
Combinatorial Pure Exploration of Multi-Armed Bandits

Anonymous Author(s)

Affiliation

Address

email

Abstract

We study the *combinatorial pure exploration (CPE)* problem in the stochastic multi-armed bandit setting, where a learner explores a set of arms with the objective of identifying the optimal member of a *decision class*, which is a collection of subsets of arms with certain combinatorial structures such as size- K subset, matching, spanning tree or path, etc. The CPE problem represents a rich class of pure exploration tasks which covers not only many existing models but also novel cases where the object of interest has a non-trivial combinatorial structure. In this paper, we provide a series of results for the general CPE problem. We present general learning algorithms that work for all decision classes that admit offline maximization oracles in both fixed confidence and fixed budget settings. We prove problem-dependent upper bounds of our algorithms. Our analysis exploits the combinatorial structures of the decision classes and introduces a new analytic tool. We also establish a general problem-dependent lower bound for the CPE problem. Our results show that the proposed algorithms achieve the optimal sample complexity (within logarithmic factors) for many decision classes. In addition, applying our results back to the problems of top- K arms identification and multiple bandit best arms identification, we recover the best available upper bounds up to constant factors and settle a conjecture on the lower bounds.

1 Introduction

Multi-armed bandit (MAB) is a predominant model for characterizing the tradeoff between exploration and exploitation in decision-making problems. Although this is an intrinsic tradeoff in many tasks, some application domains prefer a dedicated exploration procedure in which the goal is to identify an optimal object among a collection of candidates and the reward or loss incurred during exploration is irrelevant. In light of these applications, the related learning problem, called pure exploration in MABs, has received much attention. Recent advances in 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 the optimal object with certain *combinatorial structure*. 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 the initialization phase may try 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 display 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 for the framework.

In this paper, we formulate the *combinatorial pure exploration (CPE)* problem for stochastic multi-armed bandits. In the CPE 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*, which corresponds to a collection of certain combinatorial structures. During the exploration

period, in each round the learner chooses an arm to play and observes a random reward sampled from the associated distribution. The objective is when the exploration period ends, the learner outputs a member of the decision class that she believes to be optimal, in the sense that the sum of expected rewards of all arms in the output set is maximized among all members in the decision class.

The CPE 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) [18, 19, 8, 30] and the multi-bandit problem of finding the best arms simultaneously from several disjoint sets of arms (henceforth MB) [12, 8]. Further, this framework covers many more interesting cases where the decision classes correspond to collections of 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 CPE become identifying the optimal paths, spanning trees and matchings through bandit explorations, respectively. To our knowledge, there are no results available in the literature for these pure exploration tasks.

The CPE framework raises several interesting challenges to the design and analysis of pure exploration algorithms. One challenge is that, instead of solving each type of CPE task in an ad-hoc way, one requires a unified algorithm and analysis that support different decision classes. Another challenge stems from the combinatorial nature of CPE, namely that the optimal set may contain some arms with very small expected rewards (e.g. it is possible that a maximum matching contains the edge with the smallest weight); hence, arms cannot be eliminated solely based on their own rewards in the learning algorithm or ignored in the analysis. This differs from many existing approach of pure exploration MABs. Therefore, the design and analysis of algorithms for CPE demands novel techniques which take both rewards and combinatorial structures into account.

Our results. In this paper, we propose two novel learning algorithms for general CPE problem: one for the fixed confidence setting and one for the fixed budget setting. Both algorithms support a wide range of decision classes in a unified way. In the fixed confidence setting, we present Combinatorial Lower-Upper Confidence Bound (CLUCB) algorithm. The CLUCB algorithm does not need to know the definition of the decision class, as long as it has access to the decision class through a maximization oracle. We upper bound the number of samples used by CLUCB. This sample complexity bound depends on both expected reward and the structure of decision class. Our analysis relies on a novel combinatorial construction called *exchange class* which may be of independent interest for other combinatorial optimization problems. Specializing our result to TOPK and MB, we recover the best available sample complexity bounds [19, 13, 20] up to constant factors. While for other decision classes in general, our result establishes the first sample complexity upper bound. We further show that CLUCB can be easily extended to the fixed budget setting and PAC learning setting and we provide related theoretical guarantees in the supplementary material.

Moreover, we establish a problem-dependent sample complexity lower bound for the CPE problem. Our lower bound shows that the sample complexity of the proposed CLUCB algorithm is optimal (to within logarithmic factors) for many 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 CPE problems. For more general decision classes, our results show that the upper and lower bounds are within a relatively benign factor. To the best of our knowledge, there are no problem-dependent lower bounds known for pure exploration MABs besides the case of identifying the single best arm [24, 2]. We also notice that our result resolves the conjecture of Bubeck et al. [8] on the problem-dependent sample complexity lower bounds of TOPK and MB problems.

In the fixed budget setting, we present a parameter-free algorithm called Combinatorial Successive Accept Reject (CSAR) algorithm. We prove a probability of error bound of the CSAR algorithm. This bound can be shown to be equivalent to the sample complexity bound of CLUCB within logarithmic factors, although the two algorithms are based on quite different techniques. Our analysis of CSAR 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 back TOPK and MB, our bound recovers the best known result in the fixed budget setting due to Bubeck et al. [8] up to constant factors.

2 Problem Formulation

In this section, we formally define the CPE 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 φ_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 have R -sub-Gaussian tails for some known constant $R > 0$. Formally, if X is a random variable drawn according to φ_e for some $e \in [n]$, 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 the family of R -sub-Gaussian tail distributions encompasses all distributions that are supported on $[0, R]$ as well as many unbounded distributions such as Gaussian distributions with variance R^2 (cf. [27]).

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 reward¹. A learner's objective is to identify M_* from \mathcal{M} by playing the following game with the stochastic environment. 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; in each round t , the learner pulls an arm $p_t \in [n]$ and observes a reward sampled from the associated reward distribution φ_{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 literature. In the fixed confidence setting, the learner can stop the game at any round. The learner need to guarantee that $\Pr[\text{Out} = M_*] \geq 1 - \delta$ for a given confidence parameter δ . 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 Combinatorial Lower-Upper Confidence Bound (CLUCB) algorithm, a learning algorithm for the CPE 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 the widths of decision classes, which play an important role in the analysis and sample complexity bound. Furthermore, the CLUCB algorithm can be extended to the fixed budget and PAC learning settings, the discussion of which is included in the supplementary material (Appendix B).

Oracle. We allow the CLUCB 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 from a given decision class \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 for \mathcal{M} if, for all $\mathbf{v} \in \mathbb{R}^n$, we have $\text{Oracle}(\mathbf{v}) \in \arg \max_{M \in \mathcal{M}} v(M)$. It is clear that a wide range of decision classes admit such maximization oracles, including decision classes correspond to collections of matchings, paths or bases of matroids (see later for concrete examples). Besides the access to the oracle, CLUCB does not need *any* additional knowledge of the decision class \mathcal{M} .

Algorithm. Now we describe the details of CLUCB, as shown in Algorithm 1. During its execution, the CLUCB 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 $t > 0$. CLUCB begins with an initialization phase in which each arm is pulled once. Then, at round $t \geq n$, CLUCB uses the following procedure to choose an arm to play. First, CLUCB 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, CLUCB explores possible refinements of M_t . In particular, CLUCB 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

¹We denote $v(S) \triangleq \sum_{i \in S} v(i)$ for any vector $\mathbf{v} \in \mathbb{R}^n$ and any set $S \subseteq [n]$. In addition, for convenience, we will assume that M_* is unique.

Algorithm 1 CLUCB: Combinatorial Lower-Upper Confidence Bound

Require: Confidence $\delta \in (0, 1)$; Maximization oracle: $\text{Oracle}(\cdot) : \mathbb{R}^n \rightarrow \mathcal{M}$

Initialize: Play each arm $e \in [n]$ once. Initialize empirical means \bar{w}_n and set $T_n(e) \leftarrow 1$ for all e .

```
1: for  $t = n, n + 1, \dots$  do
2:    $M_t \leftarrow \text{Oracle}(\bar{w}_t)$ 
3:   Compute confidence radius  $\text{rad}_t(e)$  for all  $e \in [n]$   $\triangleright \text{rad}_t(e)$  is defined later in Theorem 1
4:   for  $e = 1, \dots, n$  do
5:     if  $e \in M_t$  then  $\tilde{w}_t(e) \leftarrow \bar{w}_t(e) - \text{rad}_t(e)$ 
6:     else  $\tilde{w}_t(e) \leftarrow \bar{w}_t(e) + \text{rad}_t(e)$ 
7:    $\tilde{M}_t \leftarrow \text{Oracle}(\tilde{w}_t)$ 
8:   if  $\tilde{w}_t(\tilde{M}_t) = \tilde{w}_t(M_t)$  then
9:     Out  $\leftarrow M_t$ 
10:    return Out
11:    $p_t \leftarrow \arg \max_{e \in (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(e)$   $\triangleright$  Break ties arbitrarily
12:   Pull arm  $p_t$  and observe the reward
13:   Update empirical means  $\bar{w}_{t+1}$  using the observed reward
14:   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  $e \neq p_t$ 
```

vector \tilde{w}_t penalizes arms belonging to the current set M_t and encourages exploring arms out of M_t . CLUCB then calls the oracle using the adjusted expectation vector \tilde{w}_t as input to compute a refined set $\tilde{M}_t = \text{Oracle}(\tilde{w}_t)$. If $\tilde{w}_t(\tilde{M}_t) = \tilde{w}_t(M_t)$ then CLUCB stops and returns $\text{Out} = M_t$. Otherwise, CLUCB 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 which has the largest confidence radius. This ends the t -th round of CLUCB. We note that CLUCB generalizes and unifies the ideas of several different fixed confidence algorithms dedicated to the TOPK and MB problems in the literature [19, 13, 20].

3.1 Sample complexity

Now we establish a problem-dependent sample complexity bound of the CLUCB algorithm. To formally state our result, we need to introduce several notions.

Gap. We begin with defining a natural hardness measure of the CPE problem. For each arm $e \in [n]$, we define its gap Δ_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 Δ_e represents the sub-optimality of the best set that includes arm e ; and, for each arm $e \in M_*$, the gap Δ_e is the sub-optimality of the best set that does not include arm e . This naturally generalizes and unifies previous definitions of gaps [2, 12, 18, 8].

Exchange class and the width of a decision class. A notable challenge of our analysis stems from the generality of CLUCB which, as we have seen, supports a wide range of decision classes \mathcal{M} . Indeed, previous algorithms for special cases including TOPK and MB require a separate analysis for each individual type of problem. Such strategy is intractable for our setting and we need a unified analysis for all decision classes. Our solution to this challenge is a novel combinatorial construction called *exchange class*, which is used as a proxy for the structure of the decision class. Intuitively, an exchange class \mathcal{B} for a decision class \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 transform M to M' by applying a series of patches of \mathcal{B} ; and each application of a patch yields a valid member of \mathcal{M} . These patches are later used by our analysis to build gadgets that interpolate between different members of the decision class and serve to bridge key quantities. Furthermore, the maximum patch size of \mathcal{B} will play an important role in our sample complexity bound.

Now we formally define the exchange class. We begin with the definition of exchange sets, which formalize the aforementioned “patches”. We define an exchange set b as an ordered pair of disjoint sets $b = (b_+, b_-)$ where $b_+ \cap b_- = \emptyset$ and $b_+, b_- \subseteq [n]$. 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. For any $M, M' \in \mathcal{M}$ such that $M \neq M'$ and for any $e \in (M \setminus M')$, there exists an exchange set $(b_+, b_-) \in \mathcal{B}$ which satisfies five constraints: **(a)** $e \in b_-$, **(b)** $b_+ \subseteq M' \setminus M$, **(c)** $b_- \subseteq M \setminus M'$, **(d)** $(M \oplus b) \in \mathcal{M}$ and **(e)** $(M' \ominus b) \in \mathcal{M}$.

Intuitively, constraints **(b)** and **(c)** resemble the concept of patches in the sense that b_+ contains only the “new” elements from M' and b_- contains only the “old” elements of M ; constraints **(d)** and **(e)** allow one to transform M one step closer to M' by applying a patch $b \in \mathcal{B}$ to yield $(M \oplus b) \in \mathcal{M}$ (and similarly for $M' \ominus b$). These transformations are the basic building blocks of our analysis. Furthermore, as we will see later in our examples, for many decision classes, there are exchange classes representing natural combinatorial structures, e.g. augmenting paths and cycles of matchings.

In our analysis, the key quantity of exchange class is called *width*, which is defined as the size of largest exchange set as follows

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

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

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

Sample complexity. Our main result of this section is a problem-dependent sample complexity bound of the CLUCB algorithm, which shows that CLUCB 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 φ_e for each arm $e \in [n]$ has mean $w(e)$ with an R -sub-Gaussian tail. Set $\text{rad}_t(e) = R\sqrt{2 \log(\frac{4nt^2}{\delta})} / T_t(e)$ for all $t > 0$ and $e \in [n]$. Then, with probability at least $1 - \delta$, the CLUCB algorithm (Algorithm 1) returns the optimal set $\text{Out} = \arg \max_{M \in \mathcal{M}} w(M)$ and*

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

where T denotes the number of samples used by Algorithm 1, \mathbf{H} is defined in Eq. (2) and $\text{width}(\mathcal{M})$ is defined in Eq. (4).

3.2 Examples of decision classes

Now we investigate several concrete types of decision classes, which correspond to different CPE tasks. We analyze the width of these decision classes and apply Theorem 1 to obtain the sample complexity bounds. A more detailed analysis and the constructions of exchange classes can be found in the supplementary material (Appendix F). 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 $\mathcal{M}_{\text{MB}(\mathcal{A})} = \{M \subseteq [n] \mid \forall i \in [m], |M \cap A_i| = 1\}$. Note that maximization oracles for these two decision classes are trivially the functions of returning the best arms or the best arm of each group.*

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 2 and 3 in the supplementary material) and therefore the sample complexity of CLUCB for solving TOPK and MB is $O(\mathbf{H} \log(\mathbf{H}/\delta))$, which matches previous results in the fixed confidence setting [19, 13, 20] up to constant factors.

Next we consider the problem of identifying the maximum matching or the problem of finding the shortest path (by negating the rewards) 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 $[n]$. Suppose that G is a bipartite graph. Let $\mathcal{M}_{\text{MATCH}(G)}$ correspond to the set of all matchings in G . Then we have $\text{width}(\mathcal{M}_{\text{MATCH}(G)}) \leq |V|$ (In fact, we construct an exchange class corresponding to the collection of augmenting cycles and augmenting paths of G ; see Fact 4).*

Next suppose that G is a directed acyclic graph and let $s, t \in V$ be two vertices. Let $\mathcal{M}_{\text{PATH}(G, s, t)}$ correspond to the set of all paths from s to t . Then we have $\text{width}(\mathcal{M}_{\text{PATH}(G, s, t)}) \leq |V|$ (In fact, we construct an exchange class corresponding to the collection of disjoint pairs of paths; see Fact 5). Therefore the sample complexity bounds of CLUCB for decision classes $\mathcal{M}_{\text{MATCH}(G)}$ and $\mathcal{M}_{\text{PATH}(G, s, t)}$ are $O(|V|^2 \mathbf{H} \log(\mathbf{H}/\delta))$.

Last, we investigate the general problem of identifying the maximum-weight basis of a matroid. Again, Theorem 1 is the first sample complexity upper bound for this pure exploration problem.

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. Assume that there is a one-to-one mapping between E and $[n]$. Recall that a basis of matroid T is a maximal independent set. Let $\mathcal{M}_{\text{MATROID}(T)}$ correspond to the set of all bases of T . Then we have $\text{width}(\mathcal{M}_{\text{MATROID}(T)}) \leq 2$ (derived from strong basis exchange property of matroids, see Fact 1) and the sample complexity of CLUCB for $\mathcal{M}_{\text{MATROID}(T)}$ is $O(\mathbf{H} \log(\mathbf{H}/\delta))$.*

The last example $\mathcal{M}_{\text{MATROID}(T)}$ is a general type of decision class which encompasses many pure exploration problems including 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, it is well-known that matchings and paths cannot be formulated as matroids since they are matroid intersections [26].

4 Lower Bound

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

We show that, for any decision class \mathcal{M} and any expected rewards \mathbf{w} , any δ -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 φ_e is given by $\varphi_e = \mathcal{N}(w(e), 1)$, where we let $\mathcal{N}(\mu, \sigma^2)$ denote the Gaussian distribution with mean μ and variance σ^2 . Then, for any $\delta \in (0, e^{-16}/4)$ and any δ -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 settles the conjecture of Bubeck et al. [8] that the lower bounds of sample complexity of TOPK and MB problem are $\Omega(\mathbf{H} \log(1/\delta))$. Moreover, in Example 1 and Example 3, we have shown that the sample complexity of CLUCB is $O(\mathbf{H} \log(n\mathbf{H}/\delta))$ for TOPK, MB and more generally the decision classes derived from matroids $\mathcal{M}_{\text{MATROID}(T)}$ (including spanning trees). Hence, we see that the CLUCB algorithm achieves the optimal sample complexity within logarithmic factors for these pure exploration tasks.

On the other hand, for general decision classes with non-constant widths, we see that there is a gap of $\tilde{\Theta}(\text{width}(\mathcal{M})^2)$ 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 CLUCB 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 (Appendix C.3), 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 Combinatorial Successive Accept Reject algorithm (CSAR), which is a parameter-free learning algorithm for the CPE problem in the fixed budget setting. Then, we upper bound the probability of error CSAR in terms of gaps and $\text{width}(\mathcal{M})$.

Constrained oracle. The CSAR algorithm requires access to a *constrained oracle*, which is a function denoted as $\text{COracle} : \mathbb{R}^n \times 2^{[n]} \times 2^{[n]} \rightarrow \mathcal{M} \cup \{\perp\}$ and satisfies

$$\text{COracle}(\mathbf{v}, A, B) = \begin{cases} \arg \max_{M \in \mathcal{M}_{A,B}} v(M) & \text{if } \mathcal{M}_{A,B} \neq \emptyset \\ \perp & \text{if } \mathcal{M}_{A,B} = \emptyset, \end{cases} \quad (7)$$

where $\mathcal{M}_{A,B} = \{M \in \mathcal{M} \mid A \subseteq M, B \cap M = \emptyset\}$ is the collection of feasible sets and \perp is a null symbol. Hence we see that $\text{COracle}(\mathbf{v}, A, B)$ returns an optimal set that includes all elements of A while excluding all elements of B ; and if there are no feasible sets, the constrained oracle $\text{COracle}(\mathbf{v}, A, B)$ returns the null symbol \perp . In the supplementary material (Appendix G), we show that constrained oracles are equivalent to maximization oracles up to a transformation on the weight vector. In addition, similar to CLUCB, CSAR does not need any additional knowledge of \mathcal{M} other than accesses to a constrained oracle for \mathcal{M} .

Algorithm. The idea of the CSAR algorithm is as follows. The CSAR algorithm divides the budget of T rounds into n phases. In the end of each phase, CSAR 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 CSAR 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 , CSAR 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, CSAR calls the constrained oracle to compute an optimal set M_t with respect to the empirical means \bar{w}_t , accepted arms A_t and rejected arms B_t , i.e. let $M_t = \text{COracle}(\bar{w}_t, A_t, B_t)$. Then, for each arm active arm e , CSAR estimates the “empirical gap” of e in the following way. If $e \in M_t$, then CSAR computes an optimal set $\tilde{M}_{t,e}$ that does not include e , i.e. $\tilde{M}_{t,e} = \text{COracle}(\bar{w}_t, A_t, B_t \cup \{e\})$. Conversely, if $e \notin M_t$, then CSAR computes an optimal $\tilde{M}_{t,e}$ which includes e , i.e. $\tilde{M}_{t,e} = \text{COracle}(\bar{w}_t, A_t \cup \{e\}, B_t)$. Then, the empirical gap of e is calculated as $\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,e})$. Finally, CSAR chooses the arm p_t with the largest empirical gap. If $p_t \in M_t$ then p_t is accepted otherwise p_t is rejected. The pseudo-code CSAR is shown in Algorithm 2. We note that CSAR can be considered as a generalization of the ideas of the two versions of SAR algorithm due to Bubeck et al. [8], which are designed specifically for the TOPK and MB problems respectively.

5.1 Probability of error

In the following theorem, we bound the probability of error of the CSAR algorithm.

Theorem 3. *Given any $T > n$, any decision class $\mathcal{M} \subseteq 2^{[n]}$ and any expected rewards $\mathbf{w} \in \mathbb{R}^n$. Assume that the reward distribution φ_e for each arm $e \in [n]$ has mean $w(e)$ with an R -sub-Gaussian tail. Let $\Delta_{(1)}, \dots, \Delta_{(n)}$ be a permutation of $\Delta_1, \dots, \Delta_n$ (defined in Eq. (1)) such that $\Delta_{(1)} \leq \dots \leq \Delta_{(n)}$. Define $\mathbf{H}_2 \triangleq \max_{i \in [n]} i \Delta_{(i)}^{-2}$. Then, the CSAR algorithm uses at most T samples and outputs a solution $\text{Out} \in \mathcal{M} \cup \{\perp\}$ 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 $\tilde{\log}(n) \triangleq \sum_{i=1}^n i^{-1}$, $M_* = \arg \max_{M \in \mathcal{M}} w(M)$ and $\text{width}(\mathcal{M})$ is defined in Eq. (4).

Algorithm 2 CSAR: Combinatorial Successive Accept Reject

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

```
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:    $\tilde{T}_t \leftarrow \left\lceil \frac{T-n}{\tilde{\log}(n)(n-t+1)} \right\rceil$ 
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  $\bar{\mathbf{w}}_t$  of each active arm
7:    $M_t \leftarrow \text{COracle}(\bar{\mathbf{w}}_t, A_t, B_t)$ 
8:   if  $M_t = \perp$  then
9:     fail: set  $\text{Out} \leftarrow \perp$  and return Out
10:   for each  $e \in [n] \setminus (A_t \cup B_t)$  do
11:     if  $e \in M_t$  then  $\tilde{M}_{t,e} \leftarrow \text{COracle}(\bar{\mathbf{w}}_t, A_t, B_t \cup \{e\})$ 
12:     else  $\tilde{M}_{t,e} \leftarrow \text{COracle}(\bar{\mathbf{w}}_t, A_t \cup \{e\}, B_t)$ 
13:    $p_t \leftarrow \arg \max_{i \in [n] \setminus (A_t \cup B_t)} \bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,i})$   $\triangleright$  Define  $\bar{w}_t(\perp) = -\infty$ ; Break ties arbitrarily.
14:   if  $p_t \in M_t$  then
15:      $A_{t+1} \leftarrow A_t \cup \{p_t\}, B_{t+1} \leftarrow B_t$ 
16:   else
17:      $A_{t+1} \leftarrow A_t, B_{t+1} \leftarrow B_t \cup \{p_t\}$ 
18:  $\text{Out} \leftarrow A_{n+1}$ 
19: return Out
```

It is well-known that \mathbf{H}_2 is equivalent to \mathbf{H} up to a logarithmic factor, i.e. $\mathbf{H}_2 \leq \mathbf{H} \leq \log(2n)\mathbf{H}_2$ (see [2]). Therefore, by setting the probability of error (the RHS of Eq. (8)) to a constant, one can see that CSAR requires a budget of $T = \tilde{O}(\text{width}(\mathcal{M})^2 \mathbf{H})$ samples. This is equivalent to the sample complexity bound of CLUCB up to logarithmic factors. In addition, applying Theorem 3 back to TOPK and MB, our bound matches the previous fixed budget algorithm due to Bubeck et al. [8].

6 Related Work

The multi-armed bandit problem has been extensively studied in both stochastic and adversarial settings [22, 3, 3]. We refer readers to [5] for a survey on recent advances. Many work in MABs focus on minimizing the cumulative regret, which is an objective known to be fundamentally different from the objective of pure exploration MABs [6]. 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) are played on each round [9, 10, 17, 25, 1, 7, 14, 23, 21]. 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 [24, 11, 6, 2, 13, 16, 15]. A flurry of recent work extend this classical problem to TOPK and MB problems and obtain algorithms, upper bounds [18, 19, 30, 8, 12, 13, 20] and worst-case lower bounds of TOPK [19, 30]. Our framework encompasses these two problems as special cases and covers a much larger class of combinatorial pure exploration problems, which are unaddressed in the current literature. Applying our results back to TOPK and MB, our upper bounds match best available problem-dependent bounds up to constant factors [13, 8, 19] in both fixed confidence and fixed budget settings; and our lower bound provides the first problem-dependent lower bounds for these two problems, which are conjectured earlier by Bubeck et al. [8].

7 Conclusion

In this paper, we proposed a general framework called combinatorial pure exploration (CPE) that can handle pure exploration tasks for many complex bandit problems with combinatorial constraints, and have 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 may be of independent interest. Our upper and lower bounds characterize the complexity of the CPE 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] C. Berge. Two theorems in graph theory. *PNAS*, 1957.
- [5] 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.
- [6] S. Bubeck, R. Munos, and G. Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412:1832–1852, 2010.
- [7] 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.
- [8] S. Bubeck, T. Wang, and N. Viswanathan. Multiple identifications in multi-armed bandits. In *ICML*, pages 258–265, 2013.
- [9] N. Cesa-Bianchi and G. Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.
- [10] W. Chen, Y. Wang, and Y. Yuan. Combinatorial multi-armed bandit: General framework and applications. In *ICML*, pages 151–159, 2013.
- [11] E. Even-Dar, S. Mannor, and Y. Mansour. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *JMLR*, 2006.
- [12] V. Gabillon, M. Ghavamzadeh, A. Lazaric, and S. Bubeck. Multi-bandit best arm identification. In *NIPS*. 2011.
- [13] V. Gabillon, M. Ghavamzadeh, and A. Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In *NIPS*, 2012.
- [14] A. Gopalan, S. Mannor, and Y. Mansour. Thompson sampling for complex online problems. In *ICML*, pages 100–108, 2014.
- [15] K. Jamieson and R. Nowak. Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. In *Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2014.
- [16] K. Jamieson, M. Malloy, R. Nowak, and S. Bubeck. lil’ucb: An optimal exploration algorithm for multi-armed bandits. *COLT*, 2014.
- [17] S. Kale, L. Reyzin, and R. E. Schapire. Non-stochastic bandit slate problems. In *NIPS*, 2010.
- [18] S. Kalyanakrishnan and P. Stone. Efficient selection of multiple bandit arms: Theory and practice. In *ICML*, pages 511–518, 2010.
- [19] S. Kalyanakrishnan, A. Tewari, P. Auer, and P. Stone. Pac subset selection in stochastic multi-armed bandits. In *ICML*, pages 655–662, 2012.
- [20] E. Kaufmann and S. Kalyanakrishnan. Information complexity in bandit subset selection. In *COLT*, pages 228–251, 2013.
- [21] B. Kveton, Z. Wen, A. Ashkan, H. Eydgahi, and B. Eriksson. Matroid bandits: Fast combinatorial optimization with learning. In *UAI*, 2014.
- [22] T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- [23] T. Lin, B. Abraham, R. Kleinberg, J. Lui, and W. Chen. Combinatorial partial monitoring game with linear feedback and its application. In *ICML*, 2014.
- [24] 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.
- [25] G. Neu, A. Györfy, and C. Szepesvári. The online loop-free stochastic shortest-path problem. In *COLT*, pages 231–243, 2010.
- [26] J. G. Oxley. *Matroid theory*. Oxford university press, 2006.
- [27] D. Pollard. Asymptopia. *Manuscript, Yale University, Dept. of Statist., New Haven, Connecticut*, 2000.
- [28] S. M. Ross. *Stochastic processes*, volume 2. John Wiley & Sons New York, 1996.
- [29] N. Spring, R. Mahajan, and D. Wetherall. Measuring isp topologies with rocketfuel. *ACM SIGCOMM Computer Communication Review*, 32(4):133–145, 2002.
- [30] Y. Zhou, X. Chen, and J. Li. Optimal pac multiple arm identification with applications to crowdsourcing. In *ICML*, 2014.

Organization

This supplementary material is organized as follows. First, we provide the deferred proofs of Theorem 1 (Appendix A). Then, we present two simple extensions of CLUCB: one to the fixed budget setting and one to the PAC learning setting; and provide related analysis (Appendix B). Next, we provide the deferred proof of our lower bound (Theorem 2) in Appendix C. In Appendix D, we provide the deferred proof of Theorem 3. We also analyze the uniform allocation algorithm as a simple benchmark in Appendix E. Afterwards, we provide the deferred discussion and constructions of example decision classes and exchange classes in Appendix F. In Appendix G, we prove the equivalence between the constrained oracles and maximization oracles. Finally, we present some preliminary experimental results in Appendix H.

A Analysis of CLUCB (Theorem 1)

In this section, we analyze the sample complexity of CLUCB and prove 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 $\text{rad}_t = (\text{rad}_t(1), \dots, \text{rad}_t(n))^T$ 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 . Finally, we let $|u| = (|u(1)|, \dots, |u(n)|)^T$ denote the element-wise absolute value of u .

A.1 Preparatory Lemmas

Let us begin with a simple lemma that characterizes the incidence vectors of exchange sets.

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 $M_2 = M_1 \setminus b_- \cup b_+$ and $b_+ \cap b_- = \emptyset$. Therefore we see that $M_2 \setminus M_1 = b_+$ and $M_1 \setminus M_2 = b_-$. We can decompose χ_{M_1} as $\chi_{M_1} = \chi_{M_1 \setminus M_2} + \chi_{M_1 \cap M_2}$. Hence, we have

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

□

The next lemma serves as a basic tool derived from exchange classes, which allows us to interpolate between different members of a decision class \mathcal{M} . It also characterizes the relationship between gaps and exchange sets.

Lemma 2 (Interpolation Lemma). *Let $\mathcal{M} \subseteq 2^{[n]}$ and let \mathcal{B} be an exchange class for \mathcal{M} . Then, for any two different members 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 Δ_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 $M_1 = \arg \max_{M' \in \mathcal{M}: e \in M'} w(M')$. Note that, by definition of gaps, one has that $w(M_*) - w(M_1) = \Delta_e$. Now we define $M_0 = M_* \ominus b$. Note that we already have that $M_0 = M_* \ominus b \in \mathcal{M}$. By combining this 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

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

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

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 $M_3 = \arg \max_{M' \in \mathcal{M}: e \notin M'} w(M')$. By definition of Δ_e , we see that $w(M_*) - w(M_3) = \Delta_e$. Now we define $M_2 = M_* \ominus b$ and notice that $M_2 \in \mathcal{M}$. By 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

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

□

Next we state two basic lemmas that help us to convert set-theoretical arguments to linear algebraic arguments.

Lemma 3. Let $M, M' \subseteq [n]$ be two sets. Let \mathbf{rad}_t be an n -dimensional vector. Then, we have

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

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 $\mathbf{rad}_t(e) \geq 0$ 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

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

□

The next lemma characterizes the property of $\tilde{\mathbf{w}}_t$ which is defined in the CLUCB algorithm.

Lemma 5. *Let M_t , $\tilde{\mathbf{w}}_t$ and \mathbf{rad}_t be defined in Algorithm 1 and Theorem 1. Let $M' \in \mathcal{M}$ be an arbitrary member of decision class. We have*

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

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

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

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

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

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 $\chi_{M'} - \chi_{M_t} \in \{-1, 0, 1\}^n$. \square

A.2 Confidence Intervals

First, we recall a standard concentration inequality of sub-Gaussian random variables.

Lemma 6 (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*

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

Next, for all $t > 0$, we define random event ξ_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 ξ_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 7. *Given any $t > 0$, assume that event ξ_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 ξ_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 ξ_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 8. Suppose that the reward distribution φ_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

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

where $T_t(e)$ 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 φ_e is a R -sub-Gaussian distribution with mean $w(e)$ and $w_t(e)$ is the empirical mean of φ_e for $T_t(e)$ samples. Using Hoeffding's inequality (see Lemma 6), we obtain

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

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

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

□

A.3 Main Lemmas

Now we state our key technical lemmas. In these lemmas, we shall use Lemma 2 to construct gadgets that interpolate between different members of a decision class. The first lemma shows that, if the confidence intervals are valid, then CLUCB always returns the correct answer when it stops.

Lemma 9. Given any $t > n$, assume that event ξ_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 $\langle \mathbf{w}, \chi_{M_*} \rangle > \langle \mathbf{w}, \chi_{M_t} \rangle$.

Let \mathcal{B} be an exchange class for \mathcal{M} . 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

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

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

$$= w(M'_t) - w(M_t) > 0, \tag{16}$$

where Eq. (14) follows from Lemma 5; and Eq. (15) follows the assumption that event ξ_t occurs and Lemma 7;

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

The next lemma shows that if the confidence interval of an arm is sufficiently small, then this arm will not be played by the algorithm.

Lemma 10. *Given any $t > 0$ and suppose that event ξ_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 $\mathcal{B} \in \arg \min_{\mathcal{B}' \in \text{Exchange}(\mathcal{M})} \text{width}(\mathcal{B}')$. Note that $\text{width}(\mathcal{B}) = \text{width}(\mathcal{M})$. Suppose, on 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): $(e \in M_* \wedge e \in c_+) \vee (e \notin M_* \wedge e \in c_-)$.

First we construct a gadget M'_t as $M'_t = \tilde{M}_t \ominus c = \tilde{M}_t \setminus c_+ \oplus c_-$ and recall that $M'_t \in \mathcal{M}$ due to the definition of exchange class.

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 $\langle w, \chi_b \rangle \geq \Delta_e > 0$. We will also use $M'_t \oplus b$ as a gadget.

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 1, 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 $\|\text{rad}_t \circ (\mathbf{d} - \chi_c)\|_\infty < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}$. Using standard set-algebraic manipulations, we have

$$\begin{aligned} M_t \setminus M'_t &= M_t \setminus (\tilde{M}_t \ominus c) \\ &= M_t \setminus (\tilde{M}_t \setminus c_+ \cup c_-) \\ &= M_t \setminus (\tilde{M}_t \setminus c_+) \cap (M_t \setminus c_-) \\ &= (M_t \cap c_+) \cup (M_t \setminus \tilde{M}_t) \cap (M_t \setminus c_-) \\ &= (M_t \setminus \tilde{M}_t) \cap (M_t \setminus c_-) \end{aligned} \tag{17}$$

$$\subseteq M_t \setminus \tilde{M}_t, \tag{18}$$

where Eq. (17) follows from $c_+ \subseteq \tilde{M}_t \setminus M_t$ and therefore $c_+ \cap M_t = \emptyset$. Similarly, we can derive $M'_t \setminus M_t$ as follows

$$\begin{aligned} M'_t \setminus M_t &= (\tilde{M}_t \ominus c) \setminus M_t = (\tilde{M}_t \setminus c_+ \cup c_-) \setminus M_t \\ &= ((\tilde{M}_t \setminus c_+) \setminus M_t) \cup (c_- \setminus M_t) \\ &= \tilde{M}_t \setminus c_+ \setminus M_t \end{aligned} \tag{19}$$

$$\subseteq \tilde{M}_t \setminus M_t, \tag{20}$$

where Eq. (19) follows from $c_- \subseteq M_t \setminus \tilde{M}_t$ and hence $c_- \setminus M_t = \emptyset$. By combining Eq. (18) and Eq. (20), we see 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

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

$$\begin{aligned}
&\leq \max_{i \in (M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)} \text{rad}_t(i) \\
&= \text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}.
\end{aligned} \tag{21}$$

Next 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

$$\begin{aligned}
\|\text{rad}_t \circ \chi_c\|_\infty &= \max_{i \in c_+ \cup c_-} \text{rad}_t(i) \\
&\leq \max_{i \in (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(i) \\
&= \text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}.
\end{aligned} \tag{22}$$

From Eq. (22), we derive

$$\langle \text{rad}_t, |\chi_c| \rangle = \langle \text{rad}_t, \chi_c^2 \rangle \tag{23}$$

$$= \langle \text{rad}_t \circ \chi_c, \chi_c \rangle \tag{24}$$

$$\leq \|\text{rad}_t \circ \chi_c\|_\infty \|\chi_c\|_1 \tag{25}$$

$$\leq \frac{\Delta_e}{3 \text{width}(\mathcal{B})} \|\chi_c\|_1 \tag{26}$$

$$\leq \frac{\Delta_e}{3}, \tag{27}$$

where Eq. (23) hold since $\chi_c \in \{-1, 0, 1\}^n$; Eq. (24) follows from Lemma 4; Eq. (25) follows from Hölder's inequality; Eq. (26) follows from Eq. (22); and Eq. (27) holds since $\|\chi_c\|_1 = |c_+| + |c_-| \leq \text{width}(\mathcal{B})$ where the inequality is due to $c \in \mathcal{B}$.

Next, we claim that $d \circ \chi_c = |\chi_c|$. Recall that $\chi_c = \chi_{c_+} - \chi_{c_-}$ and $d = \chi_{\tilde{M}_t} - \chi_{M_t} = \chi_{\tilde{M}_t \setminus M_t} - \chi_{M_t \setminus \tilde{M}_t}$. 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

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

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

Now, we bound quantity $\langle \text{rad}_t, |d_2| \rangle - \langle \text{rad}_t, |d| \rangle$ as follows

$$\langle \text{rad}_t, |d_2| \rangle - \langle \text{rad}_t, |d| \rangle = \langle \text{rad}_t, |d_2| - |d| \rangle = \langle \text{rad}_t, d_2^2 - d^2 \rangle \tag{28}$$

$$\begin{aligned}
&= \langle \text{rad}_t, (d - \chi_c + \chi_b)^2 - d^2 \rangle \\
&= \langle \text{rad}_t, \chi_b^2 + \chi_c^2 - 2\chi_b \circ \chi_c - 2d \circ \chi_c + 2d \circ \chi_b \rangle \\
&= \langle \text{rad}_t, \chi_b^2 - \chi_c^2 + 2\chi_b \circ (d - \chi_c) \rangle
\end{aligned} \tag{29}$$

$$\begin{aligned}
&= \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - 2 \langle \text{rad}_t, \chi_b \circ (d - \chi_c) \rangle \\
&= \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - 2 \langle \text{rad}_t \circ (d - \chi_c), \chi_b \rangle
\end{aligned} \tag{30}$$

$$\geq \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - 2 \|\text{rad}_t \circ (d - \chi_c)\|_\infty \|\chi_b\|_1 \tag{31}$$

$$> \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3 \text{width}(\mathcal{B})} \|\chi_b\|_1 \tag{32}$$

$$\geq \langle \text{rad}_t, |\chi_b| \rangle - \langle \text{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3}, \tag{33}$$

where Eq. (28) holds since $\mathbf{d} \in \{-1, 0, 1\}^n$ and $\mathbf{d}_2 \in \{-1, 0, 1\}^n$; Eq. (29) follows from the claim that $\mathbf{d} \circ \chi_c = |\chi_c| = \chi_c^2$; Eq. (30) and Eq. (31) follow from Lemma 4 and Hölder's inequality; Eq. (32) follows from Eq. (21); and Eq. (33) holds since $b \in \mathcal{B}$ and $\|\chi_b\|_1 = |b_+| + |b_-| \leq \text{width}(\mathcal{B})$.

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

$$\begin{aligned} \langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle + \langle \mathbf{rad}_t, |\mathbf{d}| \rangle &= \langle \bar{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \\ &= \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \\ &\geq \tilde{w}_t(M'_t \oplus b) - \tilde{w}_t(M_t) \end{aligned} \quad (34)$$

$$\begin{aligned} &= \langle \bar{\mathbf{w}}_t, \chi_{M'_t \oplus b} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{M'_t \oplus b} - \chi_{M_t}| \rangle \\ &= \langle \bar{\mathbf{w}}_t, \mathbf{d}_2 \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_2| \rangle \\ &= \langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle - \langle \bar{\mathbf{w}}_t, \chi_c \rangle + \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_2| \rangle, \end{aligned} \quad (35)$$

where Eq. (34) follows from the fact that $\tilde{w}_t(\tilde{M}_t) = \max_{M \in \mathcal{M}} \tilde{w}_t(M)$; and Eq. (35) follows from the fact that $\mathbf{d}_2 = \mathbf{d} - \chi_c + \chi_b$. Rearranging the above inequality, we obtain

$$\begin{aligned} \langle \bar{\mathbf{w}}_t, \chi_c \rangle &\geq \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_2| \rangle - \langle \mathbf{rad}_t, |\mathbf{d}| \rangle \\ &\geq \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3} \end{aligned} \quad (36)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3} \quad (37)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3} - \frac{2\Delta_e}{3} \quad (38)$$

$$= \langle \mathbf{w}, \chi_b \rangle - \Delta_e \geq 0, \quad (39)$$

where Eq. (36) uses Eq. (33); Eq. (37) follows from the assumption that event ξ_t occurs and Lemma 7; and Eq. (38) holds since Eq. (27).

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. $\bar{w}_t(M_t) \geq \bar{w}_t(M'_t)$. This is a contradiction, and therefore we have $p_t \neq e$ for this case.

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 $\tilde{M}_t \neq M_*$. On the other hand, suppose that $e \notin M_*$ 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 $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e > 0$. We will use $\tilde{M}_t \oplus b$ as a gadget for this case.

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 1 gives that $\mathbf{d}_1 = \mathbf{d} + \chi_b$.

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

$$\|\mathbf{rad}_t \circ \mathbf{d}\|_\infty \leq \max_{i: (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(i) = \text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}, \quad (40)$$

where the last inequality follows from the assumption that ξ_t occurs.

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

$$\begin{aligned} \langle \bar{\mathbf{w}}_t, \mathbf{d}_1 \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| \rangle - \langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle - \langle \mathbf{rad}_t, |\mathbf{d}| \rangle &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| - |\mathbf{d}| \rangle \\ &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, \mathbf{d}_1^2 - \mathbf{d}^2 \rangle \end{aligned} \quad (41)$$

$$= \langle \bar{w}_t, \chi_b \rangle + \langle \text{rad}_t, 2d \circ \chi_b + \chi_b^2 \rangle \quad (42)$$

$$= \langle \bar{w}_t, \chi_b \rangle + \langle \text{rad}_t, \chi_b^2 \rangle + 2 \langle \text{rad}_t \circ d, \chi_b \rangle \quad (43)$$

$$\geq \langle w, \chi_b \rangle + 2 \langle \text{rad}_t \circ d, \chi_b \rangle \quad (44)$$

$$\geq \langle w, \chi_b \rangle - 2 \|\text{rad}_t \circ d\|_\infty \|\chi_b\|_1 \quad (45)$$

$$> 0, \quad (46)$$

where Eq. (41) follows from the fact that $d_1 \in \{-1, 0, 1\}^n$ and $d \in \{-1, 0, 1\}^n$; Eq. (42) holds since $d_1 = d + \chi_b$; Eq. (43) follows from the assumption that ξ_t occurs and Lemma 7; Eq. (44) follows from Lemma 4 and Hölder's inequality; Eq. (45) is due to Eq. (40); and Eq. (46) follows from $\langle w, \chi_b \rangle \geq \Delta_e > 0$.

Therefore, we have proved that

$$\langle \bar{w}_t, d \rangle + \langle \text{rad}_t, |d| \rangle < \langle \bar{w}_t, d_1 \rangle + \langle \text{rad}_t, |d_1| \rangle. \quad (47)$$

However, we have

$$\langle \bar{w}_t, d \rangle + \langle \text{rad}_t, |d| \rangle = \langle \bar{w}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \text{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \quad (48)$$

$$= \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \quad (49)$$

$$\geq \tilde{w}_t(\tilde{M}_t \oplus b) - \tilde{w}_t(M_t)$$

$$= \langle \bar{w}_t, \chi_{\tilde{M}_t \oplus b} - \chi_{M_t} \rangle + \langle \text{rad}_t, |\chi_{\tilde{M}_t \oplus b} - \chi_{M_t}| \rangle \quad (50)$$

$$= \langle \bar{w}_t, d_1 \rangle + \langle \text{rad}_t, |d_1| \rangle,$$

where Eq. (48) follows from Lemma 5; and Eq. (49) follows from the fact that $\tilde{w}_t(\tilde{M}_t) = \max_{M \in \mathcal{M}} \tilde{w}_t(M)$. This contradicts to Eq. (47) and therefore $p_t \neq e$. \square

A.4 Proof of Theorem 1

Theorem 1 is now a straightforward corollary of Lemma 9 and Lemma 10. For the reader's 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 $w \in \mathbb{R}^n$. Assume that the reward distribution φ_e for each arm $e \in [n]$ has mean $w(e)$ with an R -sub-Gaussian tail. Set $\text{rad}_t(e) = R\sqrt{2 \log(\frac{4nt^2}{\delta})} / T_t(e)$ for all $t > 0$ and $e \in [n]$. Then, with probability at least $1 - \delta$, the CLUCB algorithm (Algorithm 1) returns the optimal set $\text{Out} = \arg \max_{M \in \mathcal{M}} w(M)$ and*

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

where T denotes the number of samples used by Algorithm 1, \mathbf{H} is defined in Eq. (2) and $\text{width}(\mathcal{M})$ is defined in Eq. (4).

Proof. Lemma 8 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 9 and the assumption on ξ , 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 $T_{t_e}(e) = T(e) - 1$. By Lemma 10, we see that $\text{rad}_{t_e}(e) \geq \frac{\Delta_e}{3 \text{width}(\mathcal{M})}$. By plugging the definition of rad_{t_e} , we have

$$\frac{\Delta_e}{3 \text{width}(\mathcal{M})} \leq R \sqrt{\frac{2 \log(4nt_e^2/\delta)}{T(e) - 1}} \leq R \sqrt{\frac{2 \log(4nT^2/\delta)}{T(e) - 1}}. \quad (51)$$

Solving Eq. (51) for $T(e)$, we obtain

$$T(e) \leq \frac{18 \text{width}(\mathcal{M})^2 R^2}{\Delta_e^2} \log(4nT^2/\delta) + 1. \quad (52)$$

Notice that $T = \sum_{i \in [n]} T(i)$. Hence the theorem follows by summing up Eq. (52) for all $e \in [n]$ and solving for T . \square

B Extensions of CLUCB

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

B.1 Fixed Budget Setting

We can extend the CLUCB algorithm to the fixed budget setting using two simple modifications: (1) requiring CLUCB to terminate after T rounds; and (2) using a different construction of confidence intervals. The first modification ensures that CLUCB 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 CLUCB 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 $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_t(e)}}$ for all arms $e \in [n]$ and all $t > 0$. Run CLUCB 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*

$$\Pr[\text{Out} \neq M_*] \leq 2Tn \exp(-2\alpha). \quad (53)$$

In particular, the right-hand side of Eq. (53) equals to $O\left(Tn \exp\left(\frac{-T}{\text{width}(\mathcal{M})^2 \mathbf{H}}\right)\right)$ when parameter $\alpha = O(T\mathbf{H}^{-1} \text{width}(\mathcal{M})^{-2})$.

Theorem 4 shows that the modified CLUCB algorithm in the fixed budget setting requires the knowledge of quantity \mathbf{H} in order to achieve the optimal performance. However \mathbf{H} is usually unknown. Therefore, although its probability of error guarantee matches the parameter-free CSAR algorithm up to logarithmic factors, this modified algorithm is considered more restricted than CSAR. Nevertheless, Theorem 4 shows that CLUCB can solve CPE in both fixed confidence and fixed budget settings and more importantly this theorem provides additional insights on the behavior CLUCB.

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 (ϵ, δ) -PAC algorithm. Formally, an algorithm \mathbb{A} is called an (ϵ, δ) -PAC algorithm if its output $\text{Out} \in \mathcal{M}$ satisfies $\Pr[w(M_*) - w(\text{Out}) > \epsilon] \leq \delta$.

We show that a simple modification on the CLUCB algorithm gives an (ϵ, δ) -PAC algorithm, with guarantees similar to Theorem 1. In fact, the only modification needed is to change the stopping condition from $\tilde{w}_t(\tilde{M}_t) = \tilde{w}_t(M_t)$ to $\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \leq \epsilon$ on line 11 of Algorithm 1. We let CLUCB-PAC denote the modified algorithm. In the following theorem, we show that CLUCB-PAC is indeed an (ϵ, δ) -PAC algorithm and has sample complexity similar to CLUCB.

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 $\text{Out} \in \mathcal{M}$ of CLUCB-PAC satisfies $w(M_*) - w(\text{Out}) \leq \epsilon$. In addition, the number of samples T used by the algorithm satisfies*

$$T \leq O\left(R^2 \sum_{e \in [n]} \min\left\{\frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2}\right\} \log\left(\frac{R^2}{\delta} \sum_{e \in [n]} \min\left\{\frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2}\right\}\right)\right), \quad (54)$$

where $K = \max_{M \in \mathcal{M}} |M|$ is the size of the largest member of decision class.

We see that if $\epsilon = 0$, the sample complexity Eq. (54) of CLUCB-PAC equals to that of CLUCB. And the sample complexity of CLUCB-PAC decreases when ϵ increases.

There are several PAC learning algorithms dedicated for the TOPK problem in the literature with different guarantees [19, 30, 13]. Zhou et al. [30] proposed an (ϵ, δ) -PAC algorithm for the TOPK problem with a problem-independent sample complexity bound of $O(\frac{K^2 n}{\epsilon^2} + \frac{Kn \log(1/\delta)}{\epsilon^2})$.² If we ignore logarithmic factors, then the sample complexity bound of CLUCB-PAC for the TOPK problem is better than theirs since $\tilde{O}(\sum_{e \in [n]} \min\{\Delta_e^{-2}, K^2 \epsilon^{-2}\}) \leq \tilde{O}(nK^2 \epsilon^{-2})$. On the other hand, the algorithms of Kalyanakrishnan et al. [19] and Gabillon et al. [13] guarantee to find K arms such that each of them is better than the K -th optimal arm within a factor of ϵ 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 CLUCB 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 11. *Fix parameter $\alpha > 0$ and the number of rounds $T > 0$. Assume that the reward distribution φ_e is a R -sub-Gaussian distribution for all $e \in [n]$. Let the confidence radius $\text{rad}_t(e)$ of arm $e \in [n]$ and round $t > 0$ be $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_t(e)}}$. Then, we have*

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

where ξ_t is the random event defined in Eq. (12).

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 9 and Lemma 10) and Lemma 11.

Proof of Theorem 4. Define random event $\xi = \bigcap_{t=1}^T \xi_t$. By Lemma 11, we see that $\Pr[\xi] \geq 1 - 2nT \exp(-2\alpha)$. In the rest of the proof, we assume that ξ happens.

Let T^* denote the round that the algorithm stops. We claim that the algorithm stops before the budget is exhausted, i.e. $T^* < T$. If the claim is true, then the algorithm stops since it meets the stopping condition on round T^* . Hence $\tilde{w}_t(\tilde{M}_{T^*}) = \tilde{w}_t(M_{T^*})$ and $\text{Out} = M_{T^*}$. By assumption on ξ and Lemma 9, 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 $T(e)$ denote the total number of pulls of arm $e \in [n]$. Let t_e be the last round that arm e is pulled. Hence $T_{t_e}(e) = T_e - 1$. By Lemma 10, 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

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

²We notice that the definition of Zhou et al. [30] allow an (ϵ', δ) -PAC algorithm to produce an output with average sub-optimality of ϵ' . This is equivalent to our definition of (ϵ, δ) -PAC algorithm with $\epsilon = K\epsilon'$ for the TOPK problem. In this paper, we translate their guarantees to our definition of PAC algorithm.

Hence we have

$$T_e \leq \frac{9R^2 \text{width}(\mathcal{B})^2}{\Delta_e^2} \cdot \alpha + 1. \quad (55)$$

By summing up Eq. (55) 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 $\alpha < \frac{1}{9}(T - n) \cdot \left(R^2 \text{width}(\mathcal{B})^2 \left(\sum_{e \in [n]} \Delta_e^{-2} \right) \right)^{-1}$.

□

B.3.2 PAC Learning (Theorem 5)

First, we prove a (ϵ, δ) -PAC counterpart of Lemma 9.

Lemma 12. *If CLUCB-PAC stops on round t and suppose that event ξ_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 CLUCB-PAC ensures that $\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \leq \epsilon$. Therefore, we have

$$\begin{aligned} \epsilon &\geq \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \\ &\geq \tilde{w}_t(M_*) - \tilde{w}_t(M_t) \end{aligned} \quad (56)$$

$$= \langle \tilde{\mathbf{w}}_t, \chi_{M_*} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{M_*} - \chi_{M_t}| \rangle \quad (57)$$

$$\begin{aligned} &\geq \langle \mathbf{w}, \chi_{M_*} - \chi_{M_t} \rangle \\ &= w(M_*) - w(M_t), \end{aligned} \quad (58)$$

where Eq. (56) follows from the definition of $\tilde{M}_t \triangleq \arg \max_{M \in \mathcal{M}} \tilde{w}_t(M)$; Eq. (57) follows from Lemma 5; Eq. (58) follows from the assumption that ξ_t occurs and Lemma 7. □

The next lemma generalizes Lemma 10. It shows that on event ξ_t 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 13. *Let $K = \max_{M \in \mathcal{M}} |M|$. For any arm $e \in [n]$ and any round $t > n$ after initialization, if $\text{rad}_t(e) \leq \max \left\{ \frac{\Delta_e}{3 \text{width}(\mathcal{M})}, \frac{\epsilon}{2K} \right\}$ and random event ξ_t occurs, 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 10 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

$$\epsilon \geq \sum_{i \in (M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)} \text{rad}_t(i) \quad (59)$$

$$= \langle \mathbf{rad}_t, |\chi_{M_t} - \chi_{\tilde{M}_t}| \rangle, \quad (60)$$

where Eq. (59) 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. (60) 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

$$\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) = \langle \tilde{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \quad (61)$$

$$\leq \langle \bar{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \epsilon \quad (62)$$

$$\begin{aligned} &= \bar{w}_t(\tilde{M}_t) - \bar{w}_t(M_t) + \epsilon \\ &\leq \epsilon, \end{aligned} \quad (63)$$

where Eq. (61) follows from Lemma 5; Eq. (62) uses Eq. (60); and Eq. (63) follows from $\bar{w}_t(M_t) \geq \bar{w}_t(\tilde{M}_t)$.

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

Using Lemma 13 and Lemma 12, we are ready to prove Theorem 5.

Proof of Theorem 5. Similar to the proof of Theorem 1, we appeal to Lemma 8, 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 ξ occurs in the rest of the proof.

By the assumption of ξ and Lemma 12, we know that $w(M_*) - w(\text{Out}) \leq \epsilon$. 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 $T_{t_e}(e) = T(e) - 1$. By Lemma 13, we see that $\text{rad}_{t_e}(e) \geq \max\{\frac{\Delta_e}{3 \text{width}(\mathcal{B})}, \frac{\epsilon}{2K}\}$. Then, by the construction of $\text{rad}_{t_e}(e)$, we have

$$\max\left\{\frac{\Delta_e}{3 \text{width}(\mathcal{B})}, \frac{\epsilon}{2K}\right\} \leq R \sqrt{\frac{2 \log(4nt_e^2/\delta)}{T(e) - 1}} \leq R \sqrt{\frac{2 \log(4nT^2/\delta)}{T(e) - 1}}. \quad (64)$$

Solving Eq. (64) for $T(e)$, we obtain

$$T(e) \leq R^2 \min\left\{\frac{18 \text{width}(\mathcal{B})^2}{\Delta_e^2}, \frac{16K^2}{\epsilon^2}\right\} \log(4nT^2/\delta) + 1. \quad (65)$$

Notice that $T = \sum_{i \in [n]} T(i)$. Hence the theorem follows by summing up Eq. (65) for all $e \in [n]$ and solving for T . \square

C Proof of Lower Bound (Theorem 2)

In this section, we prove the problem-dependent lower bound of the general CPE problem (Theorem 2). In addition, we provide evidence on the conjecture that the sample complexity should hinge on the size of exchange sets (Theorem 7), which is relevant for decision classes with non-constant widths.

Notations. In this section, we will use a notation of “next-to-optimal set” defined as follows. Fix a decision class $\mathcal{M} \subseteq 2^{[n]}$ and an expected reward vector $\mathbf{w} \in \mathbb{R}^n$. Let $M_* = \arg \min_{M \in \mathcal{M}}$ denote the optimal set. Then, for any $e \in [n]$, we define the next-to-optimal set associated with e as follows

$$M_e = \begin{cases} \arg \max_{M \in \mathcal{M}: e \in M} w(M) & \text{if } e \notin M_*, \\ \arg \max_{M \in \mathcal{M}: e \notin M} w(M) & \text{if } e \in M_*. \end{cases} \quad (66)$$

We note that, by definition of Δ_e in Eq. (1), we have $w(M_*) - w(M_e) = \Delta_e$.

C.1 Proof of Theorem 2

For reader’s convenience, we restate Theorem 2 in the following.

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 φ_e is given by $\varphi_e = \mathcal{N}(w(e), 1)$, where we let $\mathcal{N}(\mu, \sigma^2)$ denote the Gaussian distribution with mean μ and variance σ^2 . Then, for any $\delta \in (0, e^{-16}/4)$ and any δ -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).

Before stating our proof, we first introduce two technical lemmas. The first lemma is the well-known Kolmogorov's inequality.

Lemma 14. (Kolmogorov's inequality [28, Corollary 7.66]) *Let Z_1, \dots, Z_n be independent zero-mean random variables with $\text{Var}[Z_k] \leq +\infty$ for all $k \in [n]$. Then, for any $\lambda > 0$,*

$$\Pr \left[\max_{1 \leq k \leq n} |S_k| \geq \lambda \right] \leq \frac{1}{\lambda^2} \sum_{i=1}^n \text{Var}[Z_k],$$

where $S_k = X_1 + \dots + X_k$.

The second technical lemma shows that the joint likelihood of Gaussian distributions on a sequence of variables does not change much when the mean of the distribution shifts by a sufficiently small value.

Lemma 15. *Given any $d \in \mathbb{R}$ and $\theta \in (0, 1)$. Define $t = \frac{1}{16d^2} \log(1/\theta)$. Given any integer $T \leq 4t$ and any sequence s_1, \dots, s_T . Let X_1, \dots, X_T be T real numbers which satisfy the following*

$$\left| \sum_{i=1}^T X_i - \sum_{i=1}^T s_i \right| \leq \sqrt{t \log(1/\theta)}. \quad (67)$$

Then, we have

$$\prod_{i=1}^T \frac{\mathcal{N}(X_i | s_i + d, 1)}{\mathcal{N}(X_i | s_i, 1)} \geq \theta,$$

where we let $\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$ denote the probability density function of normal distribution with mean μ and variance σ^2 .

Proof. We define $v_i = X_i - s_i$ for all $i \in [T]$. Then, we have

$$\begin{aligned} \prod_{i=1}^T \frac{\mathcal{N}(X_i | s_i + d, 1)}{\mathcal{N}(X_i | s_i, 1)} &= \prod_{i=1}^T \exp\left(\frac{-(X_i - s_i - d)^2 + (X_i - s_i)^2}{2}\right) \\ &= \prod_{i=1}^T \exp\left(-v_i d - \frac{1}{2} d^2\right) \\ &= \exp\left(-\sum_{i=1}^T v_i d\right) \exp\left(-\frac{T d^2}{2}\right). \end{aligned} \quad (68)$$

We now bound each term on the right-hand side of Eq. (68) as follows

$$\begin{aligned} \exp\left(-\sum_{i=1}^T v_i d\right) &\geq \exp\left(-\left|\sum_{i=1}^T v_i\right| \cdot |d|\right) \\ &\geq \exp\left(-\sqrt{t \log(1/\theta)} d\right) \end{aligned} \quad (69)$$

$$= \exp\left(-\frac{1}{2} \log(1/\theta)\right) = \theta^{1/2}, \quad (70)$$

where Eq. (69) follows from Eq. (80); and Eq. (70) follows from the definition of t . Next we have

$$\exp\left(-\frac{T d^2}{2}\right) \geq \exp(-2t d^2) \quad (71)$$

$$= \exp\left(-\frac{1}{2} \log(1/\theta)\right) = \theta^{1/2}, \quad (72)$$

where Eq. (71) follows from $T \leq 4t$ and Eq. (72) follows from the definition of t . The lemma follows immediate by combining Eq. (68), Eq. (70) and Eq. (72). \square

Proof of Theorem 2. Fix $\delta > 0$, $\mathbf{w} = (w(1), \dots, w(n))^T$ and a δ -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

$$\mathbb{E}[T_e] \geq \frac{1}{16\Delta_e^2} \log(1/4\delta). \quad (73)$$

Notice that the theorem follows immediately by summing up Eq. (73) for all $e \in [n]$.

Now fix an arm $e \in [n]$. We define $\theta = \delta/4$ and $t_e^* = \frac{1}{16\Delta_e^2} \log(1/4\delta)$. We prove Eq. (73) by contradiction. Therefore we assume the opposite that $\mathbb{E}[T_e] < t_e^*$ in the rest of the proof.

Step (1): An alternative hypothesis. We consider two hypothesis H_0 and H_1 . Under hypothesis H_0 , all reward distributions are same with our assumption in the theorem as follows

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

On the other hand, 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 .

Now we claim that M_* is no longer the optimal set under hypothesis H_1 . Let M_e denote the next-to-optimal set defined Eq. (66). By definition of Δ_e in Eq. (1), we know that $w(M_*) - w(M_e) = \Delta_e$. Let \mathbf{w}_0 and \mathbf{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.

Step (2): Three random events. Let X_1, \dots, X_{T_e} denote the sequence of reward outcomes of arm e . Now we define three random events \mathcal{A} , \mathcal{B} and \mathcal{C} as follows

$$\mathcal{A} = \{T_e \leq 4t_e^*\}, \mathcal{B} = \{\text{Out} = M_*\} \text{ and } \mathcal{C} = \left\{ \max_{1 \leq t \leq 4t_e^*} \left| \sum_{i=1}^t X_i - t \cdot w(e) \right| < \sqrt{t_e^* \log(1/\theta)} \right\},$$

where Out is the output of algorithm \mathbb{A} .

Now we bound the probability of these events under hypothesis H_0 . First, we show that $\Pr_0[\mathcal{A}] \geq 3/4$. This can be proved by Markov inequality as follows.

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

We now show that $\Pr_0[\mathcal{C}] \geq 3/4$. Notice that $\left\{ \sum_{i=1}^t X_i - t \cdot w(e) \right\}_{t=1, \dots}$ is a martingale under H_0 . Define $K_t = \sum_{i=1}^t X_i$. Then, by Kolmogorov's inequality (Lemma 14), 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 φ_e equals to 1 and therefore $\mathbb{E}_0[(K_{4t_e^*} - 4w(e)t_e^*)^2] = 4t_e^*$; the last inequality follows since $\theta < e^{-16}$.

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

Step (3): A change of measure. 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 .

Now assume that the event \mathcal{S} occurred. We will bound the likelihood ratio $L_1(W)/L_0(W)$ under this assumption. Since H_1 and H_0 only differs on the reward distribution of arm e , we have

$$\frac{L_1(W)}{L_0(W)} = \prod_{i=1}^{T_e} \frac{\mathcal{N}(X_i|w_1(e), 1)}{\mathcal{N}(X_i|w_0(e), 1)}. \quad (74)$$

By definition of H_1 and H_0 , we see that $w_1(e) = w_0(e) \pm \Delta_e$ (where the sign depends on whether $e \in M_*$). Therefore, when event \mathcal{S} occurs, it easy to verify that we can apply Lemma 15 (by setting $d = w_1(e) - w_0(e)$, $T = T_e$ and $s_i = w_0(e)$ for all i). Hence, by Lemma 15 and Eq. (74), we have, on event \mathcal{S} ,

$$\frac{L_1(W)}{L_0(W)} \geq \theta = 4\delta. \quad (75)$$

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

$$\frac{L_1(W)}{L_0(W)} 1_{\mathcal{S}} \geq 4\delta 1_{\mathcal{S}}$$

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

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

Now we have proved that, if $\mathbb{E}_0[T_e] < t_e^*$, then $\Pr_1[\mathcal{B}] > \delta$. This means that, if $\mathbb{E}_0[T_e] < t_e^*$, algorithm \mathbb{A} will choose M_* as the output with probability at least δ , 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 δ under H_1 . This contradicts to the assumption that algorithm \mathbb{A} is a δ -correct algorithm. Hence, we must have $\mathbb{E}_0[T_e] \geq t_e^* = \frac{1}{16\Delta_e^2} \log(1/4\delta)$. \square

C.2 A lower bound in the Bernoulli reward case

NOT COMPLETED!

Our main lower bound Theorem 2 covers the cases where the reward distributions are normal distributions, which have unbounded supports. The techniques used in Theorem 2 can be adapt to establish the lower bounds for many other types of distributions with unbounded support.

In this part, we prove Theorem 6, a bounded support version of Theorem 2 accounting for Bernoulli reward distribution, using a proof idea similar with Theorem 2. Therefore, the two versions of lower bounds settle the conjecture of Bubeck et al. [8] in both bounded and unbounded support cases. On the other hand, Theorem 6 only works for decision classes with constant width and therefore is more restricted than the lower bound for the unbounded support case.

Theorem 6. Let $\mathcal{M} \subseteq 2^{[n]}$ be a decision class such that $\text{width}(\mathcal{M}) = 2$. Fix some $p_0 \in (0, 1/2)$. There exists a positive constant δ_0 , and a positive constant c_1 that only depends on p_0 which satisfying the following. For any $w \in [p_0, 1/2]^n$, suppose that for each arm $e \in [n]$, the reward distribution φ_e is given by $\varphi_e = \text{Ber}(w(e))$, where we let $\text{Ber}(p)$ denote the Bernoulli distribution with parameter p . Then, for any $\delta \in (0, \delta_0)$ and any δ -correct algorithm \mathbb{A} , we have

$$\mathbb{E}[T] \geq c_1 \mathbf{H} \log \left(\frac{1}{8\delta} \right), \quad (76)$$

where T denote the number of total samples used by algorithm \mathbb{A} .

We require two technical lemmas in order to prove Theorem 6. The first lemma upper bounds the gaps of arms in terms of the width of decision class when the expected reward is bounded.

Lemma 16. *Fix some decision class $\mathcal{M} \subseteq 2^{[n]}$ and some expected reward vector $\mathbf{w} \in [A, B]^n$. Then, for any $e \in [n]$ we have*

$$\Delta_e \leq \text{width}(\mathcal{M}) \max\{|A|, |B|\}.$$

Proof. Let M_e denote the “next-to-optimal” set as defined in Eq. XXX. By definition of gaps Eq. (1), we know that $\Delta_e = w(M_*) - w(M_e)$. We fix an exchange class $\mathcal{B} \in \text{Exchange}(\mathcal{M})$ such that $\text{width}(\mathcal{B}) = \text{width}(\mathcal{M})$. We notice that $e \in (M_* \setminus M_e) \cup (M_e \setminus M_*)$ therefore we can apply Lemma 2. Hence there exists an exchange set $b = (b_+, b_-) \in \mathcal{B}$ such that $e \in b_+ \cup b_-$, $b_- \subseteq M_e \setminus M_*$, $b_+ \subseteq M_* \setminus M_e$, $M_* \ominus b \in \mathcal{M}$ and $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e > 0$. Hence we have

$$\Delta_e \leq \langle \mathbf{w}, \chi_b \rangle = w(M_*) - w(M_* \ominus b) \quad (77)$$

$$\begin{aligned} &= \sum_{i \in M_*} w(i) - \sum_{i \in M_* \setminus b_+ \cup b_-} w(i) \\ &= \sum_{i \in b_+} w(i) - \sum_{i \in b_-} w(i) \end{aligned} \quad (78)$$

$$\begin{aligned} &\leq |b_+| \cdot |B| + |b_-| \cdot |A| \leq (|b_+| + |b_-|) \max\{|A|, |B|\} \\ &\leq \text{width}(\mathcal{M}) \max\{|A|, |B|\}, \end{aligned} \quad (79)$$

where Eq. (77) follows from Lemma 1; Eq. (78) follows from the fact that $b_+ \subseteq M_*$ and $b_- \cap M_* = \emptyset$; and Eq. (79) holds since $b \in \mathcal{B}$ and $\text{width}(\mathcal{B}) = \text{width}(\mathcal{M})$. \square

The second technical lemma is a counterpart of Lemma 15 for the Bernoulli distribution. The lemma was first discovered by Mannor and Tsitsiklis [24] for analyzing the lower bound of identifying the single best arm. In the following, we rephrase their result using the notations of our proof.

Lemma 17. *Fix some $p_0 \in (0, 1/2)$. There exists a positive constant δ_0 , and a positive constant c_1 that only depends on p_0 such that the following holds. Given any $w_0 \in [p_0, 1/2]$, $\Delta \in (0, 1/2)$ and $\theta \in (0, \delta_0/8)$. Define $t = \frac{c_1}{\Delta^2} \log(1/\theta)$. Let $X_1, \dots, X_T \in \{0, 1\}$ be T binary variables which satisfy the following*

$$T \leq 4t \quad \text{and} \quad \left| \sum_{i=1}^T X_i - T w_0 \right| \leq \sqrt{t \log(1/\theta)}. \quad (80)$$

Then, for any $w_1 \in [p_0, 1/2]$ such that $|w_0 - w_1| \leq \Delta$, we have

$$\prod_{i=1}^T \frac{\text{Ber}(X_i | w_1)}{\text{Ber}(X_i | w_0)} \geq \theta,$$

where $\text{Ber}(x|p) = p^x(1-p)^{1-x}$ denotes the probability mass function of Bernoulli distribution.

The proof of Lemma 17 can be extracted from [24, Theorem 4] and therefore is omitted.

We are now ready to prove Theorem 6 using an argument which is almost identical to Theorem 2. There are only two noticeable differences: (1) we use Lemma 17 instead of Lemma 15 for bounding the likelihood loss; (2) we use Lemma 16 to ensure the constructed alternative hypothesis still contains valid Bernoulli reward distributions. In the following, we give a proof sketch of Theorem 6 highlighting these differences and omitting some redundant arguments.

Proof Sketch of Theorem 6. aa

Step (1): An alternative hypothesis. We consider two hypotheses H_0 and H_1 . The hypothesis H_0 is identical to the assumption of the theorem, i.e. $\varphi_l = \text{Ber}(w(l))$ for all $l \in [n]$. On the other hand, under hypothesis H_1 , we change the means of reward distributions such that

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

By Lemma 16, we see that $\Delta_e \leq 2ZZZ$. \square

C.3 An exchange set size dependent lower bound

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

Theorem 7. 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 φ_e is given by $\varphi_e = \mathcal{N}(w(e), 1)$, where $\mathcal{N}(\mu, \sigma^2)$ denotes a Gaussian distribution with mean μ and variance σ^2 . Fix any $\delta \in (0, e^{-16}/4)$ and any δ -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 $\delta > 0$, $\mathbf{w} = (w(1), \dots, w(n))^T$ and a δ -correct algorithm \mathbb{A} . For each $i \in [n]$, assume that the reward distribution is given by $\varphi_i = \mathcal{N}(w(i), 1)$. For any $i \in [n]$, let T_i denote the number of trials of arm i used by algorithm \mathbb{A} .

Step (0): Setup. Fix an arm $e \in [n]$. As the first step, we construct the exchange set $b = (b_+, b_-)$ claimed in the theorem. Let M_e denote the next-to-optimal set defined Eq. (66). By definition of Δ_e in Eq. (1), we know that $w(M_*) - w(M_e) = \Delta_e$. We construct the exchange set $b = (b_+, b_-)$ where $b_+ = M_* \setminus M_e$ and $b_- = M_e \setminus M_*$. It is easy to check that $M_e \oplus b = M_*$ and $\langle \mathbf{w}, \chi_b \rangle = \Delta_e > 0$.

We have now constructed the exchange set. We define $T_{b_-} = \sum_{i \in b_-} T_i$ and $T_{b_+} = \sum_{i \in b_+} T_i$. Now we claim that

$$(a) \quad \mathbb{E}[T_{b_-}] \geq \frac{|b_-|^2}{16\Delta_e^2} \log(1/4\delta) \quad \text{and} \quad (b) \quad \mathbb{E}[T_{b_+}] \geq \frac{|b_+|^2}{16\Delta_e^2} \log(1/4\delta). \quad (81)$$

It is easy to check that theorem follows immediately from claims (a) and (b). In the rest of the proof, we focus on claim (a); the claim (b) can be proved using an almost identical similar argument.

Step (1): An alternative hypothesis. We define two hypotheses H_0 and H_1 . 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

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

Similar to the proof of Theorem 2, we let \mathbf{w}_0 and \mathbf{w}_1 denote the expected reward vectors under H_0 and H_1 respectively. One can verify that $w_1(M_*) - w_1(M_e) = -\Delta_e < 0$. This means that under H_1 , the set M_* is not the optimal set.

Step (2): Three random events. First we 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 .

Now we define three random events \mathcal{A} , \mathcal{B} and \mathcal{C} as follows

$$\mathcal{A} = \{T_{b_-} \leq 4t_{b_-}^*\}, \quad \mathcal{B} = \{\text{Out} = M_*\} \quad \text{and} \quad \mathcal{C} = \left\{ \max_{1 \leq t \leq 4t_{b_-}^*} \left| \sum_{i=1}^t X_i - \sum_{i=1}^t w(I_i) \right| < \sqrt{t_{b_-}^* \log(1/\theta)} \right\},$$

where Out is the output of algorithm \mathbb{A} . We now bound the probability of each event. First, by Markov's inequality, we have

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

Next, using Kolmogorov's inequality (Lemma 14), we obtain

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

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 $\theta < e^{-16}$.

Since \mathbb{A} is δ -correct algorithm and $\delta < 1/4$, we have $\Pr_0[\mathcal{B}] \geq 3/4$. Therefore, we have that the random event $\mathcal{S} = \mathcal{A} \cap \mathcal{B} \cap \mathcal{C}$ occurs with probability at least $1/4$ under H_0 .

Step (3): A change of measure. Similar to the proof of Theorem 2, we let L_l denote the likelihood function under hypothesis H_l for some $l \in \{0, 1\}$. Since the difference of H_0 and H_1 only lies in the arms belonging to b_- , we have

$$\frac{L_1(W)}{L_0(W)} = \prod_{i=1}^{T_l} \frac{\mathcal{N}(X_i | w_1(I_i), 1)}{\mathcal{N}(X_i | w_0(I_i), 1)},$$

where X_i and I_i is as defined in Step (2). Assume that \mathcal{S} occurs. Since $w_1(e) - w_0(e) = d$ for all $e \in b_-$ and $d = \pm \frac{2\Delta_e}{|b_-|}$, we can apply Lemma 15 here. Therefore, on event \mathcal{S} , we have

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

The rest of the proof is identical to Step (3) in the proof of Theorem 2, and one can show that $\Pr_1[\mathcal{B}] \geq \delta$, which means the probability of error of algorithm \mathbb{A} is at least δ , which contradicts to the assumption of \mathbb{A} . Therefore we have $\mathbb{E}[T_{b_-}] \geq t_{b_-}^*$ which proves claim (a) in Eq. (81). \square

D Analysis of CSAR (Theorem 3)

Notations. For convenience, we will use the following additional 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 2. Let $M \subseteq [n]$ be a set, we denote $\neg M$ to be the complement of M . We will also continue to use the notations of incidence vectors of sets and exchange sets, which are defined in Appendix A.

D.1 Confidence Intervals

First we establish the confidence bounds used for the analysis of CSAR.

Lemma 18. *Given a phase $t \in [n]$, we define random event τ_t as follows*

$$\tau_t = \left\{ \forall i \in [n] \setminus (A_t \cup B_t) \quad |\bar{w}_t(i) - w(i)| < \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \right\}. \quad (82)$$

Then, we have

$$\Pr \left[\bigcap_{t=1}^n \tau_t \right] \geq 1 - n^2 \exp \left(- \frac{2(T-n)}{9R^2 \log(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right). \quad (83)$$

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

$$\Pr \left[|\bar{w}_t(i) - w(i)| \geq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \right] \leq 2 \exp \left(-\frac{2\tilde{T}_t \Delta_{(n-t+1)}^2}{9R^2 \text{width}(\mathcal{M})^2} \right). \quad (84)$$

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. (84) can be further bounded by

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

where the last inequality follows from the definition of $\mathbf{H}_2 = \max_{i \in [n]} i \Delta_{(i)}^{-2}$. By plugging the last inequality into Eq. (84), we have

$$\Pr \left[|\bar{w}_t(i) - w(i)| \geq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \right] \leq 2 \exp \left(-\frac{2(T-n)}{9R^2 \log(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2} \right). \quad (85)$$

Now using Eq. (85) and a union bound for all $t \in [n]$ and all $i \in [n] \setminus (A_t \cup B_t)$, we have

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

□

Readers may notice that the right-hand side of Eq. (83) equals to the probability of error of CSAR claimed in Theorem 3. Indeed, we will show that the CSAR algorithm will not make a mistake if the random event $\bigcap_{t=1}^n \tau_t$ occurs.

The following lemma builds the confidence bound of inner products.

Lemma 19. Fix a phase $t \in [n]$, suppose that random event τ_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 the support of \mathbf{a} . Then, we have

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

Proof. Suppose that τ_t occurs. Then, similar to the proof of Lemma 7, we have

$$\begin{aligned} |\langle \bar{\mathbf{w}}_t, \mathbf{a} \rangle - \langle \mathbf{w}, \mathbf{a} \rangle| &= |\langle \bar{\mathbf{w}}_t - \mathbf{w}, \mathbf{a} \rangle| \\ &= \left| \sum_{i=1}^n (\bar{w}_t(i) - w(i)) a(i) \right| \\ &\leq \left| \sum_{i \in [n] \setminus (A_t \cup B_t)} (\bar{w}_t(i) - w(i)) a(i) \right| \\ &\leq \sum_{i \in [n] \setminus (A_t \cup B_t)} |(\bar{w}_t(i) - w(i)) a(i)| \\ &\leq \sum_{i \in [n] \setminus (A_t \cup B_t)} |\bar{w}_t(i) - w(i)| |a(i)| \end{aligned} \quad (86)$$

$$\begin{aligned}
&< \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \sum_{i \in [n] \setminus (A_t \cup B_t)} |a(i)| \\
&= \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\mathbf{a}\|_1,
\end{aligned} \tag{87}$$

where Eq. (86) follows from the assumption that \mathbf{a} is supported on $[n] \setminus (A_t \cup B_t)$; Eq. (87) follows from the definition of τ_t (Eq. (82)). \square

D.2 Main Lemmas

We begin with a technical lemma which characterizes several useful properties of A_t and B_t .

Lemma 20. *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

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

$$\begin{aligned}
&= (A_t \cap M) \setminus a \\
&= A_t \setminus a
\end{aligned} \tag{89}$$

$$= A_t, \tag{90}$$

where Eq. (88) holds since we have $A_t \cap b \subseteq A_t \cap (M_* \setminus M) \subseteq M \cap (M_* \setminus M) = \emptyset$; Eq. (89) follows from $A_t \subseteq M$; and Eq. (90) follows from $a \subseteq M \setminus M_*$ and $A_t \subseteq M_*$ which imply that $a \cap A_t = \emptyset$. Notice that Eq. (90) is equivalent to $A_t \subseteq (M \setminus a \cup b)$.

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

$$\begin{aligned}
B_t \cap (M \setminus a \cup b) &= (B_t \cap (M \setminus a)) \cup (B_t \cap b) \\
&= B_t \cap (M \setminus a)
\end{aligned} \tag{91}$$

$$\begin{aligned}
&= (B_t \cap M) \setminus a \\
&= \emptyset \setminus a = \emptyset,
\end{aligned} \tag{92}$$

where Eq. (91) 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. (92) 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 $A_t \subseteq M \cap M_*$. Also note that $a \subseteq M \setminus M_*$ and $b \subseteq M_* \setminus M$, we have

$$(a \cap A_t) \cup (b \cap A_t) \subseteq ((M \setminus M_*) \cap (M \cap M_*)) \cup ((M_* \setminus M) \cap (M \cap M_*)) = \emptyset. \tag{93}$$

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

$$(a \cap B_t) \cup (b \cap B_t) \subseteq ((M \setminus M_*) \cap (\neg M \cap \neg M_*)) \cup ((M_* \setminus M) \cap (\neg M \cap \neg M_*)) = \emptyset. \tag{94}$$

By combining Eq. (93) and Eq. (94), we obtain

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

\square

The next lemma provides an important insight on the correctness of CSAR. Informally speaking, suppose that the algorithm does not make an error before phase t . Then, we show that, suppose arm e has a gap Δ_e larger than the “reference” gap $\Delta_{(t-n+1)}$ of phase t , then arm e must be correctly classified by M_t , i.e. $e \in M_t$ if and only if $e \in M_*$.

Lemma 21. Fix any phase $t > 0$. Suppose that event τ_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

$$e \in (M_* \cap \neg M_t) \cup (\neg M_* \cap M_t). \quad (95)$$

Eq. (95) 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 20, 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

$$\bar{w}_t(M_t) \geq \bar{w}_t(M_t \oplus b). \quad (96)$$

On the other hand, we have

$$\bar{w}_t(M_t \oplus b) = \langle \bar{\mathbf{w}}_t, \chi_{M_t} + \chi_b \rangle \quad (97)$$

$$\begin{aligned} &= \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \bar{\mathbf{w}}_t, \chi_b \rangle \\ &> \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\chi_b\|_1 \end{aligned} \quad (98)$$

$$\begin{aligned} &\geq \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3 \text{width}(\mathcal{M})} \|\chi_b\|_1 \\ &\geq \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3} \end{aligned} \quad (99)$$

$$\geq \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle + \frac{2}{3} \Delta_e \quad (100)$$

$$\geq \langle \bar{\mathbf{w}}_t, \chi_{M_t} \rangle = \bar{w}_t(M_t), \quad (101)$$

where Eq. (97) follows from Lemma 1; Eq. (98) follows from Lemma 19 and the fact that $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$; Eq. (99) 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. (100) and Eq. (101) 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. (96). Therefore we have $e \in (M_* \cap M_t) \cup (\neg M_* \cap \neg M_t)$. \square

The next lemma takes a step further. It shows that if $\Delta_e \geq \Delta_{(t-n+1)}$ for some arm e , then the empirical gap of arm e is greater than $\frac{2}{3} \Delta_{(t-n+1)}$.

Lemma 22. Fix any phase $t > 0$. Suppose that event τ_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

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

Proof. By Lemma 21, we see that

$$e \in (M_* \cap M_t) \cup (\neg M_* \cap \neg M_t). \quad (102)$$

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. (102), 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 $e \in M_* \setminus \tilde{M}_{t,e}$. 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 20, 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

$$\bar{w}_t(M_t) \geq \bar{w}_t(M'_{t,e}). \quad (103)$$

Hence, we have

$$\begin{aligned} \bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,e}) &\geq \bar{w}_t(M'_{t,e}) - \bar{w}_t(\tilde{M}_{t,e}) \\ &= \bar{w}_t(\tilde{M}_{t,e} \oplus b) - \bar{w}_t(\tilde{M}_{t,e}) \\ &= \langle \bar{\mathbf{w}}_t, \chi_{\tilde{M}_{t,e}} + \chi_b \rangle - \langle \bar{\mathbf{w}}_t, \chi_{\tilde{M}_{t,e}} \rangle \\ &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle \end{aligned} \quad (104)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{B})} \|\chi_b\|_1 \quad (105)$$

$$\geq \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3 \text{width}(\mathcal{B})} \|\chi_b\|_1 \quad (106)$$

$$\geq \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_e}{3} \quad (107)$$

$$\geq \frac{2}{3} \Delta_e \geq \frac{2}{3} \Delta_{(n-t+1)}, \quad (108)$$

where Eq. (104) follows from Lemma 1; Eq. (105) follows from Lemma 19, the assumption on event τ_t and the fact $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$; Eq. (106) follows from the assumption that $\Delta_e \geq \Delta_{(n-t+1)}$; Eq. (107) holds since $b \in \mathcal{B}$ and therefore $\|\chi_b\|_1 \leq \text{width}(\mathcal{M})$; and Eq. (108) follows from the fact that $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e$.

□

The next lemma shows that if $\Delta_e \leq \Delta_{(t-n+1)}$ for some arm e , then the empirical gap of arm e is smaller than $\frac{1}{3} \Delta_{(t-n+1)}$.

Lemma 23. Fix any phase $t > 0$. Suppose that event τ_t occurs. Also assume that $A_t \subseteq M_*$ and $B_t \cap M_* = \emptyset$. 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

$$\bar{w}_t(\tilde{M}_{t,e}) \geq \bar{w}_t(M'_{t,e}). \quad (109)$$

From the definition of $\tilde{M}_{t,e}$ in Algorithm 2, 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 $e \in (b_- \cup b_+) = (M'_{t,e} \setminus M_t) \cup (M_t \setminus M'_{t,e})$, therefore we see that **(a)** holds. Next, we notice that **(b)** follows directly from Lemma 20 by setting $M = M_t$. Hence we have shown that Eq. (109) holds.

Hence, we have

$$\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,e}) \leq \bar{w}_t(M_t) - \bar{w}_t(M'_{t,e})$$

$$= \langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_{M_t} \rangle - \langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_{M_t} + \boldsymbol{\chi}_b \rangle \quad (110)$$

$$= -\langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_b \rangle$$

$$\leq -\langle \mathbf{w}, \boldsymbol{\chi}_b \rangle + \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\boldsymbol{\chi}_b\|_1 \quad (111)$$

$$\leq \frac{\Delta_{(n-t+1)}}{3 \text{width}(\mathcal{M})} \|\boldsymbol{\chi}_b\|_1 \leq \frac{\Delta_{(n-t+1)}}{3}, \quad (112)$$

where Eq. (110) follows from Lemma 1; Eq. (111) follows from Lemma 19, the assumption on τ_t and $(b_+ \cup b_-) \cap (A_t \cup B_t) = \emptyset$ (by Lemma 20); and Eq. (112) follows from the fact $\|\boldsymbol{\chi}_b\|_1 \leq \text{width}(\mathcal{M})$ (since $b \in \mathcal{B}$) and that $\langle \mathbf{w}, \boldsymbol{\chi}_b \rangle \geq \Delta_e \geq 0$. \square

D.3 Proof of Theorem 3

Using these technical lemmas, we are now ready to prove Theorem 3. For reader's convenience, we first restate Theorem 3 as follows.

Theorem 3. *Given any $T > n$, any decision class $\mathcal{M} \subseteq 2^{[n]}$ and any expected rewards $\mathbf{w} \in \mathbb{R}^n$. Assume that the reward distribution φ_e for each arm $e \in [n]$ has mean $w(e)$ with an R -sub-Gaussian tail. Let $\Delta_{(1)}, \dots, \Delta_{(n)}$ be a permutation of $\Delta_1, \dots, \Delta_n$ (defined in Eq. (1)) such that $\Delta_{(1)} \leq \dots \leq \Delta_{(n)}$. Define $\mathbf{H}_2 \triangleq \max_{i \in [n]} i \Delta_{(i)}^{-2}$. Then, the CSAR algorithm uses at most T samples and outputs a solution $\text{Out} \in \mathcal{M} \cup \{\perp\}$ such that*

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

where $\log(n) \triangleq \sum_{i=1}^n i^{-1}$, $M_* = \arg \max_{M \in \mathcal{M}} w(M)$ and $\text{width}(\mathcal{M})$ is defined in Eq. (4).

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

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

By Lemma 18, we know that the event $\tau \triangleq \bigcap_{t=1}^T \tau_t$ occurs with probability at least $1 - n^2 \exp\left(-\frac{2(T-n)}{9R^2 \log(n) \text{width}(\mathcal{M})^2 \mathbf{H}_2}\right)$. Therefore, we only need to prove that, under event τ , the algorithm outputs M_* . We will assume that event τ 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 .

At the beginning of phase t , there are exactly $t-1$ 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 22, we have

$$\bar{w}_t(M_t) - \bar{w}_t(M_{t,e_1}) \geq \frac{2}{3} \Delta_{(n-t+1)}. \quad (113)$$

Notice that the algorithm makes an error in 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 in 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 23, we see that

$$\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,p_t}) \leq \frac{1}{3} \Delta_{(n-t+1)}. \quad (114)$$

By combining Eq. (113) and Eq. (114), we see that

$$\bar{w}_t(M_t) - \bar{w}_t(\tilde{M}_{t,p_t}) \leq \frac{1}{3}\Delta_{(n-t+1)} < \frac{2}{3}\Delta_{(n-t+1)} \leq \bar{w}_t(M_t) - \bar{w}_t(M_{t,e_1}). \quad (115)$$

However Eq. (115) 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 $B_{n+1} \subseteq \neg M_*$ in the final phase. 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 Analysis of the Uniform Allocation Algorithm

In this section, we analyze the performance of a simple benchmark strategy UNI which plays each arm for an equal number of times and then calls a maximization oracle using the empirical means of arms as input. The pseudo-code of the UNI algorithm is listed in Algorithm 3.

Algorithm 3 UNI: Uniform Allocation

Require: Budget: $T > 0$; Maximization oracle: Oracle : $\mathbb{R}^n \rightarrow \mathcal{M}$.

- 1: Pull each arm $e \in [n]$ for $\lfloor T/n \rfloor$ times.
 - 2: Compute the empirical means $\bar{w} \in \mathbb{R}^n$ of each arm.
 - 3: $\text{Out} \leftarrow \text{Oracle}(\bar{w})$
 - 4: **return:** Out
-

The next theorem upper bounds the probability of error of UNI.

Theorem 8. *Given any $T > n$, any decision class $\mathcal{M} \subseteq 2^{[n]}$ and any expected rewards $w \in \mathbb{R}^n$. Assume that the reward distribution φ_e for each arm $e \in [n]$ has mean $w(e)$ with an R -sub-Gaussian tail. Define $\Delta_{(1)} = \min_{i \in [n]} \Delta_i$ and $\mathbf{H}_3 = n\Delta_{(1)}^{-2}$. Then, the output Out of the UNI algorithm satisfies*

$$\Pr[\text{Out} \neq M_*] \leq 2n \exp\left(-\frac{2T}{9R^2 \text{width}(\mathcal{M})^2 \mathbf{H}_3}\right), \quad (116)$$

where $M_* = \arg \max_{M \in \mathcal{M}} w(M)$.

From Theorem 8, we see that the UNI algorithm could be significantly worse than CLUCB and CSAR, since it is clear that $\mathbf{H}_3 \geq \mathbf{H} \geq \mathbf{H}_2$ and potentially one has $\mathbf{H}_3 \gg \mathbf{H} \geq \mathbf{H}_2$ for a large number of arms with heterogeneous gaps.

Now we prove Theorem 8. The proof is straightforward using tools of exchange classes.

Proof. Define $\Delta_{(1)} = \min_{i \in [n]} \Delta_i$. Define random event ξ as follows

$$\xi = \left\{ \forall i \in [n], \quad |\bar{w}(i) - w(i)| \leq \frac{\Delta_{(1)}}{3 \text{width}(\mathcal{M})} \right\}.$$

Notice that each arm is sampled for $\lfloor \frac{T}{n} \rfloor$ times. Therefore, using Hoeffding's inequality and union bound, we can bound $\Pr[\xi]$ as follows. Fix any $i \in [n]$, by Hoeffding's inequality, we have

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

Then, using a union bound, we obtain

$$\Pr[\xi] \geq 1 - 2n \exp\left(-\frac{2T\Delta_{(1)}^2}{9nR^2 \text{width}(\mathcal{M})^2}\right).$$

In addition, using an argument very similar to Lemma 19, one can show that, on event ξ , for any vector $\mathbf{a} \in \mathbb{R}^n$, it holds that

$$|\langle \bar{\mathbf{w}}, \mathbf{a} \rangle - \langle \mathbf{w}, \mathbf{a} \rangle| \leq \frac{\Delta_{(1)}}{3 \text{width}(\mathcal{M})} \|\mathbf{a}\|_1. \quad (117)$$

Now we claim that, on the event ξ , we have $\text{Out} = M_*$. Note that theorem follows immediately from the claim. Next, we prove this claim.

Suppose that, on the contrary, $\text{Out} \neq M_*$. First, we write $M = \text{Out}$. We also fix $\mathcal{B} \in \arg \min_{\mathcal{B}' \in \text{Exchange}(\mathcal{M})} \text{width}(\mathcal{B}')$. Notice that by definition $\text{width}(\mathcal{B}) = \text{width}(\mathcal{M})$.

Since $M \neq M_*$, we see that there exists $e \in (M \setminus M_*) \cup (M_* \setminus M)$. Now, by Lemma 2, we obtain that there exists $b = (b_+, b_-) \in \mathcal{B}$ such that $e \in b_+ \cup b_-$, $b_- \subseteq M \setminus M_*$, $b_+ \subseteq M_* \setminus M$, $M \oplus b \in \mathcal{M}$ and $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_e$. Also notice that $\Delta_e \geq \Delta_{(1)}$. Therefore $\langle \mathbf{w}, \chi_b \rangle \geq \Delta_{(1)}$.

Consider $M' \triangleq M \oplus b$. We have

$$\begin{aligned} \bar{w}(M') - \bar{w}(M) &= \langle \bar{\mathbf{w}}, \chi_{M'} \rangle - \langle \bar{\mathbf{w}}, \chi_M \rangle \\ &= \langle \bar{\mathbf{w}}, \chi_b \rangle \end{aligned} \quad (118)$$

$$\geq \langle \mathbf{w}, \chi_b \rangle - \frac{\Delta_{(1)}}{3 \text{width}(\mathcal{M})} \|\chi_b\|_1 \quad (119)$$

$$\geq \Delta_{(1)} - \frac{\Delta_{(1)}}{3} \quad (120)$$

$$= \frac{2}{3} \Delta_{(1)} > 0, \quad (121)$$

where Eq. (118) follows from Lemma 1; Eq. (119) follows from Eq. (117); and Eq. (120) follows from the fact that $b \in \mathcal{B}$ and hence $\|\chi_b\|_1 = |b_+| + |b_-| \leq \text{width}(\mathcal{B}) = \text{width}(\mathcal{M})$.

Hence, we have shown that $\bar{w}(M') > \bar{w}(M)$. However this contradicts to the fact that $w(M) = \max_{M_1 \in \mathcal{M}} \bar{w}(M_1)$ (by the definition of maximization oracle). Hence, by contradiction, we have proven that $\text{Out} = M_*$. \square

F Exchange Classes for Example Decision Classes

In this section, we give formal constructions of the decision classes discussed in Example 1, 2 and 3. Further, we bound the width of exchange classes for different examples. These bounds are proven using concrete constructions of exchange classes (Fact 1 through 5). The constructed exchange classes embody natural combinatorial structures. We illustrate the constructed exchange classes in Figure 1.

Notation. We need one extra notation. Let $\sigma : E \rightarrow [n]$ be a bijection from some set E with n elements to $[n]$. Let $A \subseteq E$ be an arbitrary set, we define $\sigma(A) \triangleq \{\sigma(a) \mid a \in A\}$. Conversely, for all $M \subseteq [n]$, we define $\sigma^{-1}(M) \triangleq \{\sigma^{-1}(e) \mid e \in M\}$.

Fact 1 (Matroid). Let $T = (E, \mathcal{I})$ be an arbitrary matroid, where E is the ground set of n elements and \mathcal{I} is the family of subsets of E called in the independent sets which satisfy the axioms of matroids³. Let $\sigma : E \rightarrow [n]$ be a bijection from E to $[n]$. Let $\mathcal{M}_{\text{MATROID}(T)}$ correspond to the collection of all bases of matroid T and formally we define

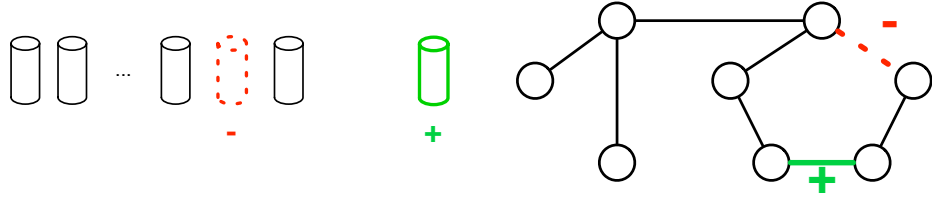
$$\mathcal{M}_{\text{MATROID}(T)} = \{M \subseteq [n] \mid \sigma^{-1}(M) \text{ is a basis of } T\}. \quad (122)$$

Define the exchange class

$$\mathcal{B}_{\text{MATROID}(n)} = \{(\{i\}, \{j\}) \mid \forall i \in [n], j \in [n]\}. \quad (123)$$

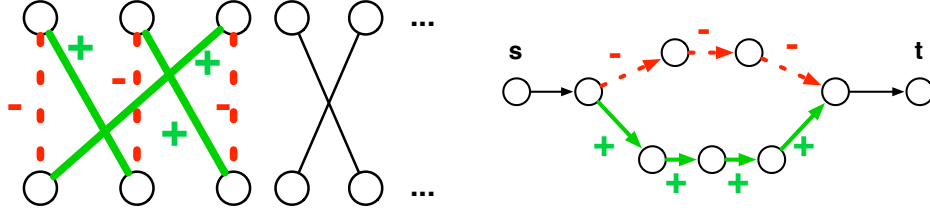
Then we have $\mathcal{B}_{\text{MATROID}(n)} \in \text{Exchange}(\mathcal{M}_{\text{MATROID}(T)})$. In addition, we have $\text{width}(\mathcal{B}_{\text{MATROID}(n)}) = 2$, which implies that $\text{width}(\mathcal{M}_{\text{MATROID}(T)}) \leq 2$.

³The three axioms of matroid are (1) $\emptyset \in \mathcal{I}$ and $\mathcal{I} \neq \{\emptyset\}$; (2) Every subsets of an independent set are independent (heredity 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 [26] for a general introduction to the matroid theory.



(a) An exchange set from $\mathcal{B}_{\text{MATROID}(n)}$ (TOPK: Fact 2; each cylinder represents an arm).

(b) An exchange set from $\mathcal{B}_{\text{MATROID}(n)}$ (Spanning trees: Fact 1; each edge corresponds to an arm).



(c) An exchange set from $\mathcal{B}_{\text{MATCH}(G)}$ (Matchings: Fact 4; each edge corresponds to an arm)

(d) An exchange set from $\mathcal{B}_{\text{PATH}(G)}$ (Paths: Fact 5; each edge corresponds to an arm).

Figure 1: Examples of exchange sets belonging to the exchange classes $\mathcal{B}_{\text{MATROID}}$ (one for TOPK and one for spanning tree), $\mathcal{B}_{\text{MATCH}(G)}$ and $\mathcal{B}_{\text{PATH}(G)}$: green-solid elements constitute the set b_+ , red-dotted elements constitute the set b_- and an example exchange set is given by $b = (b_+, b_-)$. In Figure 1a, we use TOPK as a specific instance of matroid decision class. In Figure 1b, we use spanning tree as a specific instance of matroid decision class.

To prove Fact 1, we first recall a well-known result from matroid theory which is referred as the strong basis exchange property.

Lemma 24 (Strong basis exchange [26]). *Let \mathcal{A} be the set of all bases of a matroid $T = (E, \mathcal{I})$. Let $A_1, A_2 \in \mathcal{A}$ be two bases. Then for all $x \in A_1 \setminus A_2$, there exists $y \in A_2 \setminus A_1$ such that $A_1 \setminus \{x\} \cup \{y\} \in \mathcal{A}$ and $A_2 \setminus \{y\} \cup \{x\} \in \mathcal{A}$.*

We refer readers to [26] for a proof of Lemma 24.

Proof of Fact 1. Fix a matroid $T = (E, \mathcal{I})$ where $|E| = n$ and fix the bijection $\sigma: E \rightarrow [n]$. Let $\mathcal{M}_{\text{MATROID}(T)}$ be defined as in Eq. (122) and let $\mathcal{B}_{\text{MATROID}(n)}$ be defined as in Eq. (123). Let \mathcal{A} denote the set of all bases of T . By definition, we have $\mathcal{M}_{\text{MATROID}(T)} = \{\sigma(A) \mid A \in \mathcal{A}\}$.

Now we show that $\mathcal{B}_{\text{MATROID}(n)}$ is an exchange class for $\mathcal{M}_{\text{MATROID}(T)}$. Let M, M' be two different elements of $\mathcal{M}_{\text{MATROID}(T)}$. By definition, we see that $\sigma^{-1}(M)$ and $\sigma^{-1}(M')$ are two bases of T . Consider any $e \in M \setminus M'$. Let $x = \sigma^{-1}(e)$. We see that $x \in \sigma^{-1}(M) \setminus \sigma^{-1}(M')$.

By Lemma 24, we see that there exists $y \in \sigma^{-1}(M') \setminus \sigma^{-1}(M)$ such that

$$\sigma^{-1}(M) \setminus \{x\} \cup \{y\} \in \mathcal{A} \quad \text{and} \quad \sigma^{-1}(M') \setminus \{y\} \cup \{x\} \in \mathcal{A}. \quad (124)$$

Now we define exchange set $b = (b_+, b_-)$ where $b_+ = \{\sigma(y)\}$ and $b_- = \{\sigma(x)\}$. By Eq. (124) and the fact that σ is a bijection, we see that $M \oplus b \in \mathcal{M}_{\text{MATROID}(T)}$ and $M' \ominus b \in \mathcal{M}_{\text{MATROID}(T)}$. We also have $b \in \mathcal{B}_{\text{MATROID}(n)}$. Due to M, M' and e are chosen arbitrarily, we have verified that $\mathcal{B}_{\text{MATROID}(n)}$ is an exchange class for $\mathcal{M}_{\text{MATROID}(T)}$.

To conclude, we observe that $\text{width}(\mathcal{B}_{\text{MATROID}(n)}) = 2$. \square

Now we show that TOPK and MB can be reduced to the decision classes of matroids. And therefore we can apply Fact 1 to construct exchange classes and bound the widths of these decision classes.

Fact 2 (TOPK). For all $K \in [n]$, let $\mathcal{M}_{\text{TOPK}(K)} = \{M \subseteq [n] \mid |M| = K\}$ be the collection of all subsets of size K . Then we have $\text{width}(\mathcal{M}_{\text{TOPK}(K)}) \leq 2$.

Proof. Let $U_n^K = ([n], \mathcal{I}_K)$ where \mathcal{I}_K is given by

$$\mathcal{I}_K = \{M \subseteq [n] \mid |M| \leq K\}.$$

Recall that U_n^K is a matroid (in particular, a uniform matroid of rank K) [26]. We know that a subset M of $[n]$ is basis of U_n^K if and only if $|M| = K$. Therefore, we have $\mathcal{M}_{\text{TOPK}(K)} = \mathcal{M}_{\text{MATROID}(U_n^K)}$. Then the conclusion follows immediately from Fact 1. \square

Fact 3 (MB). For any partition $\mathcal{A} = \{A_1, \dots, A_m\}$ of $[n]$, we define

$$\mathcal{M}_{\text{MB}(\mathcal{A})} = \left\{ M \subseteq [n] \mid \forall i \in [m] \quad |M \cap A_i| = 1 \right\}.$$

Then we have $\text{width}(\mathcal{M}_{\text{MB}(\mathcal{A})}) \leq 2$.

Proof. Let $P_{\mathcal{A}} = ([n], \mathcal{I}_{\mathcal{A}})$ where $\mathcal{I}_{\mathcal{A}}$ is given by

$$\mathcal{I}_{\mathcal{A}} = \{M \subseteq [n] \mid \forall i \in [m] \quad |M \cap A_i| \leq 1\}.$$

It can be shown that $P_{\mathcal{A}}$ is a matroid (known as partition matroid [26]) and each basis M of $P_{\mathcal{A}}$ satisfies $|M \cap A_i| = 1$ for all $i \in [m]$. Therefore we have $\mathcal{M}_{\text{MB}(\mathcal{A})} = \mathcal{M}_{\text{MATROID}(P_{\mathcal{A}})}$. Then the conclusion follows immediately from Fact 1. \square

Fact 4 (Matching). Let $G(V, E)$ be a bipartite graph with n edges. Let $\sigma: E \rightarrow [n]$ be a bijection. Let \mathcal{A} be the set of all valid matchings in G . We define $\mathcal{M}_{\text{MATCH}(G)}$ as follows

$$\mathcal{M}_{\text{MATCH}(G)} = \{\sigma(A) \mid A \in \mathcal{A}\}.$$

Define the exchange class

$$\mathcal{B}_{\text{MATCH}(G)} = \left\{ (\sigma(c_+), \sigma(c_-)) \mid \exists c \in \mathcal{C} \cup \mathcal{P}, \text{ the edges of } c \text{ alternate between } c_+, c_- \right\},$$

where \mathcal{C} is the set of all cycles in G and \mathcal{P} is the set of all paths in G . Then we have $\mathcal{B}_{\text{MATCH}(G)} \in \text{Exchange}(\mathcal{M}_{\text{MATCH}(G)})$. In addition, we have $\text{width}(\mathcal{B}_{\text{MATCH}(G)}) \leq |V|$, which implies that $\text{width}(\mathcal{M}_{\text{MATCH}(G)}) \leq |V|$.

To prove Fact 4, we recall a classical result on graph matching which characterizes the properties of augmenting cycles and augmenting paths [4].

Lemma 25. Let $G(V, E)$ be a bipartite graph. Let M and M' be two different matchings. Then the induced graph G' from the symmetric difference $(M \setminus M') \cup (M' \setminus M)$ consists of connected components that are one of the following

- An even cycle whose edges alternate between M and M' .
- A simple path whose edges alternate between M and M' .

Proof of Fact 4. Fix a bipartite graph $G(V, E)$ and a bijection $\sigma: E \rightarrow [n]$. Let $M, M' \in \mathcal{M}_{\text{MATCH}(G)}$ be two different elements of $\mathcal{M}_{\text{MATCH}(G)}$ and consider an arbitrary $e \in M \setminus M'$. On a high level perspective, we construct an exchange class which contains all augmenting cycles and paths of G . We know that the symmetric difference between M and M' can be decomposed into a collection of disjoint augmenting cycles and paths. And e must be on one of the augmenting cycle or path. Then, since “applying” this augmenting cycle/path on M will yield another valid matching which does not contains e . We see that this meets the requirements of an exchange class. In the rest of the proof, we carry out the technical details of this argument.

Define $A = \sigma^{-1}(M)$ and $A' = \sigma^{-1}(M')$. Let $a = \sigma^{-1}(e)$. Then A, A' are two matchings of G . Let G' be the induced graph from the symmetric difference $(A \setminus A') \cup (A' \setminus A)$. Let C be the connected component of G' which contains the edge a . Therefore, by Lemma 25, we see that C is either an even cycle or a simple path with edges alternating between A and A' . Let C_+ contains the edges

of C that belongs to $A' \setminus A$. Similarly, let C_- contains the edges of C that belongs to $A \setminus A'$. Define $b_+ = \sigma(C_+)$ and $b_- = \sigma(C_-)$. Let $b = (b_+, b_-)$ be an exchange set.

We see that $b \in \mathcal{B}_{\text{MATCH}(G)}$. Since $a \in C_-$, we obtain that $e \in b_-$. In addition, note that $C_+ \subseteq A' \setminus A$ and $C_- \subseteq A \setminus A'$. Therefore we have $b_+ \subseteq M' \setminus M$ and $b_- \subseteq M \setminus M'$.

Since C is an A -augmenting path/cycle, therefore it immediately holds that $A \setminus C_- \cup C_+$ is a valid matching. Therefore, we have $M \setminus b_- \cup b_+ \in \mathcal{M}_{\text{MATCH}(G)}$. Similarly, one can show that $M' \setminus b_+ \cup b_- \in \mathcal{M}_{\text{MATCH}(G)}$. Hence we have shown that $\mathcal{B}_{\text{MATCH}(G)}$ is an exchange class for $\mathcal{M}_{\text{MATCH}(G)}$. \square

Fact 5 (Path). Let $G(V, E)$ be a directed acyclic graph with n edges. Let $s, t \in V$ be two different vertices. Let $\sigma: E \rightarrow [n]$ be a bijection. Let $\mathcal{A}(s, t)$ be the set of all valid paths from s to t in G . We define $\mathcal{M}_{\text{PATH}(G, s, t)}$ as follows

$$\mathcal{M}_{\text{PATH}(G, s, t)} = \{\sigma(A) \mid A \in \mathcal{A}(s, t)\}.$$

Define exchange class

$$\mathcal{B}_{\text{PATH}(G)} = \{(\sigma^{-1}(p), \sigma^{-1}(q)) \mid p, q \text{ are the arcs of two disjoint paths of } G \text{ with same endpoints}\}.$$

Then, we have $\mathcal{B}_{\text{PATH}(G)} \in \text{Exchange}(\mathcal{M}_{\text{PATH}(G, s, t)})$. In addition, we have $\text{width}(\mathcal{B}_{\text{PATH}(G)}) \leq |V|$ and therefore $\text{width}(\mathcal{M}_{\text{PATH}(G, s, t)}) \leq |V|$.

Proof. Fix a directed acyclic graph $G(V, E)$ and a bijection $\sigma: E \rightarrow [n]$. Fix two vertices $s, t \in V$.

We prove that $\mathcal{B}_{\text{PATH}(G)}$ is an exchange class for $\mathcal{M}_{\text{PATH}(G, s, t)}$. Let $M, M' \in \mathcal{M}_{\text{PATH}(G, s, t)}$ be two different sets. Then $\sigma^{-1}(M), \sigma^{-1}(M')$ are two sets of arcs corresponding to two different paths from s to t . Let $P = (v_1, \dots, v_{n_1}), P' = (v'_1, \dots, v'_{n_2})$ denote two paths, respectively. Also denote $E(P) = \sigma^{-1}(M)$ and $E(P') = \sigma^{-1}(M')$.

Fix $e \in M \setminus M'$ and define $a = \sigma^{-1}(e)$. Suppose that a is an arc from u to v . Suppose $v_i = u$ and $v_{i+1} = v$. Define $j_1 = \arg \max_{j \leq i, v_j \in P'} j$ and $j_2 = \arg \min_{j \geq i+1, v_j \in P'} j$. Let $v'_{k_1} = v_{j_1}$ and $v'_{k_2} = v_{j_2}$ be the corresponding indices in P' . Then, we see that $Q_1 = (v_{j_1}, v_{j_1+1}, \dots, v_{j_2})$ and $Q_2 = (v'_{k_1}, v'_{k_1+1}, \dots, v'_{k_2})$ are two different paths from v_{j_1} to v_{j_2} . Denote the sets of arcs of Q_1 and Q_2 as $E(Q_1)$ and $E(Q_2)$.

Let $b = (b_+, b_-)$, where $b_+ = \sigma(E(Q_2)), b_- = \sigma(E(Q_1))$. We see that $b \in \mathcal{B}_{\text{PATH}(G)}$. It is clear that $a \in E(Q_1), E(Q_1) \subseteq E(P) \setminus E(P')$ and $E(Q_2) \subseteq E(P') \setminus E(P)$. Therefore $e \in b_-, b_- \subseteq M \setminus M'$ and $b_+ \subseteq M' \setminus M$.

Now it is easy to check that $E(P_1) \setminus E(Q_1) \cup E(Q_2)$ equals the set of arcs of path $(v_1, \dots, v_{j_1}, v'_{k_1+1}, \dots, v'_{k_2-1}, v_{j_2}, \dots, v_{n_1})$ (recall that $v_{j_1} = v'_{k_1}$ and $v_{j_2} = v'_{k_2}$). This means that $E(P_1) \setminus E(Q_1) \cup E(Q_2) \in \mathcal{A}(s, t)$ and therefore $M \setminus b_- \cup b_+ \in \mathcal{M}_{\text{PATH}(G, s, t)}$. Using a similar argument, one can show that $M' \setminus b_+ \cup b_- \in \mathcal{M}_{\text{PATH}(G, s, t)}$ and hence we have verified that $\mathcal{B}_{\text{PATH}(G)} \in \text{Exchange}(\mathcal{M}_{\text{PATH}(G, s, t)})$. \square

G Equivalence Between Constrained Oracles and Maximization Oracles

In this section, we present a general method to implement constrained oracles using maximization oracles. The idea of the reduction is simple: one can impose the negative constraints B by setting the corresponding weights to be sufficiently small; and one can impose the positive constraints A by setting the corresponding weights sufficiently large. The reduction method is shown in Algorithm 4. The correctness of the reduction is proved in Lemma 26. Furthermore, it is trivial to reduce from maximization oracles to constrained oracles. Therefore, Lemma 26 shows that maximization oracles are equivalent to constrained oracles up to a transformation on the weight vector.

Lemma 26. Given $\mathcal{M} \subseteq 2^{[n]}$, $\mathbf{w} \in \mathbb{R}^n$, $A \subseteq [n]$ and $B \subseteq [n]$, suppose that $A \cap B = \emptyset$. Then the output Out of Algorithm 4 satisfies $\text{Out} \in \arg \max_{M \in \mathcal{M}, A \subseteq M, B \cap M = \emptyset} w(M)$ where we use the convention that the $\arg \max$ of an empty set is \perp . Therefore Algorithm 4 is a valid constrained oracle.

Algorithm 4 COracle(w, A, B)

Require: $w \in \mathbb{R}^n$, $A \subseteq [n]$, $B \subseteq [n]$; Maximization oracle Oracle : $\mathbb{R}^n \rightarrow \mathcal{M}$

```

1:  $L_1 \leftarrow \|w\|_1$ 
2: for  $i = 1, \dots, n$  do
3:   if  $i \in A$  then
4:      $w_1(i) \leftarrow 3L_1$ 
5:   else
6:      $w_1(i) \leftarrow w(i)$ 
7:  $L_2 \leftarrow \|w_1\|_1$ 
8: for  $i = 1, \dots, n$  do
9:   if  $i \in B$  then
10:     $w_2(i) \leftarrow -3L_2$ 
11:   else
12:     $w_2(i) \leftarrow w_1(i)$ 
13:  $M \leftarrow \text{Oracle}(w_2)$ 
14: if  $B \cap M = \emptyset$  and  $A \subseteq M$  then
15:    $\text{Out} = M$ 
16: else
17:    $\text{Out} = \perp$ 
18: return: Out

```

Proof. Let w_1 and w_2 be defined as in Algorithm 4. Let $M = \text{Oracle}(w_2)$. Let $\mathcal{M}_{A,B} = \{M \in \mathcal{M} \mid A \subseteq M, B \cap M = \emptyset\}$ be the subset of \mathcal{M} which satisfies the constraints. If $\mathcal{M}_{A,B} = \emptyset$, then it is clear M cannot satisfy both of the constraints $A \subseteq M$ and $B \cap M = \emptyset$. Therefore Algorithm 4 returns \perp in this case.

In the rest of the proof, we assume that $\mathcal{M}_{A,B} \neq \emptyset$. Since $\mathcal{M}_{A,B}$ is non-empty, we can fix an arbitrary $M_0 \in \mathcal{M}_{A,B}$, which will be used later in the proof. We will also frequently use the fact that for all $v \in \mathbb{R}^n$ and all $S \subseteq [n]$, we have

$$- \|v\|_1 \leq v(S) \leq \|v\|_1. \quad (125)$$

First we claim that $B \cap M = \emptyset$. Suppose that $B \cap M \neq \emptyset$. Then there exists $i \in B \cap M$ and we fix such an i . Then we have

$$\begin{aligned} w_2(M) &= w_2(M \setminus \{i\}) + w_2(i) \\ &\leq w_2(M \setminus B) + w_2(i) \end{aligned} \quad (126)$$

$$= w_1(M \setminus B) + w_2(i) \quad (127)$$

$$\leq L_2 - 3L_2 = -2L_2, \quad (128)$$

where Eq. (126) follows from the fact that $w_2(j) = -L_2 \leq 0$ for all $j \in B \setminus \{i\}$; Eq. (127) holds since w_1 and w_2 coincide on all entries of $M \setminus B$; and Eq. (128) follows from the definition $L_2 = \|w_1\|_1$ and Eq. (125).

On the other hand, observing that $B \cap M_0 = \emptyset$, we can bound $w_2(M_0)$ as follows

$$w_2(M_0) = w_1(M_0) \geq -L_2.$$

Therefore we see that $w_2(M_0) > w_2(M)$. However, this contradicts to the definition of M since $M \in \arg \max_{M' \in \mathcal{M}} w_2(M')$. Therefore our claim $B \cap M = \emptyset$ is true. By this claim and since w_2 and w_1 coincide on entries of $[n] \setminus B$, we have

$$w_2(M) = w_1(M). \quad (129)$$

Next we claim that $A \subseteq M$. Suppose that $A \not\subseteq M$. Then we have

$$\begin{aligned} w_2(M) &= w_1(M) = w_1(M \cap A) + w_1(M \setminus A) \\ &= 3|M \cap A|L_1 + w(M \setminus A) \end{aligned} \quad (130)$$

$$\leq (3|A| - 3)L_1 + L_1 \quad (131)$$

$$= (3|A| - 2)L_1, \quad (132)$$

where Eq. (130) follows from the definition of w_1 ; and Eq. (131) follows from the assumption that $A \not\subseteq M$ and therefore $|M \cap A| \leq |A| - 1$.

On the other hand, using the fact that $A \subseteq M_0$ (since $M_0 \in \mathcal{M}_{A,B}$), we have

$$w_2(M_0) = w_1(M_0) = w_1(A) + w_1(M_0 \setminus A) \quad (133)$$

$$= 3|A|L_1 + w(M_0 \setminus A) \quad (134)$$

$$\geq 3|A|L_1 - L_1 \quad (135)$$

$$= (3|A| - 1)L_1, \quad (136)$$

where Eq. (133) follows from the fact that $M_0 \cap B = \emptyset$ and $A \subseteq M_0$; Eq. (134) follows from the definition of w_1 , which ensures that w_1 and w coincide on $M_0 \setminus A$; and Eq. (135) follows from Eq. (125).

Therefore, by combining Eq. (132) and Eq. (136), we see that $w_2(M_0) > w_2(M)$. Again this contradicts to the definition of M , which proves the claim.

Now we see that $M \in \mathcal{M}_{A,B}$. Therefore, we remain to verify that $w(M) = \max_{M' \in \mathcal{M}_{A,B}} w(M')$. Suppose that there exists $M_1 \in \mathcal{M}_{A,B}$ such that $w(M_1) > w(M)$. Notice that $B \cap M_1 = \emptyset$ and $A \subseteq M_1$, we have

$$w_2(M_1) = w_1(M_1) = w_1(M_1 \setminus A) + w_1(B) = w(M_1 \setminus A) + 3|A|L_1 = w(M_1) + 3|A|L_1 - w(A).$$

Similarly, one can show that $w_2(M) = w(M) + 3|A|L_1 - w(A)$. By combining with the assumption that $w(M_1) > w(M)$ we see that $w_2(M_1) > w_2(M)$, which contradicts to the definition of M . Hence we have verified that $w(M) = \max_{M' \in \mathcal{M}_{A,B}} w(M')$. \square

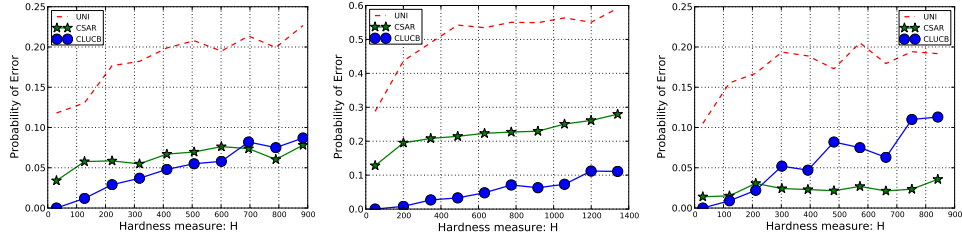
H Preliminary Experiments: Identifying the Minimum Spanning Tree

In this section, we present some preliminary experimental results of our algorithms CLUCB and CSAR. We conduct experiments on a real-world dataset with decision classes corresponding to spanning trees. We compare our algorithms with the uniform allocation benchmark UNI discussed in Appendix E. The experiment results show that the proposed algorithms are considerably more sample efficient than the UNI algorithm, which agrees with our theoretical analysis.

Setup. Our task is to identify the optimal routing tree from a networking system which has the lowest expected latency in an exploration procedure, where one can obtain noisy measurements of latencies between different nodes. We model this problem as a CPE problem where the arms correspond to edges and the decision class corresponds to the set of spanning trees (which is a special case of matroids, as we have discussed in Example 3). We use a real-world dataset called RocketFuel [29], which contains several ISP networks with routing information such as average latencies between nodes pairs. We select three medium-sized ISP networks with numbers of edges ranging from 161 to 328. For each network, we model the latency $X(e)$ of edge e as the sum of the given average latency $l(e)$ and an additive random noise $\mathcal{N}(0, 1)$. Then we model the reward of edge e as the negative latency $-X(e)$ and therefore the expected reward of e is given by $w(e) = -l(e)$. Notice that we now need to find the spanning tree that maximizes the expected reward, which is exactly an instance of CPE.

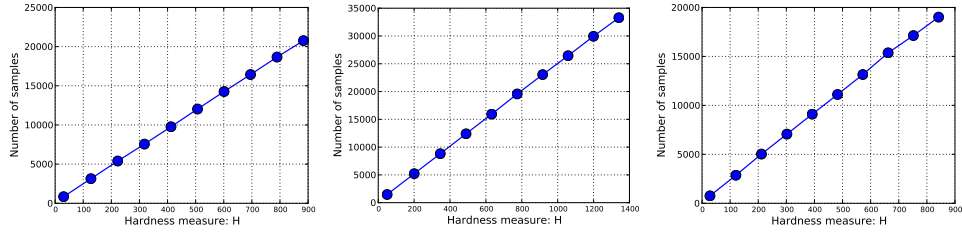
Since the ground-truth of expected reward w is known, we can compute the ground-truth of the optimal set M_* and the hardness measures \mathbf{H} . Furthermore, in order to investigate the relationship between \mathbf{H} and sample complexity empirically, we generate a number of instances with different \mathbf{H} by adjusting the expected reward of each arm $e \in M_*$ with a same additive quantity c_0 while not changing the optimality of M_* . By definition of \mathbf{H} , we see that \mathbf{H} decreases when c_0 increases.

Evaluation method. We use the following evaluation procedure to compare the sample efficiency between CLUCB, CSAR and UNI. Since CSAR and UNI are both learning algorithms in the fixed budget setting, the comparison between them is straightforward: for each given budget, we run both algorithms with this budget independently for 1000 times and compare their empirical probability of errors (the fraction of runs where a tested algorithm fails to report the ground-truth optimal set M_*). On the other hand, we use the following procedure to compare CLUCB with other fixed budget algorithms. For each instance, we run CLUCB independently for 1000 times. Suppose that the i -th run of CLUCB uses T_i samples, we also run UNI and CSAR with budget T_i . Then we compare the empirical probability of errors of the tested algorithms after the 1000 runs are completed. In this way, we see that the compared algorithms use an equal number of samples in each run, which allows us to compare their sample efficiency. Finally, we set $\delta = 0.3$ for CLUCB throughout the experiments.



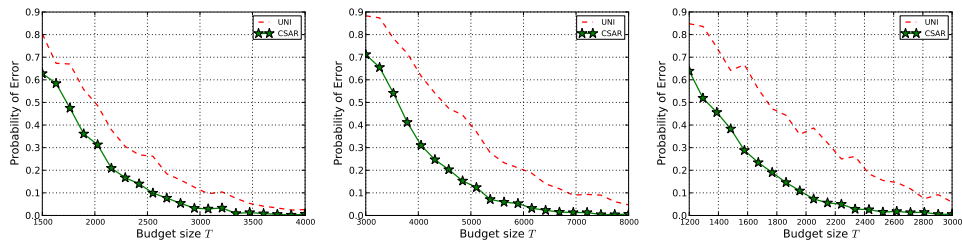
(a) Network 1755 (b) Network 3257 (c) Network 3967

Figure 2: Comparison of empirical probability of errors with respect to H .



(a) Network 1755 (b) Network 3257 (c) Network 3967

Figure 3: Empirical sample complexity of CLUCB with respect to H .



(a) Network 1755 ($H = 117.6$) (b) Network 3257 ($H = 181.5$) (c) Network 3967 ($H = 108.7$)

Figure 4: Empirical probability of error of CSAR and UNI with respect to budget size T .

Experimental results. We test all competing algorithms using the aforementioned evaluation method. The experimental results are shown in Figure 2, Figure 3 and Figure 4. From the results (Figure 2 and Figure 4), we see that both CLUCB and CSAR are consistently more sample efficient than UNI by a large margin, i.e. they incur a smaller empirical probability of error than UNI when using a same number of samples. This matches our theoretical analyses of these algorithms. We also see that the probability of error of CLUCB is always smaller than the guarantee $\delta = 0.3$ (Figure 2) and the sample complexity of CLUCB is approximately linear in \mathbf{H} (Figure 3), which agrees with our theory that the sample complexity bound for the spanning tree decision class is $\tilde{O}(\mathbf{H})$ (see Example 3).

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] C. Berge. Two theorems in graph theory. *PNAS*, 1957.
- [5] 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.
- [6] S. Bubeck, R. Munos, and G. Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412:1832–1852, 2010.
- [7] 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.
- [8] S. Bubeck, T. Wang, and N. Viswanathan. Multiple identifications in multi-armed bandits. In *ICML*, pages 258–265, 2013.
- [9] N. Cesa-Bianchi and G. Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.
- [10] W. Chen, Y. Wang, and Y. Yuan. Combinatorial multi-armed bandit: General framework and applications. In *ICML*, pages 151–159, 2013.
- [11] E. Even-Dar, S. Mannor, and Y. Mansour. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *JMLR*, 2006.
- [12] V. Gabillon, M. Ghavamzadeh, A. Lazaric, and S. Bubeck. Multi-bandit best arm identification. In *NIPS*. 2011.
- [13] V. Gabillon, M. Ghavamzadeh, and A. Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In *NIPS*, 2012.
- [14] A. Gopalan, S. Mannor, and Y. Mansour. Thompson sampling for complex online problems. In *ICML*, pages 100–108, 2014.
- [15] K. Jamieson and R. Nowak. Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. In *Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2014.
- [16] K. Jamieson, M. Malloy, R. Nowak, and S. Bubeck. lil'ucb : An optimal exploration algorithm for multi-armed bandits. *COLT*, 2014.
- [17] S. Kale, L. Reyzin, and R. E. Schapire. Non-stochastic bandit slate problems. In *NIPS*, 2010.
- [18] S. Kalyanakrishnan and P. Stone. Efficient selection of multiple bandit arms: Theory and practice. In *ICML*, pages 511–518, 2010.
- [19] S. Kalyanakrishnan, A. Tewari, P. Auer, and P. Stone. Pac subset selection in stochastic multi-armed bandits. In *ICML*, pages 655–662, 2012.
- [20] E. Kaufmann and S. Kalyanakrishnan. Information complexity in bandit subset selection. In *COLT*, pages 228–251, 2013.
- [21] B. Kveton, Z. Wen, A. Ashkan, H. Eydgahi, and B. Eriksson. Matroid bandits: Fast combinatorial optimization with learning. In *UAI*, 2014.

2214 [22] T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in*
2215 *applied mathematics*, 6(1):4–22, 1985.

2216 [23] T. Lin, B. Abrahao, R. Kleinberg, J. Lui, and W. Chen. Combinatorial partial monitoring game
2217 with linear feedback and its application. In *ICML*, 2014.

2218 [24] S. Mannor and J. N. Tsitsiklis. The sample complexity of exploration in the multi-armed bandit
2219 problem. *The Journal of Machine Learning Research*, 5:623–648, 2004.

2220 [25] G. Neu, A. György, and C. Szepesvári. The online loop-free stochastic shortest-path problem.
2221 In *COLT*, pages 231–243, 2010.

2222 [26] J. G. Oxley. *Matroid theory*. Oxford university press, 2006.

2223 [27] D. Pollard. Asymptopia. *Manuscript, Yale University, Dept. of Statist., New Haven, Connecti-*
2224 *cut*, 2000.

2225 [28] S. M. Ross. *Stochastic processes*, volume 2. John Wiley & Sons New York, 1996.

2226 [29] N. Spring, R. Mahajan, and D. Wetherall. Measuring isp topologies with rocketfuel. *ACM*
2227 *SIGCOMM Computer Communication Review*, 32(4):133–145, 2002.

2228 [30] Y. Zhou, X. Chen, and J. Li. Optimal pac multiple arm identification with applications to
2229 crowdsourcing. In *ICML*, 2014.

2230

2231

2232

2233

2234

2235

2236

2237

2238

2239

2240

2241

2242

2243

2244

2245

2246

2247

2248

2249

2250

2251

2252

2253

2254

2255

2256

2257

2258

2259

2260

2261

2262

2263

2264

2265

2266

2267