
Community Exploration: From Offline Optimization to Online Learning

Xiaowei Chen¹, Weiran Huang², Wei Chen³, John C.S. Lui¹

¹The Chinese University of Hong Kong

²Huawei Noah's Ark Lab, ³Microsoft Research

¹{xwchen, cslui}@cse.cuhk.edu.hk, ²huang.inbox@outlook.com

³weic@microsoft.com

Abstract

We introduce the community exploration problem that has many real-world applications such as online advertising. In the problem, an explorer allocates limited budget to explore communities so as to maximize the number of members he could meet. We provide a systematic study of the community exploration problem, from offline optimization to online learning. For the offline setting where the sizes of communities are known, we prove that the greedy methods for both of non-adaptive exploration and adaptive exploration are optimal. For the online setting where the sizes of communities are not known and need to be learned from the multi-round explorations, we propose an “upper confidence” like algorithm that achieves the logarithmic regret bounds. By combining the feedback from different rounds, we can achieve a constant regret bound.

1 Introduction

In this paper, we introduce the community exploration problem, which is abstracted from many real-world applications. Consider the following hypothetical scenario. Suppose that John just entered the university as a freshman. He wants to explore different student communities or study groups at the university to meet as many new friends as possible. But he only has a limited time to spend on exploring different communities, so his problem is how to allocate his time and energy to explore different student communities to maximize the number of people he would meet.

The above hypothetical community exploration scenario can also find similar counterparts in serious business and social applications. One example is online advertising. In this application, an advertiser wants to promote his products via placing advertisements on different online websites. The website would show the advertisements on webpages, and visitors to the websites may click on the advertisements when they view these webpages. The advertiser wants to reach as many unique customers as possible, but he only has a limited budget to spend. Moreover, website visitors come randomly, so it is not guaranteed that all visitors to the same website are unique customers. So the advertiser needs to decide how to spend the budget on each website to reach his customers. Of course, intuitively he should spend more budget on larger communities, but how much? And what if he does not know the user size of every website? In this case, each website is a community, consisting of all visitors to this website, and the problem can be modeled as a community exploration problem. Another example could be a social worker who wants to reach a large number of people from different communities to do social studies or improve the social welfare for a large population, while he also needs to face the budget constraint and uncertainty about the community.

In this paper, we abstract the common features of these applications and define the following community exploration problem that reflects the common core of the problem. We model the problem with m disjoint communities C_1, \dots, C_m with $C = \cup_{i=1}^m C_i$, where each community C_i has d_i

members. Each time when one explores (or visit) a community C_i , he would meet one member of C_i uniformly at random.¹ Given a budget K , the goal of community exploration is to determine the budget allocation $\mathbf{k} = (k_1, \dots, k_m) \in \mathbb{Z}_+^m$ with $\sum_{i=1}^m k_i \leq K$, such that the total number of distinct members met is maximized when each community C_i is explored k_i times.

We provide a systematic study of the above community exploration problem, from offline optimization to online learning. First, we consider the offline setting where the community sizes are known. In this setting, we further study two problem variants — the non-adaptive version and the adaptive version. The non-adaptive version requires that the complete budget allocation \mathbf{k} is decided before the exploration is started, while the adaptive version allows the algorithm to use the feedback from the exploration results of the previous steps to determine the exploration target of the next step. In both cases, we prove that the greedy algorithm provides the optimal solution. While the proof for the non-adaptive case is simple, the proof that the adaptive greedy policy is optimal is much more involved and relies on a careful analysis of transitions between system statuses. The proof techniques may be applicable in the analysis of other related problems.

Second, we consider the online setting where the community sizes are unknown in advance, which models the uncertainty about the communities in real applications. We apply the multi-armed bandit (MAB) framework to this task, in which community explorations proceed in multiple rounds, and in each round we explore communities with a budget of K , use the feedback to learn about the community size, and adjust the exploration strategy in future rounds. The reward of a round is the expected number of unique people met in the round. The goal is to maximize the cumulative reward from all rounds, or minimizing the regret, which is defined as the difference in cumulative reward between always using the optimal offline algorithm when knowing the community sizes and using the online learning algorithm. Similar to the offline case, we also consider the non-adaptive and adaptive version of exploration within each round. We provide theoretical regret bounds of $O(\log T)$ for both versions, where T is the number of rounds, which is asymptotically tight. Our analysis uses the special feature of the community exploration problem, which leads to improved coefficients in the regret bounds compared with a simple application of some existing results on combinatorial MABs. Moreover, we also discuss the possibility of using the feedback in previous round to turn the problem into the full information feedback model, which allows us to provide constant regret in this case.

In summary, our contributions include: (a) proposing the study of the community exploration problem to reflect the core of a number of real-world applications; and (b) a systematic study of the problem with rigorous theoretical analysis that covers offline non-adaptive, offline adaptive, online non-adaptive and online adaptive cases, which model the real-world situations of adapting to feedback and handling uncertainty.

2 Problem Definition

We model the problem with m disjoint communities C_1, \dots, C_m with $C = \cup_{i=1}^m C_i$, where each community C_i has d_i members. Each exploration (or visit) of one community C_i returns a member of C_i uniformly at random, and we have a total budget of K for explorations. Since we can trivially explore each community once when $K \leq m$, we assume that $K > m$.

We consider both the offline setting where the sizes of the communities d_1, \dots, d_m are known, and the online setting where the sizes of the communities are unknown. For the offline setting, we further consider two different problems: (1) non-adaptive exploration and (2) adaptive exploration. For the non-adaptive exploration, the explorer needs to predetermine the budget allocation \mathbf{k} before the exploration starts, while for the adaptive exploration, she can sequentially select the next community to explore based on previous observations (the members met in the previous community visits). Formally, we use pair (i, τ) to represent the τ -th exploration of community C_i , called an *item*. Let $\mathcal{E} = [m] \times [K]$ be the set of all possible items. A *realization* is a function $\phi: \mathcal{E} \rightarrow C$ mapping every possible item (i, τ) to a member in the corresponding community C_i , and $\phi(i, \tau)$ represents the member met in the exploration (i, τ) . We use Φ to denote a random realization, and the randomness comes from the exploration results. From the description above, Φ follows the distribution such that $\Phi(i, \tau) \in C_i$ is selected uniformly at random from C_i and is independent of all other $\Phi(i', \tau')$'s.

¹The model can be extended to meet multiple members per visit, but for simplicity, we consider meeting one member per visit in this paper.

For a budget allocation $\mathbf{k} = (k_1, \dots, k_m)$ and a realization ϕ , we define the reward R as the number of distinct members met, i.e., $R(\mathbf{k}, \phi) = \sum_{i=1}^m |\cup_{\tau=1}^{k_i} \{\phi(i, \tau)\}|$, where $|\cdot|$ is the cardinality of the set. The goal of the *non-adaptive exploration* is to find an optimal budget allocation $\mathbf{k}^* = (k_1^*, \dots, k_m^*)$ with given budget K , which maximizes the expected reward taken over all possible realizations, i.e.,

$$\mathbf{k}^* \in \arg \max_{\mathbf{k}: \|\mathbf{k}\|_1 \leq K} \mathbb{E}_{\Phi} [R(\mathbf{k}, \Phi)]. \quad (1)$$

For the adaptive exploration, the explorer sequentially picks a community to explore, meets a random member of the chosen community, then picks the next community, meets another random member of that community, and so on, until the budget is used up. After each selection, the observations so far can be represented as a *partial realization* ψ , a function from the subset of \mathcal{E} to $\mathcal{C} = \cup_{i=1}^m C_i$. Suppose that each community C_i has been explored k_i times. Then the partial realization ψ is a function mapping items in $\cup_{i=1}^m \{(i, 1), \dots, (i, k_i)\}$ (which is also called the domain of ψ , denoted as $\text{dom}(\psi)$) to members in communities. The partial realization ψ records the observation on the sequence of explored communities and the members met in this sequence. We say that a partial realization ψ is consistent with realization ϕ , denoted as $\phi \sim \psi$, if for all item (i, τ) in the domain of ψ , we have $\psi(i, \tau) = \phi(i, \tau)$. The strategy to explore the communities adaptively is encoded as a policy. The policy, denoted as π , is a function mapping ψ to an item in \mathcal{E} , specifying which community to explore next under the partial realization. Define $\pi_K(\phi) = (k_1, \dots, k_m)$, where k_i is the times the community C_i is explored via policy π under realization ϕ with budget K . More specifically, starting from the partial realization ψ_0 with empty domain, for every current partial realization ψ_s at step s , policy π determines the next community $\pi(\psi_s)$ to explore, meet the member $\phi(\pi(\psi_s))$, such that the new partial realization ψ_{s+1} is adding the mapping from $\pi(\psi_s)$ to $\phi(\pi(\psi_s))$ on top of ψ_s . This iteration continues until the communities have been explored K times, and $\pi_K(\phi) = (k_1, \dots, k_m)$ denotes the resulting exploration vector. The goal of the adaptive exploration is to find an optimal policy π^* to maximize the expected adaptive reward, i.e.,

$$\pi^* \in \arg \max_{\pi} \mathbb{E}_{\Phi} [R(\pi_K(\Phi), \Phi)]. \quad (2)$$

We next consider the online setting of community exploration. The learning process proceeds in discrete rounds. Initially, the size of communities $\mathbf{d} = (d_1, \dots, d_m)$ is unknown. In each round $t \geq 1$, the learner needs to determine an allocation or a policy (called an “*action*”) based on the previous-round observations to explore communities (non-adaptively or adaptively). When an action is played, the sets of encountered members for every community are observed as the *feedback* to the player. A learning algorithm A aims to cumulate as much reward (i.e., number of distinct members) as possible by selecting actions properly at each round. The performance of a learning algorithm is measured by the *cumulative regret*. Let Φ_t be the realization at round t . If we explore the communities with predetermined budget allocation in each round, the T -round (non-adaptive) regret of a learning algorithm A is defined as

$$\text{Reg}_{\mu}^A(T) = \mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=1}^T R(\mathbf{k}^*, \Phi_t) - R(\mathbf{k}_t^A, \Phi_t) \right], \quad (3)$$

where the budget allocation \mathbf{k}_t^A is selected by algorithm A in round t . If we explore the communities adaptively in each round, then the T -round (adaptive) regret of a learning algorithm A is defined as

$$\text{Reg}_{\mu}^A(T) = \mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=1}^T R(\pi_K^*(\Phi_t), \Phi_t) - R(\pi_K^{A,t}(\Phi_t), \Phi_t) \right], \quad (4)$$

where $\pi_K^{A,t}$ is a policy selected by algorithm A in round t . The goal of the learning problem is to design a learning algorithm A which minimizes the regret defined in (3) and (4).

3 Offline Optimization for Community Exploration

3.1 Non-adaptive Exploration Algorithms

If C_i is explored k_i times, each member in C_i is encountered at least once with probability $1 - (1 - 1/d_i)^{k_i}$. Thus we have $\mathbb{E}_{\Phi} [|\{\Phi(i, 1), \dots, \Phi(i, k_i)\}|] = d_i(1 - (1 - 1/d_i)^{k_i})$. Hence $\mathbb{E}_{\Phi} [R(\mathbf{k}, \Phi)]$ is a function of only the budget allocation \mathbf{k} and the size $\mathbf{d} = (d_1, \dots, d_m)$ of all communities.

Algorithm 1 Non-Adaptive community exploration with optimal budget allocation

```
1: procedure COMMUNITYEXPLORE( $\{\mu_1, \dots, \mu_m\}, K$ , non-adaptive)
2:   For  $i \in [m]$ ,  $k_i \leftarrow 0$   $\triangleright$  Line 2-5: budget allocation
3:   for  $s = 1, \dots, K$  do
4:      $i^* \leftarrow$  a random elements in  $\arg \max_i (1 - \mu_i)^{k_i}$   $\triangleright O(\log m)$  via using priority queue
5:      $k_{i^*} \leftarrow k_{i^*} + 1$ 
6:     For  $i \in [m]$ , explore  $C_i$  for  $k_i$  times, and put the uniformly met members in multi-set  $\mathcal{S}_i$ 
7:   end procedure
```

Algorithm 2 Adaptive community exploration with greedy policy

```
1: procedure COMMUNITYEXPLORE( $\{\mu_1, \dots, \mu_m\}, K$ , adaptive)
2:   For  $i \in [m]$ ,  $\mathcal{S}_i \leftarrow \emptyset$ ,  $c_i \leftarrow 0$   $\triangleright$  Line 2-7: adaptively explore communities with policy  $\pi^g$ 
3:   for  $s = 1, \dots, K$  do
4:      $i^* \leftarrow$  a random elements in  $\arg \max_i 1 - \mu_i c_i$ 
5:      $v \leftarrow$  a random member met when  $C_{i^*}$  is explored
6:     if  $v \notin \mathcal{S}_{i^*}$  then  $c_{i^*} \leftarrow c_{i^*} + 1$   $\triangleright v$  is not met before
7:      $\mathcal{S}_{i^*} \leftarrow \mathcal{S}_{i^*} \cup \{v\}$ 
8:   end procedure
```

Let $\mu_i = 1/d_i$, and vector $\boldsymbol{\mu} = (1/d_1, \dots, 1/d_m)$. Henceforth, we treat $\boldsymbol{\mu}$ as the parameter of the problem instance, since it is bounded with $\boldsymbol{\mu} \in [0, 1]^m$. Let $r_{\mathbf{k}}(\boldsymbol{\mu}) = \mathbb{E}_{\Phi}[R(\mathbf{k}, \Phi)]$ be the expected reward for the budget allocation \mathbf{k} . Based on the above discussion, we have

$$r_{\mathbf{k}}(\boldsymbol{\mu}) = \sum_{i=1}^m d_i (1 - (1 - 1/d_i)^{k_i}) = \sum_{i=1}^m (1 - (1 - \mu_i)^{k_i}) / \mu_i. \quad (5)$$

Since k_i must be integers, a traditional method like *Lagrange Multipliers* cannot be applied to solve the optimization problem defined in Eq. (1). We propose a *greedy method* consisting of K steps to compute the feasible \mathbf{k}^* . The greedy method is described in Line 2-5 of Algo. 1.

Theorem 1. *The greedy method obtains an optimal budget allocation.*

The time complexity of the greedy method is $O(K \log m)$, which is not efficient for large K . We find that starting from the initial allocation $k_i = \left\lceil \frac{(K-m)/\ln(1-\mu_i)}{\sum_{j=1}^m 1/\ln(1-\mu_j)} \right\rceil$, the greedy method can find the optimal budget allocation in $O(m \log m)$ ². (See Appendix A)

3.2 Adaptive Exploration Algorithms

With a slight abuse of notations, we also define $r_{\pi}(\boldsymbol{\mu}) = \mathbb{E}_{\Phi}[R(\pi_K(\Phi), \Phi)]$, since the expected reward is the function of the policy π and the vector $\boldsymbol{\mu}$. Define $c_i(\psi)$ as the number of distinct members we met in community C_i under partial realization ψ . Then $1 - c_i(\psi)/d_i$ is the probability that we can meet a new member in the community C_i if we explore community C_i one more time. A natural approach is to explore community C_{i^*} such that $i^* \in \arg \max_{i \in [m]} 1 - c_i(\psi)/d_i$ when we have partial realization ψ . We call such policy as the greedy policy π^g . The adaptive community exploration with greedy policy is described in Algo. 2. One could show that our reward function is actually an *adaptive submodular* function, for which the greedy policy is guaranteed to achieve at least $(1 - 1/e)$ of the maximized expected reward [13]. However, the following theorem shows that for our community exploration problem, our greedy policy is in fact *optimal*.

Theorem 2. *Greedy policy is the optimal policy for our adaptive exploration problem.*

Proof sketch. Note that the greedy policy chooses the next community only based on the fraction of unseen members. It does not care which members are already met. Thus, we define s_i as the percentage of members we have not met in a community C_i . We introduce the concept of *status*, denoted as $\mathbf{s} = (s_1, \dots, s_m)$. The greedy policy chooses next community based on the current

²We thank Jing Yu from School of Mathematical Sciences at Fudan University for her method to find a good initial allocation, which leads to a faster greedy method.

Algorithm 3 Combinatorial Lower Confidence Bound (CLCB) algorithm

Input budget for each round K , method (non-adaptive or adaptive)

- 1: For $i \in [m]$, $T_i \leftarrow 0$ (number of pairs), $X_i \leftarrow 0$ (collision counting), $\hat{\mu}_i \leftarrow 0$ (empirical mean)
 - 2: **for** $t = 1, 2, 3, \dots$ **do** \triangleright Line 2-8: online learning
 - 3: For $i \in [m]$, $\rho_i \leftarrow \sqrt{\frac{3 \ln t}{2T_i}}$ ($\rho_i = 0$ if $T_i = 0$) \triangleright confidence radius
 - 4: For $i \in [m]$, $\mu_i \leftarrow \max\{0, \hat{\mu}_i - \rho_i\}$ \triangleright lower confidence bound
 - 5: $\{\mathcal{S}_1, \dots, \mathcal{S}_m\} \leftarrow \text{COMMUNITYEXPLORE}(\{\mu_1, \dots, \mu_m\}, K, \text{method})$ $\triangleright \mathcal{S}_i$: set of met members
 - 6: For $i \in [m]$, $T_i \leftarrow T_i + \lfloor |\mathcal{S}_i|/2 \rfloor$ \triangleright update number of (member) pairs we observe
 - 7: For $i \in [m]$, $X_i \leftarrow X_i + \sum_{x=1}^{\lfloor |\mathcal{S}_i|/2 \rfloor} \mathbb{1}\{\mathcal{S}_i[2x-1] = \mathcal{S}_i[2x]\}$ $\triangleright \mathcal{S}_i[x]$: x -th element in \mathcal{S}_i
 - 8: For $i \in [m]$ and $|\mathcal{S}_i| > 1$, $\hat{\mu}_i \leftarrow X_i/T_i$ \triangleright update empirical mean
-

status. In the proof, we further extend the definition of reward with a non-decreasing function f as $R(\mathbf{k}, \phi) = f\left(\sum_{i=1}^m \left|\bigcup_{\tau=1}^{k_i} \{\phi(i, \tau)\}\right|\right)$. Note that the reward function corresponding to the original community exploration problem is simply the identity function $f(x) = x$. Let $F_\pi(\psi, t)$ denote the expected *marginal gain* when we further explore communities for t steps with policy π starting from a partial realization ψ . We want to prove that for all ψ, t and π , $F_{\pi^g}(\psi, t) \geq F_\pi(\psi, t)$, where π^g is the greedy policy and π is an arbitrary policy. If so, we simply take $\psi = \emptyset$, and $F_{\pi^g}(\emptyset, t) \geq F_\pi(\emptyset, t)$ for every π and t exactly shows that π^g is optimal. We prove the above result by an induction on t .

Let C_i be the community chosen by π based on the partial realization ψ . Define $c(\psi) = \sum_i c_i(\psi)$ and $\Delta_{\psi, f} = f(c(\psi) + 1) - f(c(\psi))$. We first claim that $F_{\pi^g}(\psi, 1) \geq F_\pi(\psi, 1)$ holds for all ψ and π with the fact that $F_\pi(\psi, 1) = (1 - \mu_i c_i(\psi)) \Delta_{\psi, f}$. Note that the greedy policy π^g chooses C_{i^*} with $i^* \in \arg \max_i (1 - \mu_i c_i(\psi))$. Hence, $F_{\pi^g}(\psi, 1) \geq F_\pi(\psi, 1)$.

Next we prove that $F_{\pi^g}(\psi, t+1) \geq F_\pi(\psi, t+1)$ based on the assumption that $F_{\pi^g}(\psi, t') \geq F_\pi(\psi, t')$ holds for all ψ, π , and $t' \leq t$. An important observation is that $F_{\pi^g}(\psi, t)$ has equal value for any partial realization ψ associated with the same status \mathbf{s} since the status is enough for the greedy policy to determine the choice of next community. Formally, we define $F_g(\mathbf{s}, t) = F_{\pi^g}(\psi, t)$ for any partial realization that satisfies $\mathbf{s} = (1 - c_1(\psi)/d_1, \dots, 1 - c_m(\psi)/d_m)$. Let C_{i^*} denote the community chosen by policy π^g under realization ψ , i.e., $i^* \in \arg \max_{i \in [m]} 1 - \mu_i c_i(\psi)$. Let \mathbf{I}_i be the m -dimensional unit vector with one in the i -th entry and zeros in all other entries. We show that

$$\begin{aligned} F_\pi(\psi, t+1) &\leq c_i(\psi) \cdot \mu_i F_g(\mathbf{s}, t) + (d_i - c_i(\psi)) \cdot \mu_i F_g(\mathbf{s} - \mu_i \mathbf{I}_i, t) + (1 - \mu_i c_i(\psi)) \Delta_{\psi, f} \\ &\leq \mu_{i^*} c_{i^*}(\psi) F_g(\mathbf{s}, t) + (1 - \mu_{i^*} c_{i^*}(\psi)) F_g(\mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}, t) + (1 - \mu_{i^*} c_{i^*}(\psi)) \Delta_{\psi, f} \\ &= F_g(\mathbf{s}, t+1) = F_{\pi^g}(\psi, t+1). \end{aligned}$$

The first line is derived directly from the definition and the assumption. The key is to prove the correctness of Line 2 in above inequality. It indicates that if we choose a sub-optimal community at first, and then we switch back to the greedy policy, the expected reward would be smaller. The proof is nontrivial and relies on a careful analysis based on the stochastic transitions among status vectors. Note that the reward function $r_\pi(\boldsymbol{\mu})$ is not necessary adaptive submodular if we extend the reward with the non-decreasing function f . Hence, a $(1 - 1/e)$ guarantee for adaptive submodular function [13] is not applicable in this scenario. Our analysis scheme can be applied to any adaptive problems with similar structures.

4 Online Learning for Community Exploration

The key of the learning algorithm is to estimate the community sizes. The size estimation problem is defined as inferring unknown set size d_i from random samples obtained via uniformly sampling *with replacement* from the set C_i . Various estimators have been proposed [3, 8, 10, 16] for the estimation of d_i . The core idea of estimators in [3, 16] are based on “collision counting”. Let (u, v) be an *unordered pair* of two random elements from C_i and $Y_{u,v}$ be a *pair collision* random variable that takes value 1 if $u = v$ (i.e., (u, v) is a *collision*) and 0 otherwise. It is easy to verify that $\mathbb{E}[Y_{u,v}] = 1/d_i = \mu_i$. Suppose we *independently* take T_i pairs of elements from C_i and X_i of them are collisions. Then $\mathbb{E}[X_i/T_i] = 1/d_i = \mu_i$. The size d_i can be estimated by T_i/X_i (the estimator is valid when $X_i > 0$).

We present our CLCB algorithm in Algorithm 3. In the algorithm, we maintain an unbiased estimation of μ_i instead of d_i for each community C_i for the following reasons. Firstly, T_i/X_i is not an unbiased estimator of d_i since $\mathbb{E}[T_i/X_i] \geq d_i$ according to the Jensen's inequality. Secondly, the upper confidence bound of T_i/X_i depends on d_i , which is unknown in our online learning problem. Thirdly, we need at least $(1 + \sqrt{8d_i \ln 1/\delta + 1})/2$ uniformly sampled elements in C_i to make sure that $X_i > 0$ with probability at least $1 - \delta$. We feed the lower confidence bound μ_i to the exploration process since our reward function increases as μ_i decreases. The idea is similar to CUCB algorithm [7]. The lower confidence bound is small if community C_i is not explored often (T_i is small). Small μ_i motivates us to explore C_i more times. The *feedbacks* after the exploration process at each round are the sets of encountered members $\mathcal{S}_1, \dots, \mathcal{S}_m$ in communities C_1, \dots, C_m respectively. Note that for each $i \in [m]$, all pairs of elements in \mathcal{S}_i , namely $\{(x, y) \mid x \leq y, x \in \mathcal{S}_i, y \in \mathcal{S}_i \setminus \{x\}\}$ are not mutually independent. Thus, we only use $\lfloor |\mathcal{S}_i|/2 \rfloor$ independent pairs. Therefore, T_i is updated as $T_i + \lfloor |\mathcal{S}_i|/2 \rfloor$ at each round. In each round, the community exploration could either be non-adaptive or adaptive, and the following regret analysis separately discuss these two cases.

4.1 Regret Analysis for the Non-adaptive Version

The non-adaptive bandit learning model fits into the general combinatorial multi-armed bandit (CMAB) framework of [7, 20] that deals with nonlinear reward functions. In particular, we can treat the pair collision variable in each community C_i as a base arm, and our expected reward in Eq. (5) is non-linear, and it satisfies the monotonicity and bounded smoothness properties (See Properties 1 and 2). However, directly applying the regret result from [7, 20] will give us an inferior regret bound for two reasons. First, in our setting, in each round we could have multiple sample feedback for each community, meaning that each base arm could be observed multiple times, which is not directly covered by CMAB. Second, to use the regret result in [7, 20], the bounded smoothness property needs to have a bounded smoothness constant independent of the actions, but we can have a better result by using a tighter form of bounded smoothness with action-related coefficients. Therefore, in this section, we provide a better regret result by adapting the regret analysis in [20].

We define the gap $\Delta_{\mathbf{k}} = r_{\mathbf{k}^*}(\boldsymbol{\mu}) - r_{\mathbf{k}}(\boldsymbol{\mu})$ for all action \mathbf{k} satisfying $\sum_{i=1}^m k_i = K$. For each community C_i , we define $\Delta_{\min}^i = \min_{\Delta_{\mathbf{k}} > 0, k_i > 1} \Delta_{\mathbf{k}}$ and $\Delta_{\max}^i = \max_{\Delta_{\mathbf{k}} > 0, k_i > 1} \Delta_{\mathbf{k}}$. As a convention, if there is no action \mathbf{k} with $k_i > 1$ such that $\Delta_{\mathbf{k}} > 0$, we define $\Delta_{\min}^i = \infty$ and $\Delta_{\max}^i = 0$. Furthermore, define $\Delta_{\min} = \min_{i \in [m]} \Delta_{\min}^i$ and $\Delta_{\max} = \max_{i \in [m]} \Delta_{\max}^i$. Let $K' = K - m + 1$. We have the regret for Algo. 3 as follows.

Theorem 3. *Algo. 3 with non-adaptive exploration method has regret as follows.*

$$\text{Reg}_{\boldsymbol{\mu}}(T) \leq \sum_{i=1}^m \frac{48 \binom{K'}{2} K \ln T}{\Delta_{\min}^i} + 2 \binom{K'}{2} m + \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3} m \Delta_{\max} = O \left(\sum_{i=1}^m \frac{K'^3 \log T}{\Delta_{\min}^i} \right). \quad (6)$$

The proof of the above theorem is an adaption of the proof of Theorem 4 in [20], and the full proof details as well as the detailed comparison with the original CMAB framework result are included in the supplementary materials. We briefly explain our adaption that leads to the regret improvement. We rely on the following monotonicity and 1-norm bounded smoothness properties of our expected reward function $r_{\mathbf{k}}(\boldsymbol{\mu})$, similar to the ones in [7, 20].

Property 1 (Monotonicity). *The reward function $r_{\mathbf{k}}(\boldsymbol{\mu})$ is monotonically decreasing, i.e., for any two vectors $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$ and $\boldsymbol{\mu}' = (\mu'_1, \dots, \mu'_m)$, we have $r_{\mathbf{k}}(\boldsymbol{\mu}) \geq r_{\mathbf{k}}(\boldsymbol{\mu}')$ if $\mu_i \leq \mu'_i \forall i \in [m]$.*

Property 2 (1-Norm Bounded Smoothness). *The reward function $r_{\mathbf{k}}(\boldsymbol{\mu})$ satisfies the 1-norm bounded smoothness property, i.e., for any two vectors $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$ and $\boldsymbol{\mu}' = (\mu'_1, \dots, \mu'_m)$, we have $|r_{\mathbf{k}}(\boldsymbol{\mu}) - r_{\mathbf{k}}(\boldsymbol{\mu}')| \leq \sum_{i=1}^m \binom{k_i}{2} |\mu_i - \mu'_i| \leq \binom{K'}{2} \sum_{i=1}^m |\mu_i - \mu'_i|$.*

We remark that if we directly apply the CMAB regret bound of Theorem 4 in [20], we need to revise the update procedure in Lines 6-8 of Algo. 3 so that each round we only update one observation for each community C_i if $|\mathcal{S}_i| > 1$. Then we would obtain a regret bound $O \left(\sum_i \frac{K'^4 m \log T}{\Delta_{\min}^i} \right)$, which means that our regret bound in Eq. (6) has an improvement of $O(K'm)$. This improvement is exactly due to the reason we give earlier, as we now explain with more details.

For all the random variables introduced in Algo. 3, we add the subscript t to denote their value at the end of round t . For example, $T_{i,t}$ is the value of T_i at the end of round t . First, the improvement of

the factor m comes from the use of a tighter bounded smoothness in Property 2, namely, we use the bound $\sum_{i=1}^m \binom{k_i}{2} |\mu_i - \mu'_i|$ instead of $\binom{K'}{2} \sum_{i=1}^m |\mu_i - \mu'_i|$. The CMAB framework in [20] requires the bounded smoothness constant to be independent of actions. So to apply Theorem 4 in [20], we have to use the bound $\binom{K'}{2} \sum_{i=1}^m |\mu_i - \mu'_i|$. However, in our case, when using bound $\sum_{i=1}^m \binom{k_i}{2} |\mu_i - \mu'_i|$, we are able to utilize the fact $\sum_{i=1}^m \binom{k_i}{2} \leq \binom{K'}{2}$ to improve the result by a factor of m . Second, the improvement of the $O(K')$ factor, more precisely a factor of $(K' - 1)/2$, is achieved by utilizing multiple feedback in a single round and a more careful analysis of the regret utilizing the property of the right Riemann summation. Specifically, let $\Delta_{k_t} = r_{k^*}(\mu) - r_{k_t}(\mu)$ be the reward gap. When the estimate is within the confidence radius, we have $\Delta_{k_t} \leq \sum_{i=1}^m \frac{c(k_{i,t}-1)}{2} / \sqrt{T_{i,t-1}} \leq c \sum_{i=1}^m \lfloor k_{i,t}/2 \rfloor / \sqrt{T_{i,t-1}}$, where c is a constant. In Algo. 3, we have $T_{i,t} = T_{i,t-1} + \lfloor k_{i,t}/2 \rfloor$ because we allow multiple feedback in a single round. Then $\sum_{t \geq 1, T_{i,t} \leq L_i(T)} \lfloor k_{i,t}/2 \rfloor / \sqrt{T_{i,t-1}}$ is the form of a right Riemann summation, which achieves the maximum value when $k_{i,t} = K'$. Here $L_i(T)$ is a $\ln T$ function with some constants related with community C_i . Hence the regret bound $\sum_{t=1}^T \Delta_{k_t} \leq c \sum_{i=1}^m \sum_{t \geq 1, T_{i,t} \leq L_i(T)} \lfloor \frac{k_{i,t}}{2} \rfloor / \sqrt{T_{i,t-1}} \leq 2c \sum_{i=1}^m \sqrt{L_i(T)}$. However, if we use the original CMAB framework, we need to set $T_{i,t} = T_{i,t-1} + \mathbb{1}\{k_{i,t} > 1\}$. In this case, we can only bound the regret as $\sum_{t=1}^T \Delta_{k_t} = c \sum_{i=1}^m \sum_{t \geq 1, T_{i,t} \leq L_i(T)} (k_{i,t} - 1)/2 \sqrt{T_{i,t-1}} \leq 2c \frac{K'-1}{2} \sum_{i=1}^m \sqrt{L_i(T)}$, leading to an extra factor of $(K' - 1)/2$.

Justification for Algo. 3. In Algo. 3, we only use the members in current round to update the estimator. This is practical for the situation where the member identifiers are changing in different rounds for privacy protection. Privacy gains much attention these days. Consider the online advertising scenario we explain in the introduction. Whenever a user clicks an advertisement, the advertiser would store the user information (e.g. Facebook ID, IP address etc.) to identify the user and correlated with past visits of the user. If such user identifiers are fixed and do not change, the advertiser could easily track user behavior, which may result in privacy leak. A reasonable protection for users is to periodically change user IDs (e.g. Facebook can periodically change user hash IDs, or users adopt dynamic IP addresses, etc.), so that it is difficult for the advertiser to track the same user over a long period of time. Under such situation, it may be likely that our learning algorithm can still detect ID collisions within the short period of each learning round, but cross different rounds, collisions may not be detectable due to ID changes.

Full information feedback. Now we consider the scenario where the member identifiers are fixed over all rounds, and design an algorithm with a constant regret bound. Our idea is to ensure that we can observe at least one pair of members in every community C_i in each round t . We call such guarantee as *full information feedback*. If we only use members revealed in current round, we cannot achieve this goal since we have no observation of new pairs for a community C_i when $k_i = 1$. To achieve full information feedback, we use at least one sample from the previous round to form a pair with a sample in the current round to generate a valid pair collision observation. In particular, we revise the Line 3, 6, and 7 as follows. Here we use u_0 to represent the last member in \mathcal{S}_i in the previous round (let $u_0 = \text{null}$ when $t = 1$) and $u_x (x > 0)$ to represent the x -th members in \mathcal{S}_i in the current round. The revision of Line 3 implies that we use the empirical mean $\hat{\mu}_i = X_i/T_i$ instead of the lower confidence bound in the function COMMUNITYEXPLORE.

$$\begin{aligned} \text{Line 3: } & \text{For } i \in [m], \rho_i = 0; \quad \text{Line 6: } \text{For } i \in [m], T_i \leftarrow T_i + |\mathcal{S}_i| - \mathbb{1}\{t = 1\}, \\ \text{Line 7: } & \text{For } i \in [m], X_i \leftarrow X_i + \sum_{x=0}^{|\mathcal{S}_i|-1} \mathbb{1}\{u_x = u_{x+1}\}. \end{aligned} \tag{7}$$

Theorem 4. *With the full information feedback revision in Eq. (7), Algo. 3 with non-adaptive exploration method has a constant regret bound. Specifically,*

$$\text{Reg}_\mu(T) \leq (2 + 2me^2 K'^2 (K' - 1)^2 / \Delta_{\min}^2) \Delta_{\max}.$$

Note that we cannot apply the Hoeffding bound in [14] directly since the random variables $\mathbb{1}\{u_x = u_{x+1}\}$ we obtain during the online learning process are not mutually independent. Instead, we apply a concentration bound in [9] that is applicable to variables that have local dependence relationship.

4.2 Regret Analysis for the Adaptive Version

For the adaptive version, we feed the lower confidence bound $\underline{\mu}_t$ into the adaptive community exploration procedure, namely $\text{COMMUNITYEXPLORE}(\{\underline{\mu}_1, \dots, \underline{\mu}_m\}, K, \text{adaptive})$ in round t . We denote the policy implemented by this procedure as π^t . Note that both π^g and π^t are based on the greedy procedure $\text{COMMUNITYEXPLORE}(\cdot, K, \text{adaptive})$. The difference is that π^g uses the true parameter μ while π^t uses the lower bound parameter $\underline{\mu}_t$. More specifically, given a partial realization ψ , the community chosen by π^t is C_{i^*} where $i^* \in \arg \max_{i \in [m]} 1 - c_i(\psi) \underline{\mu}_{i,t}$. Recall that $c_i(\psi)$ is the number of distinct encountered members in community C_i under partial realization ψ .

Similar to π^g , the policy π^t also chooses next community to explore based on current *status*. Let $\mathbf{s} = (s_1, \dots, s_m) = (1 - c_1(\psi) \underline{\mu}_1, \dots, 1 - c_m(\psi) \underline{\mu}_m)$ be the corresponding status to the partial realization ψ . Here s_i is the percentage of unmet members in the community C_i . For any partial realization ψ having status \mathbf{s} , the policy π^t choose C_{i^*} to explore, where $i^* \in \arg \max_{i \in [m]} (\underline{\mu}_{i,t} / \mu_i) s_i + (\mu_i - \underline{\mu}_{i,t}) / \mu_i$. When $\underline{\mu}_{i,t} \leq \mu_i$, we have $(\underline{\mu}_{i,t} / \mu_i) s_i + (\mu_i - \underline{\mu}_{i,t}) / \mu_i \geq s_i$, which means that the percentage of unmet members in C_i is overestimated by π^t .

We first properly define the metrics $\Delta_{\min}^{i,k}$ and $\Delta_{\max}^{(k)}$ used in the regret bound as follows. Consider a specific full realization ϕ where $\{\phi(i, 1), \dots, \phi(i, d_i)\}$ are d_i distinct members in C_i for $i \in [m]$. The realization ϕ indicates that we will obtain a new member in the first d_i exploration of community C_i . Let $U_{i,k}$ denote the number of times community C_i is selected by policy π^g in the first $k-1$ ($k > m$) steps under the special full realization ϕ we define previously. We define $\Delta_{\min}^{i,k} = (\mu_i U_{i,k} - \min_{j \in [m]} \mu_j U_{j,k}) / U_{i,k}$. Conceptually, the value $\mu_i U_{i,k} - \min_{j \in [m]} \mu_j U_{j,k}$ is gap in the expected reward of the next step between selecting a community by π^g (the optimal policy) and selecting community C_i , when we already meet $U_{j,k}$ distinct members in C_j for $j \in [m]$. When $\mu_i U_{i,k} = \min_{j \in [m]} \mu_j U_{j,k}$, we define $\Delta_{\min}^{i,k} = \infty$. Let π be another policy that chooses the same sequence of communities as π^g when the number of met members in C_i is no more than $U_{i,k}$ for all $i \in [m]$. Note that policy π chooses the same communities as π^g in the first $k-1$ steps under the special full realization ϕ . Actually, the policy π is the same as π^g for at least $k-1$ steps. We use Π_k to denote the set of all such policies. We define $\Delta_{\max}^{(k)}$ as the maximum reward gap between the policy $\pi \in \Pi_k$ and the optimal policy π^g , i.e., $\Delta_{\max}^{(k)} = \max_{\pi \in \Pi_k} r_{\pi^g}(\mu) - r_{\pi}(\mu)$. Let $D = \sum_{i=1}^m d_i$.

Theorem 5. *Algo. 3 with adaptive exploration method has regret as follows.*

$$\text{Reg}_{\mu}(T) \leq \left(\sum_{i=1}^m \sum_{k=m+1}^{\min\{K, D\}} \frac{6\Delta_{\max}^{(k)}}{(\Delta_{\min}^{i,k})^2} \right) \ln T + \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3} \sum_{i=1}^m \sum_{k=m+1}^{\min\{K, D\}} \Delta_{\max}^{(k)}. \quad (8)$$

Theorem 6. *With the full information feedback revision in Eq. (7), Algo. 3 with adaptive exploration method has a constant regret bound. Specifically,*

$$\text{Reg}_{\mu}(T) \leq \sum_{i=1}^m \sum_{k=m+1}^{\min\{K, D\}} (2/\varepsilon_{i,k}^4 + 1) \Delta_{\max}^{(k)}.$$

where $\varepsilon_{i,k}$ is defined as (here $i_k^* \in \arg \min_{i \in [m]} \mu_i U_{i,k}$)

$$\varepsilon_{i,k} \triangleq (\mu_i U_{i,k} - \mu_{i_k^*} U_{i_k^*,k}) / (U_{i,k} + U_{i_k^*,k}) \text{ for } i \neq i_k^* \text{ and } \varepsilon_{i,k} = \infty \text{ for } i = i_k^*.$$

Gabillon et al. [11] analyzes a general adaptive submodular function maximization in bandit setting. We have a regret bound in similar form as (8) if we directly apply Theorem 1 in [11]. However, their version of $\Delta_{\max}^{(k)}$ is an upper bound on the expected reward of policy π^g from k steps forward, which is larger than our $\Delta_{\max}^{(k)}$. Their version of $\Delta_{\min}^{i,k}$ is the minimum $(\mu_i c_i(\psi) - \min_{j \in [m]} \mu_j c_j(\psi)) / c_i(\psi)$ for all partial realization ψ obtained after policy π^g is executed for k steps, which is smaller than our $\Delta_{\min}^{i,k}$. Our regret analysis is based on counting how many times π^g and π^t choose different communities under the special full realization ϕ , while the analysis in [11] is based on counting how many times π^g and π^t choose different communities under all possible full realizations.

Discussion. In this paper, we consider the online learning problem that consists of T rounds, and during each round, we explore the communities with a budget K . Our goal is to maximize the *cumulative reward* in T rounds. Another important and natural setting is described as follows. We

start to explore communities with unknown sizes, and update the parameters every time we explore the community for *one step* (or for a few steps). Different from the setting defined in this paper, here *a member will not contribute to the reward if it has been met in previous rounds*. To differentiate the two settings, let’s call the latter one the “*interactive community exploration*”, while the former one the “*repeated community exploration*”. Both the repeated community exploration defined in this paper and the interactive community exploration we will study as the future work have corresponding applications. The former is suitable for online advertising where in each round the advertiser promotes different products. Hence the rewards in different rounds are additive. The latter corresponds to the adaptive online advertising for the same product, and thus the rewards in different rounds are dependent.

5 Related Work

Golovin and Krause [13] show that a greedy policy could achieve at least $(1 - 1/e)$ approximation for the adaptive submodular function. The result could be applied to our offline adaptive problem, but by an independent analysis we show the better result that the greedy policy is optimal. Multi-armed bandit (MAB) problem is initiated by Robbins [18] and extensively studied in [2, 4, 19]. Our online learning algorithm is based on the extensively studied *Upper Confidence Bound* approach [1]. The non-adaptive community exploration problem in the online setting fits into the general combinatorial multi-armed bandit (CMAB) framework [6, 7, 12, 17, 20], where the reward is a set function of base arms. The CMAB problem is first studied in [12], and its regret bound is improved by [7, 17]. We leverage the analysis framework in [7, 20] and prove a tighter bound for our algorithm. Gabillon et al. [11] define an adaptive submodular maximization problem in bandit setting. Our online adaptive exploration problem is a instance of the problem defined in [11]. We prove a tighter bound than the one in [11] by using the properties of our problem.

Our model bears similarities to the optimal discovery problem proposed in [5] such as we both have disjoint assumption, and both try to maximize the number of target elements. However, there are also some differences: (a) We use different estimators for our critical parameters, because our problem setting is different. (b) Their online model is closer to the interactive community exploration we explained in 4.2, while our online model is on repeated community exploration. As explained in 4.2, the two online models serve different applications and have different algorithms and analyses. (c) We also have more comprehensive studies on the offline cases.

6 Future Work

In this paper, we systematically study the community exploration problems. In the offline setting, we propose the greedy methods for both of non-adaptive and adaptive exploration problems. The optimality of the greedy methods are rigorously proved. We also analyze the online setting where the community sizes are unknown initially. We provide a CLCB algorithm for the online community exploration. The algorithm has $O(\log T)$ regret bound. If we further allow the full information feedback, the CLCB algorithm with some minor revisions has a constant regret.

Our study opens up a number of possible future directions. For example, we can consider various extensions to the problem model, such as more complicated distributions of member meeting probabilities, overlapping communities, or even graph structures between communities. We could also study the gap between non-adaptive and adaptive solutions.

Acknowledgments

We thank Jing Yu from School of Mathematical Sciences at Fudan University for her insightful discussion on the offline problems, especially, we thank Jing Yu for her method to find a good initial allocation, which leads to a faster greedy method. Wei Chen is partially supported by the National Natural Science Foundation of China (Grant No. 61433014). The work of John C.S. Lui is supported in part by the GRF Grant 14208816.

References

- [1] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [2] Donald A Berry and Bert Fristedt. Bandit problems: sequential allocation of experiments. *Chapman and Hall*, 5:71–87, 1985.
- [3] Marco Bressan, Enoch Peserico, and Luca Pretto. Simple set cardinality estimation through random sampling. *arXiv preprint arXiv:1512.07901*, 2015.
- [4] Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- [5] Sébastien Bubeck, Damien Ernst, and Aurélien Garivier. Optimal discovery with probabilistic expert advice: finite time analysis and macroscopic optimality. *JMLR*, 14(Feb):601–623, 2013.
- [6] Wei Chen, Wei Hu, Fu Li, Jian Li, Yu Liu, and Pinyan Lu. Combinatorial multi-armed bandit with general reward functions. In *NIPS*, pages 1659–1667, 2016.
- [7] Wei Chen, Yajun Wang, Yang Yuan, and Qinshi Wang. Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. *Journal of Machine Learning Research*, 17(50):1–33, 2016. A preliminary version appeared as Chen, Wang, and Yuan, “Combinatorial multi-armed bandit: General framework, results and applications”, ICML’2013.
- [8] Mary C Christman and Tapan K Nayak. Sequential unbiased estimation of the number of classes in a population. *Statistica Sinica*, pages 335–352, 1994.
- [9] Devdatt Dubhashi and Alessandro Panconesi. *Concentration of Measure for the Analysis of Randomized Algorithms*. Cambridge University Press, 1st edition, 2009.
- [10] Mark Finkelstein, Howard G. Tucker, and Jerry Alan Veeh. Confidence intervals for the number of unseen types. *Statistics & Probability Letters*, pages 423 – 430, 1998.
- [11] Victor Gabillon, Branislav Kveton, Zheng Wen, Brian Eriksson, and S Muthukrishnan. Adaptive submodular maximization in bandit setting. In *NIPS*, pages 2697–2705, 2013.
- [12] Yi Gai, Bhaskar Krishnamachari, and Rahul Jain. Combinatorial network optimization with unknown variables: Multi-armed bandits with linear rewards and individual observations. *IEEE/ACM Trans. Netw.*, 20(5):1466–1478, 2012.
- [13] Daniel Golovin and Andreas Krause. Adaptive submodularity: Theory and applications in active learning and stochastic optimization. *Journal of Artificial Intelligence Research*, 42: 427–486, 2011.
- [14] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American statistical association*, 58(301):13–30, 1963.
- [15] Svante Janson. Large deviations for sums of partly dependent random variables. *Random Structures & Algorithms*, 24(3):234–248, 2004.
- [16] Liran Katzir, Edo Liberty, and Oren Somekh. Estimating sizes of social networks via biased sampling. In *WWW*, 2011.
- [17] Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight regret bounds for stochastic combinatorial semi-bandits. In *Artificial Intelligence and Statistics*, pages 535–543, 2015.
- [18] Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*, pages 169–177. Springer, 1985.
- [19] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [20] Qinshi Wang and Wei Chen. Improving regret bounds for combinatorial semi-bandits with probabilistically triggered arms and its applications. In *NIPS*, pages 1161–1171, 2017.

Supplementary Materials

A Improved Budget Allocation Algorithm

Theorem 1. *The greedy method obtains an optimal budget allocation.*

Proof. Let $r_i(j) = \mathbb{E}_\Phi[\{\Phi(i, 1), \dots, \Phi(i, j)\}] = d_i(1 - (1 - 1/d_i)^j)$ denote the expected reward when the community i is explored j times. Then we have that the marginal gain $r_i(j+1) - r_i(j) = (1 - \mu_i)^j$. Define a matrix $\mathbf{X} \in \mathbb{R}^{m \times K}$, where the (i, j) -th entry $X_{i,j}$ is $(1 - \mu_i)^{j-1}$. When the budget allocation is $\mathbf{k} = (k_1, \dots, k_m)$, the expected reward $r_{\mathbf{k}}(\boldsymbol{\mu})$ can be written as the sum of elements in \mathbf{X} , i.e., $r_{\mathbf{k}}(\boldsymbol{\mu}) = \sum_{i=1}^m \sum_{j=1}^{k_i} X_{i,j}$. A key property of \mathbf{X} is that the value in each row is decreasing with respect to the column index j . Hence, for every $s \geq 1$, the s -th step of the greedy method chooses the s -th largest value in \mathbf{X} . At step $s = K$, the greedy method finds the largest K values in matrix \mathbf{X} . We can conclude that the greedy method obtains a budget allocation that maximizes the reward $r_{\mathbf{k}}(\boldsymbol{\mu})$. ■

We propose a budget allocation algorithm which has time complexity $O(m \log m)$ in Algo. 4. The basic idea is to find a good initial allocation that is not far from the optimal allocation. Then starting from the initial allocation, we run our original greedy method.

Algorithm 4 Budget allocation algorithm

Input parameters $\boldsymbol{\mu}$, budget $K > m$
1: For $i \in [m]$, $k_i = \lceil ((K - m) / \ln(1 - \mu_i)) / (\sum_{j=1}^m 1 / \ln(1 - \mu_j)) \rceil$ \triangleright A good initial allocation
2: **while** $\sum_{i=1}^m k_i < K$ **do**
3: $i^* \leftarrow \arg \max_i (1 - \mu_i)^{k_i}$ $\triangleright O(\log m)$ via using priority queue
4: $k_{i^*} \leftarrow k_{i^*} + 1$

Lemma 1 (Basic property of optimal allocation). *Let \mathbf{k}^* be the optimal budget allocation when the parameter of the community is $\boldsymbol{\mu}$. For $i, j \in [m]$, we have*

$$(1 - \mu_i)^{(k_i^* - 1)} \geq (1 - \mu_j)^{k_j^*}.$$

Proof. We define budget allocation \mathbf{k}' which is the same as \mathbf{k}^* except that $k'_i = k_i^* - 1$ and $k'_j = k_j^* + 1$. If $(1 - \mu_i)^{(k_i^* - 1)} < (1 - \mu_j)^{k_j^*}$ and $i \neq j$, then we have

$$r_{\mathbf{k}'}(\boldsymbol{\mu}) = r_{\mathbf{k}^*}(\boldsymbol{\mu}) - (1 - \mu_i)^{(k_i^* - 1)} + (1 - \mu_j)^{k_j^*} > r_{\mathbf{k}^*}(\boldsymbol{\mu}),$$

which is contradict with the fact that \mathbf{k}^* is the optimal solution. This proves the lemma. ■

Lemma 2 (Allocation lower bound). *Let \mathbf{k}^* be the optimal budget allocation when the parameter of the communities is $\boldsymbol{\mu}$. Define $\mathbf{k}^- = (k_1^-, \dots, k_m^-)$ where*

$$k_i^- = \frac{(K - m) / \ln(1 - \mu_i)}{\sum_{j=1}^m 1 / \ln(1 - \mu_j)}.$$

We have $k_i^ \geq k_i^-$.*

Proof. According to the definition of \mathbf{k}^- , we have $k_i^- \ln(1 - \mu_i) = k_j^- \ln(1 - \mu_j)$ for $i, j \in [m]$. If we can find i such that $k_i^- + 1 \leq k_i^*$, then

$$(1 - \mu_j)^{k_j^-} = (1 - \mu_i)^{k_i^-} \geq (1 - \mu_i)^{k_i^* - 1} \geq (1 - \mu_j)^{k_j^*}.$$

Hence $k_j^- \leq k_j^*$. On the other hand, we can always find $k_i^- + 1 \leq k_i^*$ since $\sum_{i=1}^m (k_i^- + 1) = K$. ■

In Algo. 4, we start with the lower bound \mathbf{k}^- of the optimal allocation. Since $\sum_{i=1}^m k_i^- = K - m$, we have $\sum_{i=1}^m |k_i^- - k_i^*| \leq \sum_{i=1}^m |k_i^- - k_i^*| = m$, which indicates Algo. 4 obtains the optimal budget allocation within m steps. We also provide an upper bound \mathbf{k}^+ in the following. The upper bound is also close to the optimal budget since $\sum_{i=1}^m |k_i^+ - k_i^*| \leq \sum_{i=1}^m |k_i^+ - k_i^*| = m$.

Lemma 3 (Allocation upper bound). *Let \mathbf{k}^* be the optimal budget allocation when the parameter of the communities is $\boldsymbol{\mu}$. Define $\mathbf{k}^+ = (k_1^+, \dots, k_m^+)$ where*

$$k_i^+ = \frac{K / \ln(1 - \mu_i)}{\sum_{j=1}^m 1 / \ln(1 - \mu_j)} + 1.$$

We have $k_i^ \leq k_i^+$.*

Proof. According to the definition of \mathbf{k}^+ , we have $(k_i^+ - 1) \ln(1 - \mu_i) = (k_j^+ - 1) \ln(1 - \mu_j)$ for $i, j \in [m]$. If we can find i such that $k_i^+ - 1 \geq k_i^*$, then

$$(1 - \mu_j)^{k_j^+ - 1} = (1 - \mu_i)^{k_i^+ - 1} \leq (1 - \mu_i)^{k_i^*} \leq (1 - \mu_j)^{k_j^* - 1}.$$

Hence $k_j^+ \geq k_j^*$. On the other hand, we can always find $k_i^+ - 1 \geq k_i^*$ since $\sum_{i=1}^m (k_i^+ - 1) = K$. ■

B Properties of Greedy Policy

In the following, we show some important properties of the greedy policy. We further extend the definition of reward with a non-decreasing function f as $R(\mathbf{k}, \phi) = f\left(\sum_{i=1}^m \left| \bigcup_{\tau=1}^{k_i} \{\phi(i, \tau)\} \right| \right)$.

B.1 Optimality of greedy policy

In this part, we prove that the greedy policy is the optimal policy for our adaptive community exploration problem. To prove the optimality, we first rewrite the proof sketch of Theorem 2, and then provide the supporting Lemma 4&5.

Theorem 2. *Greedy policy is the optimal policy for our adaptive exploration problem.*

Proof. Let $F_\pi(\psi, t)$ denote the expected marginal gain when we further explore communities for t steps with policy π starting from a partial realization ψ . We want to prove that for all ψ , t and π , $F_{\pi^g}(\psi, t) \geq F_\pi(\psi, t)$, where π^g is the greedy policy and π is an arbitrary policy. If so, we simply take $\psi = \emptyset$, and $F_{\pi^g}(\emptyset, t) \geq F_\pi(\emptyset, t)$ for every π and t exactly shows that π^g is optimal. We prove the above result by an induction on t . Recall that $c_i(\psi)$ is the number of distinct members met in community C_i under the partial realization ψ . Define $c(\psi) = \sum_i c_i(\psi)$ and $\Delta_{\psi, f} = f(c(\psi) + 1) - f(c(\psi))$.

For all ψ and π , we first claim that $F_{\pi^g}(\psi, 1) \geq F_\pi(\psi, 1)$ holds. Suppose that policy π chooses community C_i to explore based on the partial realization ψ . Since the exploration will return a new member with probability $1 - \mu_i c_i(\psi)$, the expected marginal gain $F_\pi(\psi, 1)$ is $(1 - \mu_i c_i(\psi)) [f(c(\psi) + 1) - f(c(\psi))]$. Note that the greedy policy π^g chooses community C_{i^*} to explore with $i^* \in \arg \max_j (1 - \mu_j c_j(\psi))$, and $\Delta_{\psi, f}$ does not depend on the policy. Hence, $F_{\pi^g}(\psi, 1) \geq F_\pi(\psi, 1)$.

Assume $F_{\pi^g}(\psi, t') \geq F_\pi(\psi, t')$ holds for all ψ , π , and $t' \leq t$. Our goal is to prove that $F_{\pi^g}(\psi, t + 1) \geq F_\pi(\psi, t + 1)$. Suppose that in the first step after ψ , policy π chooses C_i to explore based on partial realization ψ , and let $\pi(\psi) = (i, \tau)$. Define E_ψ as the event that the member $\Phi(i, \tau)$ is not met in partial realization ψ , for $\Phi \sim \psi$. In the following, we represent partial realization ψ equivalently as a relation $\{((i, \tau), \psi(i, \tau)) \mid (i, \tau) \in \text{dom}(\psi)\}$, so we could use $\psi \cup \{((i, \tau), \Phi(i, \tau))\}$ to represent the new partial realization extended from ψ by one step with (i, τ) added to the domain and $\Phi(i, \tau)$ as the member met for this exploration of C_i . Then we have

$$\begin{aligned} F_\pi(\psi, t + 1) &= \sum_{v \in C_i} \Pr(\Phi(i, \tau) = v) \mathbb{E}_\Phi[F_\pi(\psi, t + 1) \mid \Phi \sim \psi, \Phi(i, \tau) = v] \\ &= \sum_{v \in C_i} \mu_i \mathbb{E}_\Phi[F_\pi(\psi \cup \{((i, \tau), \Phi(i, \tau))\}, t) + f(c(\psi) + \mathbb{1}\{E_\psi\}) - f(c(\psi)) \mid \Phi \sim \psi, \Phi(i, \tau) = v] \\ &\leq \sum_{v \in C_i} \mu_i \mathbb{E}_\Phi[F_{\pi^g}(\psi \cup \{((i, \tau), \Phi(i, \tau))\}, t) \mid \Phi \sim \psi, \Phi(i, \tau) = v] + (1 - \mu_i c_i(\psi)) \Delta_{\psi, f}. \end{aligned}$$

The 2nd line above is derived directly from the definition of $F_\pi(\psi, t)$. The 3rd line is based on the induction hypothesis that $F_\pi(\psi', t) \leq F_{\pi^g}(\psi', t)$ holds for all ψ' . An important observation is that $F_{\pi^g}(\psi, t)$ has equal value for any partial realization ψ associated with the same status s since the

status is enough for the greedy policy to determine the choice of next community. Formally, we define $F_g(\mathbf{s}, t) = F_{\pi^g}(\psi, t)$ for any partial realization that satisfies $\mathbf{s} = (1 - c_1(\psi)/d_1, \dots, 1 - c_m(\psi)/d_m)$. Let C_{i^*} denote the community chosen by policy π^g under realization ψ , i.e., $i^* \in \arg \max_{i \in [m]} 1 - c_i(\psi)\mu_i$. Let \mathbf{I}_i be the m -dimensional unit vector with 1 in the i -th entry and 0 in all other entries. Therefore,

$$\begin{aligned} F_{\pi}(\psi, t+1) &\leq c_i(\psi) \cdot \mu_i F_g(\mathbf{s}, t) + (d_i - c_i(\psi)) \cdot \mu_i F_g(\mathbf{s} - \mu_i \mathbf{I}_i, t) + (1 - \mu_i c_i(\psi)) \Delta_{\psi, f} \\ &\leq \mu_{i^*} c_{i^*}(\psi) F_g(\mathbf{s}, t) + (1 - \mu_{i^*} c_{i^*}(\psi)) F_g(\mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}, t) + (1 - \mu_{i^*} c_{i^*}(\psi)) \Delta_{\psi, f} \\ &= F_g(\mathbf{s}, t+1) = F_{\pi^g}(\psi, t+1). \end{aligned} \quad \begin{array}{l} \text{(Lemma 5)} \\ \text{(Lemma 4)} \end{array}$$

The key is to prove the correctness of Line 2 in above equation. It indicates that if we choose a sub-optimal community at first, and then we switch back to the greedy policy, the expected reward would be smaller. The proof is nontrivial and relies on a careful analysis based on the stochastic transitions among status vectors. The above result completes the induction step for $t+1$. Thus the theorem holds. \blacksquare

Lemma 4. Let $\mathbf{s} = (s_1, \dots, s_m)$ be a status where each entry $s_i \in [0, 1]$. We have

$$F_g(\mathbf{s}, t+1) = (1 - s_{i^*}) F_g(\mathbf{s}, t) + s_{i^*} F_g(\mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}, t) + s_{i^*} (f(c(\psi) + 1) - f(c(\psi))),$$

where $i^* = \arg \max_{i \in [m]} s_i$. Here ψ is any partial realization corresponding to status \mathbf{s} .

Proof. For any partial realization ψ associated with status \mathbf{s} , π^g would choose community i^* . With probability $\mu_{i^*} c_{i^*}(\psi) = 1 - s_{i^*}$, we will obtain a member that is already met. If so, the communities stay at the same status. Hence, with probability $1 - s_{i^*}$, the expected extra reward is $F_g(\mathbf{s}, t)$ after the first step exploration. With probability $1 - \mu_{i^*} c_{i^*}(\psi) = s_{i^*}$, we will obtain an unseen member in C_{i^*} . The communities will transit to next status $\mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}$. Therefore, with probability s_{i^*} , the expected extra reward is $F_g(\mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}, t) + f(c(\psi) + 1) - f(c(\psi))$ after the first step exploration. \blacksquare

Lemma 5. Let $\mathbf{s} = (s_1, \dots, s_m)$ be a status where each entry $s_i \in [0, 1]$ and ψ be any partial realization corresponding to \mathbf{s} . We have

$$\begin{aligned} &(1 - s_i) F_g(\mathbf{s}, t) + s_i F_g(\mathbf{s} - \mu_i \mathbf{I}_i, t) + s_i \Delta_c \\ &\leq (1 - s_{i^*}) F_g(\mathbf{s}, t) + s_{i^*} F_g(\mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}, t) + s_{i^*} \Delta_c, \end{aligned} \quad (9)$$

where $i^* \in \arg \max_{i \in [m]} s_i$, $s_i < s_{i^*}$ and $\Delta_c = f(c(\psi) + 1) - f(c(\psi))$.

Proof. Let $A(\mathbf{s}, i, t)$ denote the first line of Eq. (9), i.e.,

$$A(\mathbf{s}, i, t) = (1 - s_i) F_g(\mathbf{s}, t) + s_i F_g(\mathbf{s} - \mu_i \mathbf{I}_i, t) + s_i \Delta_c.$$

Note that $A(\mathbf{s}, i, t)$ is the expected reward of the following adaptive process.

1. At the first step, choose an arbitrary community C_i (different from C_{i^*}) to explore.
2. From the second step to the $(t+1)$ -th step, explore communities with the greedy policy π^g .

Similarly, $A(\mathbf{s}, i^*, t)$ is the expected reward of the $t+1$ step community exploration via the greedy policy, i.e., $A(\mathbf{s}, i^*, t) = F_g(\mathbf{s}, t+1)$. Eq. (9) can be written as $A(\mathbf{s}, i, t) \leq F_g(\mathbf{s}, t+1)$. We prove this inequality by induction. When $t = 0$, we have $A(\mathbf{s}, i, t) = s_i \Delta_c$, and $A(\mathbf{s}, i^*, t) = s_{i^*} \Delta_c$. Hence, $A(\mathbf{s}, i, t) \leq A(\mathbf{s}, i^*, t) = F_g(\mathbf{s}, t+1)$ when $t = 0$. Assume $A(\mathbf{s}, i, t') \leq F_g(\mathbf{s}, t'+1)$ holds for any $0 \leq t' \leq t$, and any status \mathbf{s} . Our goal is to prove that $A(\mathbf{s}, i, t+1) \leq A(\mathbf{s}, i^*, t+1) = F_g(\mathbf{s}, t+2)$. We expand $A(\mathbf{s}, i, t+1)$ as follows.

$$\begin{aligned} A(\mathbf{s}, i, t+1) &= (1 - s_i) F_g(\mathbf{s}, t+1) + s_i F_g(\mathbf{s} - \mu_i \mathbf{I}_i, t+1) + s_i \Delta_c \\ &= (1 - s_i) ((1 - s_{i^*}) F_g(\mathbf{s}, t) + s_{i^*} F_g(\mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}, t) + s_{i^*} \Delta_c) \\ &\quad + s_i ((1 - s_{i^*}) F_g(\mathbf{s} - \mu_i \mathbf{I}_i, t) + s_{i^*} F_g(\mathbf{s} - \mu_i \mathbf{I}_i - \mu_{i^*} \mathbf{I}_{i^*}, t) + s_{i^*} \Delta_{c+1}) \\ &\quad + s_i \Delta_c. \end{aligned}$$

Here $\Delta_{c+1} = f(c(\psi) + 2) - f(c(\psi) + 1)$. Above expansion of $A(i, t + 1)$ is based on Lemma 4. We expand $A(s, i^*, t + 1)$ as follows.

$$\begin{aligned}
A(s, i^*, t + 1) &= (1 - s_{i^*})F_g(s, t + 1) + s_{i^*}F_g(s - \mu_{i^*}\mathbf{I}_{i^*}, t + 1) + s_{i^*}\Delta_c \\
&\geq (1 - s_{i^*})((1 - s_i)F_g(s, t) + s_iF_g(s - \mu_i\mathbf{I}_i, t) + s_i\Delta_c) \\
&\quad \text{(assumption } A(s, i, t) \leq F_g(s, t + 1)) \\
&\quad + s_{i^*}((1 - s_i)F_g(s - \mu_{i^*}\mathbf{I}_{i^*}, t) + s_iF_g(s - \mu_{i^*}\mathbf{I}_{i^*} - \mu_i\mathbf{I}_i, t) + s_i\Delta_{c+1}) \\
&\quad \text{(assumption } A(s - \mu_{i^*}\mathbf{I}_{i^*}, i, t) \leq F_g(s - \mu_{i^*}\mathbf{I}_{i^*}, t + 1)) \\
&\quad + s_{i^*}\Delta_c \\
&= A(i, t + 1).
\end{aligned}$$

This completes the proof. \blacksquare

Remarks. During the rebuttal of this paper, we realized that Bubeck et al. [5] applied similar inductive reasoning techniques to prove the optimality of the greedy policy for their optimal discovery problem (Lemma 2 of [5]). To quantitatively measure how good is the greedy policy, we also give a formula to show the exact difference between $A(s, i, t)$ and $A(s, i^*, t)$ in Sec. B.3.

B.2 Computation of expected reward

Lemma 4 indicates $r_{\pi^g}(\mu)$ can be computed in a recursive way. However, the recursive method has time complexity $O(2^K)$. It is impractical when K is large. In the following we show that the expected reward of policy π^g can be computed in polynomial time.

B.2.1 Transition probability list of greedy policy

Assume we explore the communities via the greedy policy when the communities already have partial realization ψ . Define $s_{i,0} = 1 - \mu_i c_i(\psi)$ and $\mathbf{s}_0 = (s_{1,0}, \dots, s_{m,0})$. The greedy policy will choose community i_0^* to explore, where $i_0^* \in \arg \max_i s_{i,0}$. After one step exploration, the communities stay at the same status \mathbf{s}_0 with probability $q_0 := 1 - s_{i_0^*}$. The communities transit to next status $\mathbf{s}_1 := \mathbf{s}_0 - \mu_{i_0^*}\mathbf{I}_{i_0^*}$ with probability $p_0 := s_{i_0^*}$. We recursively define \mathbf{s}_{t+1} as $\mathbf{s}_t - \mu_{i_t^*}\mathbf{I}_{i_t^*}$, where $i_t^* \in \arg \max_i s_{i,t}$. We call $p_t := \max_i s_{i,t}$ the *transition probability* and $q_t := 1 - p_t$ the *loop probability*. Each time the communities transit to next status, a new member will be met. During the exploration, the number of different statuses the communities can stay is at most $1 + \sum_i d_i - c_i(\psi)$ since there are $D := \sum_i d_i - c_i(\psi)$ unseen members in total. Based on above discussion, we define a *transition probability list* $\mathcal{P}(\pi^g, \psi) := (p_0, \dots, p_D)$, where $p_D \equiv 0$. The list $\mathcal{P}(\pi^g, \psi)$ is unique for any initial partial realization ψ . Fig. 1 gives an example to demonstrate statuses and the list $\mathcal{P}(\pi^g, \psi)$.

Corollary 1. Let ψ be any partial realization corresponding to the status $\mathbf{s} = (s_1, \dots, s_m)$. The number of unseen members $\sum_i d_i - c_i(\psi)$ is denoted as D . The probability list $\mathcal{P}(\pi^g, \psi) = (p_0, \dots, p_D)$ can be obtained by sorting $\cup_{i=1}^m \{s_i, s_i - \mu_i, \dots, \mu_i\} \cup \{0\}$ in descending order.

Corollary 1 is an important observation based on the definition of transition probability list.

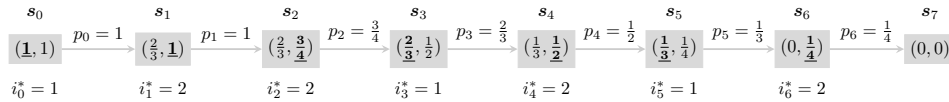


Figure 1: Illustration with $\mathbf{d} = (3, 4)$ and empty partial realization. The initial status is $(1, 1)$. The list $\mathcal{P}(\pi^g, \emptyset) = (1, 1, 3/4, 2/3, 1/2, 1/3, 1/4, 0)$.

B.2.2 Compute the expected reward efficiently

Lemma 6. Let ψ be a partial realization and \mathbf{s}_0 be the corresponding status. The number of unseen members is denoted as $D = \sum_i d_i - c_i(\psi)$. The transition probability list is $\mathcal{P}(\pi^g, \psi) = (p_0, \dots, p_D)$. Then

$$F_{\pi^g}(\psi, t) = F_g(\mathbf{s}_0, t) = \sum_{j=0}^{\min\{t, D\}} (f(j + c(\psi)) - f(c(\psi))) \times \left(\prod_{l=0}^{j-1} p_l \right) \times \left(\sum_{I \in \mathcal{I}(j, t-j)} \prod_{l \in I} q_l \right),$$

where $q_l = 1 - p_l$ and $\mathcal{I}(j, t - j)$ consists of subsets of multi-set $\{0, \dots, j\}^{t-j}$ with fixed size $t - j$.

Proof. When the communities ends at status s_j , we meet j distinct members. Let $\Pr(s_j \square)$ be the probability for this event. We can the *transition step* as the communities transit to a new status, and the *loop step* as the communities stay at the same status. When the communities ends at status s_j , we have j transition steps and $t - j$ loop steps. The communities takes loops at statuses $\{s_0, \dots, s_j\}$. Hence,

$$\Pr(s_j \square) = \sum_{I \in \mathcal{I}(j, t-j)} \prod_{l=0}^{j-1} p_l \cdot \prod_{l \in I} (1 - p_l) = \prod_{l=0}^{j-1} p_l \times \sum_{I \in \mathcal{I}(j, t-j)} \prod_{l \in I} q_l.$$

The reward $F_{\pi^g}(\psi, t) = \sum_{j=1}^{\min\{t, D\}} (f(j + c(\psi)) - f(c(\psi))) \times \Pr(s_j \square)$. ■

For later analysis, we define the *loop probability*

$$L(\{q_0, \dots, q_j\}, t) := \sum_{I \in \mathcal{I}(j, t)} \prod_{l \in I} q_l$$

since $\sum_{I \in \mathcal{I}(j, t)} \prod_{l \in I} q_l$ is just a function of $\{q_0, \dots, q_j\}$ and t ($t \geq 1$). Actually, $L(\{q_0, \dots, q_j\}, t)$ aggregates the product of all possible t elements in $\{q_0, \dots, q_j\}$. Note that each element in $\{q_0, \dots, q_j\}$ can be chosen multiple times. W.l.o.g, we define $L(\{q_0, \dots, q_j\}, t) = 1$ and $\prod_{l=0}^{t-1} p_l = 1$ when $t = 0$. Based on the definition, we can write $L(\{q_0, \dots, q_j\}, t)$ in a recursive way as follows.

$$L(\{q_0, \dots, q_j\}, t) = \sum_{s=0}^t q_a^s L(\{q_0, \dots, q_j\} \setminus \{q_a\}, t - s). \quad (10)$$

Here $a \in \{0, \dots, j\}$. According to Eq. 10, the probability $\sum_{I \in \mathcal{I}(j, t-j)} \prod_{l \in I} q_l$ can be computed in $O((t - j)j^2)$ via *dynamic programming*. Hence $r_{\pi^g}(\mu) = F_g((1, \dots, 1), K)$ can be computed in $O(K \min\{K, D\}^2)$ according to Lemma 6.

B.3 Reward gap between optimal policy and sub-optimal policy

Recall that $A(s, i, t)$ is the expected reward of the following adaptive process.

1. At the first step, choose an arbitrary community C_i (different from C_{i^*}) to explore.
2. From the second step to the $(t + 1)$ -th step, explore communities with the greedy policy π^g .

Here s is the initial status of the communities. Lemma 5 only proves that $A(s, i, t) \leq F_g(s, t + 1)$. In the following, we aim to answer the following question:

- How much is $F_g(s, t + 1)$ larger than $A(s, i, t)$?

B.3.1 Analysis of loop probability

The following two corollaries show the basic properties of the *loop probability*.

Corollary 2. For a transition probability list $\mathcal{P}(\pi^g, \psi) = (p_0, \dots, p_D)$, we have

$$\sum_{j=0}^M p_0 \times \dots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t - j) = 1,$$

where $q_j = 1 - p_j$ and $M = \min\{t, D\}$.

Corollary 2 says the probabilities that the communities ends at status $\{s_0, \dots, s_D\}$ sums up to 1.

Corollary 3. For a transition probability list $\mathcal{P}(\pi^g, \psi) = (p_0, \dots, p_D)$ and $a, b \in \{0, \dots, j\}$ ($j \leq D, t \geq 1$), we have

$$L(\{q_0, \dots, q_j\} \setminus \{q_a\}, t) - L(\{q_0, \dots, q_j\} \setminus \{q_b\}, t) = (q_b - q_a)L(\{q_0, \dots, q_j\}, t - 1),$$

where $D = \sum_i d_i - c_i(\psi)$ and $q_j = 1 - p_j$.

Proof. We prove the corollary according to Eq. (10).

$$\begin{aligned}
& L(\{q_0, \dots, q_j\} \setminus \{q_a\}, t) - L(\{q_0, \dots, q_j\} \setminus \{q_b\}, t) \\
&= \sum_{s=0}^t (q_b^s - q_a^s) L(\{q_0, \dots, q_j\} \setminus \{q_a, q_b\}, t-s) \quad (\text{by Eq. (10)}) \\
&= \sum_{s=0}^{t-1} (q_b^{s+1} - q_a^{s+1}) L(\{q_0, \dots, q_j\} \setminus \{q_a, q_b\}, t-s-1) \quad (\text{replace } s-1 \text{ as } s') \\
&= (q_b - q_a) \sum_{s=0}^{t-1} \sum_{m=0}^s q_b^{s-m} q_a^m L(\{q_0, \dots, q_j\} \setminus \{q_a, q_b\}, t-1-s) \quad (\text{sum of geometric sequence}) \\
&= (q_b - q_a) L(\{q_0, \dots, q_j\}, t-1). \quad (\text{by definition or expanding Eq. (10)})
\end{aligned}$$

This completes the proof. \blacksquare

B.3.2 Pseudo reward

Lemma 7. For a transition probability list $\mathcal{P}(\pi^g, \psi) = (p_0, \dots, p_D)$ and a non-decreasing function $f(x)$, a pseudo reward $R(k)$ is defined as

$$\begin{aligned}
R(k) &= q_k \sum_{j=0}^M f(j) \times p_0 \times \dots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t-j) \\
&\quad + p_k \sum_{j=0}^{k-1} f(j+1) \times p_0 \times \dots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t-j) \\
&\quad + \sum_{j=k}^{M'} f(j+1) \times p_0 \times \dots \times p_j \times L(\{q_0, \dots, q_{j+1}\} \setminus \{q_k\}, t-j),
\end{aligned}$$

where $M = \min\{D, t\}$ and $M' = \{D-1, t\}$. We claim that for $0 \leq k \leq M-1$,

$$R(k) - R(k+1) = (p_k - p_{k+1}) \left(\sum_{j=0}^k (f(j+1) - f(j)) p_0 \times \dots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t-j) \right).$$

Proof. We expand $R(k) - R(k+1)$ as follows using the definition.

$$\begin{aligned}
& R(k) - R(k+1) \\
&= -(p_k - p_{k+1}) \sum_{j=0}^M f(j) \times p_0 \times \dots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t-j) \\
&\quad + (p_k - p_{k+1}) \sum_{j=0}^{k-1} f(j+1) \times p_0 \times \dots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t-j) \\
&\quad + f(k+1) \times p_0 \times \dots \times p_{k-1} \times p_k \times L(\{q_0, \dots, q_{k+1}\} \setminus \{q_k\}, t-k) \quad (\text{from } R(k)) \\
&\quad - f(k+1) \times p_0 \times \dots \times p_{k-1} \times p_{k+1} \times L(\{q_0, \dots, q_k\}, t-k) \quad (\text{from } R(k+1)) \\
&\quad + \sum_{j=k+1}^{M'} f(j+1) \times p_0 \times \dots \times p_j \times (L(\{q_0, \dots, q_{j+1}\} \setminus \{q_k\}, t-j) \\
&\quad \quad \quad \underbrace{- L(\{q_0, \dots, q_{j+1}\} \setminus \{q_{k+1}\}, t-j))}_{(p_k - p_{k+1}) \sum_{j=k+1}^{M-1} f(j+1) \times p_0 \times \dots \times p_j \times L(\{q_0, \dots, q_{j+1}\}, t-j-1)}.
\end{aligned}$$

The last line of above equation can be rewritten with the Corollary 3. The summation from $j = k+1$ to $j = M-1$ in the last line cancels out with the second line when $j = k+2$ to $j = M$. The

summation from $j = 0$ to $j = k$ in the second line can be combined with the third line. We continue the computation of $R(k) - R(k + 1)$ by rearranging its expansion.

$$\begin{aligned}
& R(k) - R(k + 1) \\
&= -(p_k - p_{k+1})f(k + 1) \times p_0 \times \cdots \times p_k \times L(\{q_0, \dots, q_{k+1}\}, t - k - 1) \\
&\quad - (p_k - p_{k+1})f(k + 1) \times p_0 \times \cdots \times p_{k-1} \times L(\{q_0, \dots, q_k\}, t - k) \\
&\quad + (p_k - p_{k+1}) \sum_{j=0}^k (f(j + 1) - f(j)) \times p_0 \times \cdots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t - j) \\
&\quad + f(k + 1) \times p_0 \times \cdots \times p_{k-1} \times p_k \times L(\{q_0, \dots, q_{k+1}\} \setminus \{q_k\}, t - k) \\
&\quad - f(k + 1) \times p_0 \times \cdots \times p_{k-1} \times p_{k+1} \times L(\{q_0, \dots, q_k\}, t - k).
\end{aligned}$$

Define Δ_k as the sum of the 2nd, 3rd, 5th, 6th line in above equation. We have

$$\begin{aligned}
& R(k) - R(k + 1) \\
&= \left. \begin{aligned} & -(p_k - p_{k+1}) \times f(k + 1) \times p_0 \times \cdots \times p_{k-1} \times L(\{q_0, \dots, q_k\}, t - k) \\ & -(p_k - p_{k+1}) \times f(k + 1) \times p_0 \times \cdots \times p_k \times L(\{q_0, \dots, q_{k+1}\}, t - k - 1) \\ & + f(k + 1) \times p_0 \times \cdots \times p_{k-1} \times p_k \times L(\{q_0, \dots, q_{k+1}\} \setminus \{q_k\}, t - k) \\ & - f(k + 1) \times p_0 \times \cdots \times p_{k-1} \times p_{k+1} \times L(\{q_0, \dots, q_k\}, t - k) \end{aligned} \right\} \triangleq \Delta_k \\
&\quad + (p_k - p_{k+1}) \sum_{j=0}^k (f(j + 1) - f(j)) \times p_0 \times \cdots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t - j).
\end{aligned}$$

We rewrite Δ_k as follows.

$$\begin{aligned}
\Delta_k / f(k + 1) &= p_0 \times \cdots \times p_{k-1} \times p_k \times L(\{q_0, \dots, q_{k+1}\} \setminus \{q_k\}, t - k) \\
&\quad - p_0 \times \cdots \times p_{k-1} \times p_k \times L(\{q_0, \dots, q_k\}, t - k) \Big\} \text{cancel each other} \\
&\quad + p_0 \times \cdots \times p_{k-1} \times p_k \times L(\{q_0, \dots, q_k\}, t - k) \\
&\quad - (p_k - p_{k+1}) \times p_0 \times \cdots \times p_{k-1} \times L(\{q_0, \dots, q_k\}, t - k) \\
&\quad - (p_k - p_{k+1}) \times p_0 \times \cdots \times p_k \times L(\{q_0, \dots, q_{k+1}\}, t - k - 1) \\
&\quad - p_0 \times \cdots \times p_{k-1} \times p_{k+1} \times L(\{q_0, \dots, q_k\}, t - k).
\end{aligned}$$

According to Corollary 3, the first line and the second line of above equation equals to $(p_k - p_{k+1}) \times p_0 \times \cdots \times p_k \times L(\{q_0, \dots, q_{k+1}\}, t - k - 1)$, which cancels out with the fifth line. Hence, we have

$$\begin{aligned}
\Delta_k / f(k + 1) &= p_0 \times \cdots \times p_{k-1} \times p_k \times L(\{q_0, \dots, q_k\}, t - k) \\
&\quad - (p_k - p_{k+1}) \times p_0 \times \cdots \times p_{k-1} \times L(\{q_0, \dots, q_k\}, t - k) \\
&\quad - p_0 \times \cdots \times p_{k-1} \times p_{k+1} \times L(\{q_0, \dots, q_k\}, t - k) \\
&= 0.
\end{aligned}$$

With above result of $\Delta_k = 0$, we prove that

$$\begin{aligned}
& R(k) - R(k + 1) \\
&= (p_k - p_{k+1}) \left(\sum_{j=0}^k (f(j + 1) - f(j)) p_0 \times \cdots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t - j) \right) \geq 0.
\end{aligned}$$

This completes the proof. ■

B.3.3 Reward gap

Let ψ, ψ', ψ'' be any partial realization corresponding to the status $\mathbf{s}, \mathbf{s} - \mu_{i^*} \mathbf{I}_{i^*}, \mathbf{s} - \mu_i \mathbf{I}_i$ respectively. Define $\mathcal{P}(\pi^g, \psi) = (p_0, \dots, p_D)$, where $D = \sum_i d_i - c_i(\psi)$. Recalling Corollary 1, we know that (p_0, \dots, p_{D-1}) can be obtained by sorting $\cup_{i=1}^m \{s_i, s_i - \mu_i, \dots, \mu_i\}$. Assume the first time s_i appear in (p_0, \dots, p_D) is the k -th entry, i.e., $k = \min\{k' : 0 \leq k' \leq D, p_{k'} = s_i\}$. According to Corollary 1, we have the following.

$$\begin{aligned}
\mathcal{P}(\pi^g, \psi') &= (p_1, \dots, p_D), \\
\mathcal{P}(\pi^g, \psi'') &= (p_0, \dots, p_{k-1}, p_{k+1}, \dots, p_D).
\end{aligned}$$

Note that $p_0 = s_{i^*}$ and $p_k = s_i$. Let $M = \min\{D, t\}$, $M' = \min\{D - 1, t\}$, and $f'(j) = f(j + c(\psi)) - f(c(\psi))$. The second line of Eq. (9) is

$$\begin{aligned} R_1 &= q_0 \sum_{j=0}^M f'(j) \times p_0 \times \cdots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t - j) && ((1 - s_{i^*})F_g(\mathbf{s}, t)) \\ &+ p_0 \sum_{j=0}^{M'} f'(j + 1) \times p_1 \times \cdots \times p_j \times L(\{q_1, \dots, q_{j+1}\}, t - j). && (s_{i^*}F_g(\mathbf{s} - \mu_{i^*}\mathbf{I}_{i^*}, t) + s_{i^*}\Delta_c) \end{aligned}$$

In fact, $R_1 = F_g(\mathbf{s}, t + 1)$ based on Lemma 4. The first line of Eq. (9) is

$$\begin{aligned} R_2 &= q_k \sum_{j=0}^M f'(j) \times p_0 \times \cdots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t - j) \\ &+ p_k \sum_{j=0}^{k-1} f'(j + 1) \times p_0 \times \cdots \times p_{j-1} \times L(\{q_0, \dots, q_j\}, t - j) \\ &+ \sum_{j=k}^{M'} f'(j + 1) \times p_0 \times \cdots \times p_j \times L(\{q_0, \dots, q_{j+1}\} \setminus \{q_k\}, t - j). \end{aligned}$$

Our goal is to measure the gap $R_1 - R_2$. Let $\text{Prob}_{\mathbf{s}, t}(i)$ be the probability we can meet i distinct members if we explore communities (whose initial status is \mathbf{s}) with greedy policy for t steps. According to Lemma 7, we have

$$\begin{aligned} F_g(\mathbf{s}, t + 1) - A(\mathbf{s}, i, t) &= \sum_{j=0}^{k-1} (R(j) - R(j + 1)) \\ &= \sum_{j=0}^{k-1} (p_j - p_{j+1}) \left(\sum_{o=0}^j (f'(o + 1) - f'(o)) \text{Prob}_{\mathbf{s}, t}(o) \right) \\ &= \sum_{o=0}^{k-1} (f'(o + 1) - f'(o)) \text{Prob}_{\mathbf{s}, t}(o) \left(\sum_{j=o}^{k-1} p_j - p_{j+1} \right) \\ &= \sum_{j=0}^{k-1} (f'(j + 1) - f'(j)) (p_j - p_k) \text{Prob}_{\mathbf{s}, t}(j) \end{aligned}$$

When the reward equals to the number of distinct members, we have

$$F_g(\mathbf{s}, t + 1) - A(\mathbf{s}, i, t) = \sum_{j=0}^{k-1} (p_j - p_k) \text{Prob}_{\mathbf{s}, t}(j).$$

Besides, the gap $F_g(\mathbf{s}, t + 1) - A(\mathbf{s}, i, t)$ increases as k increases, which means the worse choice we have at first, the larger reward gap we have at end.

C Basics of online learning problems

C.1 Set size estimation by collision counting

Suppose we have a set $C_i = \{u_1, \dots, u_{d_i}\}$ whose population d_i is unknown. Let u, v be two elements selected with replacement from C_i , and $Y_{u,v}$ denote a random variable that takes value 1 if $u = v$ (a *collision*) and 0 otherwise. The expectation of $Y_{u,v}$ equals to $\frac{1}{d_i}$, i.e., $\mathbb{E}[Y_{u,v}] = \frac{1}{d_i}$. Assume we sample k_i elements *with replacement* uniformly at random from set C_i . Let S_i be the set of samples. With the sample S_i , we compute the estimator for d_i as

$$\hat{d}_i = \frac{k_i(k_i - 1)}{2X_i},$$

here $X_i = \sum_{u \in \mathcal{S}_i, v \in \mathcal{S}_i \setminus \{u\}} Y_{u,v}$ is the number of collisions in \mathcal{S}_i . According to the Jensen's inequality³, we have $d_i \leq \mathbb{E}[\hat{d}_i]$, i.e., \hat{d}_i is a biased estimator. The estimator is invalid when $X_i = 0$. Since the equality only occurs when $\text{Var}[X_i] = 0$, which is not the case Here. We have $d_i < \mathbb{E}[\hat{d}_i]$.

Independence. Let $\mathcal{S}_i = \{v_1, \dots, v_{k_i}\}$. For the two random variable Y_{v_x, v_y} ($1 \leq x < y \leq k_i$) and $Y_{v_{x'}, v_{y'}}$ ($1 \leq x' < y' \leq k_i$), we consider three difference cases.

1. There are $\binom{k_i}{2}$ occurrences when $x = x', y = y'$. Here $\mathbb{E}[Y_{v_x, v_y} Y_{v_{x'}, v_{y'}}] = 1/d_i$.
2. There are $6\binom{k_i}{3}$ occurrences when $x = x', y \neq y'$ or $x \neq x', y = y'$. $\mathbb{E}[Y_{v_x, v_y} Y_{v_{x'}, v_{y'}}] = 1/d_i^2$.
3. There are $6\binom{k_i}{4}$ occurrences when $x \neq x', y \neq y'$. Here $\mathbb{E}[Y_{v_x, v_y} Y_{v_{x'}, v_{y'}}] = 1/d_i^2$.

We say that pairs (v_x, v_y) and $(v_{x'}, v_{y'})$ are different if $x \neq x'$ or $y \neq y'$. When (v_x, v_y) and $(v_{x'}, v_{y'})$ are different, we have $\mathbb{E}[Y_{v_x, v_y} Y_{v_{x'}, v_{y'}}] = \mathbb{E}[Y_{v_x, v_y}] \mathbb{E}[Y_{v_{x'}, v_{y'}}] = 1/d_i^2$. Above discussion indicates that the $\binom{k_i}{2}$ pairs of random variables obtained from \mathcal{S}_i are 2-wise independent.

Variance. We compute the variance $\text{Var}[X_i] = \mathbb{E}[X_i^2] - \mathbb{E}^2[X_i]$ in the following.

$$\begin{aligned} \text{Var}[X_i] &= \frac{k_i(k_i-1)}{2d_i} + \frac{k_i(k_i-1)(k_i-2)}{d_i^2} + \frac{k_i(k_i-1)(k_i-2)(k_i-3)}{4d_i^2} - \frac{k_i^2(k_i-1)^2}{4d_i^2} \\ &= \binom{k_i}{2} \frac{1}{d_i} \left(1 - \frac{1}{d_i}\right) = \binom{k_i}{2} \text{Var}[Y_{u,v}]. \end{aligned}$$

Collision Since the estimator is based on the collision counting, we need to ensure that $X_i > 0$ with high probability. Let B_{k_i} denote the event that the k_i samples $\{v_1, \dots, v_{k_i}\}$ are distinct. We have

$$\begin{aligned} \Pr\{B_k\} &= 1 \cdot \left(1 - \frac{1}{d_i}\right) \left(1 - \frac{2}{d_i}\right) \dots \left(1 - \frac{k_i-1}{d_i}\right) \leq e^{-1/d_i} e^{-2/d_i} \dots e^{-(k_i-1)/d_i} \\ &= e^{-\sum_{j=1}^{k_i-1} j/d_i} = e^{-k_i(k_i-1)/2d_i}. \end{aligned}$$

To ensure that $X_i > 0$ with probability no less than $1 - \delta$, we have

$$k_i \geq \left(1 + \sqrt{8d_i \ln \frac{1}{\delta} + 1}\right) / 2.$$

C.2 Concentration bound for variables with local dependence

Note that the pairs $Y_{u,v}$ and $Y_{u',v'}$ are not mutually independent. Actually, their dependence can be described with a *dependence graph* [9, 15]. The Chernoff-Hoeffding bound in [14] can not be used directly for our estimator of μ_i . In the following, we present a concentration bound that is applicable to our problem.

Definition 1 (U-statistics). Let ξ_1, \dots, ξ_n be independent random variables, and let

$$X := \sum_{1 \leq i_1 \leq \dots \leq i_d} f_{i_1, \dots, i_d}(\xi_{i_1}, \dots, \xi_{i_d}).$$

Lemma 8 (Chapter 3.2 [9]). If $a \leq f_{i_1, \dots, i_d}(\xi_{i_1}, \dots, \xi_{i_d}) \leq b$ for every i_1, \dots, i_d for some reals $a \leq b$, we have

$$\Pr \left\{ |X - \mathbb{E}[X]| \geq \epsilon \binom{n}{d} \right\} \leq 2 \exp \left(\frac{-2 \lfloor n/d \rfloor \epsilon^2}{(b-a)^2} \right).$$

In our problem, if we get k_i samples from set C_i , then the number of collisions satisfies

$$\Pr \left\{ |X_i - \mathbb{E}[X_i]| \geq \epsilon \binom{k_i}{2} \right\} \leq 2 \exp \left(-2 \lfloor k_i/2 \rfloor \epsilon^2 \right).$$

Above inequality indicates that the actual number of independent pairs is $\lfloor k_i/2 \rfloor$ when using collisions in k_i samples to estimate μ_i .

³If X is a random variable, and φ is a convex function, then $\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$.

D Regret Analysis for Non-Adaptive Problem

D.1 Supporting Corollaries

Corollary 4. For action k with $\sum_{i=1}^m k_i = K$ and $k_i \geq 1$, we have $\sum_{i=1}^m \binom{k_i}{2} \leq \binom{K-m+1}{2}$.

Proof. We prove the corollary by simple calculation.

$$\begin{aligned} \sum_{i=1}^m \binom{k_i}{2} - \binom{K-m+1}{2} &= \frac{1}{2} \left(\sum_{i=1}^m k_i(k_i - 1) - \left(1 + \sum_{i=1}^m (k_i - 1) \right) \left(\sum_{i=1}^m (k_i - 1) \right) \right) \\ &= \frac{1}{2} \left(\sum_{i=1}^m (k_i - 1)^2 - \left(\sum_{i=1}^m (k_i - 1) \right)^2 \right) \leq 0. \quad \blacksquare \end{aligned}$$

D.2 Basics

To compare with the CUCB algorithm introduced in [20] for general CMAB problem, we propose an revised Algo. 3 that is consistent with the CUCB algorithm in [20]. We revise the Line 6-8 in Algo. 3 as follows.

$$\begin{aligned} \text{Line 6: } & \text{For } i \in [m], T_i \leftarrow T_i + \mathbb{1}\{|\mathcal{S}_i| > 1\}, \\ \text{Line 7: } & \text{For } i \in [m] \text{ and } |\mathcal{S}_i| > 1, X_{i,t} \leftarrow \sum_{x=1}^{|\mathcal{S}_i|/2} \mathbb{1}\{u_{2x-1} = u_{2x}\} / \lfloor |\mathcal{S}_i|/2 \rfloor, \quad (11) \\ \text{Line 8: } & \text{For } i \in [m] \text{ and } |\mathcal{S}_i| > 1, \hat{\mu}_i \leftarrow \hat{\mu}_i + (X_{i,t} - \hat{\mu}_i) / T_i. \end{aligned}$$

Note that $\hat{\mu}_i$ in Eq. (11) is also an unbiased estimator of μ_i . Then we can obtain the regret bound of the revised Algo. 3 by applying the Theorem 4 in the extended version of [20] directly.

$$\text{Reg}_\mu(T) \leq \sum_{i=1}^m \frac{48 \binom{K-m+1}{2} m \ln T}{\Delta_{\min}^i} + 2 \binom{K-m+1}{2} m + \frac{\pi^2}{3} \cdot m \cdot \Delta_{\max}.$$

We add superscript r to differentiate the corresponding random variables in the revised Algo. 3 from the original ones. E.g., $T_{i,t}^r$ is the value of T_i^r in the revised Algo. 3 at the end of round t . Recall that $K' = K - m + 1$, which is the maximum exploration times for a community in each round.

D.3 Proof framework

We first introduce a definition which describes the event that $\hat{\mu}_{i,t-1}$ ($\hat{\mu}_{i,t-1}^r$) is accurate at the beginning of round t .

Definition 2. We say that the sampling is nice at the beginning of round t if for every community $i \in [m]$, $|\hat{\mu}_{i,t-1} - \mu_i| \leq \rho_{i,t}$ (resp. $|\hat{\mu}_{i,t-1}^r - \mu_i| \leq \rho_{i,t}^r$), where $\rho_{i,t} = 2\sqrt{\frac{3 \ln t}{2T_{i,t-1}^r}}$ (resp. $\rho_{i,t}^r = 2\sqrt{\frac{3 \ln t}{2T_{i,t-1}^r}}$) in round t . Let \mathcal{N}_t (resp. \mathcal{N}_t^r) be such event.

Lemma 9. For each round $t \geq 1$, $\Pr \{\neg \mathcal{N}_t\} \leq 2m \lfloor K'/2 \rfloor t^{-2}$ (resp. $\Pr \{\neg \mathcal{N}_t^r\} \leq 2mt^{-2}$).

Proof. For each round $t \geq 1$, we have

$$\begin{aligned}
\Pr \{\neg \mathcal{N}_t\} &= \Pr \left\{ \exists i \in [m], |\hat{\mu}_{i,t-1} - \mu_i| \geq \sqrt{\frac{3 \ln t}{2T_{i,t-1}}} \right\} \\
&\leq \sum_{i \in [m]} \Pr \left\{ |\hat{\mu}_{i,t-1} - \mu_i| \geq \sqrt{\frac{3 \ln t}{2K_{i,t-1}}} \right\} \\
&= \sum_{i \in [m]} \sum_{k=1}^{(t-1)\lfloor K'/2 \rfloor} \Pr \left\{ T_{i,t-1} = k, |\hat{\mu}_{i,t-1} - \mu_i| \geq \sqrt{\frac{3 \ln t}{2T_{i,t-1}}} \right\} \\
&\leq \sum_{i \in [m]} \sum_{k=1}^{(t-1)\lfloor K'/2 \rfloor} \frac{2}{t^3} < 2m \lfloor K'/2 \rfloor t^{-2}. \quad (\text{Hoeffding's inequality [14]})
\end{aligned}$$

When $T_{i,t-1} = k$, $\hat{\mu}_{i,t}$ is the average of k i.i.d. random variables $Y_i^{[1]}, \dots, Y_i^{[k]}$, where $Y_i^{[j]}$ is a random variable that indicates whether two members selected with replacement from C_i are the same. Since each community is explored at most K' times in each round, $T_{i,t-1} \leq (t-1)\lfloor K'/2 \rfloor$. The last line leverages the Hoeffding's inequality [14]. By replacing the summation range $k \in [1, (t-1)\lfloor K'/2 \rfloor]$ with $k \in [1, (t-1)]$ in the 3rd line of above equation, we have $\Pr \{\neg \mathcal{N}_t\} \leq 2mt^{-2}$. ■

Secondly, we use the monotonicity and bounded smoothness properties to bound the reward gap $\Delta_{\mathbf{k}_t} = r_{\mathbf{k}^*}(\boldsymbol{\mu}) - r_{\mathbf{k}_t}(\boldsymbol{\mu})$ between our action \mathbf{k}_t and the optimal action \mathbf{k}^* .

Lemma 10. *If the event \mathcal{N}_t holds in round t , we have*

$$\Delta_{\mathbf{k}_t} \leq \sum_{i=1}^m \binom{k_{i,t}}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1}).$$

Here the function $\kappa_T(M, s)$ is defined as

$$\kappa_T(M, s) = \begin{cases} 2 & \text{if } s = 0, \\ 2\sqrt{\frac{6 \ln t}{s}} & \text{if } 1 \leq s \leq l_T(M), \\ 0 & \text{if } s \geq l_T(M) + 1, \end{cases}$$

where

$$l_T(M) = \frac{24 \binom{K'}{2}^2 \ln T}{M^2}.$$

Proof. By \mathcal{N}_t (i.e., $\underline{\boldsymbol{\mu}}_t \leq \boldsymbol{\mu}$) and the monotonicity of $r_{\mathbf{k}}(\boldsymbol{\mu})$, we have

$$r_{\mathbf{k}_t}(\underline{\boldsymbol{\mu}}_t) \geq r_{\mathbf{k}^*}(\underline{\boldsymbol{\mu}}_t) \geq r_{\mathbf{k}^*}(\boldsymbol{\mu}) = r_{\mathbf{k}_t}(\boldsymbol{\mu}) + \Delta_{\mathbf{k}_t}.$$

Then by the *bounded smoothness* properties of reward function, we have

$$\Delta_{\mathbf{k}_t} \leq r_{\mathbf{k}_t}(\underline{\boldsymbol{\mu}}_t) - r_{\mathbf{k}_t}(\boldsymbol{\mu}) \leq \sum_{i=1}^m \binom{k_{i,t}}{2} (\mu_i - \underline{\mu}_{i,t}).$$

We intend to bound $\Delta_{\mathbf{k}_t}$ by bounding $\mu_i - \underline{\mu}_{i,t}$. Before doing so, we perform a transformation. Let $M_{\mathbf{k}_t} = \max_{i \in [m], k_{i,t} > 1} \Delta_{\min}^i$. Since the action \mathbf{k}_t always satisfies $\Delta_{\mathbf{k}_t} \geq \max_{i \in [m], k_{i,t} > 1} \Delta_{\min}^i$,

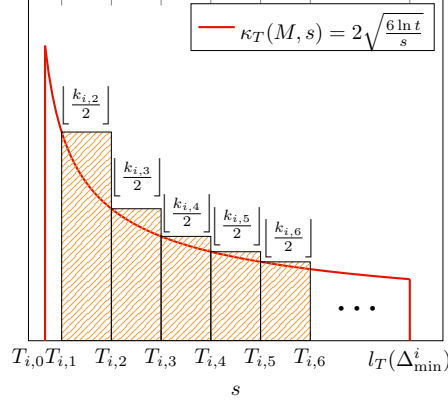


Figure 2: Demonstration of the regret summation $\sum_{t=2}^T \lfloor k_{i,t}/2 \rfloor \kappa_T(\Delta_{\min}^i, T_{i,t-1})$. It is obvious that when $k_{i,t} = K'$, then the shaded area (colored with orange) covered by the rectangles is maximized.

we have $\Delta_{\mathbf{k}_t} \geq M_{\mathbf{k}_t}$. So $\sum_i \binom{k_{i,t}}{2} (\mu_i - \underline{\mu}_{i,t}) \geq \Delta_{\mathbf{k}_t} \geq M_{\mathbf{k}_t}$. Therefore,

$$\begin{aligned}
\Delta_{\mathbf{k}_t} &\leq \sum_{i=1}^m \binom{k_{i,t}}{2} (\mu_i - \underline{\mu}_{i,t}) \leq -M_{\mathbf{k}_t} + 2 \sum_{i=1}^m \binom{k_{i,t}}{2} (\mu_i - \underline{\mu}_{i,t}) \\
&\leq -\frac{\sum_{i=1}^m \binom{k_{i,t}}{2}}{\binom{K'}{2}} M_{\mathbf{k}_t} + 2 \sum_{i=1}^m \binom{k_{i,t}}{2} (\mu_i - \underline{\mu}_{i,t}) \quad (\text{Corollary 4: } \sum_{i=1}^m \binom{k_{i,t}}{2} \leq \binom{K'}{2}) \\
&= 2 \sum_{i=1}^m \binom{k_{i,t}}{2} \left[(\mu_i - \underline{\mu}_{i,t}) - \frac{M_{\mathbf{k}_t}}{K'(K'-1)} \right] \\
&\leq 2 \sum_i \binom{k_{i,t}}{2} \left[(\mu_i - \underline{\mu}_{i,t}) - \frac{\Delta_{\min}^i}{K'(K'-1)} \right]. \quad (\text{by definition of } M_{\mathbf{k}_t})
\end{aligned}$$

By \mathcal{N}_t , we have $\mu_i - \underline{\mu}_{i,t} \leq \min\{2\rho_{i,t}, 1\}$. So

$$\mu_i - \underline{\mu}_{i,t} - \frac{\Delta_{\min}^i}{K'(K'-1)} \leq \min\{2\rho_{i,t}, 1\} - \frac{\Delta_{\min}^i}{K'(K'-1)} \leq \min\left\{\sqrt{\frac{6 \ln t}{T_{i,t-1}}}, 1\right\} - \frac{\Delta_{\min}^i}{K'(K'-1)}.$$

If $T_{i,t-1} \leq l_T(\Delta_{\min}^i)$, we have $\mu_i - \underline{\mu}_{i,t} - \frac{\Delta_{\min}^i}{K'(K'-1)} \leq \min\left\{\sqrt{\frac{6 \ln t}{T_{i,t-1}}}, 1\right\} \leq \frac{1}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1})$.

If $T_{i,t-1} > l_T(\Delta_{\min}^i) + 1$, then $\sqrt{\frac{6 \ln t}{T_{i,t-1}}} \leq \frac{\Delta_{\min}^i}{K'(K'-1)}$, so $(\mu_i - \underline{\mu}_{i,t}) - \frac{\Delta_{\min}^i}{K'(K'-1)} \leq 0 = \kappa_T(\Delta_{\min}^i, T_{i,t-1})$. In conclusion, we have

$$\Delta_{\mathbf{k}_t} \leq \sum_{i=1}^m \binom{k_{i,t}}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1}). \quad \blacksquare$$

Above result is also valid for the revised Algo. 3, i.e., $\Delta_{\mathbf{k}_t}^r \leq \sum_{i=1}^m \binom{k_{i,t}}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1}^r)$. Our third step is to prove that when \mathcal{N}_t (resp. \mathcal{N}_t^r) holds, the regret is bounded in $O(\ln T)$.

Theorem 3. Algo. 3 with non-adaptive exploration method has regret as follows.

$$\text{Reg}_{\mu}(T) \leq \sum_{i=1}^m \frac{48 \binom{K'}{2} K \ln T}{\Delta_{\min}^i} + 2 \binom{K'}{2} m + \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3} m \Delta_{\max} = O\left(\sum_{i=1}^m \frac{K'^3 \log T}{\Delta_{\min}^i}\right). \quad (6)$$

Proof. We first prove the regret when the event \mathcal{N}_t holds. In each run, we have

$$\begin{aligned} \sum_{t=1}^T \mathbf{1}(\{\Delta_{\mathbf{k}_t} \wedge \mathcal{N}_t\}) \cdot \Delta_{\mathbf{k}_t} &\leq \sum_{t=1}^T \sum_{i=1}^m \binom{k_{i,t}}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1}) \\ &= \sum_{i=1}^m \sum_{t' \in \{t \mid 1 \leq t \leq T, k_{i,t} > 1\}} \binom{k_{i,t'}}{2} \kappa_T(\Delta_{\min}^i, T_{i,t'-1}). \end{aligned}$$

Hence, we just assume $k_{i,t} > 1$ for $t > 0$.

$$\begin{aligned} \sum_{t=1}^T \mathbf{1}(\{\Delta_{\mathbf{k}_t} \wedge \mathcal{N}_t\}) \cdot \Delta_{\mathbf{k}_t} &\leq \sum_{i=1}^m \sum_{t=1}^T \binom{k_{i,t}}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1}) \\ &\leq \sum_{i=1}^m 2 \binom{k_{i,1}}{2} + K' \sum_{i=1}^m \sum_{t=2}^T \frac{(k_{i,t} - 1)}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1}) \\ &\leq 2m \binom{K'}{2} + K' \sum_{i=1}^m \sum_{t=2}^T \left\lfloor \frac{k_{i,t}}{2} \right\rfloor \kappa_T(\Delta_{\min}^i, T_{i,t-1}). \quad (\text{Fig. 2}) \end{aligned}$$

To maximize the summation $\sum_{t=2}^T \left\lfloor \frac{k_{i,t}}{2} \right\rfloor \kappa_T(\Delta_{\min}^i, T_{i,t-1})$, we just need to let $k_{i,t} = K'$ when $t > 1$.

$$\begin{aligned} \sum_{t=1}^T \mathbf{1}(\{\Delta_{\mathbf{k}_t} \wedge \mathcal{N}_t\}) \cdot \Delta_{\mathbf{k}_t} &\leq 2m \binom{K'}{2} + K' \sum_{d=0}^{l_T(\Delta_{\min}^i)/\lfloor K'/2 \rfloor} \left\lfloor \frac{K'}{2} \right\rfloor \kappa_T(\Delta_{\min}^i, 1 + d \lfloor K'/2 \rfloor) \\ &\leq 2m \binom{K'}{2} + K' \sum_{i=1}^m \sum_{d=0}^{l_{T,K}} \frac{\sqrt{24 \ln T} \lfloor K'/2 \rfloor}{\sqrt{1 + d \lfloor K'/2 \rfloor}} \quad (l_{T,K} := \frac{l_T(\Delta_{\min}^i)}{\lfloor K'/2 \rfloor}) \\ &\leq 2m \binom{K'}{2} + K' \sum_{i=1}^m \int_{x=0}^{l_{T,K}} \frac{\sqrt{24 \lfloor K'/2 \rfloor \ln T}}{\sqrt{x}} dx \\ &= 2m \binom{K'}{2} + K' \sum_{i=1}^m \sqrt{96 l_T(\Delta_{\min}^i, T) \ln T} \\ &= 2m \binom{K'}{2} + \sum_{i=1}^m \frac{48 \binom{K'}{2} K' \ln T}{\Delta_{\min}^i}. \end{aligned}$$

On the other hand, when \mathcal{N}_t does not hold, we can bound the regret as Δ_{\max} . Hence,

$$\mathbb{E} \left[\sum_{t=1}^T \mathbf{1}(\{\Delta_{\mathbf{k}_t} \wedge \neg \mathcal{N}_t\}) \cdot \Delta_{\mathbf{k}_t} \right] \leq \Delta_{\max} \sum_{t=1}^T 2m \lfloor K'/2 \rfloor t^{-2} \leq \frac{m \lfloor K'/2 \rfloor \pi^2}{3} \Delta_{\max}.$$

Based on above discussion, we have

$$\text{Reg}_{\mu}(T) \leq \frac{m \lfloor K'/2 \rfloor \pi^2}{3} \Delta_{\max} + 2m \binom{K'}{2} + \sum_{i=1}^m \frac{48 \binom{K'}{2} K' \ln T}{\Delta_{\min}^i}. \quad \blacksquare$$

Theorem 7. The revised Algo. 3 has regret as follows.

$$\text{Reg}_{\mu}^r(T) \leq \sum_{i=1}^m \frac{48 \binom{K'}{2}^2 \ln T}{\Delta_{\min}^i} + 2 \binom{K'}{2} m + \frac{\pi^2}{3} \cdot m \cdot \Delta_{\max}. \quad (12)$$

Proof. We prove the regret when the event \mathcal{N}_t^r holds. In each run, we have

$$\begin{aligned}
\sum_{t=1}^T \mathbf{1}(\{\Delta_{\mathbf{k}_t}^r \wedge \mathcal{N}_t^r\}) \cdot \Delta_{\mathbf{k}_t}^r &\leq \sum_{t=1}^T \sum_{i=1}^m \binom{k_{i,t}}{2} \kappa_T(\Delta_{\min}^i, T_{i,t-1}^r) \\
&= \sum_{i=1}^m \sum_{s=0}^{T_{i,T}^r} \binom{k_{i,s}}{2} \kappa_T(\Delta_{\min}^i, s) \\
&\leq 2m \binom{K'}{2} + \binom{K'}{2} \sum_{i=1}^m \sum_{s=1}^{l_T(\Delta_{\min}^i)} \sqrt{\frac{24 \ln T}{s}} \\
&\leq 2m \binom{K'}{2} + \sum_{i=1}^m \frac{48 \binom{K'}{2}^2 \ln T}{\Delta_{\min}^i}.
\end{aligned}$$

On the other hand, $\Pr\{\neg \mathcal{N}_t^r\} \leq 2mt^{-2}$. Hence we have

$$\begin{aligned}
\text{Reg}_{\mu}^r(T) &= \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}(\{\Delta_{\mathbf{k}_t} \wedge \neg \mathcal{N}_t^r\}) \cdot \Delta_{\mathbf{k}_t} \right] + \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}(\{\Delta_{\mathbf{k}_t} \wedge \mathcal{N}_t^r\}) \cdot \Delta_{\mathbf{k}_t} \right] \\
&\leq \frac{m\pi^2}{3} \Delta_{\max} + 2m \binom{K'}{2} + \sum_{i=1}^m \frac{48 \binom{K'}{2}^2 \ln T}{\Delta_{\min}^i}. \quad \blacksquare
\end{aligned}$$

The bound in Eq. (12) is tighter than the one obtained by directly applying [20].

D.4 Comparison

Estimator. Let $\hat{\mu}_{i,t}$ be the estimator computed in Algo. 3 by end of round t and $\hat{\mu}_{i,t}^r$ be the estimator computed with revision in Eq. (11) by end of round t . Both of $\hat{\mu}_{i,t}$ and $\hat{\mu}_{i,t}^r$ are unbiased estimator of μ_i . However, $\hat{\mu}_{i,t}$ is a *more efficient* estimator than $\hat{\mu}_{i,t}^r$. More specifically, $\text{Var}[\hat{\mu}_{i,t}] = \mu_i(1 - \mu_i)/(\sum_{t'=1}^t \lfloor k_{i,t'}/2 \rfloor)$ and $\text{Var}[\hat{\mu}_{i,t}^r] = \mu_i(1 - \mu_i) \cdot (\sum_{t'=1}^t 1/\lfloor k_{i,t'}/2 \rfloor)/(T_{i,t}^r)^2$. Here $k_{i,t}$ is the size of \mathcal{S}_i in round t , and $T_{i,t}^r = \sum_{t'=1}^t \mathbb{1}\{k_{i,t'} > 1\}$. Since the harmonic mean is always not larger than arithmetic mean, i.e., $T_{i,t}^r/(\sum_{t'=1}^t 1/\lfloor k_{i,t'}/2 \rfloor) \leq (\sum_{t'=1}^t \lfloor k_{i,t'}/2 \rfloor)/T_{i,t}^r$, we conclude that $\text{Var}[\hat{\mu}_{i,t}] \leq \text{Var}[\hat{\mu}_{i,t}^r]$.

Regret Bound. The regret bound in Eq. (6) is tighter than the one in Eq. (12) up to $(K' - 1)/2$ factor in the $O(\ln T)$ term. The bound in Eq. (6) has a larger constant term. That's because we use a smaller confidence radius, which leads to earlier exploitation of Algo. 3 than the revised one.

D.5 Full information feedback

In the following, we prove the constant regret bound of the Algo. 3 with feeding the empirical mean in COMMUNITYEXPLORE and making revision defined in Eq. (7).

Proof. We first bound $\Delta_{\mathbf{k}_t}$ by $\sum_{i=1}^m |\mu_{i,t} - \mu_i|$.

$$\begin{aligned}
\Delta_{\mathbf{k}_t} &= r_{\mathbf{k}^*}(\mu) - r_{\mathbf{k}_t}(\mu) = r_{\mathbf{k}^*}(\mu) - r_{\mathbf{k}_t}(\hat{\mu}) + r_{\mathbf{k}_t}(\hat{\mu}) - r_{\mathbf{k}_t}(\mu) \\
&\leq r_{\mathbf{k}^*}(\mu) - r_{\mathbf{k}^*}(\hat{\mu}) + r_{\mathbf{k}_t}(\hat{\mu}) - r_{\mathbf{k}_t}(\mu) \quad (r_{\mathbf{k}^*}(\hat{\mu}) \leq r_{\mathbf{k}_t}(\hat{\mu})) \\
&\leq |r_{\mathbf{k}^*}(\mu) - r_{\mathbf{k}^*}(\hat{\mu})| + |r_{\mathbf{k}_t}(\hat{\mu}) - r_{\mathbf{k}_t}(\mu)| \\
&\leq \sum_{i=1}^m \left(\binom{k_i^*}{2} + \binom{k_{i,t}}{2} \right) |\hat{\mu}_{i,t-1} - \mu_i|.
\end{aligned}$$

Leverage the fact that $\sum_{i=1}^m \binom{k_{i,t}}{2} \leq \binom{K'}{2}$. If $|\hat{\mu}_{i,t-1} - \mu_i| < \frac{\Delta_{\min}}{K'(K'-1)}$, then

$$\Delta_{\mathbf{k}_t} \leq \sum_{i=1}^m \left(\binom{k_i^*}{2} + \binom{k_{i,t}}{2} \right) \frac{\Delta_{\min}}{K'(K'-1)} < \Delta_{\min},$$

which means $\Delta_{\mathbf{k}_t} = 0$. Hence,

$$\begin{aligned} \Pr(\Delta_{\mathbf{k}_t} > 0) &\leq \sum_{i=1}^m \Pr\left(|\hat{\mu}_{i,t-1} - \mu_i| \geq \frac{\Delta_{\min}}{K'(K'-1)}\right) \\ &\leq \sum_{i=1}^m 2e^{-2(T_{i,t-1}/2)\Delta_{\min}^2/(K'(K'-1))^2}. \end{aligned} \quad (\text{Theorem 3.2 in [9]})$$

The second line of above inequality using Theorem 3.2 in [9]. Note that the $T_{i,t-1}$ member pairs using for collision counting are not independent with each other. We need to construct a *dependence graph* G to model their dependence. The dependence graph here is just a line with $T_{i,t-1}$ nodes. Since the fractional chromatic number of the dependence graph is 2, we have a $1/2$ factor for $T_{i,t-1}$ in the exponential. The regret is bounded as

$$\begin{aligned} \text{Reg}_{\mu}(T) &\leq \sum_{t=1}^T \sum_{i=1}^m \Delta_{\mathbf{k}_t} 2e^{-T_{i,t-1}\Delta_{\min}^2/(K'(K'-1))^2} \\ &\leq 2\Delta_{\max} + \sum_{i=1}^m \sum_{t=3}^T \Delta_{\mathbf{k}_t} 2e^{-(t-2)\Delta_{\min}^2/(K'(K'-1))^2} \quad (T_{i,t-1} \geq t-2) \\ &\leq 2\Delta_{\max} + 2m\Delta_{\max} \int_{t=0}^{\infty} e^{-t\Delta_{\min}^2/(K'(K'-1))^2} dt \\ &\leq \left(2 + 8me^2 \binom{K'}{2} / \Delta_{\min}^2\right) \Delta_{\max}. \quad \blacksquare \end{aligned}$$

E Regret Analysis for Adaptive Problem

E.1 Transition probability list of policy π^t

Similar to the discussion in Section B.2.1, we define a transition probability list $\mathcal{P}(\pi^t, \psi)$ for the policy π^t and write the reward function $r_{\pi^t}(\mu)$ with $\mathcal{P}(\pi^t, \emptyset)$.

Definition. Assume the initial partial realization is ψ . Let \mathbf{s}_0 be the status corresponding to ψ . Recall that $\mathbf{s}_0 = (s_{1,0}, \dots, s_{m,0}) = (1 - \mu_1 c_1(\psi), \dots, 1 - \mu_m c_m(\psi))$. At the first step, policy π^t chooses community $i_0^* = \arg \max_{i \in [m]} 1 - c_i(\psi) \mu_{i,t}$. With probability $q_0^{\pi^t} := c_{i_0^*}(\psi) \mu_{i_0^*}$, the communities stay at the same status \mathbf{s}_0 . With probability $p_0^{\pi^t} := 1 - c_{i_0^*}(\psi) \mu_{i_0^*}$, the communities transit to next status $\mathbf{s}_1 := \mathbf{s}_0 - \mu_{i_0^*} \mathbf{I}$. Note that

$$1 - c_i(\psi) \mu_{i,t} = \frac{\mu_i - (1 - s_{i,0}) \mu_{i,t}}{\mu_i} = \frac{\mu_{i,t}}{\mu_i} s_{i,0} + \frac{\mu_i - \mu_{i,t}}{\mu_i}.$$

We recursively define \mathbf{s}_{k+1} as $\mathbf{s}_k - \mu_{i_k^*} \mathbf{I}_{i_k^*}$ where $i_k^* \in \max_{i \in [m]} (\mu_{i,t}/\mu_i) s_{i,k} + (\mu_i - \mu_{i,t})/\mu_i$. The transition probability $p_k^{\pi^t} := s_{i_k^*,k}$. We define the transition probability list $\mathcal{P}(\pi^t, \psi) = (p_0^{\pi^t}, \dots, p_D^{\pi^t})$ where $D = \sum_{i=1}^m (d_i - c_i(\psi))$ is the number of distinct member we haven't meet under the partial realization ψ . Note that it is possible that $p_k^{\pi^t} = 0$. In this case, there is already no unmet members in i_k^* . The communities will be stuck in status \mathbf{s}_k since the policy π^t always chooses community i_k^* to explore after the communities reach status \mathbf{s}_k . Hence, if k is the smallest index such that $p_k^{\pi^t} = 0$, we define $p_{k'}^{\pi^t} = 0$ for all $k' > k$.

Compute $\mathcal{P}(\pi^t, \psi)$. Define $\mathcal{B}_i(\psi) = \{1 - c_i(\psi) \mu_i, 1 - (1 + c_i(\psi)) \mu_i, \dots, \mu_i, 0\}$ for $i \in [m]$. Let $b_i \in \mathcal{B}_i(\psi)$, $b_j \in \mathcal{B}_j(\psi)$, $i, j \in [m]$. We define a *sorting comparator* as follows.

$$\text{less}(b_i, b_j) = \mathbb{1}\{(\mu_{i,t}/\mu_i) \cdot b_i + (\mu_i - \mu_{i,t})/\mu_i < (\mu_{j,t}/\mu_j) \cdot b_j + (\mu_j - \mu_{j,t})/\mu_j\}$$

If $b_i \geq b_j$ and $\text{less}(b_i, b_j) = 1$, we can infer that $\mu_{i,t}/\mu_i \geq \mu_{j,t}/\mu_j$, which means the size of community j is more overestimated than the size of community i . The overestimation leads to wrong order between b_i and b_j when using the comparator less. The list $\mathcal{P}(\pi^t, \psi)$ can be computed as follows. Firstly, we sort elements in $\cup_{i \in [m]} \mathcal{B}_i$ with the comparator less. Secondly, we truncate the

sorted list at the first zero elements. Thirdly, we paddle zeros at the end of list until the length is $D + 1$. All the arguments in Section B.2-B.1 about $\mathcal{P}(\pi^g, \psi)$ can be easily extended to $\mathcal{P}(\pi^t, \psi)$.

Expected reward. In the following, we still use the extended definition of reward

$$R(\mathbf{k}, \phi) = f \left(\sum_{i=1}^m \left| \bigcup_{\tau=1}^{k_i} \{\phi(i, \tau)\} \right| \right),$$

where f is a non-decreasing function. We can write the reward function $r_{\pi^t}(\boldsymbol{\mu})$ as

$$r_{\pi^t}(\boldsymbol{\mu}) = \sum_{j=0}^{\min\{K, \sum_{i=1}^m d_i\}} f(j) \times p_0^{\pi^t} \times \cdots \times p_{j-1}^{\pi^t} \times L(\{q_0^{\pi^t}, \dots, q_j^{\pi^t}\}, K - j).$$

Here $p_j^{\pi^t}$ is element in $\mathcal{P}(\pi^t, \emptyset)$, $q_j^{\pi^t} := 1 - p_j^{\pi^t}$, and K is the budget.

E.2 Proof framework

Notations. Let $D = \sum_{i=1}^m d_i$ in this part. Let $\mathcal{P}(\pi^g, \emptyset) = (p_0^{\pi^g}, \dots, p_D^{\pi^g})$ and $\mathcal{P}(\pi^t, \emptyset) = (p_0^{\pi^t}, \dots, p_D^{\pi^t})$. According to Corollary 1, we know that $\mathcal{P}(\pi^g, \emptyset)$ can be obtained by sorting $\cup_{i \in [m]} \{1, 1 - \mu_i, 1 - 2\mu_i, \dots, \mu_i\} \cup \{0\}$. Here we define another list $\tilde{\mathcal{P}}(\pi^g)$ which is obtained by sorting $\cup_{i \in [m]} \{(i, 1), (i, 1 - \mu_i), \dots, (i, \mu_i)\}$ via comparing the second value in the pair. Let $U_{i,k}$ denote how many times pair (i, \cdot) appears in the first k positions in the list $\tilde{\mathcal{P}}(\pi^g)$. The value $U_{i,k}$ satisfies that $p_k^{\pi^g} = \max_{i=1}^m 1 - U_{i,k} \mu_i$. Note that the definition of $U_{i,k}$ are equivalent to the one defined in the main text.

Theorem 5. Algo. 3 with adaptive exploration method has regret as follows.

$$\text{Reg}_{\boldsymbol{\mu}}(T) \leq \left(\sum_{i=1}^m \sum_{k=m+1}^{\min\{K, D\}} \frac{6\Delta_{\max}^{(k)}}{(\Delta_{\min}^{i,k})^2} \right) \ln T + \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3} \sum_{i=1}^m \sum_{k=m+1}^{\min\{K, D\}} \Delta_{\max}^{(k)}. \quad (8)$$

Proof. When $\boldsymbol{\mu}_t$ is close to $\boldsymbol{\mu}$, the list $\mathcal{P}(\pi^t, \emptyset)$ is similar to the list $\mathcal{P}(\pi^g, \emptyset)$, which indicates the reward gap $r_{\pi^g}(\boldsymbol{\mu}) - r_{\pi^t}(\boldsymbol{\mu})$ is small. Let $\mathbb{1}_{i,k}(\boldsymbol{\mu}_t)$ be the indicator that takes value 1 when $\mathcal{P}(\pi^g, \emptyset)$ and $\mathcal{P}(\pi^t, \emptyset)$ are the same for the first k elements, and different at the $(k + 1)$ -th elements (i.e., $p_j^{\pi^g} = p_j^{\pi^t}$ for $0 \leq j \leq k - 1$ and $p_k^{\pi^g} \neq p_k^{\pi^t}$) with condition $p_k^{\pi^t} = 1 - U_{i,k} \mu_i$. Note that the first m elements in $\mathcal{P}(\pi^t, \emptyset)$ and $\mathcal{P}(\pi^g, \emptyset)$ equal to 1. Then the reward gap at round t is

$$\Delta_{\pi^t} = r_{\pi^g}(\boldsymbol{\mu}) - r_{\pi^t}(\boldsymbol{\mu}) = \sum_{i=1}^m \sum_{k=m+1}^{\min\{K, D\}} \mathbb{1}_{i,k}(\boldsymbol{\mu}_t) \cdot \Delta_{\max}^{i,k},$$

where $\Delta_{\max}^{i,k}$ is the maximum reward gap among all possible $\boldsymbol{\mu}_t$ such that $\mathbb{1}_{i,k}(\boldsymbol{\mu}_t) = 1$, i.e.,

$$\Delta_{\max}^{i,k} = \max_{\boldsymbol{\mu}_t, \mathbb{1}_{i,k}(\boldsymbol{\mu}_t)=1} r_{\pi^g}(\boldsymbol{\mu}) - r_{\pi^t}(\boldsymbol{\mu}).$$

Note that

$$\Delta_{\max}^{i,k} \leq \sum_{j=k}^{\min\{K, D\}} f(j) \times p_0^{\pi^g} \times \cdots \times p_{j-1}^{\pi^g} \times L(\{1 - p_0^{\pi^g}, \dots, 1 - p_j^{\pi^g}\}, K - j).$$

The expected cumulative regret can be expanded as

$$\begin{aligned} \text{Reg}_{\boldsymbol{\mu}}(T) &= \mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=1}^T \Delta_{\pi^t} \right] \leq \sum_{t=1}^T \mathbb{E}_{\Phi_1, \dots, \Phi_{t-1}} \left[\sum_{k=m+1}^{\min\{K, D\}} \sum_{i=1}^m \mathbb{1}_{i,k}(\boldsymbol{\mu}_t) \times \Delta_{\max}^{i,k} \right] \\ &= \sum_{i=1}^m \sum_{k=m+1}^M \Delta_{\min}^{i,k} \mathbb{E}_{\Phi_1, \dots, \Phi_{t-1}} \left[\sum_{t=1}^T \mathbb{1}_{i,k}(\boldsymbol{\mu}_t) \right]. \end{aligned}$$

Our next step is bound $\mathbb{E}_{\Phi_1, \dots, \Phi_{t-1}} \left[\sum_{t=1}^T \mathbb{1}_{i,k}(\underline{\mu}_t) \right]$. We rewrite the indicator $\mathbb{1}_{i,k}(\underline{\mu}_t)$ as:

$$\mathbb{1}_{i,k}(\underline{\mu}_t) = \mathbb{1}_{i,k}(\underline{\mu}_t) \mathbb{1}\{T_{i,t-1} \leq l_{i,k}\} + \mathbb{1}_{i,k}(\underline{\mu}_t) \mathbb{1}\{T_{i,t-1} > l_{i,k}\},$$

where $l_{i,k}$ is a problem-specific constant. In Lemma 11, we show that the probability we choose a wrong community when community i is probed enough times (i.e., $T_{i,t-1} > l_{i,k}$) is very small. Based on the lemma, the regret corresponding to the event $\mathbb{1}\{T_{i,t-1} > l_{i,k}\}$ is bounded as follows.

$$\sum_{i=1}^m \sum_{k=m+1}^{\min\{K,D\}} \Delta_{\min}^{i,k} \mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=1}^T \mathbb{1}_{i,k}(\underline{\mu}_t) \mathbb{1}\{T_{i,t-1} > l_{i,k}\} \right] \leq \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3} \sum_{i=1}^m \sum_{k=m+1}^{\min\{K,D\}} \Delta_{\max}^{i,k}.$$

On the other hand, the regret associated with the event $\mathbb{1}\{T_{i,t-1} \leq l_{i,k}\}$ is trivially bounded by $\sum_{i=1}^m \sum_{k=m+1}^K \Delta_{\max}^{i,k} l_{i,k}$. In conclusion, the expected cumulative regret is bound as

$$\begin{aligned} \text{Reg}_{\underline{\mu}}(T) &\leq \sum_{i=1}^m \sum_{k=m+1}^K \Delta_{\max}^{i,k} \mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=1}^T \mathbb{1}_{k,t}(\underline{\mu}_t) \right] \\ &\leq \sum_{i=1}^m \sum_{k=m+1}^K \Delta_{\max}^{i,k} l_{i,k} + \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3} \sum_{i=1}^m \sum_{k=m+1}^{\min\{K,D\}} \Delta_{\max}^{i,k} \\ &\leq \left(\sum_{i=1}^m \sum_{k=m+1}^K \frac{6\Delta_{\max}^{i,k}}{(\Delta_{\min}^{i,k})^2} \right) \ln T + \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3} \sum_{i=1}^m \sum_{k=m+1}^{\min\{K,D\}} \Delta_{\max}^{i,k}. \end{aligned}$$

Note $\Delta_{\max}^{(k)} \geq \max_{i \in [m]} \Delta_{\max}^{i,k}$. This completes the proof. \blacksquare

Lemma 11. For all $k \leq \{M, \sum_{i=1}^m d_i\}$, we have

$$\mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=1}^T \mathbb{1}_{i,k}(\underline{\mu}_t) \mathbb{1}\{T_{i,t-1} > l_{i,k}\} \right] \leq \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3}, \quad (13)$$

where $l_{i,k}$ is defined as $l_{i,k} := 6 \ln T / (\Delta_{\min}^{i,k})^2$.

Proof. The following proof is similar to the that for the traditional Upper Confidence Bound (UCB) algorithm [1]. In the following, we define $i_k^* = \max_{i \in [m]} 1 - U_{i,k} \mu_i$.

$$\begin{aligned} \sum_{t=1}^T \mathbb{1}_{i,k}(\underline{\mu}_t) \mathbb{1}\{T_{i,t-1} > l_{i,k}\} &= \sum_{t=l_{i,k}+1}^T \mathbb{1}_{i,k}(\underline{\mu}_t) \mathbb{1}\{T_{i,t-1} > l_{i,k}\} \\ &\leq \sum_{t=l_{i,k}+1}^T \mathbb{1}\{(\hat{\mu}_{i,t-1} - \rho_{i,t-1})U_{i,k} < (\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1})U_{i_k^*,k}, T_{i,t-1} > l_{i,k}\}. \end{aligned}$$

When $T_{i_k^*,t-1} > l_{i,k} \triangleq \frac{6 \ln T}{(\Delta_{\min}^{i,k})^2}$, we have

$$\rho_{i,t-1} = \sqrt{\frac{3 \ln t}{2T_{i,t-1}}} < \frac{\Delta_{\min}^{i,k}}{2} \Rightarrow \underbrace{\mu_{i_k^*} U_{i_k^*,k} < (\mu_i - 2\rho_{i,t-1})U_{i,k}}_{i \text{ and } i_k^* \text{ are distinguishable with high prob.}}$$

If $i \neq i_k^*$ exists such that

$$\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1} < \mu_{i_k^*}, \text{ and } \hat{\mu}_{i,t-1} + \rho_{i,t-1} > \mu_i,$$

we have

$$(\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1})U_{i_k^*,k} < \mu_{i_k^*} U_{i_k^*,k} < (\mu_i - 2\rho_{i,t-1})U_{i,k} < (\hat{\mu}_{i,t-1} - \rho_{i,t-1})U_{i,k},$$

which contradicts with $(\hat{\mu}_{i,t-1} - \rho_{i,t-1})U_{i,k} < (\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1})U_{i_k^*,k}$. Hence when $T_{i,t-1} > l_{i,k}$, we have

$$\begin{aligned} & \{(\hat{\mu}_{i,t-1} - \rho_{i,t-1})U_{i,k} < (\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1})U_{i_k^*,k}\} \\ & \subseteq \{\hat{\mu}_{i,t-1} + \rho_{i,t-1} \leq \mu_i \text{ or } \hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1} \geq \mu_{i_k^*}\} \end{aligned}$$

Using the union bound, we have

$$\begin{aligned} & \Pr((\hat{\mu}_{i,t-1} - \rho_{i,t-1})U_{i,k} < (\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1})U_{i_k^*,k}) \\ & \leq \Pr(\hat{\mu}_{i,t-1} + \rho_{i,t-1} \leq \mu_i \text{ or } \hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1} \geq \mu_{i_k^*}) \\ & \leq \Pr(\hat{\mu}_{i,t-1} + \rho_{i,t-1} \leq \mu_i) + \Pr(\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1} \geq \mu_{i_k^*}). \end{aligned}$$

Therefore, we can conclude that

$$\begin{aligned} & \mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=1}^T \mathbb{1}_{i,k}(\underline{\mu}_t) \mathbb{1}\{T_{i,t-1} > l_{i,k}\} \right] \\ & \leq \sum_{t=l_{i,k}+1}^T \mathbb{1}\{(\hat{\mu}_{i,t-1} - \rho_{i,t-1})U_{i,k} < (\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1})U_{i_k^*,k}, T_{i,t-1} > l_{i,k}\} \\ & \leq \sum_{t=l_{i,k}+1}^T \Pr\{\hat{\mu}_{i,t-1} + \rho_{i,t-1} \leq \mu_i\} + \Pr\{\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1} \geq \mu_{i_k^*}\} \\ & \leq \sum_{t=l_{i,k}+1}^T \left(\sum_{T_{i,t-1}=l_{i,k}+1}^{\lfloor \frac{K'}{2} \rfloor} \Pr\{\hat{\mu}_{i,t-1} + \rho_{i,t-1} \leq \mu_i | T_{i,t-1}\} \right. \\ & \quad \left. + \sum_{T_{i_k^*,t-1}=1}^{\lfloor \frac{K'}{2} \rfloor} \Pr(\hat{\mu}_{i_k^*,t-1} - \rho_{i_k^*,t-1} \geq \mu_{i_k^*} | T_{i_k^*,t-1}) \right) \\ & \leq \sum_{t=1}^{\infty} 2t \left\lfloor \frac{K'}{2} \right\rfloor \times t^{-3} = 2 \left\lfloor \frac{K'}{2} \right\rfloor \sum_{t=1}^{\infty} t^{-2} = \frac{\lfloor \frac{K'}{2} \rfloor \pi^2}{3}. \quad \blacksquare \end{aligned}$$

E.3 Full information feedback

If we feed the empirical mean in the exploration oracle, then the policy π^t is determined by $\hat{\mu}_t$. Similarly, we can define the event $\mathbb{1}_{i,k}(\hat{\mu}_t)$ by replacing $\underline{\mu}_t$ with $\hat{\mu}$ in Section E.1-E.2.

Lemma 12. *If we make revisions defined in Eq. (7) to Algo. 3 and feed the empirical mean in COMMUNITYEXPLORE to explore communities adaptively, then for all community C_i and $k \leq \{K, \sum_{i=1}^m d_i\}$, we have*

$$\mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=2}^T \mathbb{1}_{i,k}(\hat{\mu}_t) \right] \leq \frac{2}{\varepsilon_{i,k}^4}, \quad (14)$$

where $\varepsilon_{i,k}$ is defined as (here $i_k^* \in \arg \min_{i \in [m]} \mu_i U_{i,k}$)

$$\varepsilon_{i,k} \triangleq \frac{\mu_i U_{i,k} - \mu_{i_k^*} U_{i_k^*,k}}{U_{i,k} + U_{i_k^*,k}} \text{ for } i \neq i_k^* \text{ and } \varepsilon_{i,k} = \infty \text{ for } i = i_k^*.$$

Proof. We first bound the probability of the following event by relating $\mathbb{1}_{i,k}(\hat{\mu}_t)$ with the event that both $\mu_{i,t-1}$ and $\mu_{i_k^*,t-1}$ in the confidence interval $\varepsilon_{i,k}$.

$$\mathbb{1}_{i,k}(\hat{\mu}_t) \leq \mathbb{1}\{\hat{\mu}_{i,t-1} U_{i,k} < \hat{\mu}_{i_k^*,t-1} U_{i_k^*,k}\}.$$

If $i \neq i_k^*$ such that

$$\hat{\mu}_{i,t-1} > \mu_i - \varepsilon_{i,k}, \text{ and } \hat{\mu}_{i_k^*,t-1} < \mu_{i_k^*} + \varepsilon_{i,k},$$

then

$$\hat{\mu}_{i,t-1} U_{i,k} > (\mu_i - \varepsilon_{i,k}) U_{i,k} = (\mu_{i_k^*} + \varepsilon_{i,k}) U_{i_k^*,k} > \hat{\mu}_{i_k^*,t-1} U_{i_k^*,k},$$

which contradicts with that $\hat{\mu}_{i,t-1}U_{i,k} < \hat{\mu}_{i_k^*,t-1}U_{i_k^*}$. Here $(\mu_i - \varepsilon_{i,k})U_{i,k} = (\mu_{i_k^*} + \varepsilon_{i^*,k})U_{i_k^*,k}$ can be derived from the definition of $\varepsilon_{i,k}$. Therefore

$$\begin{aligned} & \mathbb{1} \{ \hat{\mu}_{i,t-1}U_{i,k} < \hat{\mu}_{i_k^*,t-1}U_{i_k^*,k} \} \\ & \leq \mathbb{1} \{ \hat{\mu}_{i,t-1} \leq \mu_i - \varepsilon_{i,k} \text{ or } \hat{\mu}_{i_k^*,t-1} \geq \mu_{i_k^*} + \varepsilon_{i,k} \}. \end{aligned}$$

With above equation and the concentration bound in [9], the expectation $\mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=2}^T \mathbb{1}_{i,k}(\hat{\mu}_t) \right]$ can be bounded as

$$\begin{aligned} & \mathbb{E}_{\Phi_1, \dots, \Phi_T} \left[\sum_{t=2}^T \mathbb{1}_{i,k}(\hat{\mu}_t) \right] \\ & \leq \sum_{t=2}^T \Pr \{ \hat{\mu}_{i,t-1}U_{i,k} < \hat{\mu}_{i_k^*,t-1}U_{i_k^*,k} \} \\ & \leq \sum_{t=2}^T \Pr \{ \hat{\mu}_{i,t-1} \leq \mu_i - \varepsilon_{i,k} \} + \Pr \{ \hat{\mu}_{i_k^*,t-1} \geq \mu_{i_k^*} + \varepsilon_{i,k} \} \\ & \leq \sum_{t=2}^T \left(\sum_{T_{i,t-1}=t-1}^{t \lfloor K'/2 \rfloor} e^{-\varepsilon_{i,k}^2 T_{i,t-1}} + \sum_{T_{i_k^*,t-1}=t-1}^{t \lfloor K'/2 \rfloor} e^{-\varepsilon_{i,k}^2 T_{i_k^*,t-1}} \right) \\ & \leq 2 \sum_{t=1}^T \sum_{s=t}^{\infty} e^{-s\varepsilon_{i,k}^2} \leq 2 \sum_{t=1}^T \frac{e^{-t\varepsilon_{i,k}^2}}{\varepsilon_{i,k}^2} \leq \frac{2}{\varepsilon_{i,k}^4}. \quad \blacksquare \end{aligned}$$

F Experimental Evaluation

In this section, we conduct simulations to validate the theoretical results claimed in the main text and provide some insight for future research.

F.1 Offline Problems

In this part, we show some simulation results for the offline problems.

Performance of Algorithm 1. In Fig. 3, we show that the allocation lower bound k^- and upper bound k^+ are close to the optimal budget allocation. From Fig. 3, we observe that the $L1$ distance between k^* and k^- (or k^+) is around $m/2$, which means the average time complexity of Algorithm 1 is $\Theta((m \log m)/2)$.

Reward v.s. Budget. We show the relationship between the reward (i.e., the number of distinct members) and the given budget in Fig. 4. From Fig. 4, we can draw the following conclusions.

- The performance of the four methods are ranked as: “Adaptive Opt.”, “Non-adaptive Opt.”, “Proportional to Size”, “Random Allocation”. This validate our optimality results in Sec. 3.
- The method “Proportional to Size” and “Non-adaptive Opt.” have similar performance. It is an intuitive idea to allocate budgets proportional to the community sizes. The simulation results also demonstrate the efficiency of such budget allocation method. In the following, we analyze the reason theoretically. Recall the definition of k^- as follows.

$$k_i^- = \frac{(K-m)/\ln(1-\mu_i)}{\sum_{j=1}^m 1/\ln(1-\mu_j)}.$$

When $\mu_i \ll 1$, we have $\ln(1-\mu_i) \approx -\mu_i$. Hence,

$$k_i^- \approx \frac{(K-m)d_i}{\sum_{j=1}^m d_j}.$$

Besides, the $L1$ distance between k^* and k^- is smaller than m . We can conclude that the budget allocation proportional to size is close to the optimal budget allocation. Fig. 6 also validates this conclusion.

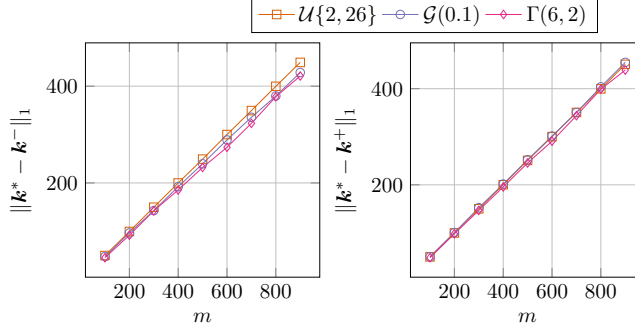


Figure 3: The $L1$ distance between \mathbf{k}^* and \mathbf{k}^- , \mathbf{k}^+ under different community size distributions. Here $\mathcal{U}\{2, 26\}$ is the discrete uniform distribution between 2 and 26. $\mathcal{G}(0.1)$ is the geometric distribution with success probability 0.1 on the support set $\{2, 3, \dots\}$. $\Gamma(\alpha, \beta)$ is the gamma distribution with shape α and rate β . We discretize the support set of the gamma distribution and add 2 to all the values in the support set to ensure that the minimum size of communities is 2. The budget K is a random number between $m + 1$ and $\sum_i d_i$. We run the simulations for 1000 times for each data point.

- The reward gap between “Non-adaptive Opt.” and “Adaptive Opt.” increases first and then decreases, as shown in Fig. 5.

Budget Allocation Comparison. Fig. 6 and Fig 5 show the budget allocation of non-adaptive optimal method and adaptive optimal method. Fig. 5 shows that the adaptive optimal method use the budget more efficiently.

F.2 Online Problems

In the following, we show the simulation results for the online, non-adaptive problem. The simulation results for online, adaptive are similar. Hence, we only present the results for online, non-adaptive problems. Fig. 7 shows the regret of three different learning methods. For illustration purpose, we set the community sizes as $bmd = (2, 3, 5, 6, 8, 10)$. From Fig. 7, we can draw the following conclusions.

- If we feed the empirical mean into the oracle, the regret grows linearly.
- The regret of CLCB algorithm is bounded logarithmically, as proved in Thm. 3.
- The regret under full information feedback setting is bounded as a problem related constant, as proved in Thm. 4.

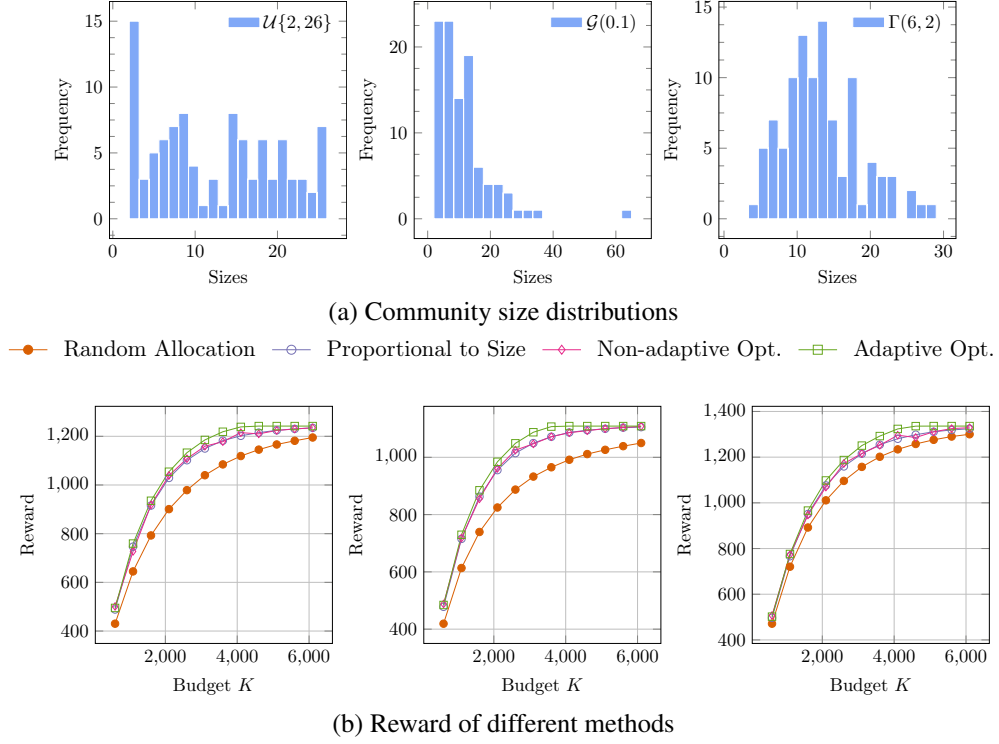


Figure 4: Reward v.s. Budget. In the first row, we show three different size distributions of $m = 100$ communities. In the second row, we show the reward of four different budget allocation methods. Here “*Random Allocation*” represents random budget allocation (sum up to K). “*Proportional to Size*” method allocates budget proportional to the community sizes. “*Non-adaptive Opt.*” corresponds to the optimal budget allocation obtained by the greedy method. “*Adaptive Opt.*” means we explore the communities with greedy adaptive policy π^g . The simulations are run for 200 times for each data point on the budget-reward curve.

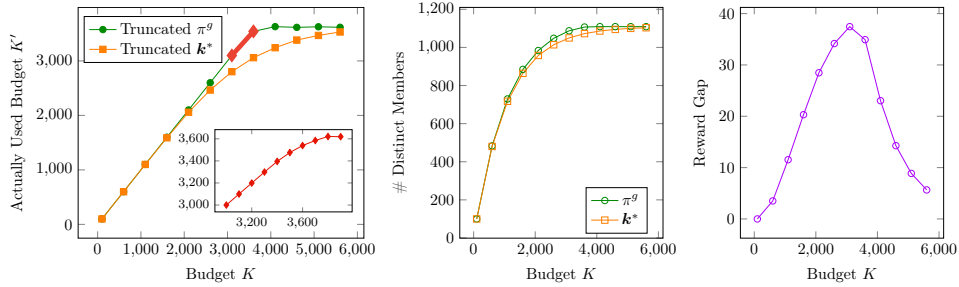


Figure 5: Actually used budget. we only show the results for the community size configuration generated by $\mathcal{G}(0.1)$, as shown in the first row of Fig. 4. The legend labels have the same meaning as in Fig. 6

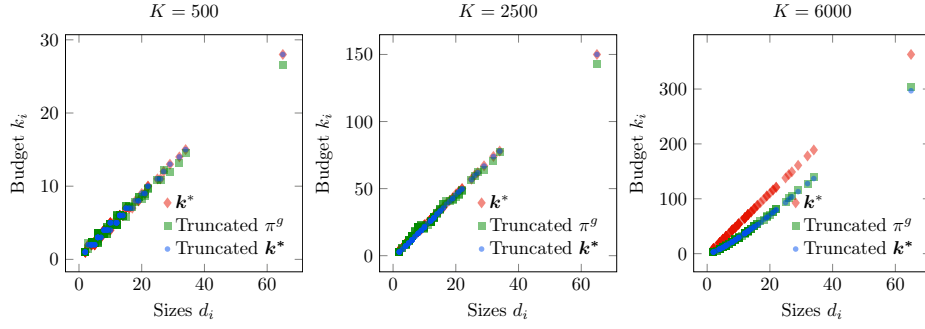
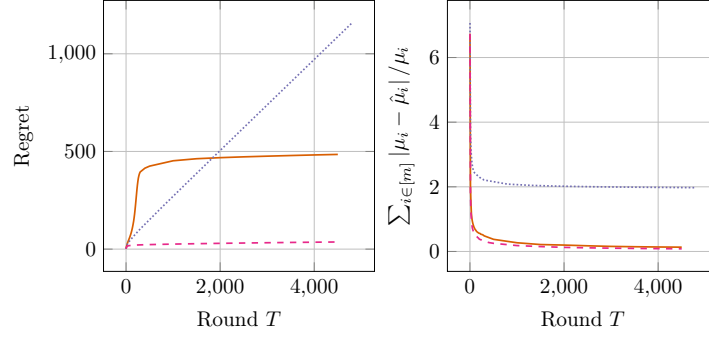
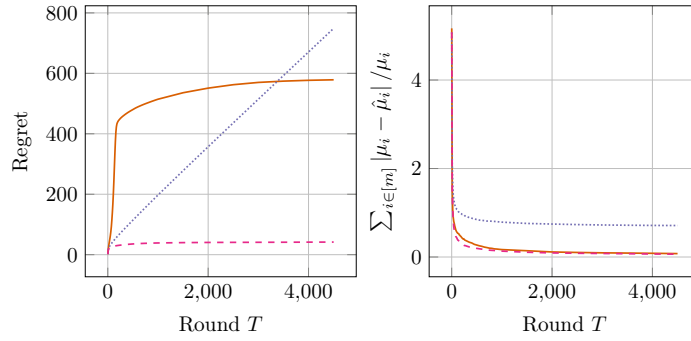


Figure 6: Comparison of different budget allocation methods. The distribution of community sizes generated by the geometric distribution with success probability 0.1, as shown in the first row of Fig. 4. The legend label “ k^* ” represents the optimal budget allocation. The “truncated π^g ” means we stop the greedy adaptive process if all the members are found. The “truncated k^* ” means we stop the non-adaptive exploration of community C_i if all the members of C_i are found. Each data point is an average of 1000 simulations.

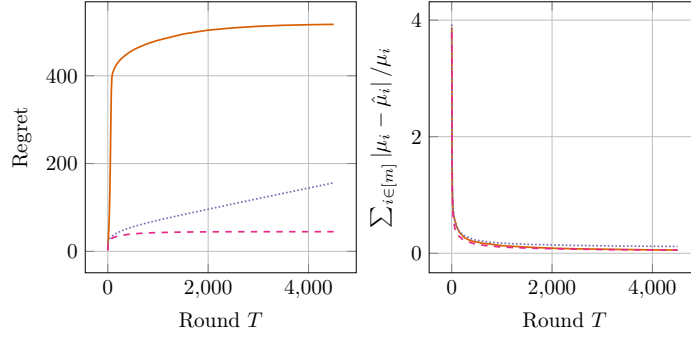
— CLCB Empirical mean - - - Full information feedback



(a) $K = 20, \mathbf{k}^* = (1, 2, 3, 3, 5, 6)$



(b) $K = 30, \mathbf{k}^* = (2, 3, 4, 5, 7, 9)$



(c) $K = 50, \mathbf{k}^* = (3, 4, 7, 9, 12, 15)$

Figure 7: Comparison of different learning algorithms. The sizes of communities are $\mathbf{d} = (2, 3, 5, 6, 8, 10)$. Here the label *Empirical mean* represents feeding the empirical mean into the oracle directly. The regret/error line plots are average of 100 simulations.