

Pure Exploration of Combinatorial Bandits

Shouyuan Chen

May 5, 2014

1 Introduction

1.1 Related Work

Notations.

2 Pure Exploration of Combinatorial Bandits

ExpCMAB: problem formulation. Suppose that there are n arms and the arms are numbered $1, 2, \dots, n$. Each arm $e \in [n]$ is associated with a reward distribution φ_e and define $w(e) = \mathbb{E}_{X \sim \varphi_e}[X]$ be the expected reward. Let $\mathbf{w} = (w(1), \dots, w(n))^T$ denote the vector of expected rewards.

Let $\mathcal{M} \subseteq 2^{[n]}$ be the family of all feasible solutions to a combinatorial problem. A learner wants to find the optimal solution of \mathcal{M} which maximizes the expected reward $M_* = \arg \max_{M \in \mathcal{M}} w(M)$ by playing the following game. At the beginning of the game, the reward distributions $\{\varphi_e\}_{e \in [n]}$ are unknown to the learner. Then, the game is played for multiple rounds; on 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 solution $\text{Out} \in \mathcal{M}$.

We consider two different stopping conditions of the game, which are known as *fixed confidence* setting and *fixed budget* setting. In the fixed confidence setting, the learner can stop the game at any point and her goal is to achieve a fixed confidence about the optimality of the returned set while uses a small number of pulls. Specifically, given a confidence parameter δ , the learner need to guarantee that $\Pr[\text{Out} = M_*] \geq 1 - \delta$. The performance is evaluated by the number of pulls used by the learner. In the fixed budget setting, the game stops after a fixed number rounds. The learner tries to minimize the probability of error $\Pr[\text{Out} \neq M_*]$ within these rounds. In this case, the learner's performance is measured by the probability of error.

Examples of combinatorial problems. The formulation of the ExpCMAB problem covers many online learning tasks. We consider the following problems as examples.

- Multi.
- Matroid.
- Match.
- Path.

We assume that all reward distributions have R -sub-Gaussian tails. Formally, for all $t \in \mathbb{R}$, we assume that $\mathbb{E}_{X \sim \varphi_e}[\exp(tX - tw(e))] \leq \exp(R^2 t^2 / 2)$. It is well known that all distributions that are supported on $[0, R]$ satisfy this property \square .

3 Algorithm and Main Result

In this section, we present **CGapExp**, a learning algorithm for the **ExpCMAB** problem in the fixed confidence setting, and analyze its sample complexity. The **CGapExp** algorithm can be extended to the fixed budget and PAC learning settings. We will discuss these extensions in Section 5.

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

Algorithm. The **CGapExp** algorithm maintains empirical mean $\bar{w}_t(e)$ and confidence radius $\text{rad}_t(e)$ for each arm $e \in [n]$ and each round t . The construction of confidence radius ensures that $|w(e) - \bar{w}_t(e)| \leq \text{rad}_t(e)$ holds with high probability for each arm $e \in [n]$ and each round $t > 0$. **CGapExp** begins with an initialization phase in which each arm is pulled once. Then, at round $t \geq n$, **CGapExp** uses the following procedure to choose an arm to play. First, **CGapExp** calls the oracle which computes the solution $M_t = \text{Oracle}(\bar{\mathbf{w}}_t)$. The solution M_t is the “best” solution with respect to the empirical means $\bar{\mathbf{w}}_t$. Then, **CGapExp** explores possible refinements of M_t . In particular, **CGapExp** uses the confidence radius to compute an adjusted expectation vector $\tilde{\mathbf{w}}_t$ in the following way: for each arm $e \in M_t$, $\tilde{w}_t(e)$ equals to the lower confidence bound $\tilde{w}_t(e) = \bar{w}_t(e) - \text{rad}_t(e)$; and for each arm $e \notin M_t$, $\tilde{w}_t(e)$ equals to the upper confidence bound $\tilde{w}_t(e) = \bar{w}_t(e) + \text{rad}_t(e)$. Intuitively, the adjusted expectation vector $\tilde{\mathbf{w}}_t$ penalizes arms belonging to the current solution M_t and encourages exploring arms out of M_t . **CGapExp** then calls the oracle using the adjusted expectation vector $\tilde{\mathbf{w}}_t$ as input to compute a refined solution $\tilde{M}_t = \text{Oracle}(\tilde{\mathbf{w}}_t)$. If $\tilde{w}_t(\tilde{M}_t) = \tilde{w}_t(M_t)$ then **CGapExp** stops and returns $\text{Out} = M_t$. Otherwise, **CGapExp** pulls the arm belonging to the symmetric difference $(\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)$ between M_t and \tilde{M}_t with the largest confidence radius in the end of round t . The pseudo-code of **CGapExp** is shown in Algorithm 1.

Algorithm 1 CGapExp: Combinatorial Gap Exploration

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

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

```

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

3.1 Analysis

Gap. We begin with defining a natural complexity measure of the ExpCMAB problem. For each arm $e \in [n]$, we define 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 \mathbf{H} as the sum of inverse squared gaps

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

From Eq. (1), we see that, for each arm $e \notin M_*$, Δ_e represents the gap between the optimal set M_* and the best set that includes arm e ; and, for each arm $e \in M_*$, Δ_e is the sub-optimality of the best set that does not include arm e . We notice that this definition resembles the definition of gaps for Multi proposed by ().

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

We begin with the definition of exchange set. We define an exchange set b as an ordered pair of disjoint sets $b = (b_+, b_-)$. Then, we define operator \oplus such that, for any set M 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 family of exchange sets \mathcal{B} an *exchange class* for \mathcal{M} if \mathcal{B} satisfies the following property. Let M and M' be two elements of \mathcal{M} . Then, for any $e \in (M \setminus M')$, there exists an exchange set $(b_+, b_-) \in \mathcal{B}$ which satisfies $e \in b_-$, $b_+ \subseteq M' \setminus M$, $b_- \subseteq M \setminus M'$, $(M \oplus b) \in \mathcal{M}$ and $(M' \ominus b) \in \mathcal{M}$. We define the *width* of exchange class \mathcal{B} as follows

$$\text{width}(\mathcal{B}) = \max_{(b_+, b_-) \in \mathcal{B}} |b_+| + |b_-|. \quad (2)$$

Intuitively, for any feasible sets M and M' , there exists an exchange set $(b_+, b_-) \in \mathcal{B}$ belonging to the exchange class \mathcal{B} which can be seen as an “operation” that transforms M one step towards M' : this operation generates a new feasible set $M \oplus b$ by removing elements (including e) from M and adding elements which belongs to M' . One can chain these operations together such that, for any $M \neq M'$, there exists a sequence of exchange sets b_1, \dots, b_k of \mathcal{B} such that $M' = M \oplus b_1 \oplus \dots \oplus b_k$.

We notice that an exchange class for \mathcal{M} can be “redundant”. It may contains some unnecessary exchange set b , such that $M \oplus b \notin \mathcal{M}$ for any $M \in \mathcal{M}$. These redundant exchange sets do not affect our analysis. But allowing them would simplify the construction and description of exchange classes for certain combinatorial problems.

It is easy to see that, for a fixed combinatorial problem \mathcal{M} , the exchange classes for \mathcal{M} are not unique. Let $\text{Exchange}(\mathcal{M})$ denote the collection of all exchange classes for \mathcal{M} . We are interested with the exchange class with small width. Formally, we define the width of a combinatorial problem \mathcal{M} as the width of the thinnest exchange class

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

Next, we construct the exchange classes for our running examples. Our constructions are summarized in Lemma 1.

Fact 1. *There exist exchange classes $\mathcal{B}_{\text{Multi}}$, $\mathcal{B}_{\text{Matroid}}$, $\mathcal{B}_{\text{Match}}$ and $\mathcal{B}_{\text{Path}}$ for $\mathcal{M}_{\text{Multi}}$, $\mathcal{M}_{\text{Matroid}}$, $\mathcal{M}_{\text{Match}}$ and $\mathcal{M}_{\text{Path}}$, respectively. These exchange classes can be constructed as follows*

1. $\mathcal{B}_{\text{Multi}} = \{(\{i\}, \{j\}) \mid \forall i \in [n], j \in [n]\}$.
2. $\mathcal{B}_{\text{Matroid}} = \{(\{i\}, \{j\}) \mid \forall i \in [n], j \in [n]\}$.
3. $\mathcal{B}_{\text{Match}} = \{(C_+, C_-) \mid C_+ \cup C_- \text{ is a cycle of } G\}$.

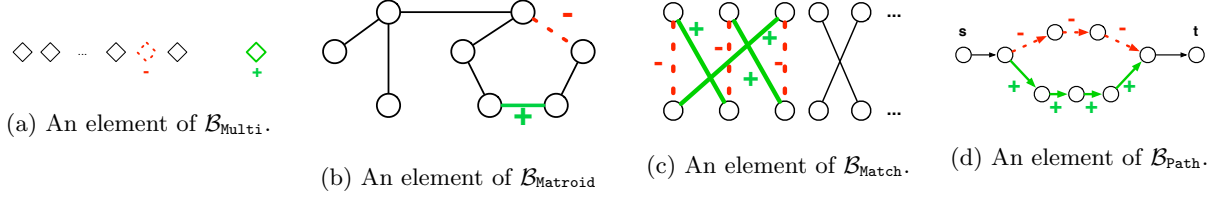


Figure 1: Examples of exchange sets belonging to the exchange classes $\mathcal{B}_{\text{Multi}}$, $\mathcal{B}_{\text{Matroid}}$, $\mathcal{B}_{\text{Match}}$ and $\mathcal{B}_{\text{Path}}$: green-solid elements constitute the set b_+ , red-dotted elements constitute the set b_- and the example exchange set is $b = (b_+, b_-)$. (In Figure 1b, we consider spanning tree as a specific instance for the **Matroid** problem.)

4. $\mathcal{B}_{\text{Path}} = \{(P_1, P_2) \mid P_1, P_2 \text{ are two disjoint paths of } G \text{ with same endpoints}\}.$

In addition, we have $\text{width}(\mathcal{B}_{\text{Multi}}) = 2$, $\text{width}(\mathcal{B}_{\text{Matroid}}) = 2$, $\text{width}(\mathcal{B}_{\text{Match}}) = |V|$ and $\text{width}(\mathcal{B}_{\text{Path}}) = |V|$. This means that $\text{width}(\text{Multi}) \leq 2$, $\text{width}(\text{Matroid}) \leq 2$, $\text{width}(\text{Match}) \leq |V|$ and $\text{width}(\text{Path}) \leq |V|$.

The construction for **Multi** problem is straightforward. For **Matroid** problem, we leverage the basis exchange property of matroids (see Lemma 13 in the appendix). And for **Match** and **Path** problems, we appeal to graph-theoretical properties of matchings and paths. We illustrate elements of these exchanges classes in Figure 1. A detailed proof of Fact 1 is deferred to the supplementary material.

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

Theorem 1. Given any $\delta \in (0, 1)$, any $\mathcal{M} \subseteq 2^{[n]}$ and any $\mathbf{w} \in \mathbb{R}^n$. Assume that the reward distribution φ_e for each arm $e \in [n]$ is R -sub-Gaussian with mean $w(e)$.

Set $\text{rad}_t(e) = R \sqrt{\frac{2 \log\left(\frac{4nt^2}{\delta}\right)}{T_e(t)}}$ for all $t > 0$ and $e \in [n]$. Then, with probability at least $1 - \delta$, the **CGapExp** algorithm (Algorithm 1) returns the optimal set $\text{Out} = M_*$ and

$$T \leq O\left(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \log\left(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \cdot n/\delta\right)\right), \quad (3)$$

where T denotes the number of samples used by Algorithm 1.

Remarks. For the **Multi** problem, we see that Fact 1 shows that $\text{width}(\text{Multi}) = O(1)$. Therefore, the sample complexity bound of **CGapExp** is $O(\mathbf{H} \log(n\mathbf{H}/\delta))$ for the **Multi** problem. This matches the previous known problem-dependent bounds for the **Multi** problem due to XXX [] within logarithmic factors. For the **Matroid** problem, we know that $\text{width}(\text{Matroid}) = O(1)$ and hence the sample complexity is also $O(\mathbf{H} \log(n\mathbf{H}/\delta))$. For **Match** and **Path** problem, we see that the sample complexity is bounded by $\tilde{O}(|V|^2 \mathbf{H})$.

4 Lower Bound

In this section, we prove a problem-dependent lower bound on the sample complexity of the **ExpCMAB** 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 algorithm \mathbb{A} .

We show that, for any combinatorial problem \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 $\mathcal{M} \subseteq 2^{[n]}$ 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 . Then, for any $\delta \in (0, e^{-16}/4)$ and any δ -correct algorithm \mathbb{A} , we have

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

where T denote the number of total samples used by algorithm \mathbb{A} and Δ_e is defined in Eq. (1).

We have already shown that, for **Multi** and **Matroid**, the sample complexity of **CGapExp** is $O(\mathbf{H} \log(n\mathbf{H}/\delta))$. Hence, for these two problems, the **CGapExp** algorithm achieves the optimal sample complexity within logarithmic factors. In addition, our result confirms the conjecture of XXX [] that the lower bound of sample complexity of **Multi** problem is $\Omega(\mathbf{H} \log(1/\delta))$.

On the other hand, for general combinatorial problems, we see that the gap between the upper bound Eq. (3) and the lower bound Eq. (4) is $\text{width}(\mathcal{M})^2$. Notice this gap only depends on the underlying combinatorial structure of \mathcal{M} and is independent of expected rewards \mathbf{w} . This suggests that the sample complexity of **CGapExp** has an optimal dependency on the gaps $\{\Delta_e\}_{e \in [n]}$ for general combinatorial problems.

We conjecture that the dependency on $\text{width}(\mathcal{M})$ of the sample complexity might be intrinsic. In the supplementary material, we provide evidence showing that the sample complexity of any δ -correct algorithm should be related to size of exchange sets.

5 Extensions

5.1 Fixed Budget Setting

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

Theorem 3. Use the same notations as in Theorem 1. Given $T > n$ and parameter $\alpha > 0$, set the confidence radius $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_e(t)}}$ for all arms $e \in [n]$ and all $t > 0$. Run **CGapExp** algorithm for at most T rounds. Then, for $0 \leq \alpha \leq \frac{1}{9}(T - n)(R^2 \text{width}(\mathcal{M})^2 \mathbf{H})^{-1}$, we have

$$\Pr[\text{Out} \neq M_*] \leq 2Tn \exp(-2\alpha). \quad (5)$$

The right-hand side of Eq. (5) equals to $O\left(Tn \exp\left(\frac{-T}{\text{width}(\mathcal{M})^2 \mathbf{H}}\right)\right)$ when parameter $\alpha = O(T\mathbf{H}^{-1} \text{width}(\mathcal{M})^{-2})$. For **Multi** problem, we see that this matches the guarantees of previous fixed budget algorithm [] up to logarithmic factors.

5.2 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 **Out** satisfies $\Pr[w(M_*) - w(\text{Out}) > \epsilon] \leq \delta$.

We show that a simple modification on the **CGapExp** 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

$\tilde{w}_t(\tilde{M}_t) \leq \tilde{w}_t(M_t)$ to $w(\tilde{M}_t) - w(M_t) \leq \epsilon$ on line 15 of Algorithm 1. We let **CGapExpPAC** denote the modified algorithm. In the following theorem, we show that **CGapExpPAC** is indeed an (ϵ, δ) -PAC algorithm and has sample complexity similar to **CGapExp**.

Theorem 4. *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 **Out** of **CGapExpPAC** satisfies $w(M_*) - w(\mathbf{Out}) \leq \epsilon$. In addition, the number of samples T used by the algorithm satisfies*

$$T \leq O \left(R^2 \sum_{e \in [n]} \min \left\{ \frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2} \right\} \log \left(\frac{R^2 n}{\delta} \sum_{e \in [n]} \min \left\{ \frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2} \right\} \right) \right), \quad (6)$$

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

We see that the sample complexity of **CGapExpPAC** decreases when ϵ increases. And if $\epsilon = 0$, the sample complexity Eq. (6) of **CGapExpPAC** equals to that of **CGapExp**. (difference??)

6 Proof of Main Result

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

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

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

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

6.1 Preparatory Lemmas

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

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

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

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

□

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

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

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

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 $(M \oplus c) \in \mathcal{M}$.

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

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

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

□

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

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

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 $\text{rad}_t(e) \geq 0$ and $\chi_{(M \setminus M') \cup (M' \setminus M)} \in \{0, 1\}^n$. □

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

Proof. We have

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

□

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

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

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

$$\begin{aligned} \langle \tilde{\mathbf{w}}_t, \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \rangle &= \langle \bar{\mathbf{w}}_t + \mathbf{rad}_t \circ (\mathbf{1}_n - 2\boldsymbol{\chi}_{M_t}), \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \rangle \\ &= \langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \rangle + \langle \mathbf{rad}_t, (\mathbf{1}_n - 2\boldsymbol{\chi}_{M_t}) \circ (\boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t}) \rangle \end{aligned} \quad (7)$$

$$\begin{aligned} &= \langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \rangle + \langle \mathbf{rad}_t, \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} - 2\boldsymbol{\chi}_{M_t} \circ \boldsymbol{\chi}_{M'} + 2\boldsymbol{\chi}_{M_t}^2 \rangle \\ &= \langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \rangle + \langle \mathbf{rad}_t, \boldsymbol{\chi}_{M'}^2 - \boldsymbol{\chi}_{M_t}^2 - 2\boldsymbol{\chi}_{M_t} \circ \boldsymbol{\chi}_{M'} + 2\boldsymbol{\chi}_{M_t}^2 \rangle \end{aligned} \quad (8)$$

$$\begin{aligned} &= \langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \rangle + \langle \mathbf{rad}_t, (\boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t})^2 \rangle \\ &= \langle \bar{\mathbf{w}}_t, \boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \rangle + \langle \mathbf{rad}_t, |\boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t}| \rangle, \end{aligned} \quad (9)$$

where Eq. (7) follows from Lemma 4; Eq. (8) holds since $\boldsymbol{\chi}_{M'} \in \{0, 1\}^n$ and $\boldsymbol{\chi}_{M_t} \in \{0, 1\}^n$ and therefore $\boldsymbol{\chi}_{M'} = \boldsymbol{\chi}_{M'}^2$ and $\boldsymbol{\chi}_{M_t} = \boldsymbol{\chi}_{M_t}^2$; and Eq. (9) follows since $\boldsymbol{\chi}_{M'} - \boldsymbol{\chi}_{M_t} \in \{-1, 0, 1\}^n$. \square

6.2 Confidence Intervals

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

$$\xi_t = \left\{ \forall i \in [n], \quad |w(i) - \bar{w}_t(i)| \leq \text{rad}_t(i) \right\}. \quad (10)$$

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. (10) to inner products as follows.

Lemma 6. Given any $t > 0$, assume that event ξ_t as defined in Eq. (10) 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 occurs. Then, we have

$$\begin{aligned} |\langle \mathbf{w}, \mathbf{a} \rangle - \langle \bar{\mathbf{w}}_t, \mathbf{a} \rangle| &= |\langle \mathbf{w} - \bar{\mathbf{w}}_t, \mathbf{a} \rangle| \\ &= \left| \sum_{i=1}^n (w(i) - \bar{w}_t(i))a(i) \right| \\ &\leq \sum_{i=1}^n |w(i) - \bar{w}_t(i)| |a(i)| \\ &\leq \sum_{i=1}^n \text{rad}_t(i) \cdot |a(i)| \\ &= \langle \mathbf{rad}_t, |\mathbf{a}| \rangle, \end{aligned} \quad (11)$$

where Eq. (11) follows the definition of event ξ_t in Eq. (10) and the assumption that it occurs. \square

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

Lemma 7. Suppose that the reward distribution 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

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

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

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

Proof. For any $t > 0$ and $e \in [n]$, notice φ_e is a R -sub-Gaussian distribution with mean $w(e)$ and $w_t(e)$ is the empirical mean of φ_e for $T_e(t)$ samples. Using Hoeffding's inequality (see Lemma 14 in Section 9), we obtain

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

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

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

□

6.3 Main Lemmas

Lemma 8. Given any $t > 0$, assume that event ξ_t (defined in Eq. (10)) occurs. Then, if Algorithm 1 terminates at round t , we have $M_t = M_*$.

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

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

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

$$\begin{aligned} \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) &\geq \tilde{w}_t(M'_t) - \tilde{w}_t(M_t) \\ &= \langle \tilde{\mathbf{w}}_t, \chi_{M'_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{M'} - \chi_{M_t}| \rangle \end{aligned} \tag{12}$$

$$\geq \langle \mathbf{w}, \chi_{M'_t} - \chi_{M_t} \rangle \tag{13}$$

$$= w(M'_t) - w(M_t) > 0, \tag{14}$$

where Eq. (12) follows from Lemma 5; and Eq. (13) follows the assumption that event ξ_t occurs and Lemma 6; Therefore Eq. (14) shows that $\tilde{w}_t(\tilde{M}_t) > \tilde{w}_t(M_t)$. However, this contradicts to the stopping condition of **CGapExp**: $\tilde{w}_t(\tilde{M}_t) \leq \tilde{w}_t(M_t)$ and the assumption that the algorithm terminates on round t . □

Lemma 9. *Given any $t > 0$ and suppose that event ξ_t (defined in Eq. (10)) occurs. For any $e \in [n]$, if $\text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{M})}$, then, arm e will not be pulled on round t , i.e. $p_t \neq e$.*

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

Now, we decompose our proof into two cases.

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

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

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

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

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

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

$$\begin{aligned} \|\text{rad}_t \circ (\mathbf{d} - \chi_c)\|_\infty &= \left\| \text{rad}_t \circ (\chi_{M'_t} - \chi_{M_t}) \right\|_\infty \\ &= \max_{i \in (M_t \setminus M'_t) \cup (M'_t \setminus M_t)} \text{rad}_t(i) \\ &\leq \max_{i \in (M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)} \text{rad}_t(i) \\ &= \text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}. \end{aligned} \quad (15)$$

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

$$\begin{aligned} \|\text{rad}_t \circ \chi_c\|_\infty &= \max_{i \in c_+ \cup c_-} \text{rad}_t(i) \\ &\leq \max_{i \in (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(i) \\ &= \text{rad}_t(e) < \frac{\Delta_e}{3 \text{width}(\mathcal{B})}. \end{aligned} \quad (16)$$

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

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

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

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

$$\langle \mathbf{rad}_t, |d_2| \rangle - \langle \mathbf{rad}_t, |d| \rangle = \langle \mathbf{rad}_t, |d_2| - |d| \rangle = \langle \mathbf{rad}_t, d_2^2 - d^2 \rangle \quad (17)$$

$$\begin{aligned} &= \langle \mathbf{rad}_t, (d - \chi_c + \chi_b)^2 - d^2 \rangle \\ &= \langle \mathbf{rad}_t, \chi_b^2 + \chi_c^2 - 2\chi_b \circ \chi_c - 2d \circ \chi_c + 2d \circ \chi_b \rangle \\ &= \langle \mathbf{rad}_t, \chi_b^2 - \chi_c^2 + 2\chi_b \circ (d - \chi_c) \rangle \end{aligned} \quad (18)$$

$$\begin{aligned} &= \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - 2 \langle \mathbf{rad}_t, \chi_b \circ (d - \chi_c) \rangle \\ &= \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - 2 \langle \mathbf{rad}_t \circ (d - \chi_c), \chi_b \rangle \end{aligned} \quad (19)$$

$$\geq \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - 2 \|\mathbf{rad}_t \circ (d - \chi_c)\|_\infty \|\chi_b\|_1 \quad (20)$$

$$> \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3 \text{width}(\mathcal{B})} \|\chi_b\|_1 \quad (21)$$

$$\geq \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3}, \quad (22)$$

where Eq. (17) holds since $d \in \{-1, 0, 1\}^n$ and $d_2 \in \{-1, 0, 1\}^n$; Eq. (18) follows from the claim that $d \circ \chi_c = |\chi_c| = \chi_c^2$; Eq. (19) and Eq. (20) follow from Lemma 4 and Hölder's inequality; Eq. (21) follows from Eq. (15); and Eq. (22) holds since $b \in \mathcal{B}$ and $\|\chi_b\|_1 = |b_+| + |b_-| \leq \text{width}(\mathcal{B})$.

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

$$\begin{aligned} \langle \bar{w}_t, d \rangle + \langle \mathbf{rad}_t, |d| \rangle &= \langle \bar{w}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \\ &= \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \\ &\geq \tilde{w}_t(M'_t) - \tilde{w}_t(M_t) \\ &= \langle \bar{w}_t, \chi_{M'_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{M'_t} - \chi_{M_t}| \rangle \\ &= \langle \bar{w}_t, d_2 \rangle + \langle \mathbf{rad}_t, |d_2| \rangle \\ &= \langle \bar{w}_t, d \rangle - \langle \bar{w}_t, \chi_c \rangle + \langle \bar{w}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |d_2| \rangle, \end{aligned}$$

where the last equality follows from the fact that $d_2 = d - \chi_c + \chi_b$. Rearranging the above inequality, we obtain

$$\begin{aligned} \langle \bar{w}_t, \chi_c \rangle &\geq \langle \bar{w}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |d_2| \rangle - \langle \mathbf{rad}_t, |d| \rangle \\ &\geq \langle \bar{w}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\chi_b| \rangle - \langle \mathbf{rad}_t, |\chi_c| \rangle - \frac{2\Delta_e}{3} \end{aligned} \quad (23)$$

$$> \langle w, \chi_b \rangle - \langle \mathbf{rad}_t, \chi_c \rangle - \frac{2\Delta_e}{3} \quad (24)$$

$$> \langle w, \chi_b \rangle - \frac{\Delta_e}{3} - \frac{2\Delta}{3} \quad (25)$$

$$= \langle w, \chi_b \rangle - \Delta_e \geq 0, \quad (26)$$

where Eq. (23) uses Eq. (22); Eq. (24) follows from the assumption that event ξ_t occurs and Lemma 6; and Eq. (24) holds since Eq. (16).

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

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

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

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

First, we claim that $\tilde{M}_t \neq M_*$. Suppose that $e \in M_*$ and $e \in c_-$. Then, we see that $e \notin \tilde{M}_t$, which implies

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

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

Next, we define vectors $\mathbf{d} = \chi_{\tilde{M}_t} - \chi_{M_t}$ and $\mathbf{d}_1 = \chi_{\tilde{M}_t \oplus b} - \chi_{M_t}$. Notice that Lemma 2 gives that $\mathbf{d}_1 = \mathbf{d} + \mathbf{b}$. Then, we apply Lemma 3 by setting $M = M_t$ and $M' = \tilde{M}_t$. This shows that

$$\|\mathbf{rad}_t \circ \mathbf{d}\|_\infty \leq \max_{i: (\tilde{M}_t \setminus M_t) \cup (M_t \setminus \tilde{M}_t)} \text{rad}_t(i) = \text{rad}_t(e) < \frac{\Delta_e}{3}. \quad (27)$$

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

$$\begin{aligned} \langle \bar{\mathbf{w}}_t, \mathbf{d}_1 \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| \rangle - \langle \bar{\mathbf{w}}_t, \mathbf{d} \rangle - \langle \mathbf{rad}_t, |\mathbf{d}| \rangle &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, |\mathbf{d}_1| - |\mathbf{d}| \rangle \\ &= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, \mathbf{d}_1^2 - \mathbf{d}^2 \rangle \end{aligned} \quad (28)$$

$$= \langle \bar{\mathbf{w}}_t, \chi_b \rangle + \langle \mathbf{rad}_t, 2\mathbf{d} \circ \chi_b + \chi_b^2 \rangle \quad (29)$$

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

$$\geq \langle \mathbf{w}, \chi_b \rangle - 2 \|\mathbf{rad}_t \circ \mathbf{d}\|_\infty \|\chi_b\|_1 \quad (31)$$

$$> \langle \mathbf{w}, \chi_b \rangle - \frac{2\Delta_e}{3} \quad (32)$$

$$\geq 0, \quad (33)$$

where Eq. (28) follows from the fact that $\mathbf{d}_1 \in \{-1, 0, 1\}^n$ and $\mathbf{d} \in \{-1, 0, 1\}^n$; Eq. (29) holds since $\mathbf{d}_1 = \mathbf{d} + \chi_b$; Eq. (30) follows from the assumption that ξ_t occurs and Lemma 6; Eq. (31) follows from Lemma 4 and Hölder's inequality; and Eq. (32) is due to Eq. (27).

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

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

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

6.4 Proof of Theorem 1

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

Theorem 1. *Given any $\delta \in (0, 1)$, any $\mathcal{M} \subseteq 2^{[n]}$ and any $\mathbf{w} \in \mathbb{R}^n$. Assume that the reward distribution φ_e for each arm $e \in [n]$ is R -sub-Gaussian with mean $w(e)$.*

*Set $\text{rad}_t(e) = R \sqrt{\frac{2 \log\left(\frac{4nt^2}{\delta}\right)}{T_e(t)}}$ for all $t > 0$ and $e \in [n]$. Then, with probability at least $1 - \delta$, the **CGapExp** algorithm (Algorithm 1) returns the optimal set $\text{Out} = M_*$ and*

$$T \leq O\left(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \log\left(R^2 \text{width}(\mathcal{M})^2 \mathbf{H} \cdot n / \delta\right)\right), \quad (3)$$

where T denotes the number of samples used by Algorithm 1.

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

By Lemma 8 and the assumption on 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 $T_e(t_e) = T_e - 1$. By Lemma 9, we see that $\text{rad}_{t_e}(e) \geq \frac{\Delta_e}{3 \text{width}(\mathcal{M})}$. Plugging in the construction radius rad , we have

$$\frac{\Delta_e}{3 \text{width}(\mathcal{M})} \leq R \sqrt{\frac{2 \log(4nt_e^2/\delta)}{T_e - 1}} \leq R \sqrt{\frac{2 \log(4nT^2/\delta)}{T_e - 1}}. \quad (34)$$

Solving Eq. (34) for T_e , we obtain

$$T_e \leq \frac{18 \text{width}(\mathcal{M})^2 R^2}{\Delta_e^2} \log(4nT^2/\delta) + 1. \quad (35)$$

Notice that $T = \sum_{i \in [n]} T_i$. Hence the theorem follows by summing up Eq. (35) for all $e \in [n]$ and solving for T . \square

7 Proof of Lower Bound

Theorem 2. 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 . Then, for any $\delta \in (0, e^{-16}/4)$ and any δ -correct algorithm \mathbb{A} , we have

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

where T denote the number of total samples used by algorithm \mathbb{A} and Δ_e is defined in Eq. (1).

Proof. 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

$$\mathbb{E}[T_e] \geq \frac{1}{16\Delta_e^2} \log(1/4\delta). \quad (36)$$

Notice that the theorem follows immediately by summing up Eq. (36) for all $e \in [n]$.

Fix an arm $e \in [n]$. We now focus on proving Eq. (36). Consider two hypothesis H_0 and H_1 . Under hypothesis H_0 , all reward distributions are same with our assumption before

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

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

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

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

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

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

By definition of 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.

Define $\theta = 4\delta$. Define

$$t_e^* = \frac{1}{16\Delta_e^2} \log\left(\frac{1}{\theta}\right). \quad (37)$$

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

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

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

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

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

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

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

where the second inequality follows from the fact that the variance of 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 $\theta < e^{-16}$.

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

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

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

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

Define the likelihood function L_l as

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

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

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

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

$$\begin{aligned} \frac{L_1(W)}{L_0(W)} &= \prod_{i=1}^{T_e} \exp\left(\frac{-(X_i - w(e) - 2\Delta_e)^2 + (X_i - w(e))^2}{2}\right) \\ &= \prod_{i=1}^{T_e} \exp(\Delta_e(2X_i - 2w(e)) - 2\Delta_e^2) \\ &= \exp(\Delta_e(2K - 2w(e)T_e) - 2\Delta_e^2T_e) \\ &= \exp(\Delta_e(2K - 2w(e)T_e)) \exp(-2\Delta_e^2T_e). \end{aligned} \quad (38)$$

Next, we bound each individual term on the right-hand side of Eq. (38). We begin with bounding the second term of Eq. (38)

$$\exp(-2\Delta_e^2T_e) \geq \exp(-8\Delta_e^2t_e^*) \quad (39)$$

$$= \exp\left(-\frac{8}{16} \log(1/\theta)\right) \quad (40)$$

$$= \theta^{1/2}, \quad (41)$$

where Eq. (39) follows from the assumption that event \mathcal{S} occurred, which implies that event \mathcal{A} occurred and therefore $T_e \leq 4t_e^*$; Eq. (40) follows from the definition of t_e^* .

Then, we bound the first term on the right-hand side of Eq. (38) as follows

$$\exp(\Delta_e(2K - 2w(e)T_e)) \geq \exp\left(-2\Delta_e\sqrt{t_e^* \log(1/\theta)}\right) \quad (42)$$

$$= \exp\left(-\frac{2}{4} \log(1/\theta)\right) \quad (43)$$

$$= \theta^{1/2}, \quad (44)$$

where Eq. (42) follows from the assumption that event \mathcal{S} occurred, which implies that event \mathcal{C} and therefore $|2K - 2w(e)T_e| \leq \sqrt{t_e^* \log(1/\theta)}$; Eq. (43) follows from the definition of t_e^* .

Combining Eq. (41) and Eq. (44), we can bound $L_1(W)/L_0(W)$ for this case as follows

$$\frac{L_1(W)}{L_0(W)} \geq \theta. \quad (45)$$

(End of Case (1).)

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

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

$$\begin{aligned}
&= \prod_{i=1}^{T_e} \exp(\Delta_e(2w(e) - 2X_i) - 2\Delta_e^2) \\
&= \exp(\Delta_e(2w(e)T_e - 2K)) \exp(-2\Delta_e^2 T_e).
\end{aligned} \tag{46}$$

Notice that the right-hand side of Eq. (46) differs from Eq. (38) only in its first term. Now, we bound the first term as follows

$$\exp(\Delta_e(2w(e)T_e - 2K)) \geq \exp(-2\Delta_e \sqrt{t_e^* \log(1/\theta)}) \tag{47}$$

$$= \exp\left(-\frac{2}{4} \log(1/\theta)\right) \tag{48}$$

$$= \theta^{1/2}, \tag{49}$$

where the inequalities hold due to reasons similar to Case (1): Eq. (47) follows from the assumption that event \mathcal{S} occurred, which implies that event \mathcal{C} and therefore $|2K - 2w(e)T_e| \leq \sqrt{t_e^* \log(1/\theta)}$; Eq. (48) follows from the definition of t_e^* .

Combining Eq. (41) and Eq. (44), we can obtain the same bound of $L_1(W)/L_0(W)$ as in Eq. (45), i.e. $L_1(W)/L_0(W) \geq \theta$.

(End of Case (2).)

At this point, we have proved that, if the event \mathcal{S} occurred, then the bound of likelihood ratio Eq. (45) holds, i.e. $\frac{L_1(W)}{L_0(W)} \geq \theta$. Hence, we have

$$\begin{aligned}
\frac{L_1(W)}{L_0(W)} &\geq \theta \\
&= 4\delta.
\end{aligned} \tag{50}$$

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

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

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

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

Now we have proved that, if $\mathbb{E}_0[T_e] \leq t_e^*$, then $\Pr_1[\mathcal{B}] > \delta$. This means that, if $\mathbb{E}_0[T_e] \leq t_e^*$, algorithm \mathbb{A} will choose M_* as the output with probability at least 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 $\mathbb{E}_0[T_e] > t_e^* = \frac{1}{16\Delta_e^2} \log(1/4\delta)$. \square

7.1 Exchange set size dependent lower bound

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

in the following theorem.

Theorem 5. 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 $\delta > 0$, $\mathbf{w} \in \mathbb{R}^n$, diff-set $b = (b_+, b_-)$ and a δ -correct algorithm \mathbb{A} . Assume that $\varphi_e(e) = \mathcal{N}(w(e), 1)$ for all $e \in [n]$.

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

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

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

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

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

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

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

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

where p_l is the probability density function under H_l .

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

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

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

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

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

Define event

$$\mathcal{B} = \{\text{Out} = M_*\}. \quad (51)$$

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

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

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

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

where the second inequality follows from the fact that all reward distributions have unit variance and hence $\mathbb{E}_0 \left[\left(\sum_{i=1}^{4t_{b_-}^*} X_i - \sum_{i=1}^{4t_{b_-}^*} p_{I_i} \right)^2 \right] = 4t_{b_-}^*$; the last inequality follows since $\theta < e^{-16}$. Last, since algorithm \mathbb{A} is a δ -correct algorithm with $\delta < 1/4$. Therefore, it is easy to see that $\Pr_0[\mathcal{B}] \geq 3/4$. And by union bound, we have

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

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

$$\begin{aligned} \frac{L_1(W)}{L_0(W)} &= \prod_{i=1}^{T_{b_-}} \exp \left(\frac{- \left(X_i - w(I_i) - \frac{2w(b)}{|b_-|} \right)^2 + (X_i - w(I_i))^2}{2} \right) \\ &= \prod_{i=1}^{T_{b_-}} \exp \left(\frac{w(b)}{|b_-|} (2X_i - 2w(I_i)) - \frac{2w(b)^2}{|b_-|^2} \right) \\ &= \exp \left(\frac{w(b)}{|b_-|} \left(\sum_{i=1}^{T_{b_-}} 2X_i - 2w(I_i) \right) - \frac{2w(b)^2}{|b_-|^2} T_{b_-} \right) \\ &= \exp \left(\frac{w(b)}{|b_-|} \left(\sum_{i=1}^{T_{b_-}} 2X_i - 2w(I_i) \right) \right) \exp \left(- \frac{2w(b)^2}{|b_-|^2} T_{b_-} \right). \end{aligned} \quad (52)$$

Then, we bound each term on the right-hand side of Eq. (52). First, we bound the second term of Eq. (52).

$$\exp \left(- \frac{2w(b)^2}{|b_-|^2} T_{b_-} \right) \geq \exp \left(- \frac{2w(b)^2}{|b_-|^2} 4t_{b_-}^* \right) \quad (53)$$

$$= \exp\left(-\frac{8}{16}\log(1/\theta)\right) \quad (54)$$

$$= \theta^{1/2}, \quad (55)$$

where Eq. (53) follows from the assumption that events \mathcal{S}_1 and \mathcal{A}_1 occurred and therefore $T_{b_-} \leq 4t_{b_-}^*$; Eq. (54) follows from the definition of $t_{b_-}^*$. Next, we bound the first term of Eq. (52) as follows

$$\exp\left(\frac{w(b)}{|b_-|} \left(\sum_{i=1}^{T_{b_-}} 2X_i - 2w(I_i)\right)\right) \geq \exp\left(-\frac{2w(b)}{|b_-|} \sqrt{t_{b_-}^* \log(1/\theta)}\right) \quad (56)$$

$$= \exp\left(-\frac{2}{4}\log(1/\theta)\right) \quad (57)$$

$$= \theta^{1/2}, \quad (58)$$

where Eq. (56) follows since event \mathcal{S}_1 and \mathcal{C}_1 occurred and therefore $|2K - 2p_e T_e| \leq \sqrt{t_e^* \log(1/\theta)}$; Eq. (57) follows from the definition of $t_{b_-}^*$.

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

$$\frac{L_1(W)}{L_0(W)} \geq \theta = 4\delta. \quad (59)$$

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

$$\begin{aligned} \Pr_1[\mathcal{B}] &\geq \Pr_1[\mathcal{S}_1] = \mathbb{E}_1[1_{\mathcal{S}_1}] \\ &= \mathbb{E}_0\left[\frac{L_1(W)}{L_0(W)} 1_{\mathcal{S}_1}\right] \\ &\geq 4\delta \mathbb{E}_0[1_{\mathcal{S}_1}] \\ &= 4\delta \Pr_0[\mathcal{S}_1] > \delta. \end{aligned} \quad (60)$$

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

$$\mathbb{E}_0[T_{b_-}] \geq t_{b_-}^* = \frac{|b_-|^2}{16w(b)^2} \log(1/4\delta). \quad (61)$$

(End of Part (1).)

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

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

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

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

Using the similar arguments, we can show that $\Pr_0[\mathcal{A}_2] \geq 3/4$ and $\Pr_0[\mathcal{C}_2] \geq 3/4$. Define event $\mathcal{S}_2 =$

$\mathcal{A}_2 \cap \mathcal{B} \cap \mathcal{C}_2$, where \mathcal{B} is defined in Eq. (51). By union bound, we see that

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

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

$$\begin{aligned} \frac{L_2(W)}{L_0(W)} &= \prod_{i=1}^{T_{b_+}} \exp \left(\frac{-\left(Y_i - w(J_i) + \frac{2w(b)}{|b_-|}\right)^2 + (Y_i - w(J_i))^2}{2} \right) \\ &= \prod_{i=1}^{T_{b_+}} \exp \left(\frac{w(b)}{|b_+|} (2w(J_i) - 2Y_i) - \frac{2w(b)^2}{|b_+|^2} \right) \\ &= \exp \left(\frac{w(b)}{|b_+|} \left(\sum_{i=1}^{T_{b_+}} 2w(J_i) - 2Y_i \right) - \frac{2w(b)^2}{|b_+|^2} T_{b_+} \right) \\ &= \exp \left(\frac{w(b)}{|b_+|} \left(\sum_{i=1}^{T_{b_+}} 2w(J_i) - 2Y_i \right) \right) \exp \left(-\frac{2w(b)^2}{|b_+|^2} T_{b_+} \right) \\ &\geq \theta \\ &= 4\delta, \end{aligned} \tag{62}$$

where Eq. (62) can be obtained using same method as in Part (1) as well as the assumption that \mathcal{S}_2 occurred. Next, similar to the derivation in Eq. (60), we see that

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

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

$$\mathbb{E}_0[T_{b_+}] \geq t_{b_+}^* = \frac{|b_+|^2}{16w(b)^2} \log(1/4\delta). \tag{63}$$

(End of Part (2))

Finally, we combine the results from both parts, i.e. Eq. (61) and Eq. (63). We obtain

$$\begin{aligned} \mathbb{E}_0[T_b] &= \mathbb{E}_0[T_{b_-}] + \mathbb{E}_0[T_{b_+}] \\ &\geq \frac{|b_+|^2 + |b_-|^2}{16w(b)^2} \log(1/4\delta) \\ &\geq \frac{|b|^2}{32w(b)^2} \log(1/4\delta). \end{aligned}$$

□

8 Proof of Extension Results

8.1 Fixed Budget Setting

In this part, we analyze the probability of error of the modified **CGapExp** algorithm in the fixed budget setting and prove Theorem 3. First, we prove a lemma which characterizes the confidence intervals constructed in Theorem 3.

Lemma 10. *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 $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_e(t)}}$. Then, we have*

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

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

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

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 3 can be obtained from the key lemmas (Lemma 8 and Lemma 9) and Lemma 10.

Theorem 3. *Use the same notations as in Theorem 1. Given $T > n$ and parameter $\alpha > 0$, set the confidence radius $\text{rad}_t(e) = R\sqrt{\frac{\alpha}{T_e(t)}}$ for all arms $e \in [n]$ and all $t > 0$. Run **CGapExp** algorithm for at most T rounds. Then, for $0 \leq \alpha \leq \frac{1}{9}(T - n) (R^2 \text{width}(\mathcal{M})^2 \mathbf{H})^{-1}$, we have*

$$\Pr [\text{Out} \neq M_*] \leq 2Tn \exp(-2\alpha). \quad (5)$$

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

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

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

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

Hence we have

$$T_e \leq \frac{9R^2 \text{width}(\mathcal{B})^2}{\Delta_e^2} \cdot \alpha + 1. \quad (64)$$

By summing up Eq. (64) for all $e \in [n]$, we have

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

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

8.2 PAC Learning

First, we prove a (ϵ, δ) -PAC counterpart of Lemma 8.

Lemma 11. *If CGapExpPAC 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 CGapExpPAC ensures that $\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \leq \epsilon$. Therefore, we have

$$\begin{aligned} \epsilon &\geq \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \\ &\geq \tilde{w}_t(M_*) - \tilde{w}_t(M_t) \end{aligned} \tag{65}$$

$$= \langle \tilde{\mathbf{w}}_t, \chi_{M_*} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{M_*} - \chi_{M_t}| \rangle \tag{66}$$

$$\begin{aligned} &\geq \langle \mathbf{w}, \chi_{M_*} - \chi_{M_t} \rangle \\ &= w(M_*) - w(M_t), \end{aligned} \tag{67}$$

where Eq. (65) follows from the definition of $\tilde{M}_t \triangleq \arg \max_{M \in \mathcal{M}} \tilde{w}_t(M)$; Eq. (66) follows from Lemma 5; Eq. (67) follows from the assumption that ξ_t occurs and Lemma 6. \square

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

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

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

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

$$\epsilon \geq \sum_{i \in (M_t \setminus \tilde{M}_t) \cup (\tilde{M}_t \setminus M_t)} \text{rad}_t(i) \tag{68}$$

$$= \langle \mathbf{rad}_t, |\chi_{M_t} - \chi_{\tilde{M}_t}| \rangle, \tag{69}$$

where Eq. (68) 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. (69) 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

$$\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) = \langle \tilde{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \langle \mathbf{rad}_t, |\chi_{\tilde{M}_t} - \chi_{M_t}| \rangle \tag{70}$$

$$\leq \langle \tilde{\mathbf{w}}_t, \chi_{\tilde{M}_t} - \chi_{M_t} \rangle + \epsilon \tag{71}$$

$$= \tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) + \epsilon$$

$$\leq \epsilon, \quad (72)$$

where Eq. (70) follows from Lemma 5; Eq. (71) uses Eq. (69); and Eq. (72) follows from $\bar{w}_t(M_t) \geq \bar{w}_t(\tilde{M}_t)$. Therefore, we see that $\tilde{w}_t(\tilde{M}_t) - \tilde{w}_t(M_t) \leq \epsilon$. By the stopping condition of **CGapExpPAC**, the algorithm must terminate on round t . This contradicts to the assumption that $p_t = e$. \square

Using Lemma 12 and Lemma 11, we are ready to prove Theorem 4.

Theorem 4. *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 **Out** of **CGapExpPAC** satisfies $w(M_*) - w(\text{Out}) \leq \epsilon$. In addition, the number of samples T used by the algorithm satisfies*

$$T \leq O \left(R^2 \sum_{e \in [n]} \min \left\{ \frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2} \right\} \log \left(\frac{R^2 n}{\delta} \sum_{e \in [n]} \min \left\{ \frac{\text{width}(\mathcal{M})^2}{\Delta_e^2}, \frac{K^2}{\epsilon^2} \right\} \right) \right), \quad (6)$$

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

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

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

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

$$\min \left\{ \frac{\Delta_e}{3 \text{width}(\mathcal{B})}, \frac{\epsilon}{2K} \right\} \leq R \sqrt{\frac{2 \log(4nt_e^2/\delta)}{T_e - 1}} \leq R \sqrt{\frac{2 \log(4nT^2/\delta)}{T_e - 1}}. \quad (73)$$

Solving Eq. (73) for T_e , we obtain

$$T_e \leq R^2 \min \left\{ \frac{18 \text{width}(\mathcal{B})^2}{\Delta_e^2}, \frac{16K^2}{\epsilon^2} \right\} \log(4nT^2/\delta) + 1. \quad (74)$$

Notice that $T = \sum_{i \in [n]} T_i$. Hence the theorem follows by summing up Eq. (74) for all $e \in [n]$ and solving for T . \square

9 Technical Lemmas

Lemma 13 (Basis exchange property). *AA*

Lemma 14 (Hoeffding's inequality). *Let X_1, \dots, X_n be n independent R -sub-Gaussian random variables. Let $\bar{X} = \frac{1}{n} \sum X_i$ be the average of these random variables. Then, we have*

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

References