

---

# Pure Exploration of Combinatorial Bandits

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

We study the structured pure exploration problem in the stochastic multi-armed bandit setting. In particular, we formulate the ExpCMAB problem where a learner’s objective is to identify the optimal set of arms from a collection of sets of arms which we called decision class. The decision class could be a collection of size  $k$  subsets, matchings, spanning trees or paths. The generality of decision classes allows ExpCMAB to represent a rich class of structured pure exploration problems. We present two algorithms for the ExpCMAB problem: one for the fixed confidence setting and one for the fixed budget setting. We prove problem-dependent upper bounds of our algorithms. Our analysis exploits the combinatorial structures of the decision classes and introduces a new analysis tool. We also establish a general problem-dependent lower bound for the ExpCMAB problem. Our results show that the proposed algorithms achieve optimal sample complexity (within logarithmic factors) for many decision classes. In addition, when applying our results back to top- $K$  arms identification and multiple bandit best arms identification, we recover the best known upper bounds and settles two open conjectures on the lower bounds.

## 1 Introduction

For more than fifty years, multi-armed bandit (MAB) has been a predominant model for characterizing the tradeoff between exploration and exploitation in decision-making problems. Although this kind of tradeoff is intrinsic in many tasks, some application domains prefer a dedicated exploration procedure in which the goal is to identify an optimal object among several candidates and the profit/cost incurred during exploration is irrelevant. In light of these applications, the related learning problem, called pure exploration in MABs, has received much attention. Recent results of pure exploration MABs have found potential applications in many domains including crowdsourcing, communication network and online advertising.

In many of these application domains, a recurring problem is to identify an optimal object with certain *combinatorial structures*. For example, a crowdsourcing application may want to find the best assignment from workers to tasks such that overall productivity of workers are maximized. A network routing system during an initialize phase may tries to build a spanning tree that minimizes the delay of links, or attempts to identify the shortest path between two sites. An online advertising system may be interested to find the best matching between ads and displaying slots. The literature of pure exploration MAB problems lacks a framework that encompasses these kinds of problems where the object of interest has a non-trivial combinatorial structure. Our paper contributes such a framework which accounts for general combinatorial structures, and develops a series of results, including algorithms, upper bounds and lower bounds.

In this paper, we formulate the Pure Exploration in Combinatorial Bandits (ExpCMAB) problem. In an ExpCMAB problem, a learner has a fixed set of arms and each arm is associated with an unknown reward distribution. The learner is also given a collection of sets of arms called *decision class*. The objective is to find the optimal set of arms, which maximizes the sum of expected reward, from the given decision class. Over a sequence of rounds, the learner chooses an arm and observes a random reward sampled from the associated distribution. In the end, the learner is asked to output a member of the decision class which she believes to be optimal.

The ExpCMAB framework represents a rich class of pure exploration problems. The conventional pure exploration problem in MAB, where the objective is to find the single best arm, clearly fits into this framework, in which the decision class is the collection of all singletons. This framework also naturally encompasses several recent extensions, including the problem of finding the top  $K$  arms (henceforth, TOPK) [15, 16, 7, 21] and the problem of finding the best arms simultaneously from several disjoint sets of arms (henceforth, MB) [10, 7], by constructing proper decision classes. Furthermore, this framework also covers many more interesting cases where the decision classes correspond to non-trivial combinatorial structures. For example, suppose that the arms represent the edges in a graph. Then a decision class could be the set of all paths between two vertices, all spanning trees or all matchings of the graph. And, in these cases, the objectives of ExpCMAB correspond to identifying the optimal paths, spanning trees and matchings, respectively. As we have explained earlier, these kind of structured pure exploration problems admit applications in diverse domains.

The generality of ExpCMAB framework raises several interesting challenges to the design and analysis of pure exploration algorithms. One challenge is that the arms with the largest mean rewards may not belong to the optimal set. For example, consider the case where the decision class is the set of all matchings in a bipartite graph. In this case, a matching consisting of edges with relatively small weights may turn out to be optimal. On the other hand, in many existing algorithms for pure exploration MABs, arms would no longer be considered once their expected reward are proven sub-optimal during the learning process. Therefore, the design and analysis of algorithms for ExpCMAB demands different techniques which take both rewards and structures into account.

**Our results.** We present Combinatorial Gap Exploration (CGapExp) algorithm, a learning algorithm for the ExpCMAB problem in the fixed confidence setting that supports a very wide range of decision classes. The proposed CGapExp algorithm does not need to know anything about the details of the decision class's definition, as long as it has access to the decision class through a maximization oracle. We prove a sample complexity bound of CGapExp. The sample complexity bound depends on both expected reward and the structure of decision class. When specializing our result into TOPK and MB, we recover previous sample complexity bounds due to Kalyanakrishnan et al. [16] and Gabillon et al. [11]. While for other decision classes in general, our result establishes the first sample complexity upper bound. Our analysis relies on a novel combinatorial construction called *exchange class* which we believed may be of independent interest for other combinatorial optimization problems. We further show that CGapExp can be easily extended to the fixed budget setting and PAC learning setting and we provide related theoretical guarantees.

Moreover, we prove a problem-dependent sample complexity lower bound for the ExpCMAB problem. Our lower bound shows that the sample complexity of the proposed CGapExp algorithm is optimal (to within logarithmic factors) for a large class of decision classes, including TOPK, MB and the decision classes derived from matroids (e.g. spanning tree). Therefore our upper and lower bounds provide a near full characterization of the sample complexity of these pure exploration problems. In addition, for general decision classes, our results find out the upper and lower bounds are within a relatively benign factor. To the best of our knowledge, there are few problem-dependent lower bounds known for pure exploration MABs besides the case of identifying the single best arm [18, 2]. We also notice that our result resolves the open conjectures of Kalyanakrishnan et al. [16] and Bubeck et al. [7] on the sample complexity lower bounds of TOPK and MB problem.

In supplement to our main algorithm CGapExp, we develop a parameter-free algorithm called Combinatorial Gap-based Elimination (CGapKill) algorithm, for the fixed budget setting. We prove a probability of error bound of the CGapKill algorithm. This bound can be shown to be equivalent to the sample complexity bound of CGapExp within logarithmic factors, although the two algorithms are based on quite different techniques. Our analysis of CGapKill re-uses exchange classes as tools. This suggests that exchange class may be useful for the analysis of similar problems. In addition, when applying the algorithm to TOPK and MB, our bound recovers a recent result due to Bubeck et al. [7].

**Notations.**

## 2 Problem Formulation

In this section, we formally define the ExpCMAB problem. Suppose that there are  $n$  arms and the arms are numbered  $1, 2, \dots, n$ . Assume that each arm  $e \in [n]$  is associated with a reward distribution  $\varphi_e$ . Let  $\mathbf{w} = (w(1), \dots, w(n))^T$  denote the vector of expected rewards, where each entry  $w(e) = \mathbb{E}_{X \sim \varphi_e}[X]$  denote the expected reward of arm  $e$ . Following standard assumptions of stochastic MABs, we assume that all reward distributions are  $R$ -sub-Gaussian for some constant  $R > 0$ . Formally, if  $X$  is a random variable drawn according to  $\varphi_e$ , then, for all  $t \in \mathbb{R}$ , one has  $\mathbb{E}[\exp(tX - t\mathbb{E}[X])] \leq \exp(R^2 t^2 / 2)$  and  $\mathbb{E}[\exp(t\mathbb{E}[X] - tX)] \leq \exp(R^2 t^2 / 2)$ . It is well known that all distributions that are supported on  $[0, R]$  satisfy this property [].

We define a *decision class*  $\mathcal{M} \subseteq 2^{[n]}$  as a collection of sets of arms. Let  $M_* = \arg \max_{M \in \mathcal{M}} w(M)$  denote the optimal set belonging to the decision class  $\mathcal{M}$  which maximizes the sum of expected reward<sup>1</sup>. A learner's objective is to identify  $M_*$  from  $\mathcal{M}$  by playing the following game. At the beginning of the game, the decision class  $\mathcal{M}$  is revealed to the learner while the reward distributions  $\{\varphi_e\}_{e \in [n]}$  are unknown to the learner. Then, the learner plays the game over a sequence of rounds; on each round  $t$ , the learner pulls an arm  $p_t \in [n]$  and observes a reward sampled from the associated reward distribution  $\varphi_{p_t}$ . The game continues until certain stopping condition is satisfied. After the game finishes, the learner need to output a set  $\text{Out} \in \mathcal{M}$ .

We consider two different stopping conditions of the game, which are known as *fixed confidence* setting and *fixed budget* setting. In the fixed confidence setting, the learner can stop the game at any point. The learner need to guarantee that  $\Pr[\text{Out} = M_*] \geq 1 - \delta$  for a given confidence parameter  $\delta$ . The learner's performance is evaluated by her *sample complexity*, i.e. the number of pulls used by the learner. In the fixed budget setting, the game stops after a fixed number  $T$  of rounds, where  $T$  is given before the game. The learner tries to minimize the *probability of error*, which is formally  $\Pr[\text{Out} \neq M_*]$ , within  $T$  rounds. In this case, the learner's performance is measured by the probability of error.

## 3 Algorithm, Exchange Class and Sample Complexity

In this section, we present CGapExp, a learning algorithm for the ExpCMAB problem in the fixed confidence setting, and analyze its sample complexity. En route to our sample complexity bound, we introduce the notions of exchange class and width of combinatorial problems, which characterize the exchange properties of combinatorial structures.

The CGapExp algorithm can be extended to the fixed budget and PAC learning settings. We will discuss these extensions in Section B.

**Oracle.** We allow the CGapExp algorithm to access a *maximization oracle*. A maximization oracle takes a weight vector  $\mathbf{v} \in \mathbb{R}^n$  as input and finds an optimal set within  $\mathcal{M}$  with respect to the weight vector  $\mathbf{v}$ . Formally, we call a function  $\text{Oracle}: \mathbb{R}^n \rightarrow \mathcal{M}$  a maximization oracle if, for all  $\mathbf{v} \in \mathbb{R}^n$ , we have  $\text{Oracle}(\mathbf{v}) \in \arg \max_{M \in \mathcal{M}} \mathbf{v}(M)$ . It is clear that a broad class of decision classes admit such maximization oracles, including all example decision classes we considered in this paper. Besides the access to the oracle, CGapExp does not need *any* additional knowledge of the decision class  $\mathcal{M}$ .

**Algorithm.** The CGapExp algorithm maintains empirical mean  $\bar{w}_t(e)$  and confidence radius  $\text{rad}_t(e)$  for each arm  $e \in [n]$  and each round  $t$ . The construction of confidence radius ensures that  $|w(e) - \bar{w}_t(e)| \leq \text{rad}_t(e)$  holds with high probability for each arm  $e \in [n]$  and each round  $t > 0$ . CGapExp begins with an initialization phase in which each arm is pulled once. Then, at round  $t \geq n$ , CGapExp uses the following procedure to choose an arm to play. First, CGapExp calls the oracle which finds the set  $M_t = \text{Oracle}(\bar{\mathbf{w}}_t)$ . The set  $M_t$  is the “best” set with respect to the empirical means  $\bar{\mathbf{w}}_t$ . Then, CGapExp explores possible refinements of  $M_t$ . In particular, CGapExp uses the confidence radius to compute an adjusted expectation vector  $\tilde{\mathbf{w}}_t$  in the following way: for each arm  $e \in M_t$ ,  $\tilde{w}_t(e)$  equals to the lower confidence bound  $\tilde{w}_t(e) = \bar{w}_t(e) - \text{rad}_t(e)$ ; and for each arm  $e \notin M_t$ ,  $\tilde{w}_t(e)$  equals to the upper confidence bound  $\tilde{w}_t(e) = \bar{w}_t(e) + \text{rad}_t(e)$ . Intuitively, the adjusted expectation vector  $\tilde{\mathbf{w}}_t$  penalizes arms belonging to the current set  $M_t$  and encourages exploring arms out of  $M_t$ . CGapExp then calls the oracle using the adjusted expectation vector  $\tilde{\mathbf{w}}_t$  as input to compute a

<sup>1</sup>For convenience, we will assume that  $M_*$  is unique throughout the paper.

refined set  $\tilde{M}_t = \text{Oracle}(\tilde{\mathbf{w}}_t)$ . If  $\tilde{w}_t(\tilde{M}_t) = \tilde{w}_t(M_t)$  then CGapExp stops and returns  $\text{Out} = M_t$ . Otherwise, CGapExp pulls the arm belonging to the symmetric difference  $(\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)$  between  $M_t$  and  $\tilde{M}_t$  with the largest confidence radius in the end of round  $t$ . The pseudo-code of CGapExp is shown in Algorithm 1.

---

**Algorithm 1** CGapExp: Combinatorial Gap Exploration

---

**Require:** Confidence parameter:  $\delta \in (0, 1)$ ; Maximization oracle:  $\text{Oracle}(\cdot) : \mathbb{R}^n \rightarrow \mathcal{M}$ .  
**Initialize:** Play each arm  $e \in [n]$  once. Initialize empirical means  $\bar{\mathbf{w}}_n$  and set  $T_n(e) \leftarrow 1$  for all  $e$ .

```

1: for  $t = n, n + 1, \dots$  do
2:    $M_t \leftarrow \text{Oracle}(\bar{\mathbf{w}}_t)$ 
3:   for  $e = 1, \dots, n$  do
4:     if  $e \in M_t$  then
5:        $\tilde{w}_t(e) \leftarrow \bar{w}_t(e) - \text{rad}_t(e)$ 
6:     else
7:        $\tilde{w}_t(e) \leftarrow \bar{w}_t(e) + \text{rad}_t(e)$ 
8:     end if
9:   end for
10:   $\tilde{M}_t \leftarrow \text{Oracle}(\tilde{\mathbf{w}}_t)$ 
11:  if  $\tilde{w}_t(\tilde{M}_t) = \tilde{w}_t(M_t)$  then
12:     $\text{Out} \leftarrow M_t$ 
13:    return  $\text{Out}$ 
14:  end if
15:   $p_t \leftarrow \arg \max_{e \in (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(e)$ 
16:  Pull arm  $p_t$  and observe the reward
17:  Update empirical means  $\bar{\mathbf{w}}_{t+1}$  using the observed reward
18:  Update number of pulls:  $T_{t+1}(p_t) \leftarrow T_t(p_t) + 1$  and  $T_{t+1}(e) \leftarrow T_t(e)$  for all  $e \neq p_t$ 
19: end for

```

---

### 3.1 Analysis

Now we prove a problem-dependent sample complexity bound of the CGapExp algorithm. Our sample complexity bound depends on the combinatorial properties of  $\mathcal{M}$ . Therefore, to formally state our result, we need to introduce several definitions.

**Gap.** We begin with defining a natural hardness measure of the ExpCMAB problem. For each arm  $e \in [n]$ , we define its gap  $\Delta_e$  as

$$\Delta_e = \begin{cases} w(M_*) - \max_{M \in \mathcal{M}: e \in M} w(M) & \text{if } e \notin M_*, \\ w(M_*) - \max_{M \in \mathcal{M}: e \notin M} w(M) & \text{if } e \in M_*, \end{cases} \quad (1)$$

where we use the convention that the maximum value of an empty set is  $-\infty$ . We also define the hardness  $\mathbf{H}$  as the sum of inverse squared gaps

$$\mathbf{H} = \sum_{e \in [n]} \Delta_e^{-2}. \quad (2)$$

From Eq. (1), we see that, for each arm  $e \notin M_*$ , the gap  $\Delta_e$  represents sub-optimality of the best set that includes arm  $e$ ; and, for each arm  $e \in M_*$ , the gap  $\Delta_e$  is the sub-optimality of the best set that does not include arm  $e$ . When specializing to TOPK and MB, our definition resembles the previous definition of gaps due to Kalyanakrishnan et al. [16] and Gabillon et al. [11].

**Exchange class and the width of a decision class.** The analysis of our algorithm depends on certain exchange properties of combinatorial structures. To capture these properties, we introduce notions of *exchange set* and *exchange class* as tools for our analysis. We present their definitions in the following.

We begin with the definition of exchange set. We define an exchange set  $b$  as an ordered pair of disjoint sets  $b = (b_+, b_-)$  where  $b_+ \cap b_- = \emptyset$ . Then, we define operator  $\oplus$  such that, for any set  $M \subseteq [n]$  and any exchange set  $b = (b_+, b_-)$ , we have  $M \oplus b \triangleq M \setminus b_- \cup b_+$ . Similarly, we also define operator  $\ominus$  such that  $M \ominus b \triangleq M \setminus b_+ \cup b_-$ .

We call a collection of exchange sets  $\mathcal{B}$  an *exchange class* for  $\mathcal{M}$  if  $\mathcal{B}$  satisfies the following property. Let  $M$  and  $M'$  be two elements of  $\mathcal{M}$ . Then, for any  $e \in (M \setminus M')$ , there exists an exchange set  $(b_+, b_-) \in \mathcal{B}$  which satisfies  $e \in b_+$ ,  $b_+ \subseteq M' \setminus M$ ,  $b_- \subseteq M \setminus M'$ ,  $(M \oplus b) \in \mathcal{M}$  and  $(M' \ominus b) \in \mathcal{M}$ . We define the *width* of exchange class  $\mathcal{B}$  to be the size of largest exchange set as follows

$$\text{width}(\mathcal{B}) = \max_{(b_+, b_-) \in \mathcal{B}} |b_+| + |b_-|. \quad (3)$$

Our analysis uses exchange classes to build gadgets that interpolate between different members of a decision class. Intuitively an exchange class  $\mathcal{B}$  for  $\mathcal{M}$  can be seen as a collection of “patches” (borrowing concepts from software engineering) such that, for any two different sets  $M, M' \in \mathcal{M}$ , one can convert from  $M$  to  $M'$  by applying a series of patches from  $\mathcal{B}$ , i.e.  $M' = M \oplus b_1 \oplus \dots \oplus b_k$ . And  $\text{width}(\mathcal{B})$  reflects the granularity of these patches.

Let  $\text{Exchange}(\mathcal{M})$  denote the family of all possible exchange classes for  $\mathcal{M}$ . We define the width of a decision class  $\mathcal{M}$  as the width of the thinnest exchange class

$$\text{width}(\mathcal{M}) = \min_{\mathcal{B} \in \text{Exchange}(\mathcal{M})} \text{width}(\mathcal{B}), \quad (4)$$

where  $\text{width}(\mathcal{B})$  is defined in Eq. (3).

**Sample complexity.** Our main result is a problem-dependent sample complexity bound of the CGapExp algorithm which shows that CGapExp returns the optimal set with high probability and uses at most  $\tilde{O}(\text{width}(\mathcal{M})^2 \mathbf{H})$  samples.

**Theorem 1.** *Given any  $\delta \in (0, 1)$ , any decision class  $\mathcal{M} \subseteq 2^{[n]}$  and any expected rewards  $\mathbf{w} \in \mathbb{R}^n$ . Assume that the reward distribution  $\varphi_e$  for each arm  $e \in [n]$  is  $R$ -sub-Gaussian with mean  $w(e)$ . Set  $\text{rad}_t(e) = R\sqrt{2 \log(\frac{4nt^2}{\delta})} / T_e(t)$  for all  $t > 0$  and  $e \in [n]$ . Then, with probability at least  $1 - \delta$ , the CGapExp algorithm (Algorithm 1) returns the optimal set  $\text{Out} = M_*$  and*

$$T \leq O(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \log(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \cdot n / \delta)), \quad (5)$$

where  $T$  denotes the number of samples used by Algorithm 1 and  $\mathbf{H}$  is defined in Eq. (2).

### 3.2 Examples of decision classes

Now we investigate several concrete types of decision classes, which correspond to different structured pure exploration problems. We analyze the width of these decision classes and apply Theorem 1 to obtain the sample complexity bounds. We begin with the problem of top- $K$  arm identification (TOPK) and multi-bandit best arms identification (MB).

**Example 1** (TOPK and MB). *For any  $K \in [n]$ , the problem of finding the top  $K$  arms with the largest expected reward can be modeled by decision class  $\mathcal{M}_{\text{TOPK}(K)} = \{M \subseteq [n] \mid |M| = K\}$ . Let  $\mathcal{A} = \{A_1, \dots, A_m\}$  be a partition of  $[n]$ . The problem of identifying the best arms from each group of arms  $A_1, \dots, A_m$  can be modeled by decision class  $\mathcal{M}_{\text{MB}(\mathcal{A})} = \{M \subseteq [n] \mid \forall i \in [m], |M \cap A_i| = 1\}$ .*

*Then we have  $\text{width}(\mathcal{M}_{\text{TOPK}(K)}) \leq 2$  and  $\text{width}(\mathcal{M}_{\text{MB}(\mathcal{A})}) \leq 2$  (see Fact 1 in the appendix) and therefore the sample complexity of CGapExp for solving TOPK and MB is  $O(\mathbf{H} \log(n\mathbf{H}/\delta))$ , which matches previous results in the fixed confidence setting [16, 11].*

Next we consider the problem of identifying the maximum matching or the shortest path in a setting where arms correspond to edges. For these problems, Theorem 1 establishes the first known sample complexity bound.

**Example 2** (Matchings and Paths). *Let  $G(V, E)$  be a graph with  $n$  edges and assume there is a one-to-one mapping between edges  $E$  and arms  $[n]$ . First let us assume that  $G$  is a bipartite graph. Let  $\mathcal{M}_{\text{MATCH}(G)}$  denote the set of all matchings in  $G$ . Then we have  $\text{width}(\mathcal{M}_{\text{MATCH}(G)}) \leq |V|$  (see Fact 1).*

*Next suppose that  $G$  is a direct acyclic graph and let  $s, t \in V$  be two vertices. Let  $\mathcal{M}_{\text{PATH}(G, s, t)}$  denote the set of all paths from  $s$  to  $t$ . Then we have  $\text{width}(\mathcal{M}_{\text{PATH}(G, s, t)}) \leq |V|$  (see Fact 1). Therefore the sample complexity bounds of CGapExp for decision classes  $\mathcal{M}_{\text{MATCH}(G)}$  and  $\mathcal{M}_{\text{PATH}(G, s, t)}$  are  $O(|V|^2 \mathbf{H} \log(n\mathbf{H}/\delta))$ .*

Last, we investigate the general problem of identifying the maximum-weight basis of a matroid, which encompasses a wide range of pure exploration problems.

**Example 3 (Matroids).** Let  $T = (E, \mathcal{I})$  be a finite matroid, where  $E$  is a set of size  $n$  (called ground set) and  $\mathcal{I}$  is a family of subsets of  $E$  (called independent sets) which satisfies the axioms of matroids<sup>2</sup>. Assume that there is a one-to-one mapping between  $E$  and  $[n]$ . And recall that a basis of matroid  $T$  is a maximal independent set. Let  $\mathcal{M}_{\text{MATROID}(T)}$  denote the set of all bases of  $T$ . Then we have  $\text{width}(\mathcal{M}_{\text{MATROID}(T)}) \leq 2$  and sample complexity of CGapExp for  $\mathcal{M}_{\text{MATROID}(T)}$  is  $O(\mathbf{H} \log(n\mathbf{H}/\delta))$ .

In our last example, we see that  $\mathcal{M}_{\text{MATROID}(T)}$  is a general type of decision class which encompasses TOPK and MB as special cases, where TOPK corresponds to uniform matroids of rank  $K$  and MB corresponds to partition matroids. It is easy to see that  $\mathcal{M}_{\text{MATROID}(T)}$  also covers the decision class that contains all spanning trees of a graph. On the other hand, the family of matchings and paths cannot be formulated as matroids and in fact they are matroid intersections (cf. [20]).

## 4 Lower Bound

In this section, we present a problem-dependent lower bound on the sample complexity of the ExpCMAB problem. To state our results, we first define the notion of  $\delta$ -correct algorithm as follows. For any  $\delta \in (0, 1)$ , we call an algorithm  $\mathbb{A}$  a  $\delta$ -correct algorithm if, for any expected reward  $\mathbf{w} \in \mathbb{R}^n$ , the probability of error of  $\mathbb{A}$  is at most  $\delta$ , i.e.  $\Pr[M_* \neq \text{Out}] \leq \delta$ , where Out is the output of algorithm  $\mathbb{A}$ .

We show that, for any decision class  $\mathcal{M}$  and any expected rewards  $\mathbf{w}$ , any  $\delta$ -correct algorithm  $\mathbb{A}$  must use at least  $\Omega(\mathbf{H} \log(1/\delta))$  samples in expectation.

**Theorem 2.** Fix any decision class  $\mathcal{M} \subseteq 2^{[n]}$  and any vector  $\mathbf{w} \in \mathbb{R}^n$ . Suppose that, for each arm  $e \in [n]$ , the reward distribution  $\varphi_e$  is given by  $\varphi_e = \mathcal{N}(w(e), 1)$ , where  $\mathcal{N}(\mu, \sigma^2)$  denotes a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ . Then, for any  $\delta \in (0, e^{-16}/4)$  and any  $\delta$ -correct algorithm  $\mathbb{A}$ , we have

$$\mathbb{E}[T] \geq \frac{1}{16} \mathbf{H} \log \left( \frac{1}{4\delta} \right), \quad (6)$$

where  $T$  denote the number of total samples used by algorithm  $\mathbb{A}$  and  $\mathbf{H}$  is defined in Eq. (2).

Theorem 2 resolves the open conjectures of Kalyanakrishnan et al. [16] and Bubeck et al. [7] that the lower bounds of sample complexity of TOPK and MB problem are  $\Omega(\mathbf{H} \log(1/\delta))$ . In addition, in Example 1 and Example 3, we have shown that the sample complexity of CGapExp is  $O(\mathbf{H} \log(n\mathbf{H}/\delta))$  for TOPK, MB and more generally the decision classes derived from matroids  $\mathcal{M}_{\text{MATROID}(T)}$ . Hence, we see that the CGapExp algorithm achieves the optimal sample complexity within logarithmic factors for these pure exploration problems.

On the other hand, for general decision classes with non-constant widths, we see that there is a gap of  $\tilde{O}(\text{width}(\mathcal{M})^2)$  between the upper bound Eq. (5) and the lower bound Eq. (6). Notice that we have  $\text{width}(\mathcal{M}) \leq n$  for any decision class  $\mathcal{M}$  and therefore the gap is relatively benign. Our lower bound also suggests that the dependency on  $\mathbf{H}$  of the sample complexity of CGapExp cannot be improved up to logarithmic factors. Furthermore, we conjecture that the sample complexity lower bound might inherently depend on the size of exchange sets. In the supplementary material, we provide evidence on this conjecture which is a lower bound on the sample complexity of exploration of exchange sets.

## 5 Fixed Budget Algorithm

In this section, we present CGapKill, a parameter-free learning algorithm for the ExpCMAB problem in the fixed budget setting. We analyze the probability of error CGapKill using the tools of exchange

<sup>2</sup>The three axioms of matroid are (1)  $\emptyset \in \mathcal{I}$  and  $\mathcal{I} \neq \{\emptyset\}$ ; (2) Every subsets of an independent set are independent (hereditary property); (3) For all  $A, B \in \mathcal{I}$  such that  $|B| = |A| + 1$  there exists an element  $e \in B \setminus A$  such that  $A \cup \{e\} \in \mathcal{I}$  (augmentation property). We refer interested readers to [20] for a general introduction to the matroid theory.



classes. Recall that, in the fixed budget setting, an algorithm is given a budget  $T > 0$  such that the algorithm can use at most  $T$  pulls. The goal of the algorithm is to minimize the probability of error.

**Constrained oracle.** The CGapKill algorithm requires access to a *constrained oracle*, which is a function denoted as  $\text{COracle} : \mathbb{R}^n \times 2^{[n]} \times 2^{[n]} \rightarrow \mathcal{M}$ . The function COracle takes three inputs: (1) a weight vector  $\mathbf{v}$ , (2) a set  $A$  of positive constraints and (3) a set  $B$  of negative constraints and returns a solution such that

$$\text{COracle}(\mathbf{v}, A, B) \in \arg \max_{M: M \in \mathcal{M}, A \subseteq M, B \cap M = \emptyset} v(M). \quad (7)$$

Hence we see that  $\text{COracle}(\mathbf{v}, A, B)$  computes an optimal solution that includes all elements of  $A$  while excluding all elements of  $B$ . In the supplementary material, we show that constrained oracles can be easily reduced to maximization oracles by modifying the weight vector according to the constraints. In addition, similar to CGapExp, CGapKill does not need any additional knowledge of  $\mathcal{M}$  other than accesses to a constrained oracle for  $\mathcal{M}$ .

**Algorithm.** The idea of the CGapKill algorithm is as follows. The CGapKill algorithm divides the budget of  $T$  rounds into  $n$  phases. In the end of each phase, CGapKill either accepts or rejects a single arm. If an arm is accepted, then it is included into the final output. Conversely, if an arm is rejected, then it is excluded from the final output. The arms that are neither accepted nor rejected are sampled for a equal number of times in the next phase.

Now we describe the procedure of the CGapKill algorithm for choosing an arm to accept/reject. Let  $A_t$  denote the set of accepted arms before phase  $t$  and let  $B_t$  denote the set of rejected arms before phase  $t$ . We call an arm  $e$  to be active if  $e \notin A_t \cup B_t$ . Then, in phase  $t$ , CGapKill samples each active arm for  $\tilde{T}_t - \tilde{T}_{t-1}$  times, where the definition of  $\tilde{T}_t$  is given in Algorithm 2. Next, CGapKill calls the constrained oracle to compute an optimal solution  $M_t$  with respect to the empirical means  $\bar{\mathbf{w}}_t$ , accepted arms  $A_t$  and rejected arms  $B_t$ , i.e. let  $M_t = \text{COracle}(\bar{\mathbf{w}}_t, A_t, B_t)$ . Then, for each arm active arm  $e$ , CGapKill estimate the gap of  $e$  in the following way. If  $e \in M_t$ , then CGapKill computes an optimal solution  $\tilde{M}_{t,e}$  that does not include  $e$ , i.e.  $\tilde{M}_{t,e} = \text{COracle}(\bar{\mathbf{w}}_t, A_t, B_t \cup \{e\})$ . Conversely, if  $e \notin M_t$ , then CGapKill computes an optimal  $\tilde{M}_{t,e}$  which includes  $e$ , i.e.  $\tilde{M}_{t,e} = \text{COracle}(\bar{\mathbf{w}}_t, A_t \cup \{e\}, B_t)$ . Then, the gap of  $e$  is calculated as  $\bar{\mathbf{w}}_t(M_t) - \bar{\mathbf{w}}_t(\tilde{M}_{t,e})$ . Finally, CGapKill chooses the arm  $p_t$  with the largest gap. If  $p_t \in M_t$  then  $p_t$  is accepted otherwise  $p_t$  is rejected. The pseudo-code of CGapKill is shown in Algorithm 2.

## 5.1 Analysis

In the following theorem, we show that the probability of error of the CGapKill algorithm is at most  $\tilde{O}(\exp(-T \text{width}(\mathcal{M})^{-2} \mathbf{H}^{-1}))$ .

**Theorem 3.** *Use the same notations as in Theorem 1. Let  $\Delta_{(1)}, \dots, \Delta_{(n)}$  be a permutation of  $\Delta_1, \dots, \Delta_n$  such that  $\Delta_{(1)} \leq \dots \leq \Delta_{(n)}$ . Define  $\mathbf{H}_2 \triangleq \max_{i \in [n]} i \Delta_{(i)}^{-2}$ . Then, given any budget  $T > n$ , the CGapKill algorithm uses at most  $T$  samples and outputs a solution  $\text{Out} \in \mathcal{M}$  such that*

$$\Pr[\text{Out} \neq M_*] \leq n^2 \exp \left( - \frac{2(T - n)}{9R^2 \tilde{\log}(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right), \quad (8)$$

where  $\tilde{\log}(n) \triangleq \sum_{i=1}^n \frac{1}{i}$ .

The quantity  $\mathbf{H}_2$  defined in Theorem 3 can be considered as a surrogate of  $\mathbf{H}$  since these two quantities are equivalent up to a logarithmic factor:  $\mathbf{H}_2 \leq \mathbf{H} \leq \log(2n) \mathbf{H}_2$ . Applying Theorem 3 to the TOPK problem, we see that our bound matches previous fixed budget algorithm due to Bubeck et al. [7].

We see that  $\text{width}(\mathcal{M})^2$  also appears in the probability of error Eq. (8) in Theorem 3, which indicates that  $\text{width}(\mathcal{M})$  may reflect the inherent hardness associate with a decision class  $\mathcal{M}$  for the ExpCMAB problem. We also notice that the CGapKill algorithm uses solely the empirical means of arms to support its decision making; while our main algorithm CGapExp uses explicitly constructed confidence intervals of the arms. From this viewpoint, the CGapKill algorithm is quite different from our main algorithm CGapExp. This actually leads to an argument substantially different from the proof of Theorem 1.

---

**Algorithm 2** CGapKill: Combinatorial Gap-based Elimination

---

**Require:** Budget:  $T > 0$ ; Constrained oracle:  $\text{COracle} : \mathbb{R}^n \times 2^{[n]} \times 2^{[n]} \rightarrow \mathcal{M}$ .

```
1: Define  $\tilde{\log}(n) \triangleq \sum_{i=1}^n \frac{1}{i}$ 
2:  $\tilde{T}_0 \leftarrow 0, A_1 \leftarrow \emptyset, B_1 \leftarrow \emptyset$ 
3: for  $t = 1, \dots, n$  do
4:    $\tilde{T}_t \leftarrow \left\lceil \frac{T-n}{\log(n)(n-t+1)} \right\rceil$ 
5:   Pull each arm  $e \in [n] \setminus (A_t \cup B_t)$  for  $\tilde{T}_t - \tilde{T}_{t-1}$  times
6:   Update the empirical means  $\bar{\mathbf{w}}_t \in \mathbb{R}^n$  of each arm
7:    $M_t \leftarrow \text{COracle}(\bar{\mathbf{w}}_t, A_t, B_t)$ 
8:   for each  $e \in [n] \setminus (A_t \cup B_t)$  do
9:     if  $e \in M_t$  then
10:       $\tilde{M}_{t,e} \leftarrow \text{COracle}(\bar{\mathbf{w}}_t, A_t, B_t \cup \{e\})$ 
11:     else
12:       $\tilde{M}_{t,e} \leftarrow \text{COracle}(\bar{\mathbf{w}}_t, A_t \cup \{e\}, B_t)$ 
13:     end if
14:   end for
15:    $p_t \leftarrow \arg \max_{i \in [n] \setminus (A_t \cup B_t)} \bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,i})$ 
16:   if  $p_t \in M_t$  then
17:      $A_{t+1} \leftarrow A_t \cup \{p_t\}, B_{t+1} \leftarrow B_t$ 
18:   else
19:      $A_{t+1} \leftarrow A_t, B_{t+1} \leftarrow B_t \cup \{p_t\}$ 
20:   end if
21: end for
22: Out  $\leftarrow A_{n+1}$ 
23: return Out
```

---

## 6 Related Work

The multi-armed bandit problem has been extensively studied in both stochastic and adversarial settings [17, 3, 3]. We refer readers to [4] for a survey on recent results. Many work on MABs focus on minimizing the cumulative regret, which is an objective known to be fundamentally different from the objective of pure exploration MABs [5]. Among these work, a recent line of research considers a generalized setting called combinatorial bandits in which a set of arms (satisfying certain combinatorial constraint) is played on each round [8, 9, 14, 19, 1, 6]. Note that the objective of these work is to minimize the cumulative regret, which differs from ours.

In the literature of pure exploration MABs, the classical problem of identifying the single best arm has been well-studied in both fixed confidence and fixed budget settings [18, 5, 2, 11, 13, 12]. A flurry of recent work extend this classical problem to TOPK and MB problems and provide algorithms and upper bounds [15, 16, 21, 7, 10, 11]. Our framework encompasses these two problems as special cases and covers a much larger class of structured pure exploration problems, which are unaddressed in the current literature. Applying our results back to TOPK and MB, our upper bounds match the best known problem-dependent bounds due to Gabillon et al. [11] and Bubeck et al. [7]; and our lower bound provides the first problem-dependent lower bounds for these two problems, which are conjectured earlier by Kalyanakrishnan et al. [16] and Bubeck et al. [7].

## 7 Conclusion

In this paper, we proposed a general framework called ExpCMAB which represents a rich class of structured pure exploration problems and admits potential applications in various domains. We have shown a number of results for the framework, including two novel learning algorithms, their related upper bounds and a novel lower bound. The proposed algorithms support a wide range of decision classes in a unifying way and our analysis introduced a novel tool called exchange class which maybe of independent interest. Our upper and lower bounds characterize the complexity of the ExpCMAB problem: the sample complexity of our algorithm is optimal (up to a logarithmic factor) for the decision classes derived from matroids (including TOPK and MB), while for general decision classes, our upper and lower bounds are within a relatively benign factor.



## References

- [1] J.-Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, 2009.
- [2] J.-Y. Audibert, S. Bubeck, and R. Munos. Best arm identification in multi-armed bandits. In *COLT*, 2010.
- [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- [4] S. Bubeck and N. Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5:1–122, 2012.
- [5] S. Bubeck, R. Munos, and G. Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412:1832–1852, 2010.
- [6] S. Bubeck, N. Cesa-bianchi, S. M. Kakade, S. Mannor, N. Srebro, and R. C. Williamson. Towards minimax policies for online linear optimization with bandit feedback. In *COLT*, 2012.
- [7] S. Bubeck, T. Wang, and N. Viswanathan. Multiple identifications in multi-armed bandits. In *ICML*, pages 258–265, 2013.
- [8] N. Cesa-Bianchi and G. Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.
- [9] W. Chen, Y. Wang, and Y. Yuan. Combinatorial multi-armed bandit: General framework and applications. In *Proceedings of The 30th International Conference on Machine Learning*, pages 151–159, 2013.
- [10] V. Gabillon, M. Ghavamzadeh, A. Lazaric, and S. Bubeck. Multi-bandit best arm identification. In *NIPS*. 2011.
- [11] V. Gabillon, M. Ghavamzadeh, and A. Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In *NIPS*, 2012.
- [12] K. Jamieson and R. Nowak. Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. In *Information Sciences and Systems (CISS), 2014 48th Annual Conference on*, pages 1–6. IEEE, 2014.
- [13] K. Jamieson, M. Malloy, R. Nowak, and S. Bubeck. lil’ucb: An optimal exploration algorithm for multi-armed bandits. *COLT*, 2014.
- [14] S. Kale, L. Reyzin, and R. E. Schapire. Non-stochastic bandit slate problems. In *NIPS*, pages 1054–1062, 2010.
- [15] S. Kalyanakrishnan and P. Stone. Efficient selection of multiple bandit arms: Theory and practice. In *ICML*, pages 511–518, 2010.
- [16] S. Kalyanakrishnan, A. Tewari, P. Auer, and P. Stone. Pac subset selection in stochastic multi-armed bandits. In *ICML*, pages 655–662, 2012.
- [17] T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- [18] S. Mannor and J. N. Tsitsiklis. The sample complexity of exploration in the multi-armed bandit problem. *The Journal of Machine Learning Research*, 5:623–648, 2004.
- [19] G. Neu, A. Györfgy, and C. Szepesvári. The online loop-free stochastic shortest-path problem. In *COLT*, pages 231–243, 2010.
- [20] J. G. Oxley. *Matroid theory*. Oxford university press, 2006.
- [21] Y. Zhou, X. Chen, and J. Li. Optimal pac multiple arm identification with applications to crowdsourcing. In *ICML*, 2014.

## A Proof of Main Result

In this section, we prove our main result: Theorem 1.

**Notations.** We need some additional notations for our analysis. For any set  $a \subseteq [n]$ , let  $\chi_a \in \{0, 1\}^n$  denote the incidence vector of set  $a \subseteq [n]$ , i.e.  $\chi_a(e) = 1$  if and only if  $e \in a$ . For an exchange set  $b = (b_+, b_-)$ , we define  $\chi_b \triangleq \chi_{b_+} - \chi_{b_-}$  as the incidence vector of  $b$ . We notice that  $\chi_b \in \{-1, 0, 1\}^n$ .

For each round  $t$ , we define vector  $\mathbf{rad}_t = (\text{rad}_t(1), \dots, \text{rad}_t(n))^T$  and recall that  $\bar{w}_t \in \mathbb{R}^n$  is the empirical mean rewards of arms up to round  $t$ .

Let  $u \in \mathbb{R}^n$  and  $v \in \mathbb{R}^n$  be two vectors. Let  $\langle u, v \rangle$  denote the inner product of  $u$  and  $v$ . We define  $u \circ v \triangleq (u(1) \cdot v(1), \dots, u(n) \cdot v(n))^T$  as the element-wise product of  $u$  and  $v$ . For any  $s \in \mathbb{R}$ , we also define  $u^s \triangleq (u(1)^s, \dots, u(n)^s)^T$  as the element-wise exponentiation of  $u$ . Let  $|u| = (|u(1)|, \dots, |u(n)|)^T$  denote the element-wise absolute value of  $u$ .

### A.1 Preparatory Lemmas

**Lemma 1.** Let  $M_1 \subseteq [n]$  be a set. Let  $b = (b_+, b_-)$  be an exchange set such that  $b_- \subseteq M_1$  and  $b_+ \cap M_1 = \emptyset$ . Define  $M_2 = M_1 \oplus b$ . Then, we have

$$\chi_{M_1} + \chi_b = \chi_{M_2}.$$

*Proof.* Recall that  $M_2 = M_1 \setminus b_- \oplus b_+$  and  $b_+ \cap b_- = \emptyset$ . Therefore we see that  $M_2 \setminus M_1 = b_+$  and  $M_1 \setminus M_2 = b_-$ . Then, we decompose  $\chi_{M_1}$  as  $\chi_{M_1} = \chi_{M_1 \setminus M_2} + \chi_{M_1 \cap M_2}$ . Hence, we have

$$\begin{aligned} \chi_{M_1} + \chi_b &= \chi_{M_1 \setminus M_2} + \chi_{M_1 \cap M_2} + \chi_{b_+} - \chi_{b_-} \\ &= \chi_{M_1 \cap M_2} + \chi_{M_2 \setminus M_1} \\ &= \chi_{M_2}. \end{aligned}$$

□

**Lemma 2.** Let  $\mathcal{M} \subseteq 2^{[n]}$  and  $\mathcal{B}$  be an exchange class for  $\mathcal{M}$ . Then, for any two different elements  $M, M'$  of  $\mathcal{M}$  and any  $e \in (M \setminus M') \cup (M' \setminus M)$ , there exists an exchange set  $b = (b_+, b_-) \in \mathcal{B}$  such that  $e \in (b_+ \cup b_-)$ ,  $b_- \subseteq (M \setminus M')$ ,  $b_+ \subseteq (M' \setminus M)$ ,  $(M \oplus b) \in \mathcal{M}$  and  $(M' \ominus b) \in \mathcal{M}$ . Moreover, if  $M' = M_*$ , then we have  $\langle w, \chi_b \rangle \geq \Delta_e > 0$ , where  $\Delta_e$  is the gap defined in Eq. (1).

*Proof.* We decompose our proof into two cases.

**Case (1):**  $e \in M \setminus M'$ .

By the definition of exchange class, we know that there exists  $b = (b_+, b_-) \in \mathcal{B}$  which satisfies that  $e \in b_+$ ,  $b_- \subseteq (M \setminus M')$ ,  $b_+ \subseteq (M' \setminus M)$ ,  $(M \oplus b) \in \mathcal{M}$  and  $(M' \ominus b) \in \mathcal{M}$ .

Next, if  $M' = M_*$ , we see that  $e \notin M_*$ . Let us consider the set  $M_1 = \arg \max_{M': M' \in \mathcal{M} \wedge e \in M'} w(M')$ . Also define  $M_0 = M_* \ominus b$ . We have already proved that  $M_0 \in \mathcal{M}$ . Combining with the fact that  $e \in M_0$ , we see that  $w(M_0) \leq w(M_1)$ . Therefore, we obtain that  $w(M_*) - w(M_0) \geq w(M_*) - w(M_1) = \Delta_e$ . Notice that the left-hand side of the former inequality can be rewritten using Lemma 1 as follows

$$w(M_*) - w(M_0) = \langle w, \chi_{M_*} \rangle - \langle w, \chi_{M_0} \rangle = \langle w, \chi_{M_*} - \chi_{M_0} \rangle = \langle w, \chi_b \rangle.$$

Therefore, we obtain  $\langle w, \chi_b \rangle \geq \Delta_e$ .

**Case (2):**  $e \in M' \setminus M$ .

Using the definition of exchange class, we see that there exists  $c = (c_+, c_-) \in \mathcal{B}$  such that  $e \in c_-$ ,  $c_- \subseteq (M' \setminus M)$ ,  $c_+ \subseteq (M \setminus M')$ ,  $(M' \oplus c) \in \mathcal{M}$  and  $(M \ominus c) \in \mathcal{M}$ .

We construct  $b = (b_+, b_-)$  by setting  $b_+ = c_-$  and  $b_- = c_+$ . Notice that, by the construction of  $b$ , we have  $M \oplus b = M \ominus c$  and  $M' \ominus b = M' \oplus c$ . Therefore, it is clear that  $b$  satisfies the requirement of the lemma.

Now, suppose that  $M' = M_*$ . In this case, we have  $e \in M_*$ . Consider the set  $M_3 = \arg \max_{M': M' \in \mathcal{M} \wedge e \notin M'} w(M')$ . We see that  $w(M_*) - w(M_3) = \Delta_e$ . Define  $M_2 = M_* \ominus b$  and notice that  $M_2 \in \mathcal{M}$ . Combining with the fact that  $e \notin M_2$ , we obtain that  $w(M_2) \leq w(M_3)$ . Hence, we have  $w(M_*) - w(M_2) \geq w(M_*) - w(M_3) = \Delta_e$ . Similar to Case (1), applying Lemma 1 again, we have

$$\langle w, \chi_b \rangle = w(M_*) - w(M_2) \geq \Delta_e.$$

□

**Lemma 3.** Let  $M$  and  $M'$  be two sets. Then, we have

$$\max_{e \in (M \setminus M') \cup (M' \setminus M)} \text{rad}_t(e) = \|\text{rad}_t \circ |\chi_{M'} - \chi_M|\|_\infty.$$

*Proof.* Notice that  $\chi_{M'} - \chi_M = \chi_{M' \setminus M} - \chi_{M \setminus M'}$ . In addition, since  $(M' \setminus M) \cap (M \setminus M') = \emptyset$ , we have  $\chi_{M' \setminus M} \circ \chi_{M \setminus M'} = \mathbf{0}_n$ . Also notice that  $\chi_{M' \setminus M} - \chi_{M \setminus M'} \in \{-1, 0, 1\}^n$ . Therefore, we have

$$\begin{aligned} |\chi_{M' \setminus M} - \chi_{M \setminus M'}| &= (\chi_{M' \setminus M} - \chi_{M \setminus M'})^2 \\ &= \chi_{M' \setminus M}^2 + \chi_{M \setminus M'}^2 + 2\chi_{M' \setminus M} \circ \chi_{M \setminus M'} \\ &= \chi_{M' \setminus M} + \chi_{M \setminus M'} \\ &= \chi_{(M' \setminus M) \cup (M \setminus M')}, \end{aligned}$$

where the third equation follows from the fact that  $\chi_{M \setminus M'} \in \{0, 1\}^n$  and  $\chi_{M' \setminus M} \in \{0, 1\}^n$ . The lemma follows immediately from the fact that  $\text{rad}_t(e) \geq 0$  and  $\chi_{(M \setminus M') \cup (M' \setminus M)} \in \{0, 1\}^n$ . □

**Lemma 4.** Let  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^n$  be three vectors. Then, we have  $\langle \mathbf{a}, \mathbf{b} \circ \mathbf{c} \rangle = \langle \mathbf{a} \circ \mathbf{b}, \mathbf{c} \rangle$ .

*Proof.* We have

$$\langle \mathbf{a}, \mathbf{b} \circ \mathbf{c} \rangle = \sum_{i=1}^n a(i)(b(i)c(i)) = \sum_{i=1}^n (a(i)b(i))c(i) = \langle \mathbf{a} \circ \mathbf{b}, \mathbf{c} \rangle.$$

□

**Lemma 5.** Let  $M_t$  and  $\tilde{w}_t$  be defined in Algorithm 1. Let  $M' \in \mathcal{M}$  be a feasible set. We have

$$\tilde{w}_t(M') - \tilde{w}_t(M_t) = \langle \tilde{w}_t, \chi_{M'} - \chi_{M_t} \rangle = \langle \bar{w}_t, \chi_{M'} - \chi_{M_t} \rangle + \langle \text{rad}_t, |\chi_{M'} - \chi_{M_t}| \rangle.$$

*Proof.* We begin with proving the first part. It is easy to verify that  $\tilde{w}_t = \bar{w}_t + \text{rad}_t \circ (\mathbf{1}_n - 2\chi_{M_t})$ . Then, we have

$$\begin{aligned} \langle \tilde{w}_t, \chi_{M'} - \chi_{M_t} \rangle &= \langle \bar{w}_t + \text{rad}_t \circ (\mathbf{1}_n - 2\chi_{M_t}), \chi_{M'} - \chi_{M_t} \rangle \\ &= \langle \bar{w}_t, \chi_{M'} - \chi_{M_t} \rangle + \langle \text{rad}_t, (\mathbf{1}_n - 2\chi_{M_t}) \circ (\chi_{M'} - \chi_{M_t}) \rangle \end{aligned} \quad (9)$$

$$\begin{aligned} &= \langle \bar{w}_t, \chi_{M'} - \chi_{M_t} \rangle + \langle \text{rad}_t, \chi_{M'} - \chi_{M_t} - 2\chi_{M_t} \circ \chi_{M'} + 2\chi_{M_t}^2 \rangle \\ &= \langle \bar{w}_t, \chi_{M'} - \chi_{M_t} \rangle + \langle \text{rad}_t, \chi_{M'}^2 - \chi_{M_t}^2 - 2\chi_{M_t} \circ \chi_{M'} + 2\chi_{M_t}^2 \rangle \end{aligned} \quad (10)$$

$$\begin{aligned} &= \langle \bar{w}_t, \chi_{M'} - \chi_{M_t} \rangle + \langle \text{rad}_t, (\chi_{M'} - \chi_{M_t})^2 \rangle \\ &= \langle \bar{w}_t, \chi_{M'} - \chi_{M_t} \rangle + \langle \text{rad}_t, |\chi_{M'} - \chi_{M_t}| \rangle, \end{aligned} \quad (11)$$

where Eq. (9) follows from Lemma 4; Eq. (10) holds since  $\chi_{M'} \in \{0, 1\}^n$  and  $\chi_{M_t} \in \{0, 1\}^n$  and therefore  $\chi_{M'} = \chi_{M'}^2$  and  $\chi_{M_t} = \chi_{M_t}^2$ ; and Eq. (11) follows since  $\chi_{M'} - \chi_{M_t} \in \{-1, 0, 1\}^n$ . □

## A.2 Confidence Intervals

For all  $t > 0$ , we define random event  $\xi_t$  as follows

$$\xi_t = \left\{ \forall i \in [n], \quad |w(i) - \bar{w}_t(i)| \leq \text{rad}_t(i) \right\}. \quad (12)$$

We notice that random event  $\xi_t$  characterizes the event that the confidence bounds of all arms are valid at round  $t$ .

If the confidence bounds are valid, we can generalize Eq. (12) to inner products as follows.

**Lemma 6.** *Given any  $t > 0$ , assume that event  $\xi_t$  as defined in Eq. (12) occurs. Then, for any vector  $\mathbf{a} \in \mathbb{R}^n$ , we have*

$$|\langle \mathbf{w}, \mathbf{a} \rangle - \langle \bar{\mathbf{w}}_t, \mathbf{a} \rangle| \leq \langle \mathbf{rad}_t, |\mathbf{a}| \rangle.$$

*Proof.* Suppose that  $\xi_t$  occurs. Then, we have

$$\begin{aligned} |\langle \mathbf{w}, \mathbf{a} \rangle - \langle \bar{\mathbf{w}}_t, \mathbf{a} \rangle| &= |\langle \mathbf{w} - \bar{\mathbf{w}}_t, \mathbf{a} \rangle| \\ &= \left| \sum_{i=1}^n (w(i) - \bar{w}_t(i)) a(i) \right| \\ &\leq \sum_{i=1}^n |w(i) - \bar{w}_t(i)| |a(i)| \\ &\leq \sum_{i=1}^n \text{rad}_t(i) \cdot |a(i)| \\ &= \langle \mathbf{rad}_t, |\mathbf{a}| \rangle, \end{aligned} \quad (13)$$

where Eq. (13) follows the definition of event  $\xi_t$  in Eq. (12) and the assumption that it occurs.  $\square$

Next, we construct the high probability confidence intervals for the fixed confidence setting.

**Lemma 7.** *Suppose that the reward distribution  $\varphi_e$  is a  $R$ -sub-Gaussian distribution for all  $e \in [n]$ . And if, for all  $t > 0$  and all  $e \in [n]$ , the confidence radius  $\text{rad}_t(e)$  is given by*

$$\text{rad}_t(e) = R \sqrt{\frac{2 \log \left( \frac{4nt^2}{\delta} \right)}{T_e(t)}},$$

where  $T_e(t)$  is the number of samples of arm  $e$  up to round  $t$ . Then, we have

$$\Pr \left[ \bigcap_{t=1}^{\infty} \xi_t \right] \geq 1 - \delta.$$

*Proof.* For any  $t > 0$  and  $e \in [n]$ , notice  $\varphi_e$  is a  $R$ -sub-Gaussian distribution with mean  $w(e)$  and  $w_t(e)$  is the empirical mean of  $\varphi_e$  for  $T_e(t)$  samples. Using Hoeffding's inequality (see Lemma 20 in Section E), we obtain

$$\Pr \left[ |\bar{w}_t(e) - w(e)| \geq R \sqrt{\frac{2 \log \left( \frac{4nt^2}{\delta} \right)}{T_e(t)}} \right] \leq \frac{\delta}{2nt^2}.$$

By union bound over all  $e \in [n]$ , we see that  $\Pr[\xi_t] \geq 1 - \frac{\delta}{2t^2}$ . Using a union bound again over all  $t > 0$ , we have

$$\begin{aligned} \Pr \left[ \bigcap_{t=1}^{\infty} \xi_t \right] &\geq 1 - \sum_{t=1}^{\infty} \Pr[\neg \xi_t] \\ &\geq 1 - \sum_{t=1}^{\infty} \frac{\delta}{2t^2} \\ &= 1 - \frac{\pi^2}{12} \delta \geq 1 - \delta. \end{aligned}$$

$\square$

### A.3 Main Lemmas

**Lemma 8.** *Given any  $t > 0$ , assume that event  $\xi_t$  (defined in Eq. (12)) occurs. Then, if Algorithm 1 terminates at round  $t$ , we have  $M_t = M_*$ .*

*Proof.* Suppose that  $M_t \neq M_*$ . By definition, we have  $w(M_*) > w(M_t)$ . Rewriting the former inequality, we obtain that  $\langle \mathbf{w}, \chi_{M_*} \rangle > \langle \mathbf{w}, \chi_{M_t} \rangle$ .

Applying Lemma 2 by setting  $M = M_t$  and  $M' = M_*$ , we see that there exists  $b = (b_+, b_-) \in \mathcal{B}$  such that  $(M_t \oplus b) \in \mathcal{M}$ .

Now define  $M'_t = M_t \oplus b$ . Recall that  $\tilde{M}_t = \arg \max_{M \in \mathcal{M}} \tilde{w}_t(M)$  and therefore  $\tilde{w}_t(\tilde{M}_t) \geq \tilde{w}_t(M'_t)$ . Hence, we have

$$\begin{aligned} \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) &\geq \tilde{w}_t(M'_t) - \tilde{w}_t(M_t) \\ &= \langle \tilde{\mathbf{w}}_t, \chi_{M'_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{M'} - \chi_{M_t}| \rangle \end{aligned} \quad (14)$$

$$\geq \langle \mathbf{w}, \chi_{M'_t} - \chi_{M_t} \rangle \quad (15)$$

$$= w(M'_t) - w(M_t) > 0, \quad (16)$$

where Eq. (14) follows from Lemma 5; and Eq. (15) follows the assumption that event  $\xi_t$  occurs and Lemma 6;

Therefore Eq. (16) shows that  $\tilde{w}_t(\tilde{M}_t) > \tilde{w}_t(M_t)$ . However, this contradicts to the stopping condition of CGapExp:  $\tilde{w}_t(\tilde{M}_t) \leq \tilde{w}_t(M_t)$  and the assumption that the algorithm terminates on round  $t$ .  $\square$

**Lemma 9.** *Given any  $t > 0$  and suppose that event  $\xi_t$  (defined in Eq. (12)) occurs. For any  $e \in [n]$ , if  $\text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{M})}$ , then, arm  $e$  will not be pulled on round  $t$ , i.e.  $p_t \neq e$ .*

*Proof.* Fix an exchange class  $\mathcal{B} \in \arg \min_{\mathcal{B}' \in \text{Exchange}(\mathcal{M})} \text{width}(\mathcal{B}')$ . Suppose, in the contrary, that  $p_t = e$ . By Lemma 2, there exists an exchange set  $c = (c_+, c_-) \in \mathcal{B}$  such that  $e \in (c_+ \cup c_-)$ ,  $c_- \subseteq (M_t \setminus \tilde{M}_t)$ ,  $c_+ \subseteq (\tilde{M}_t \setminus M_t)$ ,  $(M_t \oplus c) \in \mathcal{M}$  and  $(\tilde{M}_t \ominus c) \in \mathcal{M}$ .

Now, we decompose our proof into two cases.

**Case (1):**  $(e \in M_* \wedge e \in c_+) \vee (e \notin M_* \wedge e \in c_-)$ .

Define  $M'_t = \tilde{M}_t \ominus c$  and recall that  $M'_t \in \mathcal{M}$  due to the definition of exchange class.

First, we claim that  $M'_t \neq M_*$ . Suppose that  $e \in M_*$  and  $e \in c_+$ . Then, we see that  $e \notin M'_t$  and hence  $M'_t \neq M_*$ . On the other hand, if  $e \notin M_*$  and  $e \in c_-$ , then  $e \in M'_t$  which also means that  $M'_t \neq M_*$ . Therefore we have  $M'_t \neq M_*$  in either cases.

Next, we apply Lemma 2 by setting  $M = M'_t$  and  $M' = M_*$ . We see that there exists an exchange set  $b \in \mathcal{B}$  such that,  $e \in (b_+ \cup b_-)$ ,  $(M'_t \oplus b) \in \mathcal{M}$  and  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e > 0$ .

Now, we define vectors  $\mathbf{d} = \chi_{\tilde{M}_t} - \chi_{M_t}$ ,  $\mathbf{d}_1 = \chi_{M'_t} - \chi_{M_t}$  and  $\mathbf{d}_2 = \chi_{M'_t \oplus b} - \chi_{M_t}$ . By the definition of  $M'_t$  and Lemma 2, we see that  $\mathbf{d}_1 = \mathbf{d} - \chi_c$  and  $\mathbf{d}_2 = \mathbf{d}_1 + \chi_b = \mathbf{d} - \chi_c + \chi_b$ .

Then, we claim that  $\|\mathbf{rad}_t \circ (\mathbf{d} - \chi_c)\|_\infty < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}$ . Since  $c_- \subseteq M_t$  and  $c_+ \cap M_t = \emptyset$ , using standard set theoretical manipulations, we can show that  $M_t \setminus \tilde{M}_t = (M_t \setminus M'_t) \cup c_-$ . Similarly, one can show that  $\tilde{M}_t \setminus M_t = (M'_t \setminus M_t) \cup c_+$ . This means that  $((M_t \setminus M'_t) \cup (M'_t \setminus M_t)) \subseteq ((M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t))$ . Then, applying Lemma 3, we obtain

$$\begin{aligned} \|\mathbf{rad}_t \circ (\mathbf{d} - \chi_c)\|_\infty &= \left\| \mathbf{rad}_t \circ (\chi_{M'_t} - \chi_{M_t}) \right\|_\infty \\ &= \max_{i \in (M_t \setminus M'_t) \cup (M'_t \setminus M_t)} \text{rad}_t(i) \\ &\leq \max_{i \in (M_t \setminus M_t) \cup (\tilde{M}_t \setminus M_t)} \text{rad}_t(i) \end{aligned}$$

$$= \text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}. \quad (17)$$

We claim that  $\|\text{rad}_t \circ \chi_c\|_\infty < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}$ . Recall that, by the definition of  $c$ , we have  $c_+ \subseteq (\tilde{M}_t \setminus M_t)$  and  $c_- \subseteq (M_t \setminus \tilde{M}_t)$ . Hence  $c_+ \cup c_- \subseteq (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)$ . Since  $\chi_c \in [-1, 1]^n$ , we see that

$$\begin{aligned} \|\text{rad}_t \circ \chi_c\|_\infty &= \max_{i \in c_+ \cup c_-} \text{rad}_t(i) \\ &\leq \max_{i \in (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(i) \\ &= \text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}. \end{aligned} \quad (18)$$

Next, we claim that  $\mathbf{d} \circ \chi_c = |\chi_c|$ . Recall that  $\chi_c = \chi_{c_+} - \chi_{c_-}$  and  $\mathbf{d} = \chi_{\tilde{M}_t} - \chi_{M_t} = \chi_{\tilde{M}_t \setminus M_t} - \chi_{M_t \setminus \tilde{M}_t}$ . We also notice that  $c_+ \subseteq (\tilde{M}_t \setminus M_t)$  and  $c_- \subseteq (M_t \setminus \tilde{M}_t)$ . This implies that  $c_+ \cap (M_t \setminus \tilde{M}_t) = \emptyset$  and  $c_- \cap (\tilde{M}_t \setminus M_t) = \emptyset$ . Therefore, we have

$$\begin{aligned} \mathbf{d} \circ \chi_c &= (\chi_{\tilde{M}_t \setminus M_t} - \chi_{M_t \setminus \tilde{M}_t}) \circ (\chi_{c_+} - \chi_{c_-}) \\ &= \chi_{\tilde{M}_t \setminus M_t} \circ \chi_{c_+} + \chi_{M_t \setminus \tilde{M}_t} \circ \chi_{c_-} - \chi_{\tilde{M}_t \setminus M_t} \circ \chi_{c_-} - \chi_{M_t \setminus \tilde{M}_t} \circ \chi_{c_+} \\ &= \chi_{\tilde{M}_t \setminus M_t} \circ \chi_{c_+} + \chi_{M_t \setminus \tilde{M}_t} \circ \chi_{c_-} \\ &= \chi_{c_+} + \chi_{c_-} = |\chi_c|. \end{aligned}$$

where the last equality holds since  $c_+ \cap c_- = \emptyset$ .

Now, we bound quantity  $\langle \text{rad}_t, |\mathbf{d}_2| \rangle - \langle \text{rad}_t, |\mathbf{d}| \rangle$  as follows

$$\langle \text{rad}_t, |\mathbf{d}_2| \rangle - \langle \text{rad}_t, |\mathbf{d}| \rangle = \langle \text{rad}_t, |\mathbf{d}_2| - |\mathbf{d}| \rangle = \langle \text{rad}_t, \mathbf{d}_2^2 - \mathbf{d}^2 \rangle \quad (19)$$

$$\begin{aligned} &= \langle \text{rad}_t, (\mathbf{d} - \chi_c + \chi_b)^2 - \mathbf{d}^2 \rangle \\ &= \langle \text{rad}_t, \chi_b^2 + \chi_c^2 - 2\chi_b \circ \chi_c - 2\mathbf{d} \circ \chi_c + 2\mathbf{d} \circ \chi_b \rangle \\ &= \langle \text{rad}_t, \chi_b^2 - \chi_c^2 + 2\chi_b \circ (\mathbf{d} - \chi_c) \rangle \end{aligned} \quad (20)$$

$$\begin{aligned} &= \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - 2 \langle \text{rad}_t, \chi_b \circ (\mathbf{d} - \chi_c) \rangle \\ &= \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - 2 \langle \text{rad}_t \circ (\mathbf{d} - \chi_c), \chi_b \rangle \end{aligned} \quad (21)$$

$$\geq \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - 2 \|\text{rad}_t \circ (\mathbf{d} - \chi_c)\|_\infty \|\chi_b\|_1 \quad (22)$$

$$> \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3 \text{width}(\mathcal{B})} \|\chi_b\|_1 \quad (23)$$

$$\geq \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3}, \quad (24)$$

where Eq. (19) holds since  $\mathbf{d} \in \{-1, 0, 1\}^n$  and  $\mathbf{d}_2 \in \{-1, 0, 1\}^n$ ; Eq. (20) follows from the claim that  $\mathbf{d} \circ \chi_c = |\chi_c| = \chi_c^2$ ; Eq. (21) and Eq. (22) follow from Lemma 4 and Hölder's inequality; Eq. (23) follows from Eq. (17); and Eq. (24) holds since  $b \in \mathcal{B}$  and  $\|\chi_b\|_1 = |b_+| + |b_-| \leq \text{width}(\mathcal{B})$ .

Applying Lemma 5 by setting  $M' = M'_t \oplus b$  and using the fact that  $\tilde{w}_t(\tilde{M}_t) \geq \tilde{w}_t(M'_t \oplus b)$ , we have

$$\begin{aligned} \langle \bar{w}_t, \mathbf{d} \rangle + \langle \text{rad}_t, |\mathbf{d}| \rangle &= \langle \bar{w}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \text{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \\ &= \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \\ &\geq \tilde{w}_t(M'_t \oplus b) - \tilde{w}_t(M_t) \\ &= \langle \bar{w}_t, \chi_{M'_t \oplus b} - \chi_{M_t} \rangle + \langle \text{rad}_t, |\chi_{M'_t \oplus b} - \chi_{M_t}| \rangle \\ &= \langle \bar{w}_t, \mathbf{d}_2 \rangle + \langle \text{rad}_t, |\mathbf{d}_2| \rangle \\ &= \langle \bar{w}_t, \mathbf{d} \rangle - \langle \bar{w}_t, \chi_c \rangle + \langle \bar{w}_t, \chi_b \rangle + \langle \text{rad}_t, |\mathbf{d}_2| \rangle, \end{aligned}$$



where the last equality follows from the fact that  $\mathbf{d}_2 = \mathbf{d} - \chi_c + \chi_b$ . Rearranging the above inequality, we obtain

$$\begin{aligned} \langle \bar{\mathbf{w}}_t, \chi_c \rangle &\geq \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_2| \rangle - \langle \mathbf{rad}_t, |\mathbf{d}| \rangle \\ &\geq \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3} \end{aligned} \quad (25)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \langle \mathbf{rad}_t, \chi_c \rangle - \frac{2\Delta_e}{3} \quad (26)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3} - \frac{2\Delta_e}{3} \quad (27)$$

$$= \langle \mathbf{w}, \chi_b \rangle - \Delta_e \geq 0, \quad (28)$$

where Eq. (25) uses Eq. (24); Eq. (26) follows from the assumption that event  $\xi_t$  occurs and Lemma 6; and Eq. (26) holds since Eq. (18).

We have shown that  $\langle \bar{\mathbf{w}}_t, \chi_c \rangle > 0$ . Now we can bound  $\bar{w}_t(M'_t)$  as follows

$$\bar{w}_t(M'_t) = \langle \bar{\mathbf{w}}_t, \chi_{M'_t} \rangle = \langle \bar{\mathbf{w}}_t, \chi_{M_t} + \chi_c \rangle = \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \bar{\mathbf{w}}_t, \chi_c \rangle > \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle = w_t(M_t).$$

However, the definition of  $M_t$  ensures that  $M_t = \arg \max_{M \in \mathcal{M}} \bar{w}_t(M)$ , i.e.  $\bar{w}_t(M_t) \geq \bar{w}_t(M'_t)$ . Contradiction.

**Case (2):**  $(e \in M_* \wedge e \in c_-) \vee (e \notin M_* \wedge e \in c_+)$ .

First, we claim that  $\tilde{M}_t \neq M_*$ . Suppose that  $e \in M_*$  and  $e \in c_-$ . Then, we see that  $e \notin \tilde{M}_t$ , which implies that  $\tilde{M}_t \neq M_*$ . If  $e \notin M_*$  and  $e \in c_+$ , then  $e \in \tilde{M}_t$ , which also implies that  $\tilde{M}_t \neq M_*$ . Therefore we have  $\tilde{M}_t \neq M_*$  in either cases.

Hence, by Lemma 2, there exists an exchange set  $b = (b_+, b_-) \in \mathcal{B}$  such that  $e \in (b_+ \cup b_-)$ ,  $b_- \subseteq (\tilde{M}_t \setminus M_*)$ ,  $b_+ \subseteq (M_* \setminus \tilde{M}_t)$  and  $(\tilde{M}_t \oplus b) \in \mathcal{M}$ . Lemma 2 also indicates that  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e > 0$ .

Next, we define vectors  $\mathbf{d} = \chi_{\tilde{M}_t} - \chi_{M_t}$  and  $\mathbf{d}_1 = \chi_{\tilde{M}_t \oplus b} - \chi_{M_t}$ . Notice that Lemma 2 gives that  $\mathbf{d}_1 = \mathbf{d} + b$ .

Then, we apply Lemma 3 by setting  $M = M_t$  and  $M' = \tilde{M}_t$ . This shows that

$$\|\mathbf{rad}_t \circ \mathbf{d}\|_\infty \leq \max_{i: (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(i) = \text{rad}_t(e) < \frac{\Delta_e}{3}. \quad (29)$$

Now, we bound quantity  $\langle \bar{\mathbf{w}}_t, \mathbf{d}_1 \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| \rangle - \langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle - \langle \mathbf{rad}_t, |\mathbf{d}| \rangle$  as follows

$$\begin{aligned} \langle \bar{\mathbf{w}}_t, \mathbf{d}_1 \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| \rangle - \langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle - \langle \mathbf{rad}_t, |\mathbf{d}| \rangle &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| - |\mathbf{d}| \rangle \\ &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, \mathbf{d}_1^2 - \mathbf{d}^2 \rangle \end{aligned} \quad (30)$$

$$= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, 2\mathbf{d} \circ \chi_b + \chi_b^2 \rangle \quad (31)$$

$$= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, \chi_b^2 \rangle + 2 \langle \mathbf{rad}_t \circ \mathbf{d}, \chi_b \rangle$$

$$\geq \langle \mathbf{w}, \chi_b \rangle - 2 \langle \mathbf{rad}_t \circ \mathbf{d}, \chi_b \rangle \quad (32)$$

$$\geq \langle \mathbf{w}, \chi_b \rangle - 2 \|\mathbf{rad}_t \circ \mathbf{d}\|_\infty \|\chi_b\|_1 \quad (33)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \frac{2\Delta_e}{3} \quad (34)$$

$$\geq 0, \quad (35)$$

where Eq. (30) follows from the fact that  $\mathbf{d}_1 \in \{-1, 0, 1\}^n$  and  $\mathbf{d} \in \{-1, 0, 1\}^n$ ; Eq. (31) holds since  $\mathbf{d}_1 = \mathbf{d} + \chi_b$ ; Eq. (32) follows from the assumption that  $\xi_t$  occurs and Lemma 6; Eq. (33) follows from Lemma 4 and Hölder's inequality; and Eq. (34) is due to Eq. (29).

Therefore, we have proved that  $\langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle + \langle \mathbf{rad}_t, |\mathbf{d}| \rangle < \langle \bar{\mathbf{w}}_t, \mathbf{d}_1 \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| \rangle$ . However, Lemma 5 shows that

$$\begin{aligned} \langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle + \langle \mathbf{rad}_t, |\mathbf{d}| \rangle &= \langle \bar{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \\ &= \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \end{aligned}$$

$$\begin{aligned}
&\geq \tilde{w}_t(\tilde{M}_t \oplus b) - \tilde{w}_t(M_t) \\
&= \langle \tilde{\mathbf{w}}_t, \chi_{\tilde{M}_t \oplus b} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{\tilde{M}_t \oplus b} - \chi_{M_t}| \rangle \\
&= \langle \tilde{\mathbf{w}}_t, \mathbf{d}_1 \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| \rangle.
\end{aligned}$$

This is a contradiction and therefore  $p_t \neq e$ . □

#### A.4 Proof of Theorem 1

Theorem 1 is now a straightforward corollary of Lemma 8 and Lemma 9. For the readers' convenience, we first restate Theorem 1 in the following.

**Theorem 1.** *Given any  $\delta \in (0, 1)$ , any decision class  $\mathcal{M} \subseteq 2^{[n]}$  and any expected rewards  $\mathbf{w} \in \mathbb{R}^n$ . Assume that the reward distribution  $\varphi_e$  for each arm  $e \in [n]$  is  $R$ -sub-Gaussian with mean  $w(e)$ . Set  $\text{rad}_t(e) = R\sqrt{2 \log(\frac{4nt^2}{\delta})} / T_e(t)$  for all  $t > 0$  and  $e \in [n]$ . Then, with probability at least  $1 - \delta$ , the CGapExp algorithm (Algorithm 1) returns the optimal set  $\text{Out} = M_*$  and*

$$T \leq O\left(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \log\left(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \cdot n / \delta\right)\right), \quad (5)$$

where  $T$  denotes the number of samples used by Algorithm 1 and  $\mathbf{H}$  is defined in Eq. (2).

*Proof.* Lemma 7 indicates that the event  $\xi \triangleq \bigcap_{t=1}^{\infty} \xi_t$  occurs with probability at least  $1 - \delta$ . In the rest of the proof, we shall assume that this event holds.

By Lemma 8 and the assumption on  $\xi$ , we see that  $\text{Out} = M_*$ . Next, we focus on bounding the total number  $T$  of samples.

Fix any arm  $e \in [n]$ . Let  $T_e$  denote the total number of pull of arm  $e \in [n]$ . Let  $t_e$  be the last round which arm  $e$  is pulled, i.e.  $p_{t_e} = e$ . It is easy to see that  $T_e(t_e) = T_e - 1$ . By Lemma 9, we see that  $\text{rad}_{t_e}(e) \geq \frac{\Delta_e}{3 \text{width}(\mathcal{M})}$ . By plugging the definition of  $\text{rad}_{t_e}$ , we have

$$\frac{\Delta_e}{3 \text{width}(\mathcal{M})} \leq R\sqrt{\frac{2 \log(4nt_e^2/\delta)}{T_e - 1}} \leq R\sqrt{\frac{2 \log(4nT^2/\delta)}{T_e - 1}}. \quad (36)$$

Solving Eq. (36) for  $T_e$ , we obtain

$$T_e \leq \frac{18 \text{width}(\mathcal{M})^2 R^2}{\Delta_e^2} \log(4nT^2/\delta) + 1. \quad (37)$$

Notice that  $T = \sum_{i \in [n]} T_i$ . Hence the theorem follows by summing up Eq. (37) for all  $e \in [n]$  and solving for  $T$ . □

## B Extensions of CGapExp

CGapExp is a general and flexible learning algorithm for the ExpCMAB problem. In this section, we present two extensions to CGapExp that allow it to work in the fixed budget setting and PAC learning setting.

### B.1 Fixed Budget Setting

We can extend the CGapExp algorithm to the fixed budget setting using two simple modifications: (1) requiring CGapExp to terminate after  $T$  rounds; and (2) using a different construction of confidence intervals. The first modification ensures that CGapExp uses at most  $T$  samples, which meets the requirement of the fixed budget setting. And the second modification bounds the probability that the confidence intervals are valid for all arms in  $T$  rounds. The following theorem shows that the probability of error of the modified CGapExp is bounded by  $O\left(Tn \exp\left(\frac{-T}{\text{width}(\mathcal{M})^2 \mathbf{H}}\right)\right)$ .

**Theorem 4.** Use the same notations as in Theorem 1. Given  $T > n$  and parameter  $\alpha > 0$ , set the confidence radius  $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_e(t)}}$  for all arms  $e \in [n]$  and all  $t > 0$ . Run **CGapExp** algorithm for at most  $T$  rounds. Then, for  $0 \leq \alpha \leq \frac{1}{9}(T - n)(R^2 \text{width}(\mathcal{M})^2 \mathbf{H})^{-1}$ , we have

$$\Pr [\text{Out} \neq M_*] \leq 2Tn \exp(-2\alpha). \quad (38)$$

The right-hand side of Eq. (38) equals to  $O\left(Tn \exp\left(\frac{-T}{\text{width}(\mathcal{M})^2 \mathbf{H}}\right)\right)$  when parameter  $\alpha = O(T\mathbf{H}^{-1} \text{width}(\mathcal{M})^{-2})$ . For TOPK problem, we see that this matches the guarantees of the previous fixed budget algorithm due to Gabillon et al. [11].

## B.2 PAC Learning

Now we consider a setting where the learner is only required to report an approximately optimal set of arms. More specifically, we consider the notion of  $(\epsilon, \delta)$ -PAC algorithm. Formally, an algorithm  $\mathbb{A}$  is called an  $(\epsilon, \delta)$ -PAC algorithm if its output  $\text{Out}$  satisfies  $\Pr [w(M_*) - w(\text{Out}) > \epsilon] \leq \delta$ .

We show that a simple modification on the **CGapExp** algorithm gives an  $(\epsilon, \delta)$ -PAC algorithm, with guarantees similar to Theorem 1. In fact, the only modification needed is to change the stopping condition from  $\tilde{w}_t(\tilde{M}_t) \leq \tilde{w}_t(M_t)$  to  $w(\tilde{M}_t) - w(M_t) \leq \epsilon$  on line 15 of Algorithm 1. We let **CGapExpPAC** denote the modified algorithm. In the following theorem, we show that **CGapExpPAC** is indeed an  $(\epsilon, \delta)$ -PAC algorithm and has sample complexity similar to **CGapExp**.

**Theorem 5.** Use the same notations as in Theorem 1. Fix  $\delta \in (0, 1)$  and  $\epsilon \geq 0$ . Then, with probability at least  $1 - \delta$ , the output  $\text{Out}$  of **CGapExpPAC** satisfies  $w(M_*) - w(\text{Out}) \leq \epsilon$ . In addition, the number of samples  $T$  used by the algorithm satisfies

$$T \leq O\left(R^2 \sum_{e \in [n]} \min\left\{\frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2}\right\} \log\left(\frac{R^2 n}{\delta} \sum_{e \in [n]} \min\left\{\frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2}\right\}\right)\right), \quad (39)$$

where  $K = \max_{M \in \mathcal{M}} |M|$  is the size of the largest feasible solution.

We see that the sample complexity of **CGapExpPAC** decreases when  $\epsilon$  increases. And if  $\epsilon = 0$ , the sample complexity Eq. (39) of **CGapExpPAC** equals to that of **CGapExp**.

There are several PAC algorithms for the TOPK problem in the literature with different guarantees [16, 21, 11]. Zhou et al. [21] proposed an  $(\epsilon, \delta)$ -PAC algorithm for the TOPK problem with a problem-independent sample complexity bound of  $O\left(\frac{K^2 n}{\epsilon^2} + \frac{Kn \log(1/\delta)}{\epsilon^2}\right)$ .<sup>3</sup> If we ignore logarithmic factors, then the sample complexity bound of **CGapExpPAC** for the TOPK problem is better than theirs since  $\tilde{O}\left(\sum_{e \in [n]} \min\{\Delta_e^{-2}, K^2 \epsilon^{-2}\}\right) \leq \tilde{O}(nK^2 \epsilon^{-2})$ . On the other hand, the algorithms of Kalyanakrishnan et al. [16] and Gabillon et al. [11] guarantee to find  $K$  arms such that each of them is better than the  $K$ -th optimal arm within a factor of  $\epsilon$  with probability  $1 - \delta$ . Unless  $\epsilon = 0$ , their guarantee is different from ours which concerns the optimality of the sum of  $K$  arms.

## B.3 Proof of Extension Results

### B.3.1 Fixed Budget Setting (Theorem 4)

In this part, we analyze the probability of error of the modified **CGapExp** algorithm in the fixed budget setting and prove Theorem 4. First, we prove a lemma which characterizes the confidence intervals constructed in Theorem 4.

**Lemma 10.** Fix parameter  $\alpha > 0$  and the number of rounds  $T > 0$ . Assume that the reward distribution  $\varphi_e$  is a  $R$ -sub-Gaussian distribution for all  $e \in [n]$ . Let the confidence radius  $\text{rad}_t(e)$

<sup>3</sup>We notice that Zhou et al. [21] allow an  $(\epsilon', \delta)$ -PAC algorithm to produce an output with average sub-optimality of  $\epsilon'$  in their exposition. This is equivalent to our definition of  $(\epsilon, \delta)$ -PAC algorithm with  $\epsilon = K\epsilon'$  for the TOPK problem.

of arm  $e \in [n]$  and round  $t > 0$  be  $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_e(t)}}$ . Then, we have

$$\Pr \left[ \bigcap_{t=1}^T \xi_t \right] \geq 1 - 2nT \exp(-2\alpha).$$

*Proof.* For any  $t > 0$  and  $e \in [n]$ , using Hoeffding's inequality, we have

$$\Pr [|\bar{w}_t(e) - w(e)| \geq \text{rad}_t(e)] \leq 2 \exp(-2\alpha).$$

By a union bound over all arms  $e \in [n]$ , we see that  $\Pr[\xi_t] \geq 1 - 2n \exp(-2\alpha)$ . The lemma follows immediately by using union bound again over all round  $t \in [T]$ .  $\square$

Then, Theorem 4 can be obtained from the key lemmas (Lemma 8 and Lemma 9) and Lemma 10.

**Theorem 4.** Use the same notations as in Theorem 1. Given  $T > n$  and parameter  $\alpha > 0$ , set the confidence radius  $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_e(t)}}$  for all arms  $e \in [n]$  and all  $t > 0$ . Run *CGapExp* algorithm for at most  $T$  rounds. Then, for  $0 \leq \alpha \leq \frac{1}{9}(T - n) (R^2 \text{width}(\mathcal{M})^2 \mathbf{H})^{-1}$ , we have

$$\Pr [\text{Out} \neq M_*] \leq 2Tn \exp(-2\alpha). \quad (38)$$

*Proof.* Define random event  $\xi = \bigcap_{t=1}^T \xi_t$ . By Lemma 10, we see that  $\Pr[\xi] \geq 1 - 2nT \exp(-2\alpha)$ . In the rest of the proof, we assume that  $\xi$  happens.

Let  $T^*$  denote the round that the algorithm stops. We claim that the algorithm  $T^* < T$ . If the claim is true, then the algorithm stops since it meets the stopping condition on round  $T^*$ . Hence  $\tilde{M}_{T^*} = M_{T^*}$  and  $\text{Out} = M_{T^*}$ . By assumption on  $\xi$  and Lemma 8, we know that  $M_{T^*} = M_*$ . Therefore the theorem follows immediately from this claim and the bound of  $\Pr[\xi]$ .

Next, we show that this claim is true. Let  $t_e$  be the last round that arm  $e$  is pulled. Hence  $T_e(t_e) = T_e - 1$ . By Lemma 9, we see that  $\text{rad}_{t_e}(e) \geq \frac{\Delta}{3 \text{width}(\mathcal{B})}$ . Now plugging in the definition of  $\text{rad}_{t_e}(e)$ , we have

$$\begin{aligned} \frac{\Delta}{3 \text{width}(\mathcal{B})} &\leq \text{rad}_{t_e}(e) \\ &= R\sqrt{\frac{\alpha}{T_e(t_e)}} = R\sqrt{\frac{\alpha}{T_e - 1}}. \end{aligned}$$

Hence we have

$$T_e \leq \frac{9R^2 \text{width}(\mathcal{B})^2}{\Delta_e^2} \cdot \alpha + 1. \quad (40)$$

By summing up Eq. (40) for all  $e \in [n]$ , we have

$$T^* = \sum_{e \in [n]} T_e \leq \alpha \cdot 9R^2 \text{width}(\mathcal{B})^2 \left( \sum_{e \in [n]} \Delta_e^{-2} \right) + n < T,$$

where we have used the assumption that  $\alpha < \frac{1}{9}(T - n) \cdot \left( R^2 \text{width}(\mathcal{B})^2 \left( \sum_{e \in [n]} \Delta_e^{-2} \right) \right)^{-1}$ .  $\square$

### B.3.2 PAC Learning (Theorem 5)

First, we prove a  $(\epsilon, \delta)$ -PAC counterpart of Lemma 8.

**Lemma 11.** If *CGapExpPAC* stops on round  $t$  and suppose that event  $\xi_t$  occurs. Then, we have  $w(M_*) - w(\text{Out}) \leq \epsilon$ .

*Proof.* By definition, we know that  $\text{Out} = M_t$ . Notice that the stopping condition of CGapExpPAC ensures that  $\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \leq \epsilon$ . Therefore, we have

$$\begin{aligned} \epsilon &\geq \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \\ &\geq \tilde{w}_t(M_*) - \tilde{w}_t(M_t) \end{aligned} \quad (41)$$

$$= \langle \tilde{\mathbf{w}}_t, \chi_{M_*} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{M_*} - \chi_{M_t}| \rangle \quad (42)$$

$$\begin{aligned} &\geq \langle \mathbf{w}, \chi_{M_*} - \chi_{M_t} \rangle \\ &= w(M_*) - w(M_t), \end{aligned} \quad (43)$$

where Eq. (41) follows from the definition of  $\tilde{M}_t \triangleq \arg \max_{M \in \mathcal{M}} \tilde{w}_t(M)$ ; Eq. (42) follows from Lemma 5; Eq. (43) follows from the assumption that  $\xi_t$  occurs and Lemma 6.  $\square$

The next lemma generalizes Lemma 9. It shows that, with high probability, each arm  $e \in [n]$  will not be played on round  $t$  if  $\text{rad}_t(e) \leq \max \left\{ \frac{\Delta_e}{3 \text{width}(\mathcal{M})}, \frac{\epsilon}{2K} \right\}$ .

**Lemma 12.** *Let  $K = \max_{M \in \mathcal{M}} |M|$ . For any arm  $e \in [n]$  and any round  $t > n$  after initialization, if  $\text{rad}_t(e) \leq \max \left\{ \frac{\Delta_e}{3 \text{width}(\mathcal{M})}, \frac{\epsilon}{2K} \right\}$ , then arm  $e$  will not be played on round  $t$ , i.e.  $p_t \neq e$ .*

*Proof.* If  $\text{rad}_t(e) \leq \frac{\Delta_e}{3 \text{width}(\mathcal{M})}$ , then we can apply Lemma 9 which immediately gives that  $p_t \neq e$ . Hence, we only need to prove the case that  $\frac{\Delta_e}{3 \text{width}(\mathcal{M})} \leq \text{rad}_t(e) \leq \frac{\epsilon}{2K}$ .

Now suppose that  $p_t = e$ . By the choice of  $p_t$ , we know that for each  $i \in (M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)$ , we have  $\text{rad}_t(i) \leq \text{rad}_t(e) \leq \frac{\epsilon}{2K}$ . By summing up this inequality for all  $i \in (M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)$ , we have

$$\epsilon \geq \sum_{i \in (M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)} \text{rad}_t(i) \quad (44)$$

$$= \langle \mathbf{rad}_t, |\chi_{M_t} - \chi_{\tilde{M}_t}| \rangle, \quad (45)$$

where Eq. (44) follows from the fact that  $|(M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)| \leq |M_t| + |\tilde{M}_t| \leq 2K$ ; and Eq. (45) uses the fact that  $\chi_{(M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)} = |\chi_{M_t} - \chi_{\tilde{M}_t}|$ .

Then, we have

$$\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) = \langle \tilde{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \quad (46)$$

$$\leq \langle \tilde{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \epsilon \quad (47)$$

$$\begin{aligned} &= \bar{w}_t(\tilde{M}_t) - \bar{w}_t(M_t) + \epsilon \\ &\leq \epsilon, \end{aligned} \quad (48)$$

where Eq. (46) follows from Lemma 5; Eq. (47) uses Eq. (45); and Eq. (48) follows from  $\bar{w}_t(M_t) \geq \bar{w}_t(\tilde{M}_t)$ .

Therefore, we see that  $\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \leq \epsilon$ . By the stopping condition of CGapExpPAC, the algorithm must terminate on round  $t$ . This contradicts to the assumption that  $p_t = e$ .  $\square$

Using Lemma 12 and Lemma 11, we are ready to prove Theorem 5.

**Theorem 5.** *Use the same notations as in Theorem 1. Fix  $\delta \in (0, 1)$  and  $\epsilon \geq 0$ . Then, with probability at least  $1 - \delta$ , the output  $\text{Out}$  of CGapExpPAC satisfies  $w(M_*) - w(\text{Out}) \leq \epsilon$ . In addition, the number of samples  $T$  used by the algorithm satisfies*

$$T \leq O \left( R^2 \sum_{e \in [n]} \min \left\{ \frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2} \right\} \log \left( \frac{R^2 n}{\delta} \sum_{e \in [n]} \min \left\{ \frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2} \right\} \right) \right), \quad (39)$$

where  $K = \max_{M \in \mathcal{M}} |M|$  is the size of the largest feasible solution.

*Proof.* Similar to the proof of Theorem 1, we appeal to Lemma 7, which shows that the event  $\xi \triangleq \bigcap_{t=1}^{\infty} \xi_t$  occurs with probability at least  $1 - \delta$ . And we shall assume that  $\xi$  occurs in the rest of the proof.

By the assumption of  $\xi$  and Lemma 11, we know that  $\text{Out} = M_*$ . Therefore, we only remain to bound the number of samples  $T$ .

Consider an arbitrary arm  $e \in [n]$ . Let  $T_e$  denote the total number of pull of arm  $e \in [n]$ . Let  $t_e$  be the last round which arm  $e$  is pulled, i.e.  $p_{t_e} = e$ . Hence  $T_e(t_e) = T_e - 1$ . By Lemma 12, we see that  $\text{rad}_{t_e}(e) \geq \min\{\frac{\Delta_e}{3 \text{width}(\mathcal{B})}, \frac{\epsilon}{2K}\}$ . Then, by the construction of  $\text{rad}_{t_e}(e)$ , we have

$$\min\left\{\frac{\Delta_e}{3 \text{width}(\mathcal{B})}, \frac{\epsilon}{2K}\right\} \leq R \sqrt{\frac{2 \log(4nt_e^2/\delta)}{T_e - 1}} \leq R \sqrt{\frac{2 \log(4nT^2/\delta)}{T_e - 1}}. \quad (49)$$

Solving Eq. (49) for  $T_e$ , we obtain

$$T_e \leq R^2 \min\left\{\frac{18 \text{width}(\mathcal{B})^2}{\Delta_e^2}, \frac{16K^2}{\epsilon^2}\right\} \log(4nT^2/\delta) + 1. \quad (50)$$

Notice that  $T = \sum_{i \in [n]} T_i$ . Hence the theorem follows by summing up Eq. (50) for all  $e \in [n]$  and solving for  $T$ .  $\square$

## C Proof of Lower Bound

**Theorem 2.** Fix any decision class  $\mathcal{M} \subseteq 2^{[n]}$  and any vector  $\mathbf{w} \in \mathbb{R}^n$ . Suppose that, for each arm  $e \in [n]$ , the reward distribution  $\varphi_e$  is given by  $\varphi_e = \mathcal{N}(w(e), 1)$ , where  $\mathcal{N}(\mu, \sigma^2)$  denotes a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ . Then, for any  $\delta \in (0, e^{-16}/4)$  and any  $\delta$ -correct algorithm  $\mathbb{A}$ , we have

$$\mathbb{E}[T] \geq \frac{1}{16} \mathbf{H} \log\left(\frac{1}{4\delta}\right), \quad (6)$$

where  $T$  denote the number of total samples used by algorithm  $\mathbb{A}$  and  $\mathbf{H}$  is defined in Eq. (2).

*Proof.* Fix  $\delta > 0$ ,  $\mathbf{w} = (w(1), \dots, w(n))^T$  and a  $\delta$ -correct algorithm  $\mathbb{A}$ . For each  $e \in [n]$ , assume that the reward distribution is given by  $\varphi_e = \mathcal{N}(w(e), 1)$ . For any  $e \in [n]$ , let  $T_e$  denote the number of trials of arm  $e$  used by algorithm  $\mathbb{A}$ . In the rest of the proof, we will show that for any  $e \in [n]$ , the number of trials of arm  $e$  is lower-bounded by

$$\mathbb{E}[T_e] \geq \frac{1}{16\Delta_e^2} \log(1/4\delta). \quad (51)$$

Notice that the theorem follows immediately by summing up Eq. (51) for all  $e \in [n]$ .

Fix an arm  $e \in [n]$ . We now focus on proving Eq. (51). Consider two hypothesis  $H_0$  and  $H_1$ . Under hypothesis  $H_0$ , all reward distributions are same with our assumption before

$$H_0 : \varphi_l = \mathcal{N}(w(l), 1) \quad \text{for all } l \in [n].$$

Under hypothesis  $H_1$ , we change the means of reward distributions such that

$$H_1 : \varphi_e = \begin{cases} \mathcal{N}(w(e) - 2\Delta_e, 1) & \text{if } e \in M_* \\ \mathcal{N}(w(e) + 2\Delta_e, 1) & \text{if } e \notin M_* \end{cases} \quad \text{and } \varphi_l = \mathcal{N}(w(l), 1) \quad \text{for all } l \neq e.$$

For  $l \in \{0, 1\}$ , we use  $\mathbb{E}_l$  and  $\Pr_l$  to denote the expectation and probability, respectively, under the hypothesis  $H_l$ .

Define  $M_e$  be the “next-to-optimal” set as follows

$$M_e = \begin{cases} \arg \max_{M \in \mathcal{M}: e \in M} w(M) & \text{if } e \notin M_*, \\ \arg \max_{M \in \mathcal{M}: e \notin M} w(M) & \text{if } e \in M_*. \end{cases}$$

By definition of  $\Delta_e$  in Eq. (1), we know that  $w(M_*) - w(M_e) = \Delta_e$ .



Let  $w_0$  and  $w_1$  be expected reward vectors under  $H_0$  and  $H_1$  respectively. Notice that  $w_0(M_*) - w_0(M_e) = \Delta_e > 0$ . On the other hand, we have

$$\begin{aligned} w_1(M_*) - w_1(M_e) &= w(M_*) - w(M_e) - 2\Delta_e \\ &= -\Delta_e < 0. \end{aligned}$$

This means that under  $H_1$ , the set  $M_*$  is not the optimal set.

Define  $\theta = 4\delta$ . Define

$$t_e^* = \frac{1}{16\Delta_e^2} \log\left(\frac{1}{\theta}\right). \quad (52)$$

Recall that  $T_e$  denotes the total number of samples of arm  $e$ . Define the event  $\mathcal{A} = \{T_e \leq 4t_e^*\}$ .

First, we show that  $\Pr_0[\mathcal{A}] \geq 3/4$ . This can be proved by Markov inequality as follows.

$$\begin{aligned} \Pr_0[T_e > 4t_e^*] &\leq \frac{\mathbb{E}_0[T_e]}{4t_e^*} \\ &= \frac{t_e^*}{4t_e^*} = \frac{1}{4}. \end{aligned}$$

Let  $X_1, \dots, X_{T_e}$  denote the sequence of reward outcomes of arm  $e$ . For all  $t > 0$ , we define  $K_t = \sum_{i \in [t]} X_i$  as the sum of outcomes of arm  $e$  up to round  $t$ . Next, we define the event

$$\mathcal{C} = \left\{ \max_{1 \leq t \leq 4t_e^*} |K_t - t \cdot w(e)| < \sqrt{t_e^* \log(1/\theta)} \right\}.$$

We now show that  $\Pr_0[\mathcal{C}] \geq 3/4$ . First, notice that  $\{K_t - t \cdot w(e)\}_{t=1, \dots}$  is a martingale under  $H_0$ . Then, by Kolmogorov's inequality, we have

$$\begin{aligned} \Pr_0 \left[ \max_{1 \leq t \leq 4t_e^*} |K_t - t \cdot w(e)| \geq \sqrt{t_e^* \log(1/\theta)} \right] &\leq \frac{\mathbb{E}_0[(K_{4t_e^*} - 4w(e)t_e^*)^2]}{t_e^* \log(1/\theta)} \\ &= \frac{4t_e^*}{t_e^* \log(1/\theta)} \\ &< \frac{1}{4}, \end{aligned}$$

where the second inequality follows from the fact that the variance of  $\varphi_e$  equals to 1 and therefore  $\mathbb{E}_0[(K_{4t_e^*} - 4w(e)t_e^*)^2] = 4t_e^*$ ; the last inequality follows since  $\theta < e^{-16}$ .

Then, we define the event  $\mathcal{B}$  as the event that the algorithm eventually returns  $M_*$ , i.e.

$$\mathcal{B} = \{\text{Out} = M_*\}.$$

Since the probability of error of the algorithm is smaller than  $\delta < 1/4$ , we have  $\Pr_0[\mathcal{B}] \geq 3/4$ . Define  $\mathcal{S}$  be  $\mathcal{S} = \mathcal{A} \cap \mathcal{B} \cap \mathcal{C}$ . Then, by union bound, we have  $\Pr_0[\mathcal{S}] \geq 1/4$ .

Now, we show that if  $\mathbb{E}_0[T_e] \leq t_e^*$ , then  $\Pr_1[\mathcal{B}] \geq \delta$ . Let  $W$  be the history of the sampling process until the algorithm stops (including the sequence of arms chosen at each time and the sequence of observed outcomes). Define the likelihood function  $L_l$  as

$$L_l(w) = p_l(W = w),$$

where  $p_l$  is the probability density function under hypothesis  $H_l$ . Let  $K$  be the shorthand of  $K_{T_e}$ .

Assume that the event  $\mathcal{S}$  occurred. We will bound the likelihood ratio  $L_1(W)/L_0(W)$  under this assumption. To do this, we divide our analysis into two different cases.

**Case (1):**  $e \notin M_*$ . In this case, the reward distribution of arm  $e$  under  $H_1$  is a Gaussian distribution with mean  $w(e) + 2\Delta_e$  and variance 1. Recall that the probability density function of a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$  is given by  $\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$ . Hence, we have

$$\frac{L_1(W)}{L_0(W)} = \prod_{i=1}^{T_e} \exp\left(\frac{-(X_i - w(e) - 2\Delta_e)^2 + (X_i - w(e))^2}{2}\right)$$

$$\begin{aligned}
&= \prod_{i=1}^{T_e} \exp(\Delta_e(2X_i - 2w(e)) - 2\Delta_e^2) \\
&= \exp(\Delta_e(2K - 2w(e)T_e) - 2\Delta_e^2 T_e) \\
&= \exp(\Delta_e(2K - 2w(e)T_e)) \exp(-2\Delta_e^2 T_e).
\end{aligned} \tag{53}$$

Next, we bound each individual term on the right-hand side of Eq. (53). We begin with bounding the second term of Eq. (53)

$$\exp(-2\Delta_e^2 T_e) \geq \exp(-8\Delta_e^2 t_e^*) \tag{54}$$

$$= \exp\left(-\frac{8}{16} \log(1/\theta)\right) \tag{55}$$

$$= \theta^{1/2}, \tag{56}$$

where Eq. (54) follows from the assumption that event  $\mathcal{S}$  occurred, which implies that event  $\mathcal{A}$  occurred and therefore  $T_e \leq 4t_e^*$ ; Eq. (55) follows from the definition of  $t_e^*$ .

Then, we bound the first term on the right-hand side of Eq. (53) as follows

$$\exp(\Delta_e(2K - 2w(e)T_e)) \geq \exp\left(-2\Delta_e \sqrt{t_e^* \log(1/\theta)}\right) \tag{57}$$

$$= \exp\left(-\frac{2}{4} \log(1/\theta)\right) \tag{58}$$

$$= \theta^{1/2}, \tag{59}$$

where Eq. (57) follows from the assumption that event  $\mathcal{S}$  occurred, which implies that event  $\mathcal{C}$  and therefore  $|2K - 2w(e)T_e| \leq \sqrt{t_e^* \log(1/\theta)}$ ; Eq. (58) follows from the definition of  $t_e^*$ .

Combining Eq. (56) and Eq. (59), we can bound  $L_1(W)/L_0(W)$  for this case as follows

$$\frac{L_1(W)}{L_0(W)} \geq \theta. \tag{60}$$

(End of Case (1).)

**Case (2):**  $e \in M_*$ . In this case, we know that the mean reward of arm  $e$  under  $H_1$  is  $w(e) - 2\Delta$ . Therefore, the likelihood ratio  $L_1(W)/L_0(W)$  is given by

$$\begin{aligned}
\frac{L_1(W)}{L_0(W)} &= \prod_{i=1}^{T_e} \exp\left(\frac{-(X_i - w(e) + 2\Delta_e)^2 + (X_i - w(e))^2}{2}\right) \\
&= \prod_{i=1}^{T_e} \exp(\Delta_e(2w(e) - 2X_i) - 2\Delta_e^2) \\
&= \exp(\Delta_e(2w(e)T_e - 2K)) \exp(-2\Delta_e^2 T_e).
\end{aligned} \tag{61}$$

Notice that the right-hand side of Eq. (61) differs from Eq. (53) only in its first term. Now, we bound the first term as follows

$$\exp(\Delta_e(2w(e)T_e - 2K)) \geq \exp\left(-2\Delta_e \sqrt{t_e^* \log(1/\theta)}\right) \tag{62}$$

$$= \exp\left(-\frac{2}{4} \log(1/\theta)\right) \tag{63}$$

$$= \theta^{1/2}, \tag{64}$$

where the inequalities hold due to reasons similar to Case (1): Eq. (62) follows from the assumption that event  $\mathcal{S}$  occurred, which implies that event  $\mathcal{C}$  and therefore  $|2K - 2w(e)T_e| \leq \sqrt{t_e^* \log(1/\theta)}$ ; Eq. (63) follows from the definition of  $t_e^*$ .

Combining Eq. (56) and Eq. (59), we can obtain the same bound of  $L_1(W)/L_0(W)$  as in Eq. (60), i.e.  $L_1(W)/L_0(W) \geq \theta$ .

(End of Case (2).)

At this point, we have proved that, if the event  $\mathcal{S}$  occurred, then the bound of likelihood ratio Eq. (60) holds, i.e.  $\frac{L_1(W)}{L_0(W)} \geq \theta$ . Hence, we have

$$\begin{aligned} \frac{L_1(W)}{L_0(W)} &\geq \theta \\ &= 4\delta. \end{aligned} \tag{65}$$

Define  $1_S$  as the indicator variable of event  $\mathcal{S}$ , i.e.  $1_S = 1$  if and only if  $\mathcal{S}$  occurs and otherwise  $1_S = 0$ . Then, we have

$$\frac{L_1(W)}{L_0(W)} 1_S \geq 4\delta 1_S$$

holds regardless the occurrence of event  $\mathcal{S}$ . Therefore, we can obtain

$$\begin{aligned} \Pr_1[\mathcal{B}] &\geq \Pr_1[\mathcal{S}] = \mathbb{E}_1[1_S] \\ &= \mathbb{E}_0 \left[ \frac{L_1(W)}{L_0(W)} 1_S \right] \\ &\geq 4\delta \mathbb{E}_0[1_S] \\ &= 4\delta \Pr_0[\mathcal{S}] > \delta. \end{aligned}$$

Now we have proved that, if  $\mathbb{E}_0[T_e] \leq t_e^*$ , then  $\Pr_1[\mathcal{B}] > \delta$ . This means that, if  $\mathbb{E}_0[T_e] \leq t_e^*$ , algorithm  $\mathbb{A}$  will choose  $M_*$  as the output with probability at least  $\delta$ , under hypothesis  $H_1$ . However, under  $H_1$ , we have shown that  $M_*$  is not the optimal set since  $w_1(M_e) > w_1(M_*)$ . Therefore, algorithm  $\mathbb{A}$  has a probability of error at least  $\delta$  under  $H_1$ . This contradicts to the assumption that algorithm  $\mathbb{A}$  is a  $\delta$ -correct algorithm. Hence, we must have  $\mathbb{E}_0[T_e] > t_e^* = \frac{1}{16\Delta_e^2} \log(1/4\delta)$ .  $\square$

### C.1 Exchange set size dependent lower bound

We show that, for any arm  $e \in [n]$ , there exists an exchange set  $b$  which contains  $e$  such that a  $\delta$ -correct algorithm must spend  $\tilde{\Omega}\left(\left(|b_+| + |b_-|\right)^2 / \Delta_e^2\right)$  samples on the arms belonging to  $b$ . This result is formalized in the following theorem.

**Theorem 6.** Fix any  $\mathcal{M} \subseteq 2^{[n]}$  and any vector  $\mathbf{w} \in \mathbb{R}^n$ . Suppose that, for each arm  $e \in [n]$ , the reward distribution  $\varphi_e$  is given by  $\varphi_e = \mathcal{N}(w(e), 1)$ , where  $\mathcal{N}(\mu, \sigma^2)$  denotes a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ . Fix any  $\delta \in (0, e^{-16}/4)$  and any  $\delta$ -correct algorithm  $\mathbb{A}$ .

Then, for any  $e \in [n]$ , there exists an exchange set  $b = (b_+, b_-)$ , such that  $e \in b_+ \cup b_-$  and

$$\mathbb{E} \left[ \sum_{i \in b_+ \cup b_-} T_i \right] \geq \frac{(|b_+| + |b_-|)^2}{32\Delta_e^2} \log(1/4\delta),$$

where  $T_i$  is the number of samples of arm  $i$ .

*Proof.* Fix  $\delta > 0$ ,  $\mathbf{w} \in \mathbb{R}^n$ , diff-set  $b = (b_+, b_-)$  and a  $\delta$ -correct algorithm  $\mathbb{A}$ . Assume that  $\varphi_e(e) = \mathcal{N}(w(e), 1)$  for all  $e \in [n]$ .

We define three hypotheses  $H_0$ ,  $H_1$  and  $H_2$ . Under hypothesis  $H_0$ , the reward distribution

$$H_0 : \varphi_l = \mathcal{N}(w(l), 1) \quad \text{for all } l \in [n].$$

Under hypothesis  $H_1$ , the mean reward of each arm is given by

$$H_1 : \varphi_e = \begin{cases} \mathcal{N}\left(w(e) + 2\frac{w(b)}{|b_-|}, 1\right) & \text{if } e \in b_-, \\ \mathcal{N}(w(e), 1) & \text{if } e \notin b_-. \end{cases}$$

And under hypothesis  $H_2$ , the mean reward of each arm is given by

$$H_2 : \varphi_e = \begin{cases} \mathcal{N}\left(w(e) - 2\frac{w(b)}{|b_-|}, 1\right) & \text{if } e \in b_+, \\ \mathcal{N}(w(e), 1) & \text{if } e \notin b_+. \end{cases}$$

Since  $b \in \mathcal{B}_{\text{opt}}$ , it is clear that  $\neg b \prec M_*$ . Hence we define  $M = M_* \ominus b$ . Let  $w_0, w_1$  and  $w_2$  be the expected reward vectors under  $H_0, H_1$  and  $H_2$  respectively. It is easy to check that  $w_1(M_*) - w_1(M) = -w(b) < 0$  and  $w_2(M_*) - w_2(M) = -w(b) < 0$ . This means that under  $H_1$  or  $H_2$ ,  $M_*$  is not the optimal set. Further, for  $l \in \{0, 1, 2\}$ , we use  $\mathbb{E}_l$  and  $\Pr_l$  to denote the expectation and probability, respectively, under the hypothesis  $H_l$ . In addition, let  $W$  be the history of the sampling process until algorithm  $\mathbb{A}$  stops. Define the likelihood function  $L_l$  as

$$L_l(w) = p_l(W = w),$$

where  $p_l$  is the probability density function under  $H_l$ .

Define  $\theta = 4\delta$ . Let  $T_{b_-}$  and  $T_{b_+}$  denote the number of trials of arms belonging to  $b_-$  and  $b_+$ , respectively. In the rest of the proof, we will bound  $\mathbb{E}_0[T_{b_-}]$  and  $\mathbb{E}_0[T_{b_+}]$  individually.

**Part (1): Lower bound of  $\mathbb{E}_0[T_{b_-}]$ .** In this part, we will show that  $\mathbb{E}_0[T_{b_-}] \geq t_{b_-}^*$ , where we define  $t_{b_-}^* = \frac{|b_-|^2}{16w(b)^2} \log(1/\theta)$ .

Consider the complete sequence of sampling process by algorithm  $\mathbb{A}$ . Formally, let  $W = \{(\tilde{I}_1, \tilde{X}_1), \dots, (\tilde{I}_T, \tilde{X}_T)\}$  be the sequence of all trials by algorithm  $\mathbb{A}$ , where  $\tilde{I}_i$  denotes the arm played in  $i$ -th trial and  $\tilde{X}_i$  be the reward outcome of  $i$ -th trial. Then, consider the subsequence  $W_1$  of  $W$  which consists all the trials of arms in  $b_-$ . Specifically, we write  $W = \{(I_1, X_1), \dots, (I_{T_{b_-}}, X_{T_{b_-}})\}$  such that  $W_1$  is a subsequence of  $W$  and  $I_i \in b_-$  for all  $i$ .

Next, we define several random events in a way similar to the proof of Theorem 2. Define event  $\mathcal{A}_1 = \{T_{b_-} \leq 4t_{b_-}^*\}$ . Define event

$$\mathcal{C}_1 = \left\{ \max_{1 \leq t \leq 4t_{b_-}^*} \left| \sum_{i=1}^t X_i - \sum_{i=1}^t w(I_i) \right| < \sqrt{t_{b_-}^* \log(1/\theta)} \right\}.$$

Define event

$$\mathcal{B} = \{\text{Out} = M_*\}. \quad (66)$$

Define event  $\mathcal{S}_1 = \mathcal{A}_1 \cap \mathcal{B} \cap \mathcal{C}_1$ . Then, we bound the probability of events  $\mathcal{A}_1, \mathcal{B}, \mathcal{C}_1$  and  $\mathcal{S}_1$  under  $H_0$  using methods similar to Theorem 2. First, we show that  $\Pr_0[\mathcal{A}_1] \geq 3/4$ . This can be proved by Markov inequality as follows.

$$\begin{aligned} \Pr_0[T_{b_-} > 4t_{b_-}^*] &\leq \frac{\mathbb{E}_0[T_{b_-}]}{4t_{b_-}^*} \\ &= \frac{t_{b_-}^*}{4t_{b_-}^*} = \frac{1}{4}. \end{aligned}$$

Next, we show that  $\Pr_0[\mathcal{C}_1] \geq 3/4$ . Notice that the sequence  $\left\{ \sum_{i=1}^t X_i - \sum_{i=1}^t p_{I_i} \right\}_{t \in [4t_{b_-}^*]}$  is a martingale. Hence, by Kolmogorov's inequality, we have

$$\begin{aligned} \Pr_0 \left[ \max_{1 \leq t \leq 4t_{b_-}^*} \left| \sum_{i=1}^t X_i - \sum_{i=1}^t w(I_i) \right| \geq \sqrt{t_{b_-}^* \log(1/\theta)} \right] &\leq \frac{\mathbb{E}_0 \left[ \left( \sum_{i=1}^{4t_{b_-}^*} X_i - \sum_{i=1}^{4t_{b_-}^*} w(I_i) \right)^2 \right]}{t_{b_-}^* \log(1/\theta)} \\ &= \frac{4t_{b_-}^*}{t_{b_-}^* \log(1/\theta)} \\ &< \frac{1}{4}, \end{aligned}$$

where the second inequality follows from the fact that all reward distributions have unit variance and hence  $\mathbb{E}_0 \left[ \left( \sum_{i=1}^{4t_{b_-}^*} X_i - \sum_{i=1}^{4t_{b_-}^*} p_{I_i} \right)^2 \right] = 4t_{b_-}^*$ ; the last inequality follows since  $\theta < e^{-16}$ .

Last, since algorithm  $\mathbb{A}$  is a  $\delta$ -correct algorithm with  $\delta < 1/4$ . Therefore, it is easy to see that  $\Pr_0[\mathcal{B}] \geq 3/4$ . And by union bound, we have

$$\Pr_0[\mathcal{S}_1] \geq 1/4.$$

Now, we show that if  $\mathbb{E}_0[T_{b_-}] \leq t_{b_-}^*$ , then  $\Pr_1[\mathcal{B}] \geq \delta$ . Assume that the event  $\mathcal{S}_1$  occurred. We bound the likelihood ratio  $L_1(W)/L_0(W)$  under this assumption as follows

$$\begin{aligned}
\frac{L_1(W)}{L_0(W)} &= \prod_{i=1}^{T_{b_-}} \exp \left( \frac{-\left(X_i - w(I_i) - \frac{2w(b)}{|b_-|}\right)^2 + (X_i - w(I_i))^2}{2} \right) \\
&= \prod_{i=1}^{T_{b_-}} \exp \left( \frac{w(b)}{|b_-|} (2X_i - 2w(I_i)) - \frac{2w(b)^2}{|b_-|^2} \right) \\
&= \exp \left( \frac{w(b)}{|b_-|} \left( \sum_{i=1}^{T_{b_-}} 2X_i - 2w(I_i) \right) - \frac{2w(b)^2}{|b_-|^2} T_{b_-} \right) \\
&= \exp \left( \frac{w(b)}{|b_-|} \left( \sum_{i=1}^{T_{b_-}} 2X_i - 2w(I_i) \right) \right) \exp \left( -\frac{2w(b)^2}{|b_-|^2} T_{b_-} \right). \tag{67}
\end{aligned}$$

Then, we bound each term on the right-hand side of Eq. (67). First, we bound the second term of Eq. (67).

$$\exp \left( -\frac{2w(b)^2}{|b_-|^2} T_{b_-} \right) \geq \exp \left( -\frac{2w(b)^2}{|b_-|^2} 4t_{b_-}^* \right) \tag{68}$$

$$= \exp \left( -\frac{8}{16} \log(1/\theta) \right) \tag{69}$$

$$= \theta^{1/2}, \tag{70}$$

where Eq. (68) follows from the assumption that events  $\mathcal{S}_1$  and  $\mathcal{A}_1$  occurred and therefore  $T_{b_-} \leq 4t_{b_-}^*$ ; Eq. (69) follows from the definition of  $t_{b_-}^*$ . Next, we bound the first term of Eq. (67) as follows

$$\exp \left( \frac{w(b)}{|b_-|} \left( \sum_{i=1}^{T_{b_-}} 2X_i - 2w(I_i) \right) \right) \geq \exp \left( -\frac{2w(b)}{|b_-|} \sqrt{t_{b_-}^* \log(1/\theta)} \right) \tag{71}$$

$$= \exp \left( -\frac{2}{4} \log(1/\theta) \right) \tag{72}$$

$$= \theta^{1/2}, \tag{73}$$

where Eq. (71) follows since event  $\mathcal{S}_1$  and  $\mathcal{C}_1$  occurred and therefore  $|2K - 2p_e T_e| \leq \sqrt{t_e^* \log(1/\theta)}$ ; Eq. (72) follows from the definition of  $t_{b_-}^*$ .

Hence, if event  $\mathcal{S}_1$  occurred, we can bound the likelihood ratio as follows

$$\frac{L_1(W)}{L_0(W)} \geq \theta = 4\delta. \tag{74}$$

Let  $1_{\mathcal{S}_1}$  denote the indicator variable of event  $\mathcal{S}_1$ . Then, we have  $\frac{L_1(W)}{L_0(W)} 1_{\mathcal{S}_1} \geq 4\delta 1_{\mathcal{S}_1}$ . Therefore, we can bound  $\Pr_1[\mathcal{B}]$  as follows

$$\begin{aligned}
\Pr_1[\mathcal{B}] &\geq \Pr_1[\mathcal{S}_1] = \mathbb{E}_1[1_{\mathcal{S}_1}] \\
&= \mathbb{E}_0 \left[ \frac{L_1(W)}{L_0(W)} 1_{\mathcal{S}_1} \right] \\
&\geq 4\delta \mathbb{E}_0[1_{\mathcal{S}_1}] \\
&= 4\delta \Pr_0[\mathcal{S}_1] > \delta. \tag{75}
\end{aligned}$$

This means that, if  $\mathbb{E}_0[T_{b_-}] \leq t_{b_-}^*$ , then, under  $H_1$ , the probability of algorithm  $\mathbb{A}$  returning  $M_*$  as output is at least  $\delta$ . But  $M_*$  is not the optimal set under  $H_1$ . Hence this contradicts to the assumption that  $\mathbb{A}$  is a  $\delta$ -correct algorithm. Hence we have proved that

$$\mathbb{E}_0[T_{b_-}] \geq t_{b_-}^* = \frac{|b_-|^2}{16w(b)^2} \log(1/4\delta). \tag{76}$$

(End of Part (1).)

**Part (2): Lower bound of  $\mathbb{E}_0[T_{b_+}]$ .** In this part, we will show that  $\mathbb{E}_0[T_{b_+}] \geq t_{b_+}^*$ , where we define  $t_{b_+}^* = \frac{|b_+|^2}{16w(b)^2} \log(1/\theta)$ . The arguments used in this part are similar to that of Part (1). Hence, we will omit the redundant parts and highlight the differences.

Recall that we have defined that  $W$  to be the history of all trials by algorithm  $\mathbb{A}$ . We define  $W$  be the subsequence of  $\tilde{S}$  which contains the trials of arms belonging to  $b_+$ . We write  $S_2 = \{(J_1, Y_1), \dots, (J_{T_{b_+}}, Y_{T_{b_+}})\}$ , where  $J_i$  is  $i$ -th played arm in sequence  $S_2$  and  $Y_i$  is the associated reward outcome.

We define the random events  $\mathcal{A}_2$  and  $\mathcal{C}_2$  similar to Part (1). Specifically, we define

$$\mathcal{A}_2 = \{T_{b_+} \leq 4t_{b_+}^*\} \quad \text{and} \quad \mathcal{C}_2 = \left\{ \max_{1 \leq t \leq 4t_{b_+}^*} \left| \sum_{i=1}^t Y_i - \sum_{i=1}^t w(J_i) \right| < \sqrt{t_{b_+}^* \log(1/\theta)} \right\}.$$

Using the similar arguments, we can show that  $\Pr_0[\mathcal{A}_2] \geq 3/4$  and  $\Pr_0[\mathcal{C}_2] \geq 3/4$ . Define event  $\mathcal{S}_2 = \mathcal{A}_2 \cap \mathcal{B} \cap \mathcal{C}_2$ , where  $\mathcal{B}$  is defined in Eq. (66). By union bound, we see that

$$\Pr_0[\mathcal{S}_2] \geq 1/4.$$

Then, we show that if  $\mathbb{E}_0[T_{b_+}] \leq t_{b_+}^*$ , then  $\Pr_2[\mathcal{B}] \geq \delta$ . We bound likelihood ratio  $L_2(W)/L_0(W)$  under the assumption that  $\mathcal{S}_2$  occurred as follows

$$\begin{aligned} \frac{L_2(W)}{L_0(W)} &= \prod_{i=1}^{T_{b_+}} \exp \left( \frac{-\left(Y_i - w(J_i)\right) + \frac{2w(b)}{|b_-|} + (Y_i - w(J_i))^2}{2} \right) \\ &= \prod_{i=1}^{T_{b_+}} \exp \left( \frac{w(b)}{|b_+|} (2w(J_i) - 2Y_i) - \frac{2w(b)^2}{|b_+|^2} \right) \\ &= \exp \left( \frac{w(b)}{|b_+|} \left( \sum_{i=1}^{T_{b_+}} 2w(J_i) - 2Y_i \right) - \frac{2w(b)^2}{|b_+|^2} T_{b_+} \right) \\ &= \exp \left( \frac{w(b)}{|b_+|} \left( \sum_{i=1}^{T_{b_+}} 2w(J_i) - 2Y_i \right) \right) \exp \left( -\frac{2w(b)^2}{|b_+|^2} T_{b_+} \right) \\ &\geq \theta \\ &= 4\delta, \end{aligned} \tag{77}$$

where Eq. (77) can be obtained using same method as in Part (1) as well as the assumption that  $\mathcal{S}_2$  occurred.

Next, similar to the derivation in Eq. (75), we see that

$$\Pr_2[\mathcal{B}] \geq \Pr_2[\mathcal{S}_2] = \mathbb{E}_2[1_{\mathcal{S}_2}] = \mathbb{E}_0 \left[ \frac{L_2(W)}{L_0(W)} 1_{\mathcal{S}_2} \right] \geq 4\delta \mathbb{E}_0[1_{\mathcal{S}_2}] > \delta,$$

where  $1_{\mathcal{S}_2}$  is the indicator variable of event  $\mathcal{S}_2$ . Therefore, we see that if  $\mathbb{E}_0[T_{b_+}] \leq t_{b_+}^*$ , then, under  $H_2$ , the probability of algorithm  $\mathbb{A}$  returning  $M_*$  as output is at least  $\delta$ , which is not the optimal set under  $H_2$ . This contradicts to the assumption that algorithm  $\mathbb{A}$  is a  $\delta$ -correct algorithm. In sum, we have proved that

$$\mathbb{E}_0[T_{b_+}] \geq t_{b_+}^* = \frac{|b_+|^2}{16w(b)^2} \log(1/4\delta). \tag{78}$$

(End of Part (2))

Finally, we combine the results from both parts, i.e. Eq. (76) and Eq. (78). We obtain

$$\mathbb{E}_0[T_b] = \mathbb{E}_0[T_{b_-}] + \mathbb{E}_0[T_{b_+}]$$



$$\begin{aligned}
&\geq \frac{|b_+|^2 + |b_-|^2}{16w(b)^2} \log(1/4\delta) \\
&\geq \frac{|b|^2}{32w(b)^2} \log(1/4\delta).
\end{aligned}$$

□

## D Analysis of CGapKill

**Notations.** For convenience, we will use the following notations in the rest of this section. Let  $\mathbf{w} \in \mathbb{R}^n$  be the vector expected rewards of arms. Let  $M_* = \arg \max_{M \in \mathcal{M}} w(M)$  be the optimal solution. Let  $T$  be the budget of samples. Let  $\Delta_{(1)}, \dots, \Delta_{(n)}$  be a permutation of  $\Delta_1, \dots, \Delta_n$  such that  $\Delta_{(1)} \leq \dots \leq \Delta_{(n)}$ . Let  $A_1, \dots, A_n$  and  $B_1, \dots, B_n$  be two sequence of sets which are defined in Algorithm X.

### D.1 Confidence Intervals

**Lemma 13.** *Given a phase  $t \in [n]$ , we define random event  $\tau_t$  as follows*

$$\tau_t = \left\{ \forall i \in [n] \setminus (A_t \cup B_t) \quad |\bar{w}_t(i) - w(i)| < \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \right\}. \quad (79)$$

*Then, we have*

$$\Pr \left[ \bigcap_{t=1}^T \tau_t \right] \geq 1 - n^2 \exp \left( - \frac{2(T-n)}{9R^2 \tilde{\log}(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right). \quad (80)$$

*Proof.* Let us consider an arbitrary phase  $t \in [n]$  and an arbitrary active arm  $i \in [n] \setminus (A_t \cup B_t)$  of phase  $t$ .

Notice that the arm  $e$  has been pulled for  $\tilde{T}_t$  times during phases  $1, \dots, t$ . Therefore, by Hoeffding's inequality, we have

$$\Pr \left[ |\bar{w}_t(i) - w(i)| \geq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \right] \leq 2 \exp \left( - \frac{2\tilde{T}_t \Delta_{(n-t+1)}^2}{9R^2 \text{width}(\mathcal{M})^2} \right). \quad (81)$$

By plugging the definition of  $\tilde{T}_t$ , the quantity  $\tilde{T}_t \Delta_{(n-t+1)}^2$  on the right-hand side of Eq. (81) can be further bounded by

$$\begin{aligned}
\tilde{T}_t \Delta_{(n-t+1)}^2 &\geq \frac{T-n}{\tilde{\log}(n)(n-t+1)} \Delta_{(n-t+1)}^2 \\
&\geq \frac{T-n}{\tilde{\log}(n) \mathbf{H}_2},
\end{aligned}$$

where the last inequality follows from the definition of  $\mathbf{H}_2$ . By plugging the last inequality into Eq. (81), we have

$$\Pr \left[ |\bar{w}_t(i) - w(i)| \geq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \right] \leq 2 \exp \left( - \frac{2(T-n)}{9R^2 \tilde{\log}(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right). \quad (82)$$

Now using Eq. (82) and a union bound for all  $t \in [n]$  and all  $i \in [n] \setminus (A_t \cup B_t)$ , we have

$$\begin{aligned}
\Pr \left[ \bigcap_{t=1}^n \tau_t \right] &\geq 1 - 2 \sum_{t=1}^n (n-t+1) \exp \left( - \frac{2(T-n)}{9R^2 \tilde{\log}(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right) \\
&\geq 1 - n^2 \exp \left( - \frac{2(T-n)}{9R^2 \tilde{\log}(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right).
\end{aligned}$$

□

**Lemma 14.** Fix a phase  $t \in [n]$ , suppose that random event  $\tau_t$  occurs. For any vector  $\mathbf{a} \in \mathbb{R}^n$ , suppose that  $\text{supp}(\mathbf{a}) \cap (A_t \cup B_t) = \emptyset$ , where  $\text{supp}(\mathbf{a}) \triangleq \{i \mid a(i) \neq 0\}$  is support of  $\mathbf{a}$ . Then, we have

$$|\langle \bar{\mathbf{w}}_t, \mathbf{a} \rangle - \langle \mathbf{w}, \mathbf{a} \rangle| \leq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\mathbf{a}\|_1.$$

*Proof.* Suppose that  $\tau_t$  occurs. Then, similar to the proof of Lemma 6, we have

$$\begin{aligned} |\langle \bar{\mathbf{w}}_t, \mathbf{a} \rangle - \langle \mathbf{w}, \mathbf{a} \rangle| &= |\langle \bar{\mathbf{w}}_t - \mathbf{w}, \mathbf{a} \rangle| \\ &= \left| \sum_{i=1}^n (\bar{w}_t(i) - w(i)) a(i) \right| \\ &\leq \left| \sum_{i \in [n] \setminus (A_t \cup B_t)} (\bar{w}_t(i) - w(i)) a(i) \right| \end{aligned} \quad (83)$$

$$\begin{aligned} &\leq \sum_{i \in [n] \setminus (A_t \cup B_t)} |(\bar{w}_t(i) - w(i)) a(i)| \\ &\leq \sum_{i \in [n] \setminus (A_t \cup B_t)} |\bar{w}_t(i) - w(i)| |a(i)| \\ &\leq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \sum_{i \in [n] \setminus (A_t \cup B_t)} |a(i)| \\ &= \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\mathbf{a}\|_1, \end{aligned} \quad (84)$$

where Eq. (83) follows from the assumption that  $\mathbf{a}$  supported on  $[n] \setminus (A_t \cup B_t)$ ; Eq. (84) follows from the definition of  $\tau_t$  (Eq. (79)).  $\square$

### D.1.1 Main Lemmas

**Lemma 15.** Fix a phase  $t \in [n]$ . Suppose that  $A_t \subseteq M_*$  and  $B_t \cap M_* = \emptyset$ . Let  $M$  be a set such that  $A_t \subseteq M$  and  $B_t \cap M = \emptyset$ . Let  $a$  and  $b$  be two sets satisfying that  $a \subseteq M \setminus M_*$ ,  $b \subseteq M_* \setminus M$  and  $a \cap b = \emptyset$ . Then, we have

$$A_t \subseteq (M \setminus a \cup b) \quad \text{and} \quad B_t \cap (M \setminus a \cup b) = \emptyset \quad \text{and} \quad (a \cup b) \cap (A_t \cup B_t) = \emptyset.$$

*Proof.* We prove the first part as follows

$$\begin{aligned} A_t \cap (M \setminus a \cup b) &= (A_t \cap (M \setminus a)) \cup (A_t \cap b) \\ &= A_t \cap (M \setminus a) \end{aligned} \quad (85)$$

$$\begin{aligned} &= (A_t \cap M) \setminus a \\ &= A_t \setminus a \end{aligned} \quad (86)$$

$$= A_t, \quad (87)$$

where Eq. (85) holds since we have  $A_t \cap b \subseteq A_t \cap (M_* \setminus M) \subseteq M \cap (M_* \setminus M) = \emptyset$ ; Eq. (86) follows from  $A_t \subseteq M$ ; and Eq. (87) follows from  $a \subseteq M \setminus M_*$  and  $A_t \subseteq M_*$  which imply that  $a \cap A_t = \emptyset$ . Notice that Eq. (87) is equivalent to  $A_t \subseteq (M \setminus a \cup b)$ .

Then, we proceed to prove the second part in the following

$$\begin{aligned} B_t \cap (M \setminus a \cup b) &= (B_t \cap (M \setminus a)) \cup (B_t \cap b) \\ &= B_t \cap (M \setminus a) \end{aligned} \quad (88)$$

$$\begin{aligned} &= (B_t \cap M) \setminus a \\ &= \emptyset \setminus a = \emptyset, \end{aligned} \quad (89)$$

where Eq. (88) follows from the fact that  $B_t \cap b \subseteq B_t \cap (M_* \setminus M) \subseteq \neg M_* \cap (M_* \setminus M) = \emptyset$ ; and Eq. (89) follows from the fact that  $B_t \cap M = \emptyset$ .

Last, we prove the third part. By combining the assumptions that  $A_t \subseteq M_*$  and  $A_t \subseteq M$ , we see that  $A_t \subseteq M \cap M_*$ . Therefore, we have

$$(a \cap A_t) \cup (b \cap A_t) \subseteq ((M \setminus M_*) \cap (M \cap M_*)) \cup ((M_* \setminus M) \cap (M \cap M_*)) = \emptyset. \quad (90)$$

Similarly, we have  $B_t \subseteq \neg M \cap \neg M_*$ . Hence, we derive

$$(a \cap B_t) \cup (b \cap B_t) \subseteq ((M \setminus M_*) \cap (\neg M \cap \neg M_*)) \cup ((M_* \setminus M) \cap (\neg M \cap \neg M_*)) = \emptyset. \quad (91)$$

By combining Eq. (90) and Eq. (91), we obtain

$$(a \cup b) \cap (A_t \cup B_t) = (a \cap A_t) \cup (b \cap A_t) \cup (a \cap B_t) \cup (b \cap B_t) = \emptyset.$$

□

**Lemma 16.** Fix any round  $t > 0$ . Suppose that event  $\tau_t$  occurs. Also assume that  $A_t \subseteq M_*$  and  $B_t \cap M_* = \emptyset$ . Let  $e \in [n] \setminus (A_t \cup B_t)$  be an active arm. Suppose that  $\Delta_{(t-n+1)} \leq \Delta_e$ . Then, we have  $e \in (M_* \cap M_t) \cup (\neg M_* \cap \neg M_t)$ .

*Proof.* Fix an exchange class  $\mathcal{B} \in \arg \min_{\mathcal{B}' \in \text{Exchange}(\mathcal{M})} \text{width}(\mathcal{B}')$ . Suppose that  $e \notin (M_* \cap M_t) \cup (\neg M_* \cap \neg M_t)$ . This is equivalent to the following

$$e \in (M_* \cap \neg M_t) \cup (\neg M_* \cap M_t). \quad (92)$$

Eq. (92) can be further rewritten as

$$e \in (M_* \setminus M_t) \cup (M_t \setminus M_*).$$

From this assumption, it is easy to see that  $M_t \neq M_*$ . Therefore we can apply Lemma 2. Then we know that there exists  $b = (b_+, b_-) \in \mathcal{B}$  such that  $e \in b_- \cup b_+$ ,  $b_- \subseteq M_t \setminus M_*$ ,  $b_+ \subseteq M_* \setminus M_t$ ,  $M_t \oplus b \in \mathcal{M}$  and  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e \geq 0$ .

Using Lemma 15, we see that  $(M_t \oplus b) \cap B_t = \emptyset$ ,  $A_t \subseteq (M_t \oplus b)$  and  $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$ . Now recall the definition  $M_t \in \arg \max_{M \in \mathcal{M}, A_t \subseteq M, B_t \cap M = \emptyset} \bar{w}_t(M)$  and also recall that  $M_t \oplus b \in \mathcal{M}$ . Therefore, we obtain that

$$\bar{w}_t(M_t) \geq \bar{w}_t(M_t \oplus b). \quad (93)$$

On the other hand, we have

$$\bar{w}_t(M_t \oplus b) = \langle \bar{\mathbf{w}}_t, \chi_{M_t} + \chi_b \rangle \quad (94)$$

$$= \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \bar{\mathbf{w}}_t, \chi_b \rangle$$

$$> \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\chi_b\|_1 \quad (95)$$

$$> \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3 \text{width}(\mathcal{M})} \|\chi_b\|_1$$

$$\geq \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3} \quad (96)$$

$$\geq \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \frac{2}{3} \Delta_e \quad (97)$$

$$\geq \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle = \bar{w}_t(M_t), \quad (98)$$

where Eq. (94) follows from Lemma 1; Eq. (95) follows from Lemma 14 and the fact that  $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$ ; Eq. (96) holds since  $b \in \mathcal{B}$  which implies that  $\|\chi_b\|_1 = |b_+| + |b_-| \leq \text{width}(\mathcal{B}) = \text{width}(\mathcal{M})$ ; and Eq. (97) and Eq. (98) hold since we have shown that  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e \geq 0$ .

This means that  $\bar{w}_t(M_t \oplus b) > \bar{w}_t(M_t)$ . This contradicts to Eq. (93). Therefore we have  $e \in (M_* \cap M_t) \cup (\neg M_* \cap \neg M_t)$ .

□

**Lemma 17.** Fix any round  $t > 0$ . Suppose that event  $\tau_t$  occurs. Also assume that  $A_t \subseteq M_*$  and  $B_t \cap M_* = \emptyset$ . Let  $e \in [n] \setminus (A_t \cup B_t)$  be an active arm such that  $\Delta_{(t-n+1)} \leq \Delta_e$ . Then, we have

$$\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,e}) \geq \frac{2}{3} \Delta_{(t-n+1)}.$$

*Proof.* By Lemma 16, we see that

$$e \in (M_* \cap M_t) \cup (\neg M_* \cap \neg M_t). \quad (99)$$

We claim that  $e \in (\tilde{M}_{t,e} \setminus M_*) \cup (M_* \setminus \tilde{M}_{t,e})$  and therefore  $M_* \neq \tilde{M}_{t,e}$ . By Eq. (99), we see that either  $e \in (M_* \cap M_t)$  or  $e \in (\neg M_* \cap \neg M_t)$ . First let us assume that  $e \in M_* \cap M_t$ . Then, by definition of  $\tilde{M}_{t,e}$ , we see that  $e \notin \tilde{M}_{t,e}$ . Therefore  $e \in M_t \setminus \tilde{M}_{t,e}$ . On the other hand, suppose that  $e \in \neg M_* \cap \neg M_t$ . Then, we see that  $e \in \tilde{M}_{t,e}$ . This means that  $e \in \tilde{M}_{t,e} \setminus M_*$ .

Hence we can apply Lemma 2. Then we obtain that there exists  $b = (b_+, b_-) \in \mathcal{B}$  such that  $e \in b_+ \cup b_-$ ,  $b_+ \subseteq M_* \setminus \tilde{M}_{t,e}$ ,  $b_- \subseteq \tilde{M}_{t,e} \setminus M_*$ ,  $\tilde{M}_{t,e} \oplus b \in \mathcal{M}$  and  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e$ .

Define  $M'_{t,e} \triangleq \tilde{M}_{t,e} \oplus b$ . Using Lemma 15, we have  $A_t \subseteq M'_{t,e}$ ,  $B_t \cap M'_{t,e} = \emptyset$  and  $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$ . Since  $M'_{t,e} \in \mathcal{M}$  and by definition  $M_t \in \arg \max_{M \in \mathcal{M}, A_t \subseteq M, B_t \cap M = \emptyset} \bar{w}_t(M)$ , we have

$$\bar{w}_t(M_t) \geq \bar{w}_t(M'_{t,e}). \quad (100)$$

Hence, we have

$$\begin{aligned} \bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,e}) &\geq \bar{w}_t(M'_{t,e}) - \bar{w}_t(\tilde{M}_{t,e}) \\ &= \bar{w}_t(\tilde{M}_{t,e} \oplus b) - \bar{w}_t(\tilde{M}_{t,e}) \\ &= \langle \bar{\mathbf{w}}_t, \chi_{\tilde{M}_{t,e}} + \chi_b \rangle - \langle \bar{\mathbf{w}}_t, \chi_{\tilde{M}_{t,e}} \rangle \\ &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle \end{aligned} \quad (101)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{B})} \|\chi_b\|_1 \quad (102)$$

$$\geq \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3 \text{width}(\mathcal{B})} \|\chi_b\|_1 \quad (103)$$

$$\geq \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3} \quad (104)$$

$$\geq \frac{2}{3} \Delta_e \geq \frac{2}{3} \Delta_{(n-t+1)}, \quad (105)$$

where Eq. (101) follows from Lemma 1; Eq. (102) follows from Lemma 14, the assumption on event  $\tau_t$  and the fact  $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$ ; Eq. (103) follows from the assumption that  $\Delta_e \geq \Delta_{(n-t+1)}$ ; Eq. (104) holds since  $b \in \mathcal{B}$  and therefore  $\|\chi_b\|_1 \leq \text{width}(\mathcal{M})$ ; and Eq. (105) follows from the fact that  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e$ .  $\square$

**Lemma 18.** Fix any phase  $t > 0$ . Suppose that event  $\xi_t$  occurs. Suppose an active arm  $e \in [n] \setminus (A_t \cup B_t)$  satisfies that  $e \in (M_* \cap \neg M_t) \cup (\neg M_* \cap M_t)$ . Then, we have

$$\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,e}) \leq \frac{1}{3} \Delta_{(n-t+1)}.$$

*Proof.* Fix an exchange class  $\mathcal{B} \in \arg \min_{\mathcal{B}' \in \text{Exchange}(\mathcal{M})} \text{width}(\mathcal{B}')$ .

The assumption that  $e \in (M_* \cap \neg M_t) \cup (\neg M_* \cap M_t)$  can be rewritten as  $e \in (M_* \setminus M_t) \cup (M_t \setminus M_*)$ . This shows that  $M_t \neq M_*$ . Hence Lemma 2 applies here. Therefore we know that there exists  $b = (b_+, b_-) \in \mathcal{B}$  such that  $e \in b_+ \cup b_-$ ,  $b_+ \subseteq M_* \setminus M_t$ ,  $b_- \subseteq M_t \setminus M_*$ ,  $M_t \oplus b \in \mathcal{M}$  and  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e \geq 0$ .

Define  $M'_{t,e} \triangleq M_t \oplus b$ . We claim that

$$\bar{w}_t(\tilde{M}_{t,e}) \geq \bar{w}_t(M'_{t,e}). \quad (106)$$

By definition of  $\tilde{M}_{t,e}$ , we only need to show that **(a)**:  $e \in (M'_{t,e} \setminus M_t) \cup (M_t \setminus M'_{t,e})$  and **(b)**:  $A_t \subseteq M'_{t,e}$  and  $B_t \cap M'_{t,e} = \emptyset$ . First we prove **(a)**. Notice that  $b_+ \cap b_- = \emptyset$  and  $b_- \subseteq M_t$ . Hence we see that  $M'_{t,e} \setminus M_t = (M_t \setminus b_- \cup b_+) \setminus M_t = b_+$  and  $M_t \setminus M'_{t,e} = M_t \setminus (M_t \setminus b_- \cup b_+) = b_-$ . In addition, we have that  $e \in (b_- \cup b_+)$ . Therefore we see that **(a)** holds by combining these relations. Next, we notice that **(b)** follows directly from Lemma 15. Hence we have shown that Eq. (106) holds.

Hence, we have

$$\begin{aligned} \bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,e}) &\leq \bar{w}_t(M_t) - \bar{w}_t(M'_{t,e}) \\ &= \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle - \langle \bar{\mathbf{w}}_t, \chi_{M_t} + \chi_b \rangle \\ &= -\langle \bar{\mathbf{w}}_t, \chi_b \rangle \end{aligned} \quad (107)$$

$$\leq -\langle \mathbf{w}, \chi_b \rangle + \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\chi_b\|_1 \quad (108)$$

$$\leq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\chi_b\|_1 \leq \frac{\Delta_{(n-t+1)}}{3}, \quad (109)$$

where Eq. (107) follows from Lemma 1; Eq. (108) follows from Lemma 14, the assumption on  $\tau_t$  and  $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$  (by Lemma 15); and Eq. (109) follows from the fact  $\|\chi_b\|_1 \leq \text{width}(\mathcal{M})$  (since  $b \in \mathcal{B}$ ) and that  $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e \geq 0$ .  $\square$

## D.2 Proof of Theorem 3

For reader's convenience, we first restate Theorem 3 in the following.

**Theorem 3.** *Use the same notations as in Theorem 1. Let  $\Delta_{(1)}, \dots, \Delta_{(n)}$  be a permutation of  $\Delta_1, \dots, \Delta_n$  such that  $\Delta_{(1)} \leq \dots \leq \Delta_{(n)}$ . Define  $\mathbf{H}_2 \triangleq \max_{i \in [n]} i \Delta_{(i)}^{-2}$ . Then, given any budget  $T > n$ , the CGapKill algorithm uses at most  $T$  samples and outputs a solution  $\text{Out} \in \mathcal{M}$  such that*

$$\Pr[\text{Out} \neq M_*] \leq n^2 \exp \left( -\frac{2(T-n)}{9R^2 \tilde{\log}(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right), \quad (8)$$

where  $\tilde{\log}(n) \triangleq \sum_{i=1}^n \frac{1}{i}$ .

*Proof.* First, we show that the algorithm at most  $T$  samples. It is easy to see that exactly one arm is pulled for  $\tilde{T}_1$  times, one arm is pulled for  $\tilde{T}_2$  times,  $\dots$ , and one arm is pulled for  $\tilde{T}_n$  times. Therefore, the total number samples used by the algorithm is bounded by

$$\begin{aligned} \sum_{t=1}^n \tilde{T}_t &\leq \sum_{t=1}^n \left( \frac{T-n}{\tilde{\log}(n)(n-t+1)} + 1 \right) \\ &= \frac{T-n}{\tilde{\log}(n)} \tilde{\log}(n) + n = T. \end{aligned}$$

By Lemma 13, we know that the event  $\tau \triangleq \bigcap_{t=1}^T \tau_t$  occurs with probability at least  $1 - n^2 \exp \left( -\frac{2(T-n)}{9R^2 \tilde{\log}(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right)$ . Therefore, we only need to prove that, under event  $\tau$ , the algorithm outputs  $M_*$ . We will assume that event  $\tau$  occurs in the rest of the proof.

We prove by induction. Fix a phase  $t \in [T]$ . Suppose that the algorithm does not make any error before phase  $t$ , i.e.  $A_t \subseteq M_*$  and  $B_t \cap M_* = \emptyset$ . We show that the algorithm does not err at phase  $t$ .

In the beginning phase  $t$ , there are only  $t-1$  inactive arms  $|A_t \cup B_t| = t-1$ . Therefore there must exist an active arm  $e_1 \in [n] \setminus (A_t \cup B_t)$  such that  $\Delta_{e_1} \geq \Delta_{(n-t+1)}$ . Hence, by Lemma 17, we have

$$\bar{w}_t(M_t) - \bar{w}_t(M_{t,e_1}) \geq \frac{2}{3} \Delta_{(n-t+1)}. \quad (110)$$

Notice that the algorithm makes an error on phase  $t$  if and only if it accepts an arm  $p_t \notin M_*$  or rejects an arm  $p_t \in M_*$ . On the other hand, by design, arm  $p_t$  is accepted when  $p_t \in M_t$  and is rejected when  $p_t \notin M_t$ . Therefore, we see that the algorithm makes an error on phase  $t$  if and only if  $p_t \in (M_* \cap \neg M_t) \cup (\neg M_* \cap M_t)$ .

Suppose that  $p_t \in (M_* \cap \neg M_t) \cup (\neg M_* \cap M_t)$ . Now appeal to Lemma 18, we see that

$$\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,p_t}) \leq \frac{1}{3} \Delta_{(n-t+1)}. \quad (111)$$

By combining Eq. (110) and Eq. (111), we see that

$$\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,p_t}) \leq \frac{1}{3} \Delta_{(n-t+1)} < \frac{2}{3} \Delta_{(n-t+1)} \leq \bar{w}_t(M_t) - \bar{w}_t(M_{t,e_1}). \quad (112)$$

However Eq. (112) is contradictory to the definition of  $p_t \triangleq \arg \max_{i \in [n] \setminus (A_t \cup B_t)} \bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,i})$ . Therefore we have proved that  $p_t \notin (M_* \cap \neg M_t) \cup (\neg M_* \cap M_t)$ . This means that the algorithm does not err at phase  $t$ , or equivalently  $A_{t+1} \subseteq M_*$  and  $B_{t+1} \cap M_* = \emptyset$ . By induction, we have proved that the algorithm does not err at any phase  $t \in [n]$ .

Hence we have  $A_{n+1} \subseteq M_*$  and  $B_{n+1} \subseteq \neg M_*$ . Notice that  $|A_{n+1}| + |B_{n+1}| = n$  and  $A_{n+1} \cap B_{n+1} = \emptyset$ . This means that  $A_{n+1} = M_*$  and  $B_{n+1} = \neg M_*$ . Therefore the algorithm outputs  $\text{Out} = A_{n+1} = M_*$  after phase  $n$ .  $\square$

## E Technical Lemmas

**Fact 1.** Let  $\mathcal{M} \subseteq 2^{[n]}$  be one of our example types of decision classes. Then we can construct the exchange class for  $\mathcal{M}$  and upper bound  $\text{width}(\mathcal{M})$  as follows.

- Define  $\mathcal{B}_{\text{MATROID}(n)} = \{(\{i\}, \{j\}) \mid \forall i \in [n], j \in [n]\}$ . If  $\mathcal{M} = \mathcal{M}_{\text{MATROID}(T, \sigma)}$ , then we have  $\mathcal{B}_{\text{MATROID}(n)} \in \text{Exchange}(\mathcal{M})$  and  $\text{width}(\mathcal{M}) \leq 2$ .
- Define  $\mathcal{B}_{\text{MATCH}(G, \sigma)} = \{(C_+, C_-) \mid \sigma^{-1}(C_+ \cup C_-) \text{ is a cycle of } G\}$ . If  $\mathcal{M} = \mathcal{M}_{\text{MATCH}(G, \sigma)}$ , then we have  $\mathcal{B}_{\text{MATCH}(G)} \in \text{Exchange}(\mathcal{M})$  and  $\text{width}(\mathcal{M}) \leq |V|$ .
- Define  $\mathcal{B}_{\text{PATH}(G, \sigma)} = \{(P_1, P_2) \mid \sigma^{-1}(P_1) \text{ and } \sigma^{-1}(P_2) \text{ are two disjoint paths of } G \text{ with same endpoints}\}$ . If  $\mathcal{M} = \mathcal{M}_{\text{PATH}(G, s, t, \sigma)}$ , then we have  $\mathcal{B}_{\text{PATH}(G, \sigma)} \in \text{Exchange}(\mathcal{M})$  and  $\text{width}(\mathcal{M}) \leq |V|$ .

Moreover, since MATROID encompasses both TOPK and MB types of decision classes, we see that  $\text{width}(\mathcal{M}) \leq 2$  for decision classes  $\mathcal{M}$  of these types.

**Lemma 19** (Basis exchange property). AA

**Lemma 20** (Hoeffding's inequality). Let  $X_1, \dots, X_n$  be  $n$  independent  $R$ -sub-Gaussian random variables. Let  $\bar{X} = \frac{1}{n} \sum X_i$  be the average of these random variables. Then, we have

$$\Pr \left[ |\bar{X} - \mathbb{E}[\bar{X}]| \geq t \right] \leq 2 \exp \left( -\frac{2nt^2}{R^2} \right).$$

## References

- [1] J.-Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, 2009.
- [2] J.-Y. Audibert, S. Bubeck, and R. Munos. Best arm identification in multi-armed bandits. In *COLT*, 2010.
- [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- [4] S. Bubeck and N. Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5:1–122, 2012.
- [5] S. Bubeck, R. Munos, and G. Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412:1832–1852, 2010.



- 1728 [6] S. Bubeck, N. Cesa-bianchi, S. M. Kakade, S. Mannor, N. Srebro, and R. C. Williamson.  
1729 Towards minimax policies for online linear optimization with bandit feedback. In *COLT*, 2012.
- 1730 [7] S. Bubeck, T. Wang, and N. Viswanathan. Multiple identifications in multi-armed bandits. In  
1731 *ICML*, pages 258–265, 2013.
- 1732 [8] N. Cesa-Bianchi and G. Lugosi. Combinatorial bandits. *Journal of Computer and System*  
1733 *Sciences*, 78(5):1404–1422, 2012.
- 1734 [9] W. Chen, Y. Wang, and Y. Yuan. Combinatorial multi-armed bandit: General framework and  
1735 applications. In *Proceedings of The 30th International Conference on Machine Learning*,  
1736 pages 151–159, 2013.
- 1737 [10] V. Gabillon, M. Ghavamzadeh, A. Lazaric, and S. Bubeck. Multi-bandit best arm identification.  
1738 In *NIPS*. 2011.
- 1739 [11] V. Gabillon, M. Ghavamzadeh, and A. Lazaric. Best arm identification: A unified approach to  
1740 fixed budget and fixed confidence. In *NIPS*, 2012.
- 1741 [12] K. Jamieson and R. Nowak. Best-arm identification algorithms for multi-armed bandits in  
1742 the fixed confidence setting. In *Information Sciences and Systems (CISS), 2014 48th Annual*  
1743 *Conference on*, pages 1–6. IEEE, 2014.
- 1744 [13] K. Jamieson, M. Malloy, R. Nowak, and S. Bubeck. lil’ucb: An optimal exploration algorithm  
1745 for multi-armed bandits. *COLT*, 2014.
- 1746 [14] S. Kale, L. Reyzin, and R. E. Schapire. Non-stochastic bandit slate problems. In *NIPS*, pages  
1747 1054–1062, 2010.
- 1748 [15] S. Kalyanakrishnan and P. Stone. Efficient selection of multiple bandit arms: Theory and  
1749 practice. In *ICML*, pages 511–518, 2010.
- 1750 [16] S. Kalyanakrishnan, A. Tewari, P. Auer, and P. Stone. Pac subset selection in stochastic multi-  
1751 armed bandits. In *ICML*, pages 655–662, 2012.
- 1752 [17] T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in*  
1753 *applied mathematics*, 6(1):4–22, 1985.
- 1754 [18] S. Mannor and J. N. Tsitsiklis. The sample complexity of exploration in the multi-armed bandit  
1755 problem. *The Journal of Machine Learning Research*, 5:623–648, 2004.
- 1756 [19] G. Neu, A. György, and C. Szepesvári. The online loop-free stochastic shortest-path problem.  
1757 In *COLT*, pages 231–243, 2010.
- 1758 [20] J. G. Oxley. *Matroid theory*. Oxford university press, 2006.
- 1759 [21] Y. Zhou, X. Chen, and J. Li. Optimal pac multiple arm identification with applications to  
1760 crowdsourcing. In *ICML*, 2014.
- 1761
- 1762
- 1763
- 1764
- 1765
- 1766
- 1767
- 1768
- 1769
- 1770
- 1771
- 1772
- 1773
- 1774
- 1775
- 1776
- 1777
- 1778
- 1779
- 1780
- 1781