in Fig. 8(a). We have $K = 90$ arms, where each arm is a normal distribution with mean equal to the Borda safety score of the image and standard deviation $\sigma = 0.05$. The largest gap of 0.029 is between arms 2 and 3, and the second largest gap is 0.024. In Fig. 8(b), we plot the fraction of times $\hat{C}_1 \neq \{1, 2\}$ in 120 runs as a function of the number of samples, for four algorithms, viz., random (non-adaptive) sampling, MaxGapElim, MaxGapUCB, and MaxGapTop2UCB. The error bars denote standard deviation over the runs. MaxGapUCB and MaxGapTop2UCB require 6-7x fewer samples than random sampling.

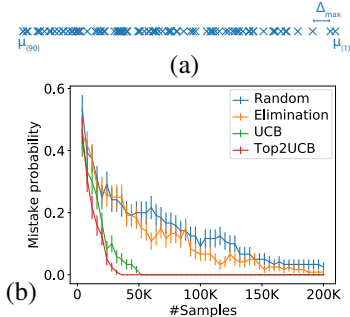


Figure 8: (a) Borda safety scores for Streetview images. (b) Probability of returning a wrong cluster.

8 Conclusion

In this paper, we proposed the MaxGap-bandit problem: a novel maximum-gap identification problem that can be used as a basic primitive for clustering and approximate ranking. Our analysis shows a novel hardness parameter for the problem, and our experiments show 6-8x gains compared to non-adaptive algorithms. We use simple Hoeffding based confidence intervals in our analysis for simplicity, but better bounds can be obtained using tighter confidence intervals [13]. Several extensions of this basic problem are possible. An ϵ -relaxation of the MaxGap Bandit is useful when the largest and second-largest gaps are close to each other. Other possibilities include identifying the largest gap within a top quantile of the arms, or clustering with a constraint that the returned clusters are of similar cardinality. All of these extensions will likely require new ideas, as it is unclear how to obtain a lower bound for the gap associated with every arm. Finding an instance-dependent lower bound for MaxGap-bandit is an intriguing problem. Finally, one way to cluster the distributions into more than two clusters is to apply the max-gap identification algorithms recursively; however it would be interesting to come up with algorithms that can perform this clustering directly.

Acknowledgments

Ardhendu Tripathy would like to thank Ervin Tanczos for helpful discussions. This work was partially supported by AFOSR/AFRL grants FA8750-17-2-0262 and FA9550-18-1-0166.

References

- [1] Arpit Agarwal, Shivani Agarwal, Sepehr Assadi, and Sanjeev Khanna. Learning with limited rounds of adaptivity: Coin tossing, multi-armed bandits, and ranking from pairwise comparisons. In *Conference on Learning Theory*, pages 39–75, 2017.
- [2] Mark Braverman, Jieming Mao, and S Matthew Weinberg. Parallel algorithms for select and partition with noisy comparisons. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 851–862. ACM, 2016.
- [3] Sebastien Bubeck, Tengyao Wang, and Nitin Viswanathan. Multiple identifications in multi-armed bandits. In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, ICML’13*, pages I–258–I–265. JMLR.org, 2013. URL <http://dl.acm.org/citation.cfm?id=3042817.3042848>.
- [4] Lijie Chen, Jian Li, and Mingda Qiao. Nearly instance optimal sample complexity bounds for top-k arm selection. In *Artificial Intelligence and Statistics*, pages 101–110, 2017.
- [5] Shouyuan Chen, Tian Lin, Irwin King, Michael R Lyu, and Wei Chen. Combinatorial pure exploration of multi-armed bandits. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 379–387. Curran Associates, Inc., 2014. URL <http://papers.nips.cc/paper/5433-combinatorial-pure-exploration-of-multi-armed-bandits.pdf>.
- [6] Susan Davidson, Sanjeev Khanna, Tova Milo, and Sudeepa Roy. Top-k and clustering with noisy comparisons. *ACM Transactions on Database Systems (TODS)*, 39(4):35, 2014.
- [7] Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7(Jun):1079–1105, 2006.
- [8] Uriel Feige, Prabhakar Raghavan, David Peleg, and Eli Upfal. Computing with noisy information. *SIAM Journal on Computing*, 23(5):1001–1018, 1994.
- [9] Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In *Advances in Neural Information Processing Systems*, pages 3212–3220, 2012.
- [10] Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In *Conference on Learning Theory*, pages 998–1027, 2016.
- [11] Weiran Huang, Jungseul Ok, Liang Li, and Wei Chen. Combinatorial pure exploration with continuous and separable reward functions and its applications. In *IJCAI*, volume 18, pages 2291–2297, 2018.
- [12] Kevin Jamieson and Robert Nowak. Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. In *2014 48th Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2014.
- [13] Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lil’ucb: An optimal exploration algorithm for multi-armed bandits. In *Conference on Learning Theory*, pages 423–439, 2014.
- [14] Kwang-Sung Jun, Kevin G Jamieson, Robert D Nowak, and Xiaojin Zhu. Top arm identification in multi-armed bandits with batch arm pulls. In *AISTATS*, pages 139–148, 2016.

- [15] Shivaram Kalyanakrishnan and Peter Stone. Efficient selection of multiple bandit arms: Theory and practice. In *ICML*, volume 10, pages 511–518, 2010.
- [16] Shivaram Kalyanakrishnan, Ambuj Tewari, Peter Auer, and Peter Stone. Pac subset selection in stochastic multi-armed bandits. In *ICML*, volume 12, pages 655–662, 2012.
- [17] Sumeet Katariya, Lalit Jain, Nandana Sengupta, James Evans, and Robert Nowak. Chicago streetview dataset. 2018. URL https://github.com/sumeetsk/coarse_ranking/.
- [18] Sumeet Katariya, Lalit Jain, Nandana Sengupta, James Evans, and Robert Nowak. Adaptive sampling for coarse ranking. In *International Conference on Artificial Intelligence and Statistics*, pages 1839–1848, 2018.
- [19] Sumeet Katariya, Ardhendu Tripathy, and Robert Nowak. Code for maxgap bandit algorithms and experiments. 2019. URL https://github.com/sumeetsk/maxgap_bandit.
- [20] 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.
- [21] Cheng Mao, Jonathan Weed, and Philippe Rigollet. Minimax rates and efficient algorithms for noisy sorting. In Firdaus Janoos, Mehryar Mohri, and Karthik Sridharan, editors, *Proceedings of Algorithmic Learning Theory*, volume 83 of *Proceedings of Machine Learning Research*, pages 821–847. PMLR, 07–09 Apr 2018. URL <http://proceedings.mlr.press/v83/mao18a.html>.
- [22] Marta Soare, Alessandro Lazaric, and Rémi Munos. Best-arm identification in linear bandits. In *Advances in Neural Information Processing Systems*, pages 828–836, 2014.
- [23] Honglei Zhuang, Chi Wang, and Yifan Wang. Identifying outlier arms in multi-armed bandit. In *Advances in Neural Information Processing Systems*, pages 5204–5213, 2017.

A Details for Section 1.2: Comparison to a Naive Algorithm

The naive algorithm first sorts the arms to determine the adjacent arms for every arm, and then runs a best-arm identification bandit algorithm on the gaps to identify the largest gap. An unbiased sample of the gap between two arms can be obtained by taking the difference of the samples from the two arms. Here we analyze the sample complexity of the naive algorithm for a general arrangement of the means.

Consider an arm $i \notin \{(m), (m+1)\}$, and let us analyze the number of times arm i is sampled by the naive algorithm. Let $\Delta_{i,j} = \mu_j - \mu_i$. Then $\Delta_i^r = \min_{j: \Delta_{i,j} > 0} \Delta_{i,j}$ is the right gap of arm i (Δ_i^l is defined analogously). In the first step of the naive algorithm, arm i needs to be sampled at least $(\Delta_i^r)^{-2}$ times to determine its right neighbor. Once the right neighbor has been determined, the best-arm identification requires at least $(\Delta_{\max} - \Delta_i^r)^{-2}$ samples to distinguish arm i 's right gap from Δ_{\max} . Since samples from the first step can be reused, the minimum number of samples required by the naive algorithm to rule out arm i 's right gap is $(\tilde{\gamma}_i^r)^{-2}$ where

$$\tilde{\gamma}_i^r = \min_{j: \Delta_{i,j} > 0} \{\Delta_{i,j}, \Delta_{\max} - \Delta_{i,j}\} \quad (13)$$

We can define $(\tilde{\gamma}_i^l)^{-2}$ analogously, and the naive algorithm collects $\Omega(\tilde{\gamma}_i^{-2})$ from arm i , where $\tilde{\gamma}_i = \min\{\tilde{\gamma}_i^r, \tilde{\gamma}_i^l\}$.

The hardness parameter of our active algorithms that is analogous to (13) is given by (4), repeated here for convenience

$$\gamma_i^r := \max_{j: \Delta_{i,j} > 0} \min\{\Delta_{i,j}, \Delta_{\max} - \Delta_{i,j}\}. \quad (14)$$

Comparing (14) to (13), we see that $\gamma_i^r > \tilde{\gamma}_i^r$.

For the toy problem discussed in Section 1.2, if we assume that $\Delta_{\min} < \Delta_{\max}/2$, we have that $\tilde{\gamma}_i = \Delta_{\min}$, while $\gamma_i = \Delta_{\max}/2 \forall i \notin \{(m), (m+1)\}$, which results in $(\Delta_{\max}/\Delta_{\min})^2$ order savings in the number of samples.

B Details for Section 4: Confidence Bounds for Gaps

We first explain the mixed integer program formulation for obtaining the upper confidence bounds on the mean gaps in Appendix B.1, and then prove the validity of Algorithm 4 in Appendix B.2.

B.1 MIP Formulation of Confidence Bounds for Gaps

Conceptually, the confidence intervals on the arm means can be used to construct upper confidence bounds on the mean gaps $\{\Delta_i\}_{i \in [K]}$ in the following manner. Consider all possible configurations of the arm means that satisfy the confidence interval constraints in (5). Each configuration fixes the gaps associated with any arm $a \in [K]$. Then the maximum gap value over all configurations is the upper confidence bound on arm a 's gap; we denote it as $\text{U}\Delta_a$.

If we focus on the right gap of arm a , the above procedure is equivalent to solving the following optimization problem.

$$\text{U}\Delta_a^r(t) \triangleq \max_{b \in [K] \setminus \{a\}} \max_{\mu'_1, \dots, \mu'_K} \mu'_b - \mu'_a \quad (15)$$

$$\text{subject to: } l_i(t) \leq \mu'_i \leq r_i(t) \quad \forall i \in [K], \text{ and} \quad (16)$$

$$\mu'_i \notin (\mu'_a, \mu'_b) \quad \forall i \in [K] \setminus \{a, b\}. \quad (17)$$

Constraint (16) ensures that μ'_i is in the confidence interval for the mean of arm i at time t , and constraint (17) ensures that arm b is the right neighbor of arm a .

The constraint (17) is a sorting constraint that can only be formulated using a binary variable. For an arm $i \in [K] \setminus \{a, b\}$ (17) can be formulated using a constant M as

$$\mu'_i \leq \mu'_a + M(1 - z_i), \quad (18a)$$

$$\mu'_i \geq \mu'_b - Mz_i, \quad (18b)$$

$$z_i \in \{0, 1\}. \quad (18c)$$

The value of M is chosen to be large number. Replacing constraint (17) by constraints (18a), (18b), (18c) for all $i \in [K] \setminus \{a, b\}$ gives an equivalent optimization problem whose optimum value is $\text{UD}_a^r(t)$. This can be seen to be true by considering the cases based on the value of z_i . If $z_i = 0$, $\mu'_i \geq \mu'_b$ and if $z_i = 1$, $\mu'_i \leq \mu'_a$. Because M is chosen to be a large number, in either case $\mu'_i \notin (\mu'_a, \mu'_b)$ and constraint (17) is satisfied. If constraint (17) is satisfied, then a similar argument allows us to choose the value of z_i that satisfies constraints (18a) and (18b).

B.2 Validity of Algorithm 4

In this section we find the value of $\text{UD}_a^r(t)$ as defined in (15) by first obtaining an upper bound to it. The proof of the upper bound is *constructive* in nature, showing that the upper bound is actually achievable. That is, (a) there is a set of real numbers $\{\mu'_i : i \in [K]\}$ which satisfy (16), (b) an index a_* which satisfies (17) with $x = \mu'_{a_*}$, such that $\text{UD}_a^r(t) = \mu'_{a_*} - \mu'_a$.

We first find an upper bound to the right gap of an arm a assuming we know its true mean μ_a , but only have confidence intervals for the means of the other arms $\mu_i \in [l_i(t), r_i(t)] \forall i \neq a$.

Lemma 2. *If all the arm means are known, the right gap associated with an arm $a \in [K]$ is $\min_{i: \mu_i > \mu_a} \mu_i - \mu_a$; if the domain is empty we say that arm a 's right gap is 0. For any $x \in \mathbb{R}$, define a function $G_a^r(\cdot)$ of the confidence intervals as follows.*

$$G_a^r(x, t) \triangleq \begin{cases} \min_{j: l_j(t) > x} r_j(t) - x & \text{if } \{j : l_j(t) > x\} \neq \emptyset, \\ \max_{j \neq a} r_j(t) - x & \text{otherwise.} \end{cases}$$

Suppose we know the value of arm a 's mean, i.e. μ_a and the confidence intervals $[l_i(t), r_i(t)] \forall i \neq a$. Then the largest possible right gap of arm a is $G_a^r(\mu_a, t)$.

Proof. We suppose that the right gap of arm a is greater than the upper bound and show a contradiction to the good event (6).

Case I: $\{j : l_j(t) > \mu_a\} \neq \emptyset$. Identify the arm $j_* = \arg \min_{j: l_j(t) > \mu_a} r_j(t)$ such that $G_a^r(\mu_a, t) = r_{j_*}(t) - \mu_a$. Let the true right gap for arm a be $\mu_k - \mu_a$. If $k = j_*$, then $\mu_k - \mu_a > r_{j_*}(t) - \mu_a$ would mean that $\mu_{j_*} > r_{j_*}(t)$, which is a contradiction. If $k \neq j_*$ and the right gap is $\mu_k - \mu_a$, then all arms $j \in [K]$ are such that $\mu_j \notin (\mu_a, \mu_k)$. But if $\mu_k - \mu_a > G_a^r(\mu_a, t)$ then $\mu_k > r_{j_*}(t)$, and from the domain in the definition of j_* , its left bound $l_{j_*}(t) > \mu_a$. Hence the confidence interval of j_* satisfies $\mu_a < l_{j_*}(t) < r_{j_*}(t) < \mu_k$. If $\mu_{j_*} \notin (\mu_a, \mu_k)$ then $\mu_{j_*} \notin [l_{j_*}(t), r_{j_*}(t)]$ and that is a contradiction.

Case II: $\{j : l_j(t) > \mu_a\} = \emptyset$. Identify the arm $j_* = \arg \max_{j \neq a} r_j(t)$ such that $G_a^r(\mu_a, t) = r_{j_*}(t) - \mu_a$. Let the true right gap for arm a be $\mu_k - \mu_a$. If $\mu_k - \mu_a > G_a^r(\mu_a, t)$ then $\mu_k > \max_{j \neq a} r_j(t)$ and that is a contradiction.

Thus the right gap of arm a is at most $G_a^r(\mu_a, t)$. We can achieve this upper bound by choosing the set of means $\{\mu'_i : i \in [K] \setminus a\}$ in the following manner. If the value of $G_a^r(\mu_a, t)$ is given by the first branch, set $\mu'_i = r_i(t) \forall i : r_i(t) > \mu_a$ and $\mu'_i = l_i(t) \forall i : l_i(t) < \mu_a$. Otherwise if the value is given by the second branch set $\mu'_{a_*} = r_{a_*}(t)$ for the arm $a_* \neq a$ which has the largest right bound, and set all other $\mu'_i = l_i(t)$ (c.f. Fig. 3 in Section 4). \square

The *left* gap analog of the above proposition can also be proved in a similar manner as above.

Lemma 3. *For any $x \in \mathbb{R}$ and arm $a \in [K]$, define a function $G_a^l(\cdot)$ of the confidence intervals as follows.*

$$G_a^l(x, t) \triangleq \begin{cases} x - \max_{j: r_j(t) < x} l_j(t) & \text{if } \{j : r_j(t) < x\} \neq \emptyset, \\ x - \min_{j \neq a} l_j(t) & \text{otherwise.} \end{cases} \quad (19)$$

Suppose we know μ_a . Using the confidence intervals $[l_i(t), r_i(t)] \forall i \neq a$, an upper bound to the left gap of arm a is $G_a^l(\mu_a, t)$.

We now replace our knowledge of the true mean value μ_a by the good event fact that $\mu_a \in [l_a(t), r_a(t)]$ at all times t . The following lemma is instrumental in arriving at an upper bound for the right gap of arm a that is consistent with the all the arms' confidence intervals.

Lemma 4. At time t , for any arm a its true mean $\mu_a \in [l_a(t), r_a(t)]$ in the good event. Define a subset of arms $\mathcal{I}_a^R(t) \triangleq \{i : l_i(t) \in [l_a(t), r_a(t)]\}$ whose left bounds lie within the confidence interval of arm a . Consider a set of K real numbers $\mathcal{P}' \triangleq \{\mu'_i \in [l_i(t), r_i(t)] : i \in [K]\}$, each associated with a corresponding arm. The largest value for the right gap of arm a if the means are \mathcal{P}' , i.e.,

$$\max\{\mu'_i - \mu'_a : \mu'_i > \mu'_a, \nexists \mu'_j \in (\mu'_a, \mu'_i), i, j \in [K] \setminus a\}$$

occurs when $\mu'_a = l_i(t)$ for some $i \in \mathcal{I}_a^R(t)$.

Proof. Suppose the largest right gap occurs when $\mu'_a \neq l_i(t)$ for any $i \in \mathcal{I}_a^R(t)$. Note that $a \in \mathcal{I}_a^R(t)$ and hence the set is not empty. We show that the right gap can be larger while still satisfying event (6). Let $l_{i_a}(t) = \max_{i \in \mathcal{I}_a^R(t)} \{l_i(t) < \mu'_a\}$. Collect all arms in the set $\mathcal{J}_a = \{j : \mu'_j \in [l_{i_a}(t), \mu'_a]\}$. Consider an alternate bandit model whose arm means are denoted by $\mathcal{Q} \triangleq \{q_i : i \in [K]\}$. We assign

$$q_i = l_{i_a}(t) \forall i \in \mathcal{J}_a \text{ and } q_i = \mu'_i \forall i \notin \mathcal{J}_a.$$

This mean assignment satisfies $q_i \in [l_i(t), r_i(t)] \forall i \in [K]$. This is because by definition of arm i_a in the original bandit model \mathcal{P}' , for all arms $j \in \mathcal{J}_a$ their left bounds satisfy $l_j(t) \leq l_{i_a}(t)$. Thus both the original \mathcal{P}' and the alternate \mathcal{Q} are possible bandit models in the good event (6) up till current time t . However, the right gap for a is larger in the alternate model \mathcal{Q} as shown next. Let arm i result in the right gap for a in the original model \mathcal{P}' , i.e., the right gap is

$$\mu'_i - \mu'_a, \text{ and } \nexists \mu'_j \in (\mu'_a, \mu'_i).$$

Then in the alternate model, $q_i = \mu'_i$, $q_a = l_{i_a}(t)$ and there is no mean $q_j \in (l_{i_a}(t), \mu'_i)$. Then the right gap of arm a is $\mu'_i - l_{i_a}(t) > \mu'_i - \mu'_a$. This contradicts the supposition that the right gap is the largest possible in the original bandit model \mathcal{P}' . \square

An analogous lemma for the *left gap* states that for any set of possible arm means \mathcal{P}' that are consistent with the current confidence intervals, the largest possible left gap of arm a occurs when $\mu'_a = r_i(t)$ for some arm $i \in \mathcal{I}_a^L(t) \triangleq \{i : r_i(t) \in [l_a(t), r_a(t)]\}$. Using the above, we can state the upper bound for the gap of an arm a in terms of all the confidence intervals as follows.

Theorem 3. At any time t , denote the upper bound to the right (resp. left) gap of arm a by $\mathcal{U}\Delta_a^r(t)$ (resp. $\mathcal{U}\Delta_a^l(t)$). The expressions for these upper bounds in terms of the confidence intervals and the functions $G_a^r(\cdot)$, $G_a^l(\cdot)$ in Lemma 2, Lemma 3 are as follows.

$$\begin{aligned} \mathcal{U}\Delta_a^r(t) &\triangleq \max\{G_a^r(l_j(t), t) : l_j(t) \in [l_a(t), r_a(t)]\}, \\ \mathcal{U}\Delta_a^l(t) &\triangleq \max\{G_a^l(r_j(t), t) : r_j(t) \in [l_a(t), r_a(t)]\}. \end{aligned} \quad (20)$$

Then an upper bound to the gap associated with arm a at time t is $\max\{\mathcal{U}\Delta_a^r(t), \mathcal{U}\Delta_a^l(t)\}$. Algorithm 4 gives pseudocode that evaluates $\mathcal{U}\Delta_a^r(t)$.

Proof. We argue for the right gap, an analogous proof gives the statement for the left gap. At any time t in the good event $\mu_i \in [l_i(t), r_i(t)] \forall i \in [K]$, in particular any number in the range $[l_a(t), r_a(t)]$ can be potentially the mean of arm a . From Lemma 4, we know that for a set of numbers \mathcal{P}' that satisfy all current confidence intervals and also maximize the right gap for arm a , the value $\mu'_a = l_i(t)$ for some left bound $l_i(t) \in [l_a(t), r_a(t)]$. If $\mu'_a = l_i(t)$ then by Lemma 2 $G_a^r(l_i(t), t)$ is the largest possible value for arm a in the bandit model \mathcal{P}' . Taking the maximum over all arms in the set $\mathcal{I}_a^R(t) = \{i \in [K] : l_i(t) \in [l_a(t), r_a(t)]\}$, we get the right gap upper bound $\mathcal{U}\Delta_a^r(t)$.

We note that the value $\mathcal{U}\Delta_a^r(t)$ is achievable by an assignment of means that satisfy the confidence bounds at time t . Without loss of generality, assume $\mathcal{U}\Delta_a(t) = \mathcal{U}\Delta_a^r(t) = G_a^r(l_{a_*}(t), t)$ for some arm a_* . One can assign $\mu_a = l_{a_*}(t)$ and other means in a way similar to that in the proof of Lemma 2 to obtain a right gap for arm a equal to the value $G_a^r(l_{a_*}(t), t)$. \square

C Details for Section 5: Accuracy

Theorem 1. With probability $1 - \delta$, MaxGapElim, MaxGapUCB and MaxGapTop2UCB cluster the arms according to the maximum gap, i.e., they satisfy (3).

Proof. Recall that the true maximum gap exists between arms (m) and $(m+1)$. The algorithms return a wrong clustering $U\Delta_{(m)}(t) < L\Delta(t)$ for any time t . We show that this leads to a contradiction if the good event (6) holds.

Assume (6) holds and $U\Delta_{(m)}(t) < L\Delta(t)$ at some time t . Recall that $L\Delta(t)$ is computed using (9), and let $(s)_t$ be the maximizer in (9). Let a be such that $a \in \{(1)_t, \dots, (s)_t\}$ and $a+1 \in \{(s+1)_t, \dots, (K)_t\}$. If (3) holds, we have that

$$\Delta_{\max} \leq U\Delta_{(m)}(t) < L\Delta(t) \stackrel{(a)}{\leq} l_a(t) - r_{a+1}(t) \leq \mu_a - \mu_{a+1},$$

where (a) holds because $L\Delta(t)$ is the minimum gap between a left confidence interval in $\{(1)_t, \dots, (s)_t\}$ and a right confidence interval in $\{(s+1)_t, \dots, (K)_t\}$. This contradicts the fact that Δ_{\max} is the largest gap. \square

D Sample Complexity: Proof of Theorem 2

To state our sample complexity bounds we use a constant α defined as follows [7].

Remark 1. *There exists constant α such that for all $x > 0$, if the number of samples $s \geq \alpha \frac{\log(K/\delta x)}{x^2}$, then $c_s \leq x$, where c_s is the confidence interval given by (5).*

D.1 Sample Complexity of MaxGapElim

Early Stopping Rule for Clustering: In the pseudocode in Algorithm 1, MaxGapElim stops when the size of the active set $|A| \leq 2$ (line 7). However, if we are only interested in clustering the arms according to the maximum gap and not interested in the identities of the arms which share the maximum gap (arms $(m), (m+1)$), we can stop earlier as follows. Assume that (9) is greater than 0 and let $(k_*)_t$ be the maximizer. This partitions the arms into the sets $\{(1)_t, \dots, (k_*)_t\}$ and $\{(k_*+1)_t, \dots, (K)_t\}$. MaxGapElim can terminate when the maximum left gap of all arms in $\{(1)_t, \dots, (k_*)_t\}$ and the maximum right gap of all arms in $\{(k_*+1)_t, \dots, (K)_t\}$ are both less than the lower bound $L\Delta(t)$. The termination condition can be expressed as $S = 1$, where

$$S = 1\{\mathcal{U}\Delta_a^r(t) < L\Delta(t), \forall a : l_a(t) \geq l_{(k_*)_t}(t) \cdot 1\{\mathcal{U}\Delta_a^l(t) < L\Delta(t), \forall a : r_a(t) \leq l_{(k_*+1)_t}(t)\}\}. \quad (21)$$

To account for the lower sample complexity as a result of the stopping rule for clustering, we modify (10) and (11) and define new parameters that yield an improved sample complexity than that stated in Theorem 2. Define

$$\rho_a^r = \max \left\{ \max_{j: \Delta_{a,j} > 0} (\min\{\Delta_{a,j}/4, ((\Delta_{\max} - \Delta_{a,j})/8)\}), ((\Delta_{\max} - \Delta_{a,1})/8) \right\}, \quad (22)$$

$$\rho_a^l = \max \left\{ \max_{j: \Delta_{a,j} < 0} (\min\{\Delta_{a,j}/4, ((\Delta_{\max} - \Delta_{j,a})/8)\}), ((\Delta_{\max} - \Delta_{a,K})/8) \right\}, \quad (23)$$

where just like in (10), the maxima assumed to be infinity if there is no j that satisfies the constraint under the inner maximization. We define $\rho_a = \min\{\rho_a^r, \rho_a^l\}$ as before and state our improved sample complexity bound for MaxGapElim next.

Theorem 4. *With probability at least $1 - \delta$, the sample complexity of MaxGapElim is bounded by*

$$H = \alpha \sum_{\substack{a \in [K]: \\ a \notin \{(m), (m+1)\}}} \frac{\log(K/\delta \rho_a)}{\rho_a^2}.$$

Proof. Arm a is eliminated in MaxGapElim when $U\Delta_a(t) < L\Delta(t)$, where $U\Delta_a(t)$ is defined as the maximum of the left and right gap upper bounds (see Section 4). Lemma 6 and Lemma 7 prove that the sufficient condition for each of these upper bounds to be less than $L\Delta(t)$ is $c_{T_a}(t) \leq \rho_a$. The result then follows by Remark 1. \square

Lemma 5. *If the good event (6) holds, then for all $a \in [K]$, for all $t \in \mathbb{N}$,*

$$l_a(t) \geq \mu_a - 2c_{T_a(t)} \text{ and } r_a(t) \leq \mu_a + 2c_{T_a(t)}$$

where $c_s = \sqrt{\frac{\beta_\delta(s)}{s}}$.

Proof. We have

$$\hat{\mu}_a(t) + c_{T_a(t)} \stackrel{(a)}{\geq} \mu_a \Rightarrow l_a(t) = \hat{\mu}_a(t) - c_{T_a(t)} \geq \mu_a - 2c_{T_a(t)}.$$

Similarly,

$$\hat{\mu}_a(t) - c_{T_a(t)} \stackrel{(a)}{\leq} \mu_a \Rightarrow r_a(t) = \hat{\mu}_a(t) + c_{T_a(t)} \leq \mu_a + 2c_{T_a(t)}.$$

In both the equations above, (a) holds by (6). \square

Lemma 6. *Assume (6) holds, and consider $a \neq m+1$. In MaxGapElim if t is such that $c_{T_a(t)} \leq \rho_a^r$, then*

$$\text{U}\Delta_a^r(t) < \text{L}\Delta(t).$$

Proof. Note that at time t in Algorithm 1, $T_a(t) = t$ and $c_{T_a(t)} = c_t$ for all arms $a \in A$. Assume (6) holds. We have $c_t < \rho_a^r < \Delta_{\max}/4$. This implies that

$$l_m(t) \stackrel{(a)}{\geq} \mu_m - 2c_t = \mu_{m+1} + \Delta_{\max} - 2c_t \stackrel{(a)}{\geq} r_{m+1}(t) + \Delta_{\max} - 4c_t \geq r_{m+1}(t). \quad (24)$$

where (a) holds by Lemma 5.

From (24) we have that

$$\text{L}\Delta(t) \geq l_m(t) - r_{m+1}(t) \geq \Delta_{\max} - 4c_t \quad (25)$$

Recall from (22) that for $a \neq 1$,

$$\rho_a^r = \max \left\{ \max_{j: \Delta_{a,j} > 0} (\min\{\Delta_{a,j}/4, ((\Delta_{\max} - \Delta_{a,j})/8)\}), ((\Delta_{\max} - \Delta_{a,1})/8) \right\}. \quad (26)$$

There are two terms in ρ_a^r and c_t could be less than either of these terms. First, suppose that

$$c_t < \max_{j: \Delta_{a,j} > 0} (\min\{\Delta_{a,j}/4, ((\Delta_{\max} - \Delta_{a,j})/8)\}),$$

and let

$$e = \arg \max_{j: \Delta_{a,j} > 0} (\min\{\Delta_{a,j}/4, ((\Delta_{\max} - \Delta_{a,j})/8)\}). \quad (27)$$

For any arm j such that $\Delta_{\max} < \Delta_{a,j}$, the inner minimum in (27) will be negative. On the other hand, since $a \neq m+1$, there must exist an arm j such that $\Delta_{\max} > \Delta_{a,j}$, and for such an arm j the inner minimum will be positive. Since e is the arm that maximizes the inner minimum, the inner minimum must be positive for e . Thus we have that $\Delta_{\max} > \Delta_{a,e}$.

From (26), (27), we have that

$$c_t < \Delta_{a,e}/4 \quad \text{and} \quad c_t < (\Delta_{\max} - \Delta_{a,e})/8. \quad (28)$$

Since $c_t < \Delta_{a,e}/4$, by following an argument similar to (24) we have that $l_e(t) \geq r_a(t)$, and hence the first branch of (7) will be used to compute $\text{U}\Delta_a^r(t)$. Hence we have

$$\text{U}\Delta_a^r(t) \stackrel{(a)}{\leq} r_e(t) - l_a(t) \stackrel{(b)}{\leq} \Delta_{a,e} + 4c_t \stackrel{(c)}{\leq} \Delta_{\max} - 4c_t \stackrel{(d)}{\leq} \text{L}\Delta(t)$$

where (a) follows from (7) and (8), (b) holds from Lemma 5, (c) follows by (28), and (d) holds by (25).

For the second case, assume

$$c_t < (\Delta_{\max} - \Delta_{a,1})/8.$$

Let $e = \arg \max_{i \neq a} r_i(t)$. From (8), we have that

$$\text{U}\Delta_a^r(t) \leq r_e(t) - l_a(t) \stackrel{(a)}{\leq} \Delta_{a,e} + 4c_t \leq \Delta_{a,1} + 4c_t \stackrel{(b)}{\leq} \Delta_{\max} - 4c_t \stackrel{(c)}{\leq} \text{L}\Delta(t),$$

where (a) holds by Lemma 5, (b) holds by the case assumption, and (c) holds by (25). \square

Lemma 7. Assume (6) holds, and consider $a \neq m$. In MaxGapElim if t is such that $c_t \leq \rho_a^l$, then

$$\mathbf{U}\Delta_a^l(t) < \mathbf{L}\Delta(t)$$

Proof. The proof is analogous to the proof of Lemma 6. \square

D.2 Sample Complexity of MaxGapUCB

For the sample complexity analysis of MaxGapUCB , we use a modified version of the left and right confidence bounds introduced in (5). We redefine

$$l'_i(t) \triangleq \max_{s \leq t} l_i(s), \quad r'_i(t) \triangleq \min_{s \leq t} r_i(s). \quad (29)$$

The nice property that these bounds have is that $[l'_i(t), r'_i(t)] \subseteq [l_i(s), r_i(s)]$ for all $t \geq s$. Lemma 9 shows that these modified bounds retain the same confidence guarantee for the arm mean values as the original confidence bounds. In what follows, we will exclusively use the modified confidence bounds (except in Lemma 9 where we show they are correct). We drop the prime symbol in their notation for brevity and henceforth $l_i(t), r_i(t)$ denote the modified confidence bounds given in (29).

We state and prove our main sample complexity result in Theorem 5.

Theorem 5. With probability at least $1 - \delta$, the number of times MaxGapUCB samples a sub-optimal arm, i.e. an arm $i \notin \{(m), (m+1)\}$, is upper bounded by $6\alpha\gamma_i^{-2} \log(K/\delta\gamma_i)$. The constant α is defined in Remark 1. Thus, the number of times MaxGapUCB samples suboptimal arms is

$$H = 6\alpha \sum_{\substack{i \in [K]: \\ i \notin \{(m), (m+1)\}}} \frac{\log(K/\delta\gamma_i)}{\gamma_i^2}.$$

Proof. We show that the result holds true as long as the confidence intervals for the means are correct (6). Let

$$\tau_r = \alpha \frac{\log(K/\delta\gamma_i^r)}{(\gamma_i^r)^2}, \quad \text{and} \quad \tau_l = \alpha \frac{\log(K/\delta\gamma_i^l)}{(\gamma_i^l)^2}, \quad (30)$$

where α is defined in Remark 1. Note that $\mathbf{U}\Delta_{(m)}(t) = \mathbf{U}\Delta_{(m+1)}(t) \geq \Delta_{\max} \forall t$. Arm i is sampled either because $\mathbf{U}\Delta_i^r$ is the largest or $\mathbf{U}\Delta_i^l$ is the largest. We prove in Lemma 8 below that when i is sampled $3\tau_r$ times due to its right gap, $\mathbf{U}\Delta_i^r < \Delta_{\max}$. Hence MaxGapUCB will not sample i due to its right gap more than $3\tau_r$ times because beyond this point $\mathbf{U}\Delta_{(m)}$ will be higher. It can similarly be proved that when i is sampled $3\tau_l$ times due to its left gap, $\mathbf{U}\Delta_i^l < \Delta_{\max}$. Thus, arm i will be sampled at most $3(\tau_r + \tau_l) \leq 6 \max\{\tau_r, \tau_l\}$ times. \square

To ease the explanation, we only focus on the right gap of i from here onwards and set

$$\tau = \alpha \frac{\log(K/\delta\gamma_i^r)}{(\gamma_i^r)^2}. \quad (31)$$

Furthermore, in the lemmas below, we only focus on samples of i drawn when $\mathbf{U}\Delta_i^r$ was the largest upper bound. With a slight overload of notation, let $t(i, s)$ denote the (random) smallest time when arm i has been sampled s times by MaxGapUCB (owing to its right gap).

Lemma 8. With probability $1 - \delta$, $\mathbf{U}\Delta_i^r(t(i, 3\tau)) < \Delta_{\max}$.

Proof. Since the proof is long and technical, we first give an outline of the entire proof.

Outline: Define i_l^t and i_r^t to be the arms that form the left and right boundaries of $\mathbf{U}\Delta_i^r(t)$ (the two arms that result in the maximum value of (8) at t). By this definition, $\mathbf{U}\Delta_i^r(t) = G_{i_l^t}^r(i_r^t(t), t) = r_{i_r^t}(t) - l_{i_l^t}(t)$. Consider the arms used in computing $\mathbf{U}\Delta_i^r(t(i, 3\tau))$, i.e. $i_l^{t(i, 3\tau)}$, $i_r^{t(i, 3\tau)}$, and denote them as i_l, i_r for brevity. Fig. 9 shows the confidence intervals of i_l in blue and those of i_r in green. Initially the confidence intervals are large, i.e., the width between the right and left bounds of arm i is greater than Δ_{\max} before time $t(i, \tau)$. After $t(i, 2\tau)$ rounds of MaxGapUCB , the confidence interval of i will have shrunk. However, note that the right gap of arm i involves either i_l and/or i_r . Since

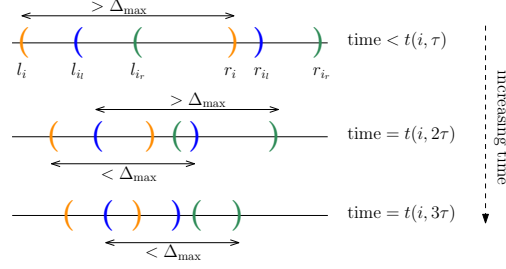


Figure 9: Illustration of left and right confidence bounds during a run of MaxGapUCB at three different times, the argument t for the bounds are omitted. Arms i_l, i_r are such that $\mathcal{U}\Delta_i^r(t(i, 3\tau)) = r_{i_r}(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau))$.

MaxGapUCB samples *all* arms that can attain the highest gap upper bound, it turns out that it will also sample i_l enough times to make $r_{i_l}(t(i, 2\tau)) - l_i(t(i, 2\tau)) < \Delta_{\max}$. If i is still sampled after $t(i, 2\tau)$ rounds due to its right gap, then its gap upper bound must involve an arm which is disjoint from i 's confidence interval, such as the arm i_r . Then from $t(i, 2\tau)$ to $t(i, 3\tau)$, MaxGapUCB samples i_l and i_r enough times to make $\mathcal{U}\Delta_i^r(t(i, 3\tau)) = r_{i_r}(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau)) < \Delta_{\max}$.

We divide the proof into four parts. In the first part, we divide all the arms into subsets (which we refer to as levels). These subsets are defined such that arms within a subset obey collective properties, that we study in some of the subsequent lemmas. In the second and third part, we prove that arms $i_r^{t(i, 3\tau)}$ and $i_l^{t(i, 3\tau)}$ are always sampled whenever i is sampled from $[t(i, 2\tau), t(i, 3\tau)]$. Finally in part four, we use these arms to argue that $\mathcal{U}\Delta_i^r(t(i, 3\tau)) < \Delta_{\max}$.

Level Sets:

At any time t , we can identify three subsets of arms with respect to arm i that we refer to as level 0, level 1, and level 2 arms respectively, and argue that the arms that define $\mathcal{U}\Delta_i^r(t)$ must lie in one of these subsets. These levels sets are defined as follows. Let

$$\mathcal{A}_i^0(t) = \{a \in [K] : l_i(t) \leq r_a(t) < r_i(t)\}, \quad (32)$$

$$\mathcal{A}_i^1(t) = \{a \in [K] : l_a(t) \leq r_i(t) \leq r_a(t)\}, \quad (33)$$

$$\mathcal{A}_i^2(t) = \left\{ a \in [K] : r_i(t) < l_a(t) \leq \min_{j: l_j(t) > r_i(t)} r_j(t) \right\}. \quad (34)$$

From their definitions the three subsets are pairwise disjoint at every $t \in \mathbb{N}$. Let

$$\mathcal{A}_i(t) = \mathcal{A}_i^0(t) \cup \mathcal{A}_i^1(t) \cup \mathcal{A}_i^2(t) \quad (35)$$

denote the union of the three levels. From the definition of $\mathcal{U}\Delta_i^r(t)$ in (8) and (32), (33), the arm

$$i_l^t \in \mathcal{A}_i^0(t) \cup \mathcal{A}_i^1(t) \forall t. \quad (36)$$

Lemma 10 proves that the arm $i_r^t \in \mathcal{A}_i(t) \forall t$. Thus at any time t , only arms in $\mathcal{A}_i(t)$ are relevant for the right gap of arm i .

Suppose $t(i, 3\tau) < \infty$, i.e., arm i is sampled at least 3τ times. To avoid clutter, we let

$$i_r = i_r^{t(i, 3\tau)} \quad \text{and} \quad i_l = i_l^{t(i, 3\tau)},$$

and use the full notation i_r^t for $t \neq t(i, 3\tau)$. We next argue that i_r and i_l must be sampled τ times before $t(i, 3\tau)$.

i_r must have been sampled at least τ times before $t(i, 3\tau)$:

By Lemma 10, $i_r \in \mathcal{A}_i(t(i, 3\tau))$. From Corollary 2, $i_r \in \mathcal{A}_i(s) \forall s \in [t(i, 2\tau), t(i, 3\tau)]$. If $i_r \in \mathcal{A}_i^0(s) \cup \mathcal{A}_i^1(s)$ for any $s \in [t(i, 2\tau), t(i, 3\tau)]$, then $r_{i_r}(t(i, 3\tau)) - l_i(t(i, 3\tau)) \leq r_{i_r}(s) - l_i(s) < \Delta_{\max}$ by Lemma 9 and Lemma 11, and we are done. Let us hence look at the case when $i_r \in \mathcal{A}_i^2(s) \forall s \in [t(i, 2\tau), t(i, 3\tau)]$. We have by Lemma 11 that $i_r^s \in \mathcal{A}_i^2(s) \forall s \in [t(i, 2\tau), t(i, 3\tau)]$, and Lemma 13 then implies that i_r must be sampled whenever i was sampled for $s \in [t(i, 2\tau), t(i, 3\tau)]$. Hence i_r is sampled at least τ times before $t(i, 3\tau)$.

i_l must have been sampled at least τ times before $t(i, 3\tau)$:

From Corollary 2, since the level of an arm cannot decrease from 2 to 1, $i_l \in \mathcal{A}_i^0(s) \cup \mathcal{A}_i^1(s)$ for all $s \in [t(i, 2\tau), t(i, 3\tau)]$. If $i_l \in \mathcal{A}_i^1(s) \forall s \in [t(i, 2\tau), t(i, 3\tau)]$, then i_l is sampled τ times whenever i is sampled by Lemma 13.

On the other hand, if $i_l \in \mathcal{A}_i^0(s)$ for some $s \in [t(i, 2\tau), t(i, 3\tau)]$, let $\text{U}\Delta_i^r(s) = r_{i_r^s}(s) - l_{i_l^s}(s)$. We consider two cases, $r_{i_l}(s) < l_{i_l^s}(s)$ and $r_{i_l}(s) \geq l_{i_l^s}(s)$. First, if $r_{i_l}(s) < l_{i_l^s}(s)$, then $l_{i_l}(s') < r_{i_l}(s') < l_{i_l^s}(s')$ for all $s' \geq s$ by Lemma 9. Since $i_l^s \in \mathcal{A}_i^0(s) \cup \mathcal{A}_i^1(s)$, we have by Lemma 9 and Lemma 11 that

$$r_{i_l^s}(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau)) \leq r_{i_l^s}(t(i, 3\tau)) - l_i(t(i, 3\tau)) \leq \Delta_{\max}. \quad (37)$$

By the definition of $\text{U}\Delta_i^r$ in (8) we have that

$$\text{U}\Delta_i^r(t(i, 3\tau)) = G_i^r(l_{i_l}(t(i, 3\tau)), t(i, 3\tau)) \leq r_{i_l^s}(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau)). \quad (38)$$

(38) and (37) imply that $\text{U}\Delta_i^r(t(i, 3\tau)) \leq \Delta_{\max}$, and we are done. For the second case, suppose $r_{i_l}(s) \geq l_{i_l^s}(s)$. Lemma 12 then gives that arm i_l is also sampled at time s . Thus, we have shown that either $\text{U}\Delta_i^r(t(i, 3\tau)) < \Delta_{\max}$, or i_l is sampled whenever i is sampled in $[t(i, 2\tau), t(i, 3\tau)]$.

We now show that $\text{U}\Delta_i^r(t(i, 3\tau)) < \Delta_{\max}$.

$\text{U}\Delta_i^r(t(i, 3\tau)) < \Delta_{\max}$:

Recall that $T_i(t(i, 3\tau)), T_{i_l}(t(i, 3\tau)), T_{i_r}(t(i, 3\tau))$ are all larger than τ . Let

$$j_* = \arg \max_{j: 0 < \Delta_{i,j} < \Delta_{\max}} \min\{\Delta_{i,j}/4, (\Delta_{\max} - \Delta_{i,j})/4\}$$

be the maximizer in (10), and note that $\mu_i < \mu_{j_*}$ by definition. Also note that τ and γ_i^r are defined in (31) and (10) respectively so that

$$4c_\tau \leq \Delta_{\max} - \Delta_{i,j_*} \quad \text{and} \quad 4c_\tau \leq \Delta_{i,j_*} \quad (39)$$

We split the proof into various cases depending on the ordering of the means $\mu_i, \mu_{i_l}, \mu_{i_r}, \mu_{j_*}$. First, note that if $\mu_{i_r} \leq \mu_{i_l}$, then

$$\text{U}\Delta_i^r(t(i, 3\tau)) = r_{i_r}(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau)) \leq \mu_{i_r} - \mu_{i_l} + 4c_\tau \leq \Delta_{\max}$$

by (39). Second, if $\max\{\mu_i, \mu_{i_l}\} < \mu_{i_r} < \mu_{j_*}$, then

$$\begin{aligned} \text{U}\Delta_i^r(t(i, 3\tau)) &\leq r_{i_r}(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau)) \leq r_{i_r}(t(i, 3\tau)) - l_i(t(i, 3\tau)) \\ &\leq \mu_{i_r} - \mu_i + 4c_\tau \leq \Delta_{i,j_*} + 4c_\tau \leq \Delta_{\max} \end{aligned}$$

by (39). Third, we show that it cannot be the case that $\mu_i < \mu_{j_*} < \mu_{i_l} < \mu_{i_r}$. Assume to the contrary. This implies that

$$l_{i_l}(t(i, 3\tau)) - r_i(t(i, 3\tau)) \geq \mu_{i_l} - \mu_i - 4c_\tau \geq \mu_{j_*} - \mu_i - 4c_\tau > 0,$$

which contradicts Eq. (36). Fourth, it cannot be the case that $\mu_{i_l} < \mu_{i_r} < \mu_i < \mu_{j_*}$, because $i_r \in \mathcal{A}_i^2(t(i, 3\tau))$ by Lemma 11. The only case that remains is $\max\{\mu_i, \mu_{i_l}\} < \mu_{j_*} < \mu_{i_r}$, which we prove next by showing that $T_{j_*}(t(i, 3\tau)) \geq \tau$.

$\max\{\mu_i, \mu_{i_l}\} < \mu_{j_*} < \mu_{i_r}$:

For any time $s \in [t(i, 2\tau), t(i, 3\tau)]$ such that $j_* \in \mathcal{A}_i^1(s) \cup \mathcal{A}_i^2(s)$, we have by Lemma 11 and Lemma 13 that j_* is sampled whenever i is sampled. Thus we only need to focus on times s when $j_* \in \mathcal{A}_i^0(s)$.

Suppose now that $j_* \in \mathcal{A}_i^0(s)$ for some $s \in [t(i, 2\tau), t(i, 3\tau)]$ when i was sampled and $\text{U}\Delta_i^r(s) = r_{i_r^s}(s) - l_{i_l^s}(s)$. Recall that $i_l^s \in \mathcal{A}_i^0(s) \cup \mathcal{A}_i^1(s)$. We consider two cases depending on whether $l_{i_l^s}(t(i, 3\tau)) > l_{i_l}(t(i, 3\tau))$ or $l_{i_l^s}(t(i, 3\tau)) \leq l_{i_l}(t(i, 3\tau))$.

- $l_{i_l^s}(t(i, 3\tau)) > l_{i_l}(t(i, 3\tau))$: We have

$$\begin{aligned} \text{U}\Delta_i^r(t(i, 3\tau)) &= G(l_{i_l}(t(i, 3\tau)), t(i, 3\tau)) \stackrel{(a)}{\leq} r_{i_l^s}(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau)) \\ &\leq r_{i_l^s}(t(i, 3\tau)) - l_i(t(i, 3\tau)) \stackrel{(b)}{\leq} r_{i_l^s}(s) - l_i(s) \stackrel{(c)}{\leq} \Delta_{\max}, \end{aligned}$$

where (a) holds by (7), (b) holds by Lemma 9, and (c) holds by Lemma 11.

- $l_{i_l^s}(t(i, 3\tau)) \leq l_{i_l}(t(i, 3\tau))$: Since $\max\{\mu_i, \mu_{i_l}\} < \mu_{j_*}$, we have $l_{i_l}(t) \leq r_{j_*}(t) \forall t$. Hence, $l_{i_l^s}(t(i, 3\tau)) < r_{j_*}(t(i, 3\tau))$, and Lemma 9 implies that $l_{i_l^s}(s) \leq r_{j_*}(s)$. Recall that s is a time such that $j_* \in \mathcal{A}_i^0(s)$, and hence

$$r_{j_*}(s) - l_{i_l^s}(s) \leq r_{j_*}(s) - l_i(s) \leq \Delta_{\max}.$$

Now, since i is sampled at time s , we have $\mathcal{U}\Delta_i^r(s) > \Delta_{\max}$, and (7) then implies that $l_{j_*}(s) < l_{i_l^s}(s)$. Hence by Lemma 12 MaxGapUCB must also sample arm j_* at time s .

This proves that $T_{j_*}^r(t(i, 3\tau)) \geq \tau$. We use this to prove that $\mathcal{U}\Delta_i^r(t(i, 3\tau)) < \Delta_{\max}$ as follows. First note that

$$l_{j_*}^r(t(i, 3\tau)) - r_i(t(i, 3\tau)) \geq p_{j_*} - p_i - 4c_\tau \geq 0.$$

Second, since arm $i_l \in \mathcal{A}_i^0(t(i, 3\tau)) \cup \mathcal{A}_i^1(t(i, 3\tau))$, and hence

$$l_{i_l}(t(i, 3\tau)) < r_i(t(i, 3\tau)) \leq l_{j_*}(t(i, 3\tau)).$$

Hence

$$\begin{aligned} \mathcal{U}\Delta_i^r(t(i, 3\tau)) &= G_i^r(l_{i_l}(t(i, 3\tau)), t(i, 3\tau)) \leq r_{j_*}^r(t(i, 3\tau)) - l_{i_l}(t(i, 3\tau)) \\ &\leq r_{j_*}^r(t(i, 3\tau)) - l_i(t(i, 3\tau)) \leq \mu_{j_*} - \mu_i + 4c_\tau \leq \Delta_{\max}. \end{aligned}$$

□

Lemma 9. *Over the sigma-algebra generated by all the arm rewards up till any time $t \in \mathbb{N}$, we have that*

$$\mathbb{P}\left(\forall t \in \mathbb{N}, \forall i \in [K], \mu_i \in \left[\max_{t' \leq t} l_i(t'), \min_{t' \leq t} r_i(t')\right]\right) = \mathbb{P}(\forall t \in \mathbb{N}, \forall i \in [K], \mu_i \in [l_i(t), r_i(t)]). \quad (40)$$

Proof. Let E' be the event in the LHS of (40) and let E be the good event. First we show that $E' \subseteq E$. The event E' implies that at any time t and for any arm i , we have that

$$\mu_i \in \left[\max_{t' \leq t} l_i(t'), \min_{t' \leq t} r_i(t')\right] \implies \mu_i \in [l_i(t'), r_i(t')] \forall t' \leq t.$$

Hence the good event is true in this case.

Now we show that $E \subseteq E'$. Suppose that E' is not true, so there is a time t and arm i such that $\mu_i \notin [\max_{t' \leq t} l_i(t'), \min_{t' \leq t} r_i(t')]$. Choose two time instants $s_l, s_r \in \mathbb{N}$ such that $s_l \in \arg \max_{t' < t} l_i(t')$, $s_r \in \arg \min_{t' < t} r_i(t')$. Then the supposition implies that either

$$\mu_i \notin [l_i(s_l), r_i(s_l)] \quad \text{or/and} \quad \mu_i \notin [l_i(s_r), r_i(s_r)].$$

Either of the above statements imply that the good event is not true. Hence $E \implies E'$. □

Corollary 1. *For any two time instants $s, t \in \mathbb{N}$ if $s < t$ then $\mathcal{U}\Delta_i^r(s) \geq \mathcal{U}\Delta_i^r(t)$.*

Proof. The quantity $\mathcal{U}\Delta_i^r(t)$ is defined in (15) as an optimization problem over a set of K real numbers $\mathcal{P}' = \{\mu'_i \in [l_i(t), r_i(t)] : i \in [K]\}$. For a time $s < t$, the $\mathcal{U}\Delta_i^r(s)$ is an optimization over $\mathcal{P}'' = \{\mu''_i \in [l_i(s), r_i(s)] : i \in [K]\}$. Lemma 9 states that $[l_i(t), r_i(t)] \subseteq [l_i(s), r_i(s)]$, hence we have that $\mathcal{U}\Delta_i^r(s) \geq \mathcal{U}\Delta_i^r(t)$. □

Corollary 2. *For all $k \in [K]$ if $k \in \mathcal{A}_i^2(t)$ then $k \in \mathcal{A}_i(s)$ at all time instants $s \leq t$. If $k \in \mathcal{A}_i^2(t)$ then $k \notin \mathcal{A}_i^0(s') \cup \mathcal{A}_i^1(s')$ at all $s' \geq t$.*

Proof. Define $\mathcal{J}(t) \triangleq \{j \in [K] : l_j(t) > r_i(t)\}$. For any $s \leq t$, if $j \in \mathcal{J}(s)$ then using Lemma 9,

$$l_j(t) \geq l_j(s) > r_i(s) \geq r_i(t) \implies j \in \mathcal{J}(t). \quad (41)$$

Hence if $k \in \mathcal{A}_i^2(t)$, from (34) we have that $l_k(t) \leq \min_{j \in \mathcal{J}(t)} r_j(t)$ and we get

$$l_k(s) \leq l_k(t) \leq \min_{j \in \mathcal{J}(t)} r_j(t) \stackrel{(a)}{\leq} \min_{j \in \mathcal{J}(t)} r_j(s) \stackrel{(b)}{\leq} \min_{j \in \mathcal{J}(s)} r_j(s),$$

where the inequality (a) is true because of Lemma 9 and inequality (b) is true as $\mathcal{J}(s) \subseteq \mathcal{J}(t)$ by (41). This implies that $k \in \mathcal{A}_i^0(s) \cup \mathcal{A}_i^1(s) \cup \mathcal{A}_i^2(s) = \mathcal{A}_i(s)$.

If $k \in \mathcal{A}_i^2(t)$ we have that $r_i(t) < l_k(t)$. At any $s' \geq t$, from Lemma 9 we have that $r_i(s') \leq r_i(t) < l_k(t) \leq l_k(s')$, i.e., the arm $k \notin \mathcal{A}_i^0(s') \cup \mathcal{A}_i^1(s')$. \square

Lemma 10. *The arms i_r^t, i_l^t are such that $\mathcal{U}\Delta_i^r(t) = r_{i_r^t}(t) - l_{i_l^t}(t)$. For the sets as defined in (32), (33), (34) the arm $i_r^t \in \mathcal{A}_i(t) \triangleq \mathcal{A}_i^0(t) \cup \mathcal{A}_i^1(t) \cup \mathcal{A}_i^2(t)$.*

Proof. Suppose arm $i_r^t \notin \mathcal{A}_i(t)$, then either $r_{i_r^t}(t) < l_i(t)$ which would give a negative value for $\mathcal{U}\Delta_i^r(t)$, or we have that $l_{i_r^t}(t) > \min_{a: l_a(t) > r_i(t)} r_a(t) \triangleq r_{a_*}(t)$. From the definition of arm i_l^t , its $l_{i_l^t}(t) \leq r_i(t)$. Using this and (7), we have that

$$G_i^r(l_{i_l^t}(t), t) = \min_{j: l_j(t) > l_{i_l^t}(t)} r_j(t) - l_{i_l^t}(t) \leq \min_{j: l_j(t) > r_i(t)} r_j(t) - l_{i_l^t}(t) = r_{a_*}(t) - l_{i_l^t}(t). \quad (42)$$

From the definition of arm i_r^t , we have that $\mathcal{U}\Delta_i^r(t) = r_{i_r^t}(t) - l_{i_l^t}(t) \leq r_{a_*}(t) - l_{i_l^t}(t)$ as argued above. That implies $r_{a_*}(t) \geq r_{i_r^t}(t) > l_{i_r^t}(t)$, which contradicts the supposition. \square

Lemma 11. *At any time $t \geq t(i, 2\tau)$, all arms $j \in \mathcal{A}_i^0(t) \cup \mathcal{A}_i^1(t)$ are such that $r_j(t) - l_i(t) \leq \Delta_{\max}$.*

Proof. Consider an arm $j \in \mathcal{A}_i^0(t) \cup \mathcal{A}_i^1(t)$, then $j \notin \mathcal{A}_i^2(s)$ for all $s \leq t$ for otherwise that would contradict corollary 2. Thus $j \in \mathcal{A}_i^0(s) \cup \mathcal{A}_i^1(s)$ for all $s \leq t$.

By choice of τ we have that $r_i(t(i, \tau)) - l_i(t(i, \tau)) = 2c_\tau \leq \Delta_{\max}$ from (39). Hence if arm $j \in \mathcal{A}_i^0(s)$ for any $s \in [t(i, \tau), t(i, 2\tau)]$, we have that $r_j(s) - l_i(s) \leq r_i(s) - l_i(s) \stackrel{(a)}{\leq} r_i(t(i, \tau)) - l_i(t(i, \tau)) \leq \Delta_{\max}$, where inequality (a) is by Lemma 9.

Hence $j \in \mathcal{A}_i^1(s)$ for all $s \in [t(i, \tau), t(i, 2\tau)]$. If $\mathcal{U}\Delta_i^r(s)$ is the largest gap upper bound then $i_R^s \notin \mathcal{A}_i^0(s)$ by the above reasoning. Then Lemma 13 states that arm j was sampled anytime arm i was sampled between $t(i, \tau)$ to $t(i, 2\tau)$. This implies that $T_j(t(i, 2\tau)) \geq \tau$, and we argue that $r_j(t(i, 2\tau)) - l_i(t(i, 2\tau)) \leq \Delta_{\max}$ in the following manner. The arm j_* is the maximizer in (10).

Case I: $\mu_i < \mu_{j_*} < \mu_j$. Here we argue that $j \notin \mathcal{A}_i^1(t(i, 2\tau))$ because $l_j(t(i, 2\tau)) \geq r_i(t(i, 2\tau))$ as shown below.

$$\begin{aligned} l_j(t(i, 2\tau)) - r_i(t(i, 2\tau)) &\geq \mu_j - 2c_{T_j(t(i, 2\tau))} - (\mu_i + 2c_{T_i(t(i, 2\tau))}) \quad (\text{Lemma 5}) \\ &\geq \mu_{j_*} - \mu_i - 4c_\tau \quad (\text{Assumption on means and monotonicity of } c(s)) \\ &\geq \mu_{j_*} - \mu_i - \Delta_{i, j_*} = 0. \quad (\text{Using (39)}) \end{aligned}$$

Case II: $\max\{\mu_i, \mu_j\} < \mu_{j_*}$. Here we argue that $r_j(t(i, 2\tau)) - l_i(t(i, 2\tau)) \leq \Delta_{\max}$ as shown below.

$$\begin{aligned} r_j(t(i, 2\tau)) - l_i(t(i, 2\tau)) &\leq \mu_j + 2c_{T_j(t(i, 2\tau))} - (\mu_i - 2c_{T_i(t(i, 2\tau))}) \quad (\text{Lemma 5}) \\ &\leq \mu_{j_*} - \mu_i + 4c_\tau \quad (\text{Assumption on means and monotonicity of } c(s)) \\ &\leq \mu_{j_*} - \mu_i + \Delta_{\max} - \Delta_{i, j_*} \leq \Delta_{\max}. \quad (\text{Using (39)}) \end{aligned}$$

\square

Lemma 12. *Suppose arm i is sampled at time t because $\mathcal{U}\Delta_i^r(t) = r_{i_r^t}(t) - l_{i_l^t}(t)$ is the largest gap upper bound. Consider an arm j whose confidence bounds satisfy any one of the following conditions.*

1. $l_j(t) < l_{i_l^t}(t) < r_j(t)$, or
2. $l_j(t) < r_{i_r^t}(t) < r_j(t)$.

Then MaxGapUCB samples arm j as well at time t .

Proof. Suppose arm j satisfies condition (1). Consider the right gap of arm j , we have that $\mathcal{U}\Delta_j^r(t) \geq G_j^r(l_{i_l^t}(t), t)$. If the value of $G_j^r(l_{i_l^t}(t), t)$ is obtained by the first branch of (7), then the value of $G_j^r(l_{i_l^t}(t), t)$ is also given by its first branch. That implies $\mathcal{U}\Delta_j^r(t) = \mathcal{U}\Delta_i^r(t)$, and hence j is sampled if i is sampled. If $j = i_r^t$, by condition (1) we have that $l_{i_r^t}(t) < l_{i_l^t}(t)$, which implies that $G_j^r(l_{i_l^t}(t), t)$ is obtained by the second branch in (7). Hence for all arms $a \neq i_l^t$ we have $l_a(t) < l_{i_l^t}(t)$ and $r_{i_r^t}(t) = r_j(t) = \max_{a \neq i} r_a(t)$. Considering the left gap of arm j , since $\{a : r_a(t) < r_j(t)\} \neq \emptyset$,

$$G_j^l(r_j(t), t) = r_j(t) - \max_{a: r_a(t) < r_j(t)} l_a(t) = r_{i_r^t}(t) - l_{i_l^t}(t) = \mathcal{U}\Delta_i^r(t),$$

and arm j is sampled if i is sampled. Finally suppose the value of $G_i^r(l_{i_l^t}(t), t)$ is obtained by the second branch in (7), and $j \neq i_r^t$. Then

$$\begin{aligned} G_i^r(l_{i_l^t}(t), t) &= \max_{a \neq i} r_a(t) - l_{i_l^t}(t) = r_{i_r^t}(t) - l_{i_l^t}(t), \\ G_j^r(l_{i_l^t}(t), t) &= \max_{a \neq j} r_a(t) - l_{i_l^t}(t) = \max\{r_i(t), r_{i_r^t}(t)\} - l_{i_l^t}(t) = r_{i_r^t}(t) - l_{i_l^t}(t), \end{aligned}$$

where the last equality is true because if not, then $\mathcal{U}\Delta_j^r(t) \geq G_j^r(l_{i_l^t}(t), t) > G_i^r(l_{i_l^t}(t), t) = \mathcal{U}\Delta_i^r(t)$, which contradicts the condition that $\mathcal{U}\Delta_i^r(t)$ is the largest. Hence arm j is sampled if i is sampled.

Suppose now that arm j satisfies condition (2). We divide the proof of this part into two cases.

Case I: Suppose $r_{i_r^t}(t) > r_{i_l^t}(t)$.

If the arm $i_l^t \neq i$, then we show that $G_i^r(l_{i_l^t}(t), t)$ cannot be the largest gap upper bound. Consider the arm $a_* \triangleq \arg \max_{a: l_a(t) < l_{i_l^t}(t)} l_a(t)$, it satisfies $l_i(t) \leq l_{a_*}(t) < l_{i_l^t}(t)$. Then $G_i^r(l_{a_*}(t), t) = \min_{a: l_a(t) > l_{a_*}(t)} r_a(t) - l_{a_*}(t)$, where the first branch of (7) is active because of arm i_l^t . But

$$\min_{a: l_a(t) > l_{a_*}(t)} r_a(t) = \min\{r_{i_l^t}(t), \min_{a: l_a(t) > l_{i_l^t}(t)} r_a(t)\} = \min\{r_{i_l^t}(t), r_{i_r^t}(t)\} = r_{i_r^t}(t).$$

That would imply

$$G_i^r(l_{a_*}(t), t) = r_{i_r^t}(t) - l_{a_*}(t) > r_{i_r^t}(t) - l_{i_r^t}(t) = G_i^r(l_{i_r^t}(t), t),$$

which contradicts the identification of arm i_l^t as the one giving the value of $\mathcal{U}\Delta_i^r(t)$. The case that remains is if the arm $i = i_l^t$. For this part consider the following two sub-cases:

Sub-case Ia: The set of arms $\{a : r_a(t) < r_{i_r^t}(t)\} = \emptyset$. Since the number of arms $K > 2$, the value $\max_{a \neq i} r_a(t) > r_{i_r^t}(t)$, and hence if $G_i^r(l_{i_l^t}(t), t) = r_{i_r^t}(t) - l_{i_l^t}(t)$, then it must be due to the first branch in (7). That implies $l_{i_r^t}(t) > l_{i_l^t}(t) = l_i(t)$. Then consider the left gap for arm j that satisfies condition (2). Since the set $\{a : r_a(t) < r_{i_r^t}(t)\} = \emptyset$, we have

$$G_j^l(r_{i_r^t}(t), t) = r_{i_r^t}(t) - \min_{a \neq j} l_a(t) = r_{i_r^t}(t) - l_{i_l^t}(t) = \mathcal{U}\Delta_i^r(t),$$

which implies that arm j will be sampled if $\mathcal{U}\Delta_i^r(t)$ is the largest.

Sub-case Ib: The set of arms $\{a : r_a(t) < r_{i_r^t}(t)\} \neq \emptyset$. Consider the arm $a_* \triangleq \arg \max_{a: l_a(t) < l_{i_l^t}(t)} l_a(t)$, the domain in the maximization is not empty because of the following. By the case assumption, there is an arm a whose $r_a(t) < r_{i_r^t}(t)$. If the left bound of this arm $l_a(t) > l_i(t)$, then $l_i(t) < l_a(t) < r_a(t) < r_{i_r^t}(t)$, which contradicts the identification of arm i_r^t for $G_i^r(l_{i_l^t}(t), t)$. Hence its left bound must satisfy $l_a(t) < l_i(t)$. Now consider the right gap of arm a_* defined above. Since $l_i(t) > l_{a_*}(t)$, we have that $G_{a_*}^r(l_{a_*}(t), t) = \min_{a: l_a(t) > l_{a_*}(t)} r_a(t) - l_{a_*}(t)$. But

$$\min_{a: l_a(t) > l_{a_*}(t)} r_a(t) = \min\{r_i(t), \min_{a: l_a(t) > l_i(t)} r_a(t)\} = \min\{r_i(t), r_{i_r^t}(t)\} = r_{i_r^t}(t),$$

which implies that $G_{a_*}^r(l_{a_*}(t), t) = r_{i_r^t}(t) - l_{a_*}(t) > r_{i_r^t}(t) - l_i(t) = \mathcal{U}\Delta_i^r(t)$, which is a contradiction. We are left with the following Case II.

Case II: Suppose $r_{i_l^t}(t) < r_{i_r^t}(t)$.

Let a_* be such that $l_{a_*}(t) \triangleq \max_{a: r_a(t) < r_{i_r^t}(t)} l_a(t)$. Then $l_{a_*}(t) \geq l_{i_r^t}(t)$. If the previous inequality is strict, then we have that

$$l_{i_r^t}(t) < l_{a_*}(t) < r_{a_*}(t) < r_{i_r^t}(t),$$

which contradicts the identification of arm i_r^t as the one giving the value of $G_i^r(l_{i_r^t}(t), t)$. Hence we have that

$$G_j^l(r_{i_r^t}(t), t) = r_{i_r^t}(t) - \max_{a: r_a(t) < r_{i_r^t}(t)} l_a(t) = r_{i_r^t}(t) - l_{i_r^t}(t) = \mathcal{U}\Delta_i^r(t),$$

and arm j is sampled if arm i is sampled because of $\mathcal{U}\Delta_i^r(t)$. \square

Lemma 13. Suppose arm i is sampled at time t when $\mathcal{U}\Delta_i^r(t) = r_{i_r^t}(t) - l_{i_r^t}(t)$. If $i_r^t \in \mathcal{A}_i^2(t)$ then all arms in the set $\mathcal{A}_i^1(t) \cup \mathcal{A}_i^2(t)$ are sampled by *MaxGapUCB*. If $i_r^t \in \mathcal{A}_i^1(t)$ then all arms in the set $\mathcal{A}_i^1(t)$ are sampled by *MaxGapUCB*.

Proof. The qualifying condition states that the arm $i_r^t \in \mathcal{A}_i^1(t) \cup \mathcal{A}_i^2(t)$, hence from definitions (33), (34) we have that $r_{i_r^t}(t) \geq r_i(t)$. By definition (8) the arm i_r^t is such that $l_{i_r^t}(t) \in [l_i(t), r_i(t)]$. We first argue that all arms in the set $\mathcal{A}_i^1(t)$ are sampled. For arm $j \in \mathcal{A}_i^1(t)$, $r_i(t) \leq r_j(t)$. If $l_j(t) \leq l_{i_r^t}(t)$, arm j satisfies condition (1) of Lemma 12 and hence is sampled if $\mathcal{U}\Delta_i^r(t)$ is the largest. If on the other hand $r_j(t) \geq r_{i_r^t}(t)$, then arm j satisfies condition (2) of Lemma 12 and hence it is sampled if i is sampled. The remaining case is if $l_{i_r^t}(t) < l_j(t) < r_j(t) < r_{i_r^t}(t)$, but that would contradict the identification of the arm i_r^t for $\mathcal{U}\Delta_i^r(t)$.

Now suppose arm $i_r^t \in \mathcal{A}_i^2(t)$, what is left to prove is that all arms in the set $\mathcal{A}_i^2(t)$ are sampled. Since $i_r^t \in \{a : l_a(t) > r_i(t) > l_{i_r^t}(t)\}$, we have that

$$G_i^r(l_{i_r^t}(t), t) = \min_{j: l_j(t) > l_{i_r^t}(t)} r_j(t) - l_{i_r^t}(t) = r_{i_r^t}(t) - l_{i_r^t}(t) = \min_{j \in \mathcal{A}_i^2(t)} r_j(t) - l_{i_r^t}(t),$$

where the last equality is true because arm $i_r^t \in \mathcal{A}_i^2(t)$ satisfies $l_{i_r^t}(t) > r_i(t) \geq l_{i_r^t}(t)$. From definition (34), any $j \in \mathcal{A}_i^2(t)$ is such that $l_j(t) \leq r_{i_r^t}(t)$ and satisfies condition (2) of Lemma 12. Hence arm j is sampled if i is sampled because of its right gap. \square

E Details for Section 6: Proof of Lemma 1

Lemma 1. Consider a model \mathcal{B} with $K = 4$ normal distributions $\mathcal{P}_i = \mathcal{N}(\mu_i, 1)$, where

$$\mu_4 = 0, \quad \mu_3 = \epsilon, \quad \mu_2 = \nu + 2\epsilon, \quad \mu_1 = 2\nu + 2\epsilon,$$

for some $\nu \gg \epsilon > 0$. Then any algorithm that is correct with probability at least $1 - \delta$ must collect $\Omega(1/\epsilon^2)$ samples of arm 4 in expectation.

Proof. The maximum gap in \mathcal{B} is $\Delta_{\max} = \Delta_{3,2} = \nu + \epsilon$. Define an alternate bandit model \mathcal{B}' with 4 normal distributions $\mathcal{P}'_i = \mathcal{N}(\mu'_i, 1)$ where

$$\mu'_i = \mu_i \quad \forall i \neq 4, \quad \mu'_4 = 2.1\epsilon.$$

Note that the ordering of the means in \mathcal{B}' does not follow the subscript indices, indeed $\mu'_3 < \mu'_4$.

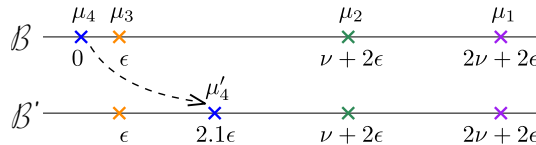


Figure 10: Changing the original bandit model \mathcal{B} to \mathcal{B}' . μ_4 is shifted to the right by 2.1ϵ . As a result, the maximum gap in \mathcal{B}' is between green and purple.

The two measures are illustrated in Fig. 10. The maximum gap in \mathcal{B}' is $\Delta'_{\max} = \Delta'_{2,1} = \nu$ and $\Delta'_{3,2}$ is no longer a valid gap between consecutive arms. Consider algorithm for identifying the maximum gap and let \hat{C}_1 denote the top-cluster returned by the algorithm when it stops at time τ . Let

$E = \{\widehat{C}_1 = \{1, 2\}\}$. Assume that $\mathbb{P}_{\mathcal{B}}(E) \geq 1 - \delta$ and $\mathbb{P}_{\mathcal{B}'}(E) \leq \delta$. Letting $d(\cdot)$ denote the binary relative entropy, Lemma 1 in Garivier and Kaufmann [10] implies that

$$\begin{aligned} \sum_{a=1}^4 \mathbb{E}_{\mathcal{B}}[T_a(\tau)] \text{KL}(\mathcal{P}_a, \mathcal{P}'_a) &\geq d(\mathbb{P}_{\mathcal{B}}(E), \mathbb{P}_{\mathcal{B}'}(E)) \geq d(1 - \delta, \delta) \\ \implies \mathbb{E}_{\mathcal{B}}[T_4(\tau)](\mu_4 - \mu'_4)^2 &\geq \log \frac{1}{2.4\delta} \implies \mathbb{E}_{\mathcal{B}}[T_4(\tau)] \geq \frac{1}{(2.1\epsilon)^2} \log \frac{1}{2.4\delta}. \end{aligned}$$

Similarly, one can show that $\mathbb{E}_{\mathcal{B}}[T_1(\tau)] \geq 1/\epsilon^2$ by creating an alternative bandit instance \mathcal{B}'' identical to \mathcal{B} except $\mu''_1 = 2\nu + 3.1\epsilon$. \square