# Sequential community mode estimation

Shubham Anand Jain[a], Shreyas Goenka[a], Divyam Bapna[a], Nikhil Karamchandani[a],
Jayakrishnan Nair[a]

*[a]Department of Electrical Engineering, IIT Bombay, India*

**Abstract**

We consider a population, partitioned into a set of communities, and study the problem of identifying the largest community within the population via sequential, random sampling of individuals. There are multiple sampling domains, referred to as *boxes*, which also partition the population. Each box may consist of individuals of different communities, and each community may in turn be spread across multiple boxes. The learning agent can, at any time, sample (with replacement) a random individual from any chosen box; when this is done, the agent learns the community the sampled individual belongs to, and also whether or not this individual has been sampled before. The goal of the agent is to minimize the probability of mis-identifying the largest community in a *fixed budget* setting, by optimizing both the sampling strategy as well as the decision rule. We propose and analyse novel algorithms for this problem, and also establish information theoretic lower bounds on the probability of error under any algorithm. In several cases of interest, the exponential decay rates of the probability of error under our algorithms are shown to be optimal up to constant factors. The proposed algorithms are further validated via simulations on real-world datasets.

*Keywords:* mode estimation, limited precision sampling, sequential algorithms, fixed budget, multi-armed bandits
*PACS:* 0000, 1111
*2000 MSC:* 0000, 1111

## 1. Introduction

Several applications in online learning involve sequential sampling/polling of an underlying population. A classical learning task in this space is *online cardinality estimation*, where the goal is to estimate the size of a set by sequential sampling of elements from the set (see, for example, [1, 2, 3]). The key idea here is to use 'collisions,' i.e., instances where the same element is sampled more than once, to estimate the size of the set. Another recent application is *community exploration*, where the goal of the learning agent is to sample as many distinct elements as possible, given a family of sampling distributions/domains to poll from (see [4, 5]).

In this paper, we focus on the related problem of *community mode estimation*. Here, the goal of the learning agent is to estimate the largest community within a population of individuals, where each individual belongs to a unique community. The agent has access to a set of sampling domains, referred to as *boxes* in this paper, which also partition the population. The agent can, at any sampling epoch, choose which box to sample from. Having chosen one such box to sample from, a random individual from this box gets revealed to the agent, along with the community

arXiv:2111.08535v1 [stat.ML] 16 Nov 2021

that individual belongs to. After a fixed budget of samples is exhausted, the learning agent reveals its estimate of the largest community (a.k.a., the community mode) in the population. The goal of the agent is in turn to minimize the probability of mis-identifying the community mode, by optimizing (i) the policy for sequential sampling of boxes, and (ii) the decision rule that determines the agent's response as a function of all observations.

One application that motivates this formulation is election polling. In this context, communities might correspond to the party/candidate an individual votes for, while boxes might correspond, for instance, to different cities/states that individuals reside in. In this case, community mode identification corresponds to predicting the winning party/candidate. A related (and contemporary) application is the detection of the dominant strain of a virus/pathogen within a population of infected individuals. Here, communities would correspond to different strains, and boxes would correspond to different regions/jurisdictions.

Another application of a different flavour is as follows. Consider a setting where an agent interacts with a database which has several entries, each with an associated label, and the agent is interested in identifying the most represented label in the database. For concreteness, consider a user who polls a movie recommendation engine which hosts a large catalogue of movies, each belonging to a particular genre, to discover the most prevalent genre in the catalogue.[1] In each round, the user might provide a genre (*community*) to the recommendation engine which then suggests a movie (*individual*) from that genre (perhaps based on other user ratings). Depending on the recommendations seen thus far, the user selects the next genre to poll and so on. Now, either due to privacy considerations or simply the lack of knowledge of all the available genres, it might not be feasible for the user to share the exact genre he/she wants to view in each round and might only provide coarser directions (*box*). For example, while there might be specific genres available such as dark comedy, romantic comedy, slapstick comedy etc., the user might only indicate its choice as 'comedy' and then let the recommendation engine suggest some movie belonging to any of the sub-genres in the broad genre. At one extreme, the user might prefer complete privacy and not suggest any genre in each round, in which case the recommendation engine will have to choose a movie over the entire database. This resembles the *mixed community setting* studied in this paper. The opposite end of the spectrum is where the user does not care about privacy and instead specifies a sub-genre in each round from which the recommendation engine can then suggest a movie. This corresponds to the *separated community setting*. We refer to the intermediate scenario where the user provides coarse directives as the *community-disjoint box setting*.

The formulation we consider here has some parallels with the classical multi-armed bandit (MAB) problem [6]; specifically, the fixed budget best arm identification formulation [7]. Indeed, one may interpret communities in our formulation as arms in an MAB problem. However, there are two crucial differences between the two formulations. The first difference lies in the stochastic behavior of the reward/observation sequence. In the classical MAB problem, each pull of an arm yields an i.i.d. reward drawn from an arm specific reward distribution. However, in the community mode detection problem, the sequence of collisions (or equivalently, the evolution of the number of distinct individuals seen) does not admit an i.i.d. description. (Indeed, whether or not a certain sample from a box results in a collision depends in a non-stationary manner on the history of observations from that box.) The second difference between the two formulations lies in the extent of sampling control on part of the agent. In the MAB setting, the agent can pull

---

any arm it chooses at any sampling epoch. However, in our formulation, the agent cannot sample directly from a community of its choice; it must instead choose a box to sample from, limiting its ability to target specific communities to explore.

In terms of the extent of sampling control that the agent has, the opposite end of the spectrum to the MAB setting is when samples are simply generated by an underlying distribution and the agent can only use these observations to estimate some property of the underlying distribution. This classical problem of property estimation from samples generated from an underlying distribution has a long and rich history. There has been a lot of work recently on characterizing the optimal sample complexity for estimating various properties of probability distributions including entropy [8, 9], support size and coverage [10, 11], and 'Lipschitz' properties [12] amongst others. Closer to the problem studied in this paper, the problem of mode estimation was originally studied in [13, 14] with the focus on statistical properties of various estimators such as consistency. More recently, the instance-optimal sample complexity of mode estimation for any discrete distribution was derived in [15]. Our formulation differs from this line of work in the non-i.i.d. nature of the observations as well as the partial ability that the agent has to control the sampling process, by being able to query any box at a given instant.

Our contributions are summarized as follows.

- We begin by considering a special case of our model where the entire population is contained within a single box; we refer to this as the *mixed community setting* (see Section 3). In this setting, the sampling process is not controlled, and the learning task involves only the decision rule. We show that a simple decision rule, based on counting the number of distinct individuals encountered from each community, is optimal, via comparison of an upper bound on the probability of error (mis-identification of the community mode) under the proposed algorithm with an information theoretic lower bound. For this setting, we also highlight the impact of being able to identify sampled individuals (i.e., determine whether or not the sampled individual has been seen before) on the achievable performance in community mode estimation.

- Next, we consider the case where each community lies in its own box; the so-called *separated community setting* (see Section 4). Here, we show that the commonly used approach of detecting pairwise collisions (see [4]) is sub-optimal. Next, a near-optimal algorithm is proposed that borrows the sampling strategy of the classical *successive rejects* policies for MABs [7], but differentiates communities based on the number of distinct individuals encountered (which is different from the classical MAB setting where arms are differentiated based on their empirical average rewards).

- Next, we consider a setting that encompasses both the mixed community as well as the separated community settings; we refer to it as the *community-disjoint box setting* (see Section 5). Here, each community is contained within a single box (though a box might contain multiple communities). For this case, we propose novel algorithms that combine elements from the mixed and separated community settings. Finally, we show how the algorithms designed for the community-disjoint box setting can be extended to the fully general case, where communities are arbitrarily spread across boxes.

- Finally, we validate the algorithms proposed on both synthetic as well as real-world datasets (see Section 6).

3

We conclude this section by making a comparison between our contributions and the literature on the fixed budget MAB problem. Near optimal algorithms for the fixed budget MAB problem (see, for example, [7, 16]) follow a sampling strategy of *successive rejection* of arms, wherein the sampling budget is split across multiple phases, and at the end of each phase, a certain number of (worst performing) arms are eliminated from further consideration. Some of our algorithms for the community mode estimation problem follow a similar sampling strategy and eliminate boxes in phases; specifically, we often use the same sampling schedule as in the classical successive rejects algorithm proposed in [7]. However, the elimination criterion we use is different: it is based on the number of distinct individuals seen (so far) from each community. Given that this statistic evolves in a non-stationary Markovian fashion over time, this distinction makes our analysis more complex.

Our information theoretic lower bounds are inspired by the framework developed in [17] for the fixed budget MAB problem. However, as before, the key distinction in our proofs stems from the difference in stochastic nature of the observation process: while reward observations for each arm in the classical MAB setup are i.i.d., the number of distinct individuals seen from each community evolves as an absorbing Markov chain in the community mode estimation problem.

## 2. Problem Formulation

Consider a population consisting of $N$ individuals. Each individual belongs to exactly one out of $m$ communities, labelled $1, 2, \cdots, m$. Additionally, the population is partitioned across $b$ *sampling domains*, also referred to as 'boxes' in this paper. The boxes are labelled $1, 2, \cdots, b$. Our learning goal is to identify, via random sequential sampling of the boxes, the largest community (a.k.a., the community mode).

We represent the partitioning of the population across communities and boxes via a $b \times m$ matrix $D$. The entry in the $i$th row and $j$th column of this matrix, denoted by $d_{ij}$, equals the number of individuals in box $i$ who are in community $j$. Throughout the paper, we refer to $D$ as the *instance* associated with the learning task. Let $d_j := \sum_i d_{ij}$ denote the size of community $j$, and $N_i := \sum_j d_{ij}$ denote the size of box $i$.

The learning agent a priori knows only the set of boxes and the set of communities. It can access the population by querying an oracle. The input to this oracle is a box number, and the response from the oracle is a (uniformly chosen) random individual from this box and the community that individual belongs to. Individuals are sampled with replacement, i.e., the same individual can be sampled multiple times. Additionally, we assume that the learning agent is able to 'identify' the sampled individual, such that it knows whether (and when) the sampled individual had been seen before.[2] For each query, the agent can decide which box to sample based on the oracle responses received thus far. At the end of a fixed budget of $t$ oracle queries, the agent outputs its estimate $\hat{h}^* \in [m]$ of the community mode $h^*(D) = \arg\max_{j \in [m]} d_j$ in the underlying instance $D$.[3] The agent makes an error if $\hat{h}^* \notin h^*(D)$, and the broad goal of this paper is to design sequential community mode estimation algorithms that minimize the probability of error.

---

[2]Note that this does not require the agent to store a unique identifier (like, say, the social security number) associated with each sampled individual. The agent can simply assign its own *pseudo-identity* to an individual the first time the individual is seen. This sampling model has been applied before in a variety of contexts, including cardinality estimation (see [1, 2]) and community exploration (see [4]).

[3]We use the notation $[a : b]$ to denote the set $\{a, a + 1, \ldots, b\}$ for any $a, b \in \mathbb{Z}, b \geq a$. For $b \in \mathbb{N}$, $[b] := [1 : b]$.

Formally, for any $k \in [t]$, a sequential algorithm $\mathcal{A}$ has to specify a box $b_k$ to sample for the $k$th query, this choice being a function of only past observations. The probability of error for an algorithm $\mathcal{A}$ under an instance $D$, with a budget of $t$ oracle queries, is given by $P_e(D, \mathcal{A}, t) \overset{\Delta}{=} \mathbb{P}(\hat{h}^* \notin h^*(D))$. An algorithm $\mathcal{A}$ is said to be *consistent* if, for any instance $D$, $\lim_{t \to \infty} P_e(D, \mathcal{A}, t) = 0$. We often suppress the dependence on the budget $t$ and also the algorithm $\mathcal{A}$ (when the algorithm under consideration is clear from the context) when expressing the probability of error, denoting it simply as $P_e(D)$.

For notational simplicity, we assume throughout that the instance $D$ is has a unique largest community, with $h^*(D)$ denoting the largest community; our results easily generalize to the case where $D$ has more than one largest community. In the following sections, for various settings of interest, we prove instance-specific upper bounds on the probability of error of our proposed algorithms. We are also able to prove information theoretic lower bounds on the probability of error under *any* algorithm (within a broad class of *reasonable* algorithms). In some cases, we show that the exponential decay rate of the information theoretic lower bound with respect to the horizon matches (up to a factor that is logarithmic in the number of boxes) the corresponding decay rate for our algorithm-specific upper bounds; this implies the near optimality of our algorithms.

**Remark:** As is also the case with algorithms for the fixed budget MAB problem, the probability of error under our proposed algorithms typically decays exponentially with respect to the budget $t$, i.e., $P_e(D) \leq \mu(D)e^{-\lambda(D)t}$, where $\mu(D)$, and $\lambda(D)$ are instance (and algorithm) dependent positive constants. Our primary goal would be to characterize and optimize the exponential decay rate $\lambda(D)$ above. With the focus thus being on the decay rate, the value of the exponential pre-factor $\mu(D)$ in our bounds will often be loose; this is also the case in the fixed budget MAB literature.

**Remark:** It is also important to note that in the classical fixed budget MAB problem, the decay rates associated with the upper bounds on the probability of error under the best known algorithms *do not* match exactly the decay rates corresponding to the best known information theoretic lower bounds: the two decay rates differ by a multiplicative factor that is logarithmic in the number of arms [18]. Given this fundamental gap in the state of the art, it is common practice to refer an algorithm as near optimal if the decay rate associated with its upper bound is a logarithmic (in the number of arms) factor away from the decay rate in the best known information theoretic lower bound. Interestingly, we observe a similar multiplicative mismatch between the decay rates in our upper and lower bound for the community mode estimation problem (as noted above).

The remainder of this paper is organized as follows. We begin by considering the *mixed community setting* in Section 3, where all individuals belong to a single box ($b = 1$); in this special case, the instance matrix $D$ has a single row. Note that in the mixed community setting, the agent has no control on the sampling process. Next, in Section 4, we study the opposite end of the spectrum with respect to sampling selectivity, where each community constitutes a unique box ($b = m$); this corresponds to $D$ being a diagonal matrix (up to row permutations). We refer to this special case as the *separated community setting.* Next, in Section 5, we consider the intermediate setting, where each community is entirely contained within a single box. This corresponds to each column of $D$ having exactly one non-zero entry. The algorithms presented in this section also extend to the most general case, where each community may be spread across multiple boxes. Finally, in Section 6, we present simulation results that compare the proposed algorithms on both synthetic data as well as several real-world datasets. We conclude this section with a summary

5

of our main results.

*Summary of main results*

In Tables 1, 2, and 3, we present a summary of our results, classified by setting. For ease of presentation, only the decay rates associated with our (upper and lower) bounds on probability of error are mentioned here.

Table 1: Summary of the mixed community setting (decay rates)

| SAMPLING MODEL | LOWER BOUND | ALGORITHM | UPPER BOUND |
|---|---|---|---|
| IDENTITYLESS | $\log\left(\frac{N}{N-\left(\sqrt{d_1}-\sqrt{d_2}\right)^2}\right)$ (THEOREM 2) | SFM | $\log\left(\frac{N}{N-\left(\sqrt{d_1}-\sqrt{d_2}\right)^2}\right)$ (THEOREM 1) |
| IDENTITY | $\log\left(\frac{N}{N-(d_1-d_2+1)}\right)$ (THEOREM 4) | DSM | $\log\left(\frac{N}{N-(d_1-d_2)}\right)$ (THEOREM 3) |

Table 1 summarizes our results for the mixed-community setting, where for simplicity, we have represented the community sizes as $d_1, d_2, \ldots, d_m$, with $d_1 > d_2 \geq d_3 \geq \cdots \geq d_m$. In this case, we consider both an *identityless* sampling model, wherein the identity of the sampled individual is not revealed to the learning agent, as well as the identity-based model described in our problem formulation. As we point out in Section 3, the decay rate corresponding to the identity-based sampling model exceeds that under the identityless model, indicating that identity information helps to improve the performance of mode identification. Note that the decay rates corresponding to our upper and lower bounds match exactly for the identity-based sampling model, and almost exactly for the identity-based model. Since the mixed-community setting consists of a single box, the multiplicative discrepancy described above between the decay rates in the upper and lower bounds does not arise here.

In Table 2, we summarize our main results for the separated community setting. Since there is a single community per box here, we once again represent the community/box sizes as $d_1, d_2, \ldots, d_b$, with $d_1 > d_2 \geq d_3 \geq \cdots \geq d_b$. The decay rate in our lower bound is expressed in terms of the instance-dependent complexity metric $H_2(D) := \sum_{i=2}^{b} \frac{1}{\log(d_1)-\log(d_i)}$, and that in our upper bound is expressed in terms of the related complexity metric $H(D)$, which is within a $\overline{log}(b) = \frac{1}{2} + \sum_{i=2}^{b} \frac{1}{i}$ factor of $H_2(D)$ (see Lemma 7).

Table 3 summarizes our main results for the community-disjoint box setting. Here, $d_{11}$ denotes the size of the largest community, which is contained in Box 1, $c_1$ denotes the size of the second largest community in Box 1, and for $i \geq 2$, $c_i$ denotes the size of the largest community in Box $i$. The remaining constants in the decay rate expressions are defined in Section 5. The decay rates corresponding to the upper and lower bounds are expressed as a minimum of two terms: the first corresponds to the (sub)task of identifying the box containing the largest community, while the second corresponds to the (sub)task of identifying the largest community within that box. As we elaborate in Section 5, for a certain class of (reasonable) instances, the two decay rates can be shown to be within constant factors of one another.

## 3. Mixed Community Setting

We first consider the mixed community setting, where $b = 1$, i.e., the instance matrix $D$ has a single row. In other words, the population is completely 'mixed' and for each query, the agent

Table 2: Summary of the separated community setting (decay rates)

| Lower Bound | Algorithm | Upper Bound |
|---|---|---|
| $\frac{3}{H_2(D)}$ (Theorem 8) | DS-SR | $\frac{1}{\overline{\log(b)}H(D)}$ (Theorem 6) |

Table 3: Summary of the community-disjoint box setting (decay rates)

| Lower Bound | Algorithms | Upper Bound |
|---|---|---|
| $\min\left(\frac{\Gamma}{H_2^b(D)}, \log\left(\frac{N_1}{N_1 - (d_{11} - c_1 + 1)}\right)\right)$ (Theorems 12, 13) | DS-SR, ENDS-SR | $\min\left(\frac{1}{\overline{\log(b)}H^b(D)}, \frac{1}{2\overline{\log(b)}}\log\left(\frac{N_1}{N_1 - d_{11} + c_1}\right)\right)$ (Theorem 10) |

obtains a uniformly random sample from the entire population. Thus, the sampling process in this case is uncontrolled, and the learning task is to simply identify the largest community based on the $t$ samples obtained.

In the mixed community setting, we also consider an *identity-less* sampling model, wherein the agent only learns the community that the sampled individual belongs to, without any other identifying information. Under this sampling model, the agent cannot tell whether or not an individual who has been sampled has been seen before. This model not only forms a benchmark for our subsequent analysis of identity-based sampling, but is also of independent interest, given its privacy-preserving property.

Throughout this section, since there is a single box, we drop the first index in $d_{ij}$, and represent the instance simply as $D = (d_1, d_1, \cdots, d_m)$. Also, without loss of generality, we order the communities as $d_1 > d_2 \geq d_3 \geq \cdots \geq d_m$.

### 3.1. Identity-less sampling

We begin by analysing the identity-less sampling model in the mixed community setting. Note that in this case, the response to each oracle query is community $i$, with a probability proportional to the size of the $i$th community. Thus, the agent receives $t$ i.i.d. samples from the discrete distribution $(p_1, p_2, \cdots, p_m)$, where $p_i = d_i/N$. Hence, the learning task boils down to the identification of the mode of this distribution, using a fixed budget of $t$ i.i.d. samples.[4]

### 3.1.1. Algorithm

We consider a natural algorithm in this setting, which we call the Sample Frequency Maximization (SFM) algorithm: return the empirical mode, i.e., the community which has produced the largest number of samples, with ties broken randomly. One would anticipate that this algorithm is optimal, since the vector $(\hat{\mu}_j(t),\ 1 \leq j \leq m)$, where $\hat{\mu}_j(t)$ denotes the number of samples from community $j$ over $t$ oracle queries, is a sufficient statistic for the distribution $D$. The probability of error under the SFM algorithm is bounded from above as follows.

---

[4]The same mode identification problem was considered in the *fixed confidence* setting recently in [15].

**Theorem 1.** *Consider the mixed community setting, under the identity-less sampling model. For any instance D, the Sample Frequency Maximization algorithm has a probability of error upper bounded as*

$$P_e(D) \leq (m-1)\left(1 - \frac{(\sqrt{d_1} - \sqrt{d_2})^2}{N}\right)^t.$$

The proof, which follows from a straightforward application of the Chernoff bound, can be found in Appendix A. Note that the probability of error under the SFM algorithm decays exponentially with the budget $t$, the decay rate being (at least) $\log\left(\frac{N}{N-(\sqrt{d_1}-\sqrt{d_2})^2}\right)$. The optimality of this decay rate is established next, via an information-theoretic lower bound on the probability of error under any consistent algorithm.

*3.1.2. Lower Bound*

The following theorem establishes an asymptotic lower bound on the probability of error under any consistent algorithm which uses identity-less sampling. Recall that under a consistent algorithm, for any underlying instance $D$ the probability of error converges to zero as $t \to \infty$.

**Theorem 2.** *In the mixed community setting, under the identity-less sampling model, any consistent algorithm on an instance D satisfies*

$$\liminf_{t\to\infty} \frac{1}{t} \log(P_e(D)) \geq -\log\left(\frac{N}{N - (\sqrt{d_1} - \sqrt{d_2})^2}\right).$$

The proof of this theorem, which uses ideas from the proof of [17, Theorem 12], can be found in Appendix B. Since the exponential decay rate in the above lower bound matches that in the upper bound corresponding to the SFM algorithm for any instance $D$, it follows that SFM is asymptotically decay-rate optimal (under identity-less sampling).

*3.2. Identity Sampling*

Having considered the case of identity-less sampling in the previous section, we now revert to the identity-based sampling model described in Section 2. We show that identity information can be used to improve the accuracy of community mode estimation. We begin by proposing and analysing a simple algorithm for community mode estimation, and then establish information-theoretic lower bounds.

*3.2.1. Algorithm*

Under identity-based sampling, we propose a simple *Distinct Samples Maximization* (DSM) algorithm: The DSM algorithm tracks the number of *distinct* individuals seen from each community, and returns the community that has produced the greatest number over the $t$ queries, with ties broken randomly. As before, this is the natural algorithm to consider under identity-based sampling, given that the vector $(S_j(t), \ 1 \leq j \leq m)$, where $S_j(t)$ denotes the number of distinct individuals from community $j$ seen over $t$ oracle queries, is a sufficient statistic for $D$ (see [2]). The probability of error under the DSM algorithm is bounded as follows.

8

**Theorem 3.** *In the mixed community setting, for any instance D, the Distinct Samples Maximization (DSM) algorithm has a probability of error upper bounded as*

$$P_e(D) \le 2(m-1)\exp\left(-\frac{t\left(d_1 - \frac{\sum_{i=2}^{m} d_i}{m-1}\right)^2}{32Nd_1}\right) \quad \text{for } t \le \min\left\{\frac{d_1 + d_m}{2d_1}N, \frac{16Nd_1}{(d_1 - d_m)^2}\right\}, \quad (1)$$

$$P_e(D) \le \binom{d_1}{d_2}\left(1 - \frac{d_1 - d_2}{N}\right)^t = \binom{d_1}{d_2}\exp\left(-t\log\left(\frac{N}{N - d_1 + d_2}\right)\right) \quad \forall t. \quad (2)$$

Theorem 3 provides two upper bounds on the probability of error. The bound (2) holds for all values of budget $t$, while the bound (1) which is only applicable for small to moderate budget values, tends to be tighter for small values of $t$. Note that (2) implies that the probability of error under the DSM algorithm decays exponentially with $t$, with decay rate (at least) $\log\left(\frac{N}{N-(d_1-d_2)}\right)$. Note that this decay rate exceeds the optimal decay rate under identity-less sampling from Theorem 2, since

$$d_1 - d_2 > (\sqrt{d_1} - \sqrt{d_2})^2 \Rightarrow \log\left(\frac{N}{N - (d_1 - d_2)}\right) > \log\left(\frac{N}{N - (\sqrt{d_1} - \sqrt{d_2})^2}\right).$$

This shows that identity information indeed improves the accuracy of community mode estimation.

*Proof.* The proof of (1) relies on an argument using McDiarmid's inequality, and is given in Appendix C. The proof of (2) is given by a coupon collector style argument. The error probability is upper bounded by the probability of the event that there exists a subset of $d_1 - d_2$ individuals in the largest community $C_1$, such that none of them are sampled in the $t$ queries. Thus we have

$$P_e(D) \le \binom{d_1}{d_2}\left(1 - \frac{d_1 - d_2}{N}\right)^t.$$

The details can be found in Appendix C. □

*3.2.2. Lower Bounds*
Next, we show that the exponential decay rate of the probability of error under the DSM algorithm is (nearly) optimal via an information-theoretic lower bound.

**Theorem 4.** *In the mixed community setting, for any consistent algorithm, the probability of error corresponding to an instance D is bounded below asymptotically as*

$$\liminf_{t \to \infty} \frac{\log(P_e(D))}{t} \ge -\log\left(\frac{N}{N - (d_1 - d_2 + 1)}\right).$$

Note that Theorem 4 implies that the DSM algorithm is nearly decay-rate optimal; the small discrepancy between the decay rate under DSM and that in the lower bound ($(d_1 - d_2)$ replaced by $(d_1 - d_2 + 1)$) stems from the discreteness of the space of alternative instances in our change of measure argument. The proof of this Theorem can be found in Appendix D.

## 4. Separated Community Setting

In this section, we consider the *separated community* setting, where each box contains a single and unique community (so that $b = m$). Compared to the mixed community setting considered in Section 3, this setting represents the opposite end of the spectrum with respect to sampling selectivity on part of the agent—the agent can now choose exactly which community to sample from at any time. Note that identity-less sampling is not meaningful in the separated community setting, since the agent can only gauge the size of a community by observing 'collisions,' which occur when the same individual is sampled again.

At a high level, the separated community setting has connections with the (fixed budget) multi-armed bandit (MAB) problem, with boxes/communities corresponding to arms. However, the reward structure in the separated community setting is different from that in a classical MAB problem; indeed, whether or not a sample taken from any community represents a collision depends on past samples from that community. Nevertheless, we show that tools from the MAB literature can still be adapted to design near-optimal algorithms for estimating the largest community in our setting.

Throughout this section, we denote the size of the community in the $b$th box by $d_b$, dropping the redundant second index since there is only one community in each box. Thus, an instance can be defined by the vector $D = (d_1, d_2, \cdots, d_b)$. WLOG, we order the communities such that $d_1 > d_2 \geq d_3... \geq d_b$.

We begin by considering a simple approach, where at each decision epoch, the agent queries a pair of samples from any chosen community, and checks whether or not a collision has occurred, i.e., the same individual has been sampled both times. Since the event of such a (pairwise, consecutive) collision is independent of past samples, and its probability is inversely proportional to the size of the community, this provides a direct mapping to the MAB setting, allowing off-the-shelf MAB algorithms to be applied.[5] However, we find that this approach, which has been used before in the literature (for example, see [4] for an application of this approach to community exploration), is sub-optimal. Next, we propose and analyse an algorithm that tracks the number of distinct individuals seen from each community, and performs a successive elimination of communities until one 'winner' remains. We show that this approach is near-optimal, by comparing its performance to an information-theoretic lower bound.

### 4.1. Algorithms

We begin by describing the successive rejects (SR) algorithm for fixed-budget MABs, proposed in [7] for best arm identification. The SR algorithm is known to be near-optimal in this setting. Our algorithms for the estimation of the largest community, which borrow the sampling framework of the SR algorithm, are described next.

**Successive rejects algorithm:** Consider an MAB problem with $b$ arms. The class of successive rejects (SR) algorithms is parameterized by natural numbers $K_1, K_2, \cdots, K_{b-1}$, satisfying $0 =: K_0 \leq K_1 \leq K_2 \leq \cdots \leq K_{b-1}$, and $\sum_{j=1}^{b-2} K_j + 2K_{b-1} \leq t$, where $t$ denotes the budget/horizon. The algorithm proceeds in $b - 1$ phases, with one arm being rejected from further consideration at the end of each phase. Specifically, in Phase $r$, the $b - r + 1$ surviving arms are each pulled $K_r - K_{r-1}$

---

[5]Note that this approach only looks for 'immediate' collisions and does not track collisions across the entire observation history.

---

**Algorithm 1** Consecutive-collision SR algorithm

---
1: Set $\mathcal{B} = [b]$                                                   ▷ Set of surviving boxes
2: Set $K_0 = 0$, $K_r = \lceil \frac{1}{\overline{\log(b)}} \frac{t/2-b}{b-r+1} \rceil$    $(1 \le r \le b - 1)$
3: **for** $r = 1, 2, ..b - 1$ **do**
4:      For each box in $\mathcal{B}$, perform $(K_r - K_{r-1})$ sample pairs
5:      Set $C_i^r$ as number of consecutive (within disjoint sample pairs) collisions in box $i \in \mathcal{B}$
6:      $\mathcal{B} = \mathcal{B} \setminus \{\arg\max_{i \in \mathcal{B}} C_i^r\}$            (ties broken randomly)
7: **Return** $\hat{h}^* = $ lone surviving box in $\mathcal{B}$

---

times. At the end of this round, the worst performing[6] surviving arm, based on the $K_r$ samples seen so far, is rejected. The output of the algorithm is the arm that survives rejection at the end of Phase $b - 1$. The original SR algorithm proposed in [7] used $K_r \propto \frac{t-b}{b-r+1}$, so that

$$K_r = \left\lceil \frac{1}{\overline{\log}(b)} \frac{t-b}{b-r+1} \right\rceil, \tag{3}$$

where $\overline{\log}(b) = \frac{1}{2} + \sum_{i=2}^{b} \frac{1}{i}$. Other SR variants, including *uniform exploration* ($K_r = \lfloor t/b \rfloor$ for $1 \le r \le b - 1$) and *successive halving* (see [19]) have also been considered in the literature. In the remainder of this paper, when we refer to the SR algorithm, we mean the specific algorithm proposed in [7], with phases defined via (3).

**Consecutive-collision SR algorithm:** In this algorithm, we map the largest community identification problem to an MAB best arm identification problem. Each community is treated as an arm, and an arm pull consists of two samples drawn from that community. The reward is binary, being 1 if the arm pull does not result in a collision, and 0 if it does. Thus, the mean reward associated with arm (community) $i$ equals $1 - \frac{1}{d_i}$, so that the best arm (the one with the highest mean reward) corresponds to the largest community. Note that since each arm pull corresponds to 2 samples, the budget of the MAB reformulation equals $t/2$. On this MAB reformulation, we apply the SR algorithm of [7] to identify the largest community; this is formalized as Algorithm 1. Adapting the proof of [7, Theorem 2] for our setting yields the following upper bound on the probability of error under the Consecutive-collision SR (CC-SR) algorithm.

**Theorem 5.** *In the separated community setting, for any instance D, the Consecutive-collision SR (CC-SR) algorithm given in Algorithm 1 has a probability of error that is upper bounded as*

$$P_e(D) \le \frac{b(b-1)}{2} \exp\left( -\frac{(t/2 - b)}{4\overline{\log}(b)H^c(D)} \right),$$

*where* $\Delta_i = \frac{1}{d_i} - \frac{1}{d_1}$, *and* $H^c(D) = \max\limits_{i \in [2:b]} \frac{i\Delta_i^{-2}}{d_i}$.

The proof of Theorem 5, which uses the Chernoff bound to concentrate the number of consecutive collisions from each community, can be found in Appendix E.

**Distinct Samples SR algorithm:** We now present an algorithm that ranks communities by the number of distinct individuals seen. Note that this involves tracking collisions across the entire

---

[6]In the classical setting where the best arm is defined as the one with the greatest mean reward, the worst performing arm would be the one with the smallest empirical mean estimate.

**Algorithm 2** Distinct Samples SR algorithm (separated community setting)

1: Set $\mathcal{B} = [b]$                                                  ▷ Set of surviving boxes
2: Set $K_0 = 0$, $K_r = \lceil \frac{1}{\overline{\log(b)}} \frac{t-b}{b-r+1} \rceil$    $(1 \le r \le b-1)$
3: **for** $r = 1, 2, ..b-1$ **do**
4:     Sample each box in $\mathcal{B}$, $K_r - K_{r-1}$ times
5:     Set $S_i^r$ as number of distinct individuals seen so far from box $i \in \mathcal{B}$
6:     $\mathcal{B} = \mathcal{B} \setminus \{\arg\min_{i \in \mathcal{B}} S_i^r\}$            (ties broken randomly)
7: Set $\hat{b}$ as lone surviving box in $\mathcal{B}$
8: **Return** $\hat{h}^* =$ lone surviving box in $\mathcal{B}$

observation history of each community. Specifically, we use the same sampling strategy as the SR algorithm, and at the end of each phase, eliminate from further consideration that community which has produced the least number of distinct individuals so far.[7] This algorithm, which we refer to as the Distinct Samples SR (DS-SR) algorithm, is stated formally as Algorithm 2.

**Theorem 6.** *In the separated community setting, for any instance D the Distinct Samples SR (DS-SR) algorithm given in Algorithm 2 has a probability of error that is upper bounded as*

$$P_e(D) \le \left( \sum_{r=1}^{b-1} \binom{d_1}{d_{b-r+1}} \right) \exp\left( -\frac{(t-b)}{\overline{log}(b)H(D)} \right),$$

*where $H(D) = \max\limits_{i \in [2:b]} \frac{i}{\log(d_1) - \log(d_i)}$.*

*Proof.* We begin by noting that $P_e(D) = \sum_r P_e^r(D)$, where $P_e^r(D)$ is the probability that box 1 is eliminated in phase $r$. Since at least one of the $r$ smallest communities is guaranteed to survive in phase $r$, box 1 will not be eliminated in the $r$th phase if the agent has seen at least $d_{b-r+1} + 1$ distinct samples from box 1. Thus, $P_e^r(D)$ is upper bounded by the probability of the event that there exists a subset of $d_1 - d_{b-r+1}$ individuals in box 1, such that none of them are sampled in the $K_r$ queries made until the end of the $r$th phase. Therefore,

$$P_e^r(D) \le \binom{d_1}{d_{b-r+1}} \left( 1 - \frac{(d_1 - d_{b-r+1})}{d_1} \right)^{K_r}$$

$$\implies P_e^r(D) \le \binom{d_1}{d_{b-r+1}} \exp\left( -K_r \log\left( \frac{d_1}{d_{b-r+1}} \right) \right)$$

Summing across $r$, we get that

$$P_e(D) \le \sum_{r=1}^{b-1} \binom{d_1}{d_{b-r+1}} \exp\left( -K_r \log\left( \frac{d_1}{d_{b-r+1}} \right) \right) \tag{4}$$

Using $K_r = \lceil \frac{1}{\overline{log}(b)} \frac{t-b}{b-r+1} \rceil$ for $1 \le r \le b-1$, we note that

$$K_r \log\left( \frac{d_1}{d_{b-r+1}} \right) \ge \frac{(t-b) \log\left( \frac{d_1}{d_{b-r+1}} \right)}{\overline{log}(b)(b-r+1)} \ge \frac{(t-b)}{\overline{log}(b)H(D)}.$$

---

[7] Note however that in the original SR algorithm for MABs, the cumulative reward from each arm has i.i.d. increments. In the present setting however, the cumulative number of distinct individuals seen from any community does not have i.i.d. increments.

Combining with (4), we have

$$P_e(D) \leq \left( \sum_{r=1}^{b-1} \binom{d_1}{d_{b-r+1}} \right) \exp\left( -\frac{(t-b)}{\overline{\log(b)}H(D)} \right).$$

$\square$

Having analysed the CC-SR algorithm and the DS-SR algorithms, it is instructive to compare the exponential decay rates corresponding to the upper bounds of the probability of error under these algorithms. From Theorems 5 and 6, this boils down to comparing the instance-dependent parameters $H^c(D)$ and $H(D)$ respectively, which encode the 'hardness' of the underlying instance. Note that the values of these parameters are larger for instances where the size of the largest community is close to the sizes of the competing communities, and hence it would be harder for an algorithm to correctly estimate the mode. Consequently, the achievable probability of error from Theorems 5 and 6 is also higher for harder instances. Furthermore, note that

$$H^c(D) = \max_{i \in [2:b]} \frac{i d_1^2 d_i}{(d_1 - d_i)^2} \overset{(a)}{>} \max_{i \in [2:b]} \frac{d_1 d_i}{d_1 - d_i} \frac{i}{\log(d_1) - \log(d_i)}$$

$$\geq \frac{d_1 d_b}{d_1 - d_b} \max_{i \in [2:b]} \frac{i}{\log(d_1) - \log(d_i)} = \frac{d_1 d_b}{d_1 - d_b} H(D).$$

Here, the bound $(a)$ follows from the fact that $\log(x) > \frac{x-1}{x}$ for $x > 1$. Since $H^c(D) > \frac{d_1 d_b}{d_1 - d_b} H(D)$, this means that $H^c(D) \gg H(D)$ for most instances of interest, which suggests that the DS-SR algorithm has a far superior performance as compared to the CC-SR algorithm (at least for large budget values). Our simulation results in Section 6 are also consistent with this observation.

Next, we establish the near optimality of the Distinct Samples SR algorithm via an information theoretic lower bound.

### 4.2. Lower Bounds

While the decay rate in the upper bound of the DS-SR algorithm was expressed in terms of the hardness parameter $H(D)$, the information theoretic lower bound for the separated community setting is expressed in terms of a related hardness parameter $H_2(D) := \sum_{i=2}^{b} \frac{1}{\log(d_1) - \log(d_i)}$. $H(D)$ and $H_2(D)$ are comparable upto a logarithmic (in the number of boxes) factor, as shown below.

**Lemma 7.** $\frac{H(D)}{2} \leq H_2(D) \leq \overline{\log}(b) H(D).$

The proof of Lemma 7 can be found in Appendix I.

We now state a lower bound on the probability of error in the separate community setting for any algorithm in a natural algorithm class. The lower bound is non-asymptotic and is expressed in terms of the maximum of the probability of error under the original instance and an alternate instance which has a lower 'hardness'. This is similar in form to the corresponding lower bound for the standard multi-armed bandit setting in [17, Theorem 16].

**Theorem 8.** *In the separated community setting, consider any algorithm that only uses the number of distinct samples from each community (box) to decide which box to sample from at each instant as well as to make the final estimate of the community mode. For any instance D, there exists an alternate instance $D^{[a]}, a \in [2:b]$, such that $H_2(D^{[a]}) \leq H_2(D)$ and*

$$\max\left( P_e(D), P_e(D^{[a]}) \right) \geq \frac{1}{4} \exp\left( -\frac{3t}{H_2(D)} \right).$$

*In the alternate instance $D^{[a]}$, only the size of community $a$ is changed from $d_a$ to $\lceil \frac{d_1^2}{d_a} \rceil$.*

The proof of Theorem 8 uses the following lemma.

**Lemma 9.** *For any algorithm $\mathcal{A}$ and instance $D$, there exists a box (community) $a \in [2 : b]$ such that $E_D[N_a(t)] \le \frac{t}{(\log(d_1) - \log(d_a))H_2(D)}$, where $N_a(t)$ denotes the number of times box $a$ is sampled in $t$ queries under $\mathcal{A}$.*

*Proof.* Assume there exists no such community. Then,

$$\sum_{a=2}^{b} E_D[N_a(t)] > \sum_{a=2}^{b} \frac{t}{(\log(d_1) - \log(d_a))H_2(D)} = t,$$

which is a contradiction. $\qquad\square$

*Proof of Theorem 8.* Consider an algorithm $\mathcal{A}$ which bases all decisions only on the number of distinct individuals seen from each community (box). In this case, $S_j$, the number of distinct samples from box (community) $j$ evolves as a Markov chain over $[0 : d_j]$, with transitions occurring each time the box is pulled. From state $s$, this chain transitions to (the same) state $s$ with probability $q_D^j(s, s) = \frac{s}{d_j}$, and to state $s + 1$ with probability $q_D^j(s, s) = \frac{d_j - s}{d_j}$.

Now, from Lemma 9 there exist a box $a \in [2 : b]$ which satisfies $E[N_a(t)] \le \frac{t}{(\log(d_1) - \log(d_a))H_2(D)}$. Consider the alternate instance $D^{[a]} = (d_1', d_2', \ldots, d_b')$ mentioned in the statement of the theorem, wherein $d_a' = \lceil d_1^2/d_a \rceil$, $d_j' = d_j$ $\forall j \ne a$. Note that the community mode under the alternate instance $D'$ is $a$, different from that under the original instance $D$. Furthermore, note that under the alternate instance $D^{[a]}$ the transition probabilities $q_{D^{[a]}}^k(u, v)$ remain the same for all $k \ne a$. For box $a$,

$$\log\left(\frac{q_D^a(s, s)}{q_{D^{[a]}}^a(s, s)}\right) = \log\left(\frac{\lceil d_1^2/d_a \rceil}{d_a}\right) \le \log\left(\frac{d_1^3}{d_a^3}\right), \tag{5}$$

$$\log\left(\frac{q_D^a(s, s + 1)}{q_{D^{[a]}}^a(s, s + 1)}\right) = \log\left(\frac{1 - s/d_a}{1 - s/\lceil d_1^2/d_a \rceil}\right). \tag{6}$$

Here, (5) because

$$\lceil d_1^2/d_a \rceil \le 1 + d_1^2/d_a = (d_a + d_1^2)/d_a \Rightarrow \frac{\lceil d_1^2/d_a \rceil}{d_a} \le \frac{d_a + d_1^2}{d_a^2} = \frac{d_a^2 + d_1^2 d_a}{d_a^3} \le \frac{d_1^3}{d_a^3}.$$

Next, let $\mathbb{P}_D, \mathbb{P}_{D^{[a]}}$ denote the probability measures induced by the algorithm under consideration by the instances $D, D^{[a]}$, respectively. Then, given a trajectory $x = (a(1), s(1), \cdots, a(t), s(t))$, where $a(k)$ denotes the box pulled on the $k$th query (action), and $s(k) = (s_j(k), j \in [b])$ is the vector of states corresponding to the arms after the $k$th query, the log-likelihood ratio is given by

$$\log \frac{\mathbb{P}_D(x)}{\mathbb{P}_{D^{[a]}}(x)} = \sum_k \sum_{u,v} N_k(u, v, 0, t) \log\left(\frac{q_D^k(u, v)}{q_{D^{[a]}}^k(u, v)}\right),$$

14

where $N_k(u, v, 0, t)$ represents the number of times the transition from state $u$ to state $v$ happens in the Markov chain corresponding to box $k$ over the $t$ queries. Combining with (5), (6), we get

$$D(\mathbb{P}_D \| \mathbb{P}_{D^{[a]}}) = E_D\left[\log \frac{\mathbb{P}_D(x)}{\mathbb{P}_{D^{[a]}}(x)}\right]$$

$$\leq \sum_s E_D[N_a(s, s, 0, t)] \log\left(\frac{d_1^3}{d_a^3}\right) + E_D[N_a(s, s+1, 0, t)] \log\left(\frac{1 - s/d_a}{1 - s/\lceil d_1^2/d_a \rceil}\right)$$

where $D(\cdot \| \cdot)$ denotes the Kullback-Leibler divergence. Note that

$$\lceil d_1^2/d_a \rceil > d_a \implies \frac{1 - s/d_a}{1 - s/\lceil d_1^2/d_a \rceil} \leq 1 \implies \log\left(\frac{1 - s/d_a}{1 - s/\lceil d_1^2/d_a \rceil}\right) \leq 0.$$

Thus, we have

$$D(\mathbb{P}_D \| \mathbb{P}_{D^{[a]}}) \leq \sum_s E_D[N_a(s, s, 0, t)] \log\left(\frac{d_1^3}{d_a^3}\right) \leq E_D[N_a(t)] \log\left(\frac{d_1^3}{d_a^3}\right)$$

Next, we use Lemma 20 from [17] (alternatively, see Lemma 21 in Appendix I) to get that

$$max\left(P_e(D), P_e(D^{[a]})\right) \geq \frac{1}{4} \exp\left(-D(\mathbb{P}_D \| \mathbb{P}_{D^{[a]}})\right) \geq \frac{1}{4} \exp\left(-E_D[N_a(t)] \log\left(\frac{d_1^3}{d_a^3}\right)\right),$$

where $P_e(D)$ is the probability of error under instance $D$. Finally, we use the bound on $E_D[N_a(t)]$ from Lemma 9 to get

$$\max\left(P_e(D), P_e(D^{[a]})\right) \geq \frac{1}{4} \exp\left(-\frac{3t}{H_2(D)}\right).$$

It now remains to show that $H_2(D^{[a]}) \leq H_2(D)$. This is equivalent to showing

$$\sum_{i \in [b], i \neq a} \frac{1}{\log(\lceil \frac{d_1^2}{d_a} \rceil) - \log(d_i)} \leq \sum_{i \in [b], i \neq 1} \frac{1}{\log(d_1) - \log(d_i)}.$$

This condition follows from the following term-by-term comparisons:

$$\frac{1}{\log(\lceil \frac{d_1^2}{d_a} \rceil) - \log(d_i)} \leq \frac{1}{\log(d_1) - \log(d_i)} \quad (i \neq 1, a)$$

$$\frac{1}{\log(\lceil \frac{d_1^2}{d_a} \rceil) - \log(d_1)} \leq \frac{1}{\log(d_1) - \log(d_a)}$$

$\square$

Comparing the upper and lower bounds on the probability of error for the separated community setting in Theorems 6 and 8, we see that the expressions for the decay rates differ (ignoring universal constants) in terms of $H(D)$ vs $H_2(D)$, which from Lemma 7, are at most a factor of $\overline{log}(b)$ apart. In other words, the decay rate under DS-SR is optimal, upto a logarithmic (in the number of boxes) factor. This is similar to the optimality guarantees available in fixed-budget MAB setting (see [7, 17]).

---

**Algorithm 3** Distinct Samples SR algorithm (community-disjoint box setting)

---

1: Set $\mathcal{B} = [b]$                                              ▷ Set of surviving boxes
2: Set $K_0 = 0$, $K_r = \lceil \frac{1}{\log(b)} \frac{t-b}{b-r+1} \rceil$     $(1 \le r \le b-1)$
3: **for** $r = 1, 2, ..b-1$ **do**
4:      Sample each box in $\mathcal{B}$, $K_r - K_{r-1}$ times
5:      Set $S_{ij}^r$ as number of distinct individuals seen so far from community $j$ in box $i \in \mathcal{B}$
6:      Set, for $i \in \mathcal{B}$, $f_i = \max_j S_{ij}^r$
7:      $\mathcal{B} = \mathcal{B} \setminus \{\arg\min_{i \in \mathcal{B}} f_i\}$           (ties broken randomly)
8: Set $\hat{b}$ as lone surviving box in $\mathcal{B}$
9: **Return** $\hat{h}^* = \arg\max_j S_{\hat{b}j}^{(b-1)}$          (ties broken randomly)

---

## 5. Community-disjoint Box Setting

In this section, we consider an intermediate setting that generalizes both the mixed and separated community settings. Specifically, we consider the case where each community exists in exactly one box; i.e, all the members of a community $j$ are present in the same box. (Though any box may contain multiple communities.) In this setting, which we refer to as the *community-disjoint box setting,* we propose algorithms that combine elements from the algorithms presented before for the mixed and separated community settings. For a class of reasonable instances, we are also able to establish the near optimality of certain algorithms. Finally, we show that the algorithms presented in this section can be generalized to handle the most general model, where communities are arbitrarily spread across boxes.

Under the community-disjoint box setting, each column of the instance matrix $D$ has exactly one non-zero entry. Without loss of generality, we assume that $d_{11}$ is the largest value in the matrix $D$; hence, box 1 contains the largest community (also labeled 1). Also without loss of generality, we order boxes by the sizes of the largest communities in them; i.e, if $g_i$, $1 \le i \le b$ is the size of the largest community in box $i$, then $d_{11} = g_1 > g_2 \ge g_3 \ge ... \ge g_b$. Additionally, we define $c_i$ to be the largest *competing* community in a box–that is, $c_i = g_i, i \ne 1$, and $c_1$ is the second largest community in the first box. We state our results in terms of $d_{11}$ and ($c_i$, $i \in [b]$).

### 5.1. Algorithms

The first algorithm we consider for this setting is a generalization of the Distinct Samples SR algorithm from Algorithm 2, where we now eliminate boxes successively. Specifically, the algorithm proceeds in $b-1$ phases; one box being eliminated from subsequent consideration in each of the phases. At the end of the final phase, the algorithm outputs the community that produced the largest number of distinct samples from the last surviving box. Since we have multiple communities in each box, our elimination criterion in each phase is based on the seemingly largest community in each surviving box. In particular, let $S_{ij}^r$ denote the number of distinct individuals encountered from community $j$ in box $i$ at the end of phase $r$. We eliminate, at the end of phase $r$, the (surviving) box that minimizes $\max_j S_{ij}^r$. This algorithm, which we continue to refer to as the Distinct Samples SR (DS-SR) algorithm (with some abuse of notation), is presented formally in Algorithm 3.

**Theorem 10.** *In the community-disjoint box setting, for any instance D, the Distinct Samples SR (DS-SR) algorithm given in Algorithm 3 has a probability of error upper bounded as*

$$P_e(D) \leq \left( \sum_{i=2}^{b} \binom{d_{11}}{c_i} \right) \exp\left( -\frac{(t-b)}{\overline{log}(b)H^b(D)} \right) + \binom{d_{11}}{c_1} \exp\left( -\frac{(t-b)\log\left(\frac{N_1}{N_1-d_{11}+c_1}\right)}{2\overline{log}(b)} \right), \qquad (7)$$

*where $H^b(D) = \max\limits_{i \in [2:b]} \frac{i}{\log(N_1) - \log(N_1 - d_{11} + c_i)}$.*

The upper bound on the probability of error under the DS-SR algorithm above is a sum of two terms. The first term in (7) bounds the probability of misidentifying the box containing the largest community, while the second term in (7) bounds the probability of misidentifying the largest community within the correct box (box 1). Not surprisingly, the second term is structurally similar to the bound (2) we obtained in Theorem 3 for the mixed community setting (restricted to box 1). The proof of Theorem 10 can be found in Appendix F.

The DS-SR algorithm works well in practice, particularly for large budget values. However, its performance can be sub-par for moderate budget values on certain types of instances; particularly instances where the largest community is contained within a very large box. In such cases, it can happen that $\mathbb{E}\left[S_{11}^r\right] < \mathbb{E}\left[S_{ij}^r\right]$ for another community $j$ in a box $i \neq 1$, making it likely that box 1 gets eliminated early. We propose modified algorithms to resolve this issue, under the additional assumption that the box sizes are known a priori to the learning agent.[8] The first modification replaces uniform exploration of boxes with a proportional exploration of the surviving boxes in each phase, resulting in a sampling process (within each phase) somewhat analogous to the mixed community setting considered in Section 3. A second class of algorithms retains uniform box exploration, but normalizes $S_{ij}^r$ to reflect the size of each box (algorithms in this class differing with respect to the specific normalization performed). This latter class of algorithm can also be extended to the original setting where the box sizes are unknown, by replacing the box size by its maximum likelihood estimator.

We begin by describing our first modification of the DS-SR algorithm, which we refer to as the Distinct Samples Proportional SR (DS-PSR) algorithm. The DS-PSR algorithm apportions the budget across phases in the same manner as DS-SR, but the queries within each phase are distributed across surviving boxes in proportion to their sizes. Formally, this corresponds to the same description as Algorithm 3, except that in Line 4, each box $i \in \mathcal{B}$ is sampled $T(\mathcal{B}, r, i)$ times, where $T(\mathcal{B}, r, i) := \lfloor \frac{N_i}{\sum_{k \in \mathcal{B}} N_k}(K_r - K_{r-1})(b - r + 1) \rfloor$. Experimentally, we find that DS-PSR performs very well. However, a tight characterization of the decay rate corresponding to the probability of error is challenging, since the number of queries available to each surviving box in phase $r$, for $1 < r \leq b - 1$, is a random quantity, that depends on the sequence of prior box eliminations.

Next, we describe the normalized variants of the DS-SR algorithm. The first, which we refer to as the Normalized Distinct Samples SR (NDS-SR) algorithm, is described by changing the definition of $f_i$ in Line 6 of Algorithm 3 to

$$f_i^{\text{NDS-SR}} = \max_j \frac{S_{ij}^r}{S_i^r} N_i,$$

---

[8]This is a natural assumption is several applications. For example, in the context of election polling, an agent might know a priori the total number of voters in each city/state.

where $S_i^r$ denotes the number of distinct individuals seen from box $i$ (across different communities) by the end of phase $r$. This normalization is justified as follows: $S_{ij}^r/S_i^r$ is an unbiased estimator of $d_{ij}/N_i$, i.e., the fraction of box $i$ that is comprised by community $j$.

The final variant we propose, referred to as the Expectation-Normalized Distinct Samples SR (ENDS-SR) algorithm, uses the following alternative normalization of $f_i$ in Line 6 of Algorithm 3:

$$f_i^{\text{ENDS-SR}} = \max_j \frac{S_{ij}^r}{\mathbb{E}\left[S_i^r\right]} N_i.$$

This normalization has a similar justification: indeed, $\frac{S_{ij}^r}{\mathbb{E}[S_i^r]}$ is another (more tractable) unbiased estimator of $d_{ij}/N_i$.

Both NDS-SR and ENDS-SR perform quite well in practice. It is challenging to analytically bound the performance of NDS-SR, due to the difficulty in concentrating the fractions $S_{ij}^r/S_i^r$. However, the probability of error under ENDS-SR admits an upper bound analogous to that under DS-SR (albeit more cumbersome). Interestingly, the exponential decay rate of the probability of error under ENDS-SR is identical to that under DS-SR.

**Theorem 11.** *In the community-disjoint box setting, for any instance D,*

$$\limsup_{t \to \infty} \frac{\log P_e(D, \text{ENDS-SR}, t)}{t} \leq -\frac{1}{\overline{log}(b)} \min\left(\frac{1}{H^b(D)}, \frac{1}{2}\log\left(\frac{N_1}{N_1 - d_{11} + c_1}\right)\right).$$

The proof of Theorem 11 can be found in Appendix G. The intuition behind Theorem 11 is that for large $t$, $\mathbb{E}\left[S_i^r\right] \approx N_i$, so that $f_i^{\text{ENDS-SR}} \approx S_{ij}^r$, making the elimination criterion under ENDS-SR nearly identical to that under DS-SR.

### 5.2. Lower Bounds

We now derive information theoretic lower bounds on the probability of error in the community-disjoint box setting, and compare the decay rates suggested by the lower bounds to the decay rate under DS-SR.

Our first lower bound captures the complexity of simply identifying the largest community from within box 1.

**Theorem 12.** *For any consistent algorithm, the probability of error corresponding to an instance D in the community-disjoint box setting is asymptotically bounded below as*

$$\liminf_{t \to \infty} \frac{P_e(D)}{t} \geq -\log\left(\frac{N_1}{N_1 - (d_{11} - c_1 + 1)}\right).$$

Note that Theorem 12 follows directly from Theorem 3 for the mixed community setting.

Our second lower bound is complementary, in that it captures the complexity of identifying the box containing the largest community. To state this bound, we define $H_2^b(D) = \sum_{i=2}^{b} \frac{1}{\log(N_1) - \log(N_1 - d_{11} + c_i)}$. Then, following along similar lines as the proof of Theorem 6, we can show that

$$\frac{H^b(D)}{2} \leq H_2^b(D) \leq \overline{log}(b)H^b(D).$$

18

**Theorem 13.** *In the community-disjoint box setting, consider any algorithm that only uses the number of distinct samples from each community to decide which box to sample from at each instant as well as to make the final estimate for the community mode. For any instance D, there exists an alternate instance $D^{[a]}$, $a \in [2:b]$, with $H_2^b(D^{[a]}) \leq H_2^b(D)$ such that*

$$\max\left(P_e(D), P_e(D^{[a]})\right) \geq \frac{1}{4} \exp\left(-\frac{t\Gamma}{H_2^b(D)}\right),$$

*where* $\Gamma = \max\left(\frac{\log\left(\lceil \frac{N_1(N_a-c_a+d_{11})}{(N_1-d_{11}+c_a)}\rceil\right)-\log(N_a)}{\log\left(\frac{N_1}{N_1-d_{11}+c_a}\right)}, \max_{i=2}^b \frac{\log\left(\lceil \frac{N_1(N_a-c_a+c_i)}{(N_1-d_{11}+c_i)}\rceil\right)-\log(N_a)}{\log\left(\frac{N_1}{N_1-d_{11}+c_a}\right)}\right)$. *The alternate instance* $D^{[a]}$ *is constructed by increasing the size of only the largest community in box a, such that the new size of box a is* $N_a' = \max\left(\lceil N_1\frac{(N_a-c_a+d_{11})}{(N_1-d_{11}+c_a)}\rceil, \max_{i=2}^b \lceil N_1\frac{(N_a-c_a+c_i)}{(N_1-d_{11}+c_i)}\rceil\right)$.

The proof of Theorem 13 follows along similar lines as the proof of Theorem 8. Details can be found in Appendix H.

Comparing the upper and lower bounds on the probability of error for the box setting in Theorems 10, 12, and 13, we see that the expressions for the exponents differ primarily in i) the presence of $H^b(D)$ vs $H_2^b(D)$, which differ by at most a factor of $\overline{log}(b)$; and ii) the presence of an additional factor $\Gamma$ in the lower bound. Note that

$$\max\left(\lceil N_1\frac{(N_a-c_a+d_{11})}{(N_1-d_{11}+c_a)}\rceil, \max_{i=2}^b \lceil N_1\frac{(N_a-c_a+c_i)}{(N_1-d_{11}+c_i)}\rceil\right) \leq \lceil\frac{N_1(N_a-c_a+d_{11})}{N_1-d_{11}+c_b}\rceil$$
$$\leq \frac{N_1(N_a-c_a+d_{11})}{N_1-d_{11}+c_b} + 1 \leq \frac{N_1(N_a-c_a+d_{11}+1)}{N_1-d_{11}+c_b}.$$

Using $\frac{x-1}{x} \leq \log(x) \leq x-1$ for all $x > 0$ and the above inequality, we get

$$\Gamma \leq \frac{\log(N_1) + \log(N_a-c_a+d_{11}+1) - \log(N_a) - \log(N_1-d_{11}+c_b)}{\log(N_1) - \log(N_1-d_{11}+c_a)}$$
$$= \frac{\log(N_1/(N_1-d_{11}+c_b)) + \log((N_a-c_a+d_{11}+1)/N_a)}{\log(N_1/(N_1-d_{11}+c_a))}$$
$$\leq \frac{(d_{11}-c_b)/(N_1-d_{11}+c_b) + (d_{11}-c_a+1)/N_a}{(d_{11}-c_a)/N_1}$$
$$\leq \frac{(d_{11}-c_b)}{(d_{11}-c_a)} \cdot \frac{N_1}{(N_1-d_{11}+c_b)} \cdot \frac{(2N_1+N_a)}{N_a}.$$

In particular, the above inequality implies that $\Gamma$ is bounded by a constant under the following natural assumptions on the class of underlying instances: i) the largest community size is at most a fraction of its corresponding box size, i.e., $d_{11} \leq (1-\delta_1)N_1$ for some $\delta_1 > 0$; ii) the size of the competing communities in other boxes is most a fraction of the largest community size, i.e., $c_a \leq (1-\delta_2)d_{11}$ for some $\delta_2 > 0$ $\forall a \neq 1$; and iii) all the box sizes are within a multiplicative constant factor $\beta$ of each other ($\beta > 1$). Under these assumptions, $\Gamma \leq \frac{2\beta+1}{\delta_1\delta_2}$.

We compare this lower bound to the first term in the upper bound given in Theorem 10. We note that these terms only differ by an order of $\overline{log}(b)\Gamma$. When $\Gamma$ is bounded from above, such as in the case described above, the DS-SR estimator matches the lower bound upto logarithmic factors for the problem of picking the correct box in the final stage of the algorithm, and is hence near-optimal. Comparing the second term in the upper bound from Theorem 10 to Theorem 12,

19

we find a similar logarithmic factor between the decay rates. Thus, the DS-SR algorithm is decay rate optimal up to logarithmic factors for the problem of picking the right community out of a box, given the correct box. This is natural and intuitive, due to its similarity with the mixed community DSM algorithm. Hence, the set of instances where DS-SR might not perform well in comparison to other algorithms can be characterized as instances where it is hard to pick the correct box containing the largest community; intuitively, these instances would produce a large value of the parameter $\Gamma$.

### 5.3. The general setting

Finally, we consider the most general setting, where communities are arbitrarily spread across boxes. From an algorithmic standpoint, the key challenge here is that it is no longer appropriate to eliminate boxes from consideration sequentially as in SR algorithms, since the largest community might be spread across multiple boxes. Accordingly, the algorithms we propose for the general setting are 'single phase' variants of the algorithms proposed in Section 5.1.

The single phase variant of Algorithm 3, which we refer to as the Distinct Samples Uniform Exploration (DS-UE) algorithm is stated as follows: sample each box $\lfloor t/b \rfloor$ times, and return the community that produces the largest number of distinct individuals. The probability of error under this algorithm can be bounded using the ideas we have used before, only the bounds are more cumbersome.

If the box sizes are known, one can also perform a single-phase proportional sampling of boxes, resulting effectively in a sampling process similar to the mixed community setting (except the budget is apportioned deterministically across boxes rather than the random allocation in the mixed community setting) . We refer to the corresponding algorithm, which outputs the community that produced the largest number of distinct individuals after $t$ queries, as the Distinct Samples Proportional Exploration (DS-PE) algorithm.

Finally, we state the normalized single phase variant of DS-UE, which we refer to as NDS-UE: Each box is sampled $\lfloor t/b \rfloor$ times, and the output of NDS-UE is the community that maximizes $\sum_i \frac{S_{ij}}{S_i} N_i$. ENDS-UE can be analogously defined.

To summarize, some of our algorithms for the disjoint box setting can indeed be applied and evaluated analytically in the general setting. However, we do not at present have a tight information theoretic lower bound for the general setting (or indeed, even for the disjoint box setting); the proof techinques we have used in the lower bounds for the mixed/separated community settings appear to be insufficient to handle the general case. So even though our algorithms for the general setting perform well in empirical evaluations (see Section 6), new methodological innovations are required to close the gap between upper and lower bounds.

## 6. Experimental Results

In this section, we present extensive simulation results comparing the performance of various algorithms discussed in the previous sections. We use both synthetic data as well as data gathered from real-world datasets for our experiments. For each experiment, we averaged the results over multiple runs (500-3000 depending on the complexity of the instance).

### 6.1. Mixed Community Mode Estimation

We begin with the mixed community setting studied in Section 3 where all individuals are placed in a single box. We demonstrate the difference in performance of the identity-less Sample

(a) Instance: [1000, 990, 600, 500, 500, 410]  (b) Instance: [1000, 900, 630, 520, 520, 430]

Figure 1: $\log(P_e(D))$ vs $t$ for mixed community setting



(a) Instance: [1000, 990, 600, 500, 500, 410]  (b) Instance: [1000, 900, 630, 520, 520, 430]

Figure 2: $\log(P_e(D))$ vs $t$ for separated community setting

Frequency Maximization (SFM) and the identity-based Distinct Samples Maximization (DSM) algorithms via simulations on synthetic data. We consider two instances, each with 4000 individuals in a single box, partitioned into communities as [1000, 990, 600, 500, 500, 410] and [1000, 900, 630, 520, 520, 430] respectively. As suggested by Theorems 2 and 4, we find that the difference in the convergence rates of the two estimators becomes more pronounced when the two largest communities are close in size. See Figure 1 where we plot the probability of error $\log(P_e)$ vs the query budget $t$ for the two instances.

## 6.2. Separated Community Mode Estimation

Next, we consider the separated community setting studied in Section 4 where each community is in a unique box. As above, we consider two instances with community sizes given by [1000, 990, 600, 500, 500, 410] and [1000, 900, 630, 520, 520, 430] respectively. We plot the performance of the Consecutive-Collision SR (CC-SR) and Distinct Samples SR (DS-SR) algorithms in Figure 2. As indicated by our results in Theorems 5 and 6, the DS-SR algorithm greatly outperforms the CC-SR algorithm.

21

*6.3. Community-Disjoint Box Mode Estimation*

Here, we look at the setting where the communities are partitioned across the boxes and thus each box can have multiple communities, as described in Section 5. We use the following two real-world datasets for comparing the performance of various estimators under this setting.

- Brazil Real Estate Dataset [20]: This dataset contains a total of 97353 apartment listings spread across 26 states and 3273 municipalities in Brazil. Mapping it to our framework, the apartments correspond to individual entities, the municipalities represent communities and the states they are located in denote the boxes. Our goal is to identify the municipality (community) with the largest number of listings by (randomly) sampling apartment listings from various states.

  Corresponding to this dataset, the four largest communities (municipalities with the most listed apartments) are of sizes [3929, 2322, 2414, 1876]. The top five box sizes are [80935, 3551, 2035, 1871, 1646], with the largest box corresponding to the state of Sao Paolo. Thus, one box has a much larger size than all others in this dataset and in fact, contains all of the the four largest communities.

- Airbnb Rental Listing Dataset [21]: This dataset contains a total of 48895 rental listings spread across 5 regions and 221 neighborhoods in New York city. Here, the apartments correspond to individual entities, the neighbourhoods represent communities and the broad regions they are located in denote the boxes.

  The top five communities (neighbourhoods) have sizes [3920, 3741, 2658, 2465, 1971]. The top 5 box sizes are [21661, 20104, 5666, 1091, 373]. Unlike the previous dataset, the two largest boxes (corresponding to Manhattan and Brooklyn respectively) are of comparable size here. Furthermore, the two boxes contain multiple competing communities of size comparable to the largest community. The largest box contains the communities with sizes 2658 and 1971, while the second largest box contains communities of sizes 3920 (mode), 3714, and 2465.

**Results** We compare the performance of the various algorithms discussed in Section 5.1 on the two datasets described above. These include the Distinct Samples-Successive Rejects (DS-SR) and its generalization Distinct Samples Proportional SR (DS-PSR) when the box sizes are known. We also consider the normalized variants of DS-SR, given by Normalized Distinct Samples SR (NDS-SR) and Expectation-Normalized Distinct Samples SR (ENDS-SR) when box sizes are known as well as Normalized Distinct Samples SR (NDS-SR (MLE)) when the box sizes are unknown, by replacing the box size by its maximum likelihood estimator.
Figure 3a shows the performance of the various algorithms on the Brazil Real Estate dataset. DS-SR which splits queries uniformly across all surviving boxes performs the worst while DS-PSR which does the division in proportion to box sizes performs the best. This is to be expected since there is one box which is much larger than all others and this box contains all of the competing largest communities. Thus, because of the uniform exploration in DS-SR, there might be fewer samples from the individual communities in the largest box in the initial rounds and it might get eliminated, which explains the poor performance for moderate query budgets. This shortcoming is addressed by DS-PSR which assign many more queries to the largest box which contains the community mode. The normalized variants NDS-SR and ENDS-SR also perform much better than DS-SR since they use the box sizes to determine the elimination criteria in each round. In

(a) Brazil Real Estate Dataset

(b) Airbnb Rental Listing Dataset

Figure 3: $\log(P_e(D))$ vs $t$ for box community setting



Figure 4: $\log(P_e(D))$ vs $t$ for the Youtube Video Dataset, General Box Setting

comparison to these, the NDS-SR (MLE) performs poorer for low query budget due to erroneous box size estimates but demonstrates similar performance for larger budgets.

Figure 3b shows the performance of the various algorithms on the Airbnb Apartment Listing dataset. Here again, DS-PSR performs the best since it allocates queries in proportion to box sizes. However, unlike the previous dataset, all the other algorithms have comparable performance. This includes DS-SR which does not use any box size information and is still able to perform better since the box sizes are relatively closer to each other for this dataset and the number of communities in each box are also fewer which makes it unlikely that the box containing the largest community is eliminated.

## 6.4. General Setting Mode Estimation

Finally, we consider the general setting where individuals in a community can be spread across multiple boxes. Section 5.3 described various single-round algorithms for this setting, namely the Distinct Samples Uniform Exploration (DS-UE) which doesn't need any box size informa-

tion and divides the query budget equally among all boxes; the Distinct Samples Proportional Exploration (DS-PE) which assigns queries in proportion to the box sizes; and the various normalized single phase variants of DS-UE, which we refer to as NDS-UE, ENDS-UE and NDS-UE (MLE). To compare the performance of these different estimators under the general setting, we use the following dataset.

- Trending Youtube Video Statistics Dataset [22]: This dataset contains the top trending videos for different regions such as Canada, US, and Japan, out of which we consider six regions. Mapped to our framework, a region corresponds to a box, a channel denotes a community, and each video represents an individual entity. The goal is to find the most popular channel which has the largest number of trending videos across the six regions. Note that a particular channel (community) can have trending videos (individuals) spread across different regions (boxes) and thus this dataset corresponds to the general setting. This dataset contains 239662 videos, each associated with one of 17773 channels. Top 5 channels have [870, 809, 752, 717, 712] top trending videos across regions. The boxes have comparable size, given by [40881, 40840, 40724, 38916, 37352, 40949].

Figure 4 shows the performance of the various algorithms on the above dataset. Note that all the estimators are able to achieve an exponential decay in the probability of error with the query budget even in this general setting. Furthermore, here the rate of decay for all the estimators is comparable since the box sizes are all similar and thus the knowledge of box sizes does not provide a distinct advantage. However, in terms of the absolute value, DS-UE performs slightly poorly as compared to the other algorithms which either use prior knowledge of box sizes or learn estimates for them using samples.

## Appendix A. Proof of Theorem 1

Let $\hat{\mu}_i(t)$ be the number of samples seen from $C_i$ over the horizon. We have

$$\hat{\mu}_i(t) = \sum_{j=1}^{t} \mathbb{1}_{\{\text{person } j \in C_i\}}$$

$$\Rightarrow E[\hat{\mu}_i(t)] = \mu_i(t) = \frac{td_i}{N}.$$

Using the union bound on $P_e(D)$, we get

$$P_e(D) \leq \sum_{i=2}^{m} P(\hat{\mu}_i(t) - \hat{\mu}_1(t) \geq 0).$$

The Chernoff bound gives us

$$P\left(\hat{\mu}_k(t) - \hat{\mu}_1(t) - (\mu_k(t) - \mu_1(t)) \geq w\right) \leq \min_{\lambda > 0} e^{-\lambda w} E\left[e^{\lambda(\hat{\mu}_k(t) - \hat{\mu}_1(t) - (\mu_k(t) - \mu_1(t)))}\right]$$

$$= \min_{\lambda > 0} e^{-\lambda[w + (\mu_k(t) - \mu_1(1))]} \left[\frac{d_k e^{\lambda}}{N} + \frac{d_1 e^{-\lambda}}{N} + \left(1 - \frac{d_1 + d_k}{N}\right)\right]^t.$$

Choosing $w = \mu_1(t) - \mu_k(t)$ and minimizing over $\lambda$,

$$P(\hat{\mu}_k(t) - \hat{\mu}_1(t) \geq 0) \leq \left[1 - \frac{(\sqrt{d_1} - \sqrt{d_k})^2}{N}\right]^t \qquad (A.1)$$

24

$$\Rightarrow P_e(D) \leq \sum_{i=2}^{m} P(\hat{\mu}_i(t) - \hat{\mu}_1(t) \geq 0) \leq \sum_{i=2}^{m} \left[ 1 - \frac{(\sqrt{d_1} - \sqrt{d_k})^2}{N} \right]^t$$

$$\leq (m-1) \left[ 1 - \frac{(\sqrt{d_1} - \sqrt{d_2})^2}{N} \right]^t.$$

## Appendix B. Proof of Theorem 2

To prove the theorem, we consider two instances $D = (d_1, d_2, \ldots, d_m)$ and $D' = (d'_1, d'_2, \ldots, d'_m)$, where the optimal community in $D$ is $C_1$ and the optimal community in $D'$ is $C_2$. We note that the mixed community setting can be modelled as a probability distribution over communities, with the probability of sampling $C_i$ under $D$ and $D'$ being $p_i = d_i/N$ and $p'_i = d'_i/N$ respectively. Let the probability distributions corresponding to instances $D$ and $D'$ be $\Theta = (p_1, p_2, \ldots p_m)$ and $\Theta' = (p'_1, p'_2, \ldots p'_m)$ respectively. Further, let the sequence of $t$ samples be denoted by $X_1, X_2, \ldots, X_t$ where $X_i$ is the index of the community that is sampled at time $i$, and let $\mathbb{P}_\Theta, \mathbb{P}_{\Theta'}$ denote the probability measures induced on the sample sequence by the instances $D, D'$. Next, we state a few lemmas which will help in the proof of the theorem.

**Lemma 14.** *For every event $\mathcal{E} \in F_t$, where $F_t = \sigma(X_1, X_2, \ldots X_t)$,*

$$\mathbb{P}_{\Theta'}(\mathcal{E}) = \mathbb{E}_\Theta[\mathbb{1}_\mathcal{E} \exp(-L_t)],$$

*where $L_t = \sum_{i=1}^{t} \log \left( \frac{p_{X_i}}{p'_{X_i}} \right)$ and $\mathbb{1}$ is the indicator random variable.*

*Proof.* This is analogous to [17, Lemma 18]. $\square$

**Lemma 15.** *For every event $\mathcal{E} \in F_t$,*

$$\mathbb{E}_\Theta[L_t|\mathcal{E}] \geq \log \frac{\mathbb{P}_\Theta(\mathcal{E})}{\mathbb{P}_{\Theta'}(\mathcal{E})}.$$

*Proof.* From Lemma 14, we know that $\mathbb{P}_{\Theta'}(\mathcal{E}) = \mathbb{E}_\Theta[\exp(-L_t)\mathbb{1}_\mathcal{E}]$. Then, using Jensen's inequality on $\exp(-x)$, we have that

$$\mathbb{P}_{\Theta'}(\mathcal{E}) = \mathbb{E}_\Theta[\exp(-L_t)\mathbb{1}_\mathcal{E}] = \mathbb{E}_\Theta[\mathbb{E}_\Theta[\exp(-L_t)|\mathbb{1}_\mathcal{E}]\mathbb{1}_\mathcal{E}] \geq \mathbb{E}_\Theta[\exp(-\mathbb{E}_\Theta[L_t|\mathcal{E}])\mathbb{1}_\mathcal{E}]$$

$$= \exp(-\mathbb{E}_\Theta[L_t|\mathcal{E}])\mathbb{P}_\Theta(\mathcal{E})$$

The last line above proves the lemma. $\square$

**Lemma 16.** *If $d(x, y) = x \log \left( \frac{x}{y} \right) + (1-x) \log \left( \frac{(1-x)}{(1-y)} \right)$, then for every event $\mathcal{E} \in F_t$,*

$$\mathbb{E}_{\Theta'}[-L_t] \geq d(\mathbb{P}_{\Theta'}(\mathcal{E}), \mathbb{P}_\Theta(\mathcal{E})).$$

*Proof.* From Lemma 15 we know that

$$\mathbb{E}_{\Theta'}[-L_t|\mathcal{E}] \geq \log \left( \frac{\mathbb{P}_{\Theta'}(\mathcal{E})}{\mathbb{P}_\Theta(\mathcal{E})} \right), \mathbb{E}_{\Theta'}[-L_t|\mathcal{E}^c] \geq \log \left( \frac{\mathbb{P}_{\Theta'}(\mathcal{E}^c)}{\mathbb{P}_\Theta(\mathcal{E}^c)} \right).$$

Using the total law of probability and the above inequality, we get

$$\mathbb{E}_{\Theta'}[-L_t] = \mathbb{E}_{\Theta'}[-L_t|\mathcal{E}]\mathbb{P}_{\Theta'}(\mathcal{E}) + \mathbb{E}_{\Theta'}[-L_t|\mathcal{E}^c]\mathbb{P}_{\Theta'}(\mathcal{E}^c) \geq d(\mathbb{P}_{\Theta'}(\mathcal{E}), \mathbb{P}_\Theta(\mathcal{E}^c)).$$

$\square$

Consider a consistent algorithm $\mathcal{A}$, and let $P_e(D)$ and $P_e(D')$ denote the probabilities of error for $\mathcal{A}$ under the instances $D$ and $D'$ respectively. Denote the community that is output by $\mathcal{A}$ as $\hat{h}^*$, and let $S$ be the event that $\hat{h}^* = 1$. Thus, $P_e(D) = 1 - \mathbb{P}_\Theta(S)$ and $P_e(D') \geq \mathbb{P}_{\Theta'}(S)$. Since algorithm $\mathcal{A}$ is consistent and thus its probability of error on both $D, D'$ goes to zero as the number of samples $t$ grows large, we have that for every $\epsilon > 0$ there exists $t_0(\epsilon)$ such that for all $t \geq t_0(\epsilon), \mathbb{P}_{\Theta'}(S) \leq \epsilon \leq \mathbb{P}_\Theta(S)$. For $t \geq t_0(\epsilon)$,

$$\mathbb{E}_{\Theta'}[-L_t] \geq d(\mathbb{P}_{\Theta'}(S), \mathbb{P}_\Theta(S)) \geq d(\epsilon, \mathbb{P}_\Theta(S)) \geq \epsilon \log\left(\frac{\epsilon}{\mathbb{P}_\Theta(S)}\right) + (1 - \epsilon) \log\left(\frac{(1-\epsilon)}{P_e(D)}\right)$$

$$\geq \epsilon \log(\epsilon) + (1 - \epsilon) \log\left(\frac{(1-\epsilon)}{P_e(D)}\right)$$

Taking the limsup, using $\mathbb{E}_{\Theta'}[-L_t] = t.D(\Theta'\|\Theta)$ where $D(\cdot\|\cdot)$ denotes the Kullback-Leibler divergence, and letting $\epsilon \to 0$, we get

$$\limsup_{t\to\infty} -\frac{1}{t} \log(P_e(D)) \leq D(\Theta'\|\Theta).$$

Consider $\Theta = (p_1, p_2, ...p_m)$ and $\Theta' = (\frac{\sqrt{p_1 p_2} - \delta}{C}, \frac{\sqrt{p_1 p_2} + \delta}{C}, \frac{p_3}{C}, ... \frac{p_m}{C})$, where $C = 1 - (\sqrt{p_1} - \sqrt{p_2})^2$ and $\delta > 0$ is sufficiently small so that $\Theta'$ is a probability distribution. Then, we get

$$\limsup_{t\to\infty} -\frac{1}{t} \log(P_e(D)) \leq \log\left(\frac{1}{C}\right) + \left(\frac{\sqrt{p_1 p_2} - \delta}{C}\right) \log\left(\frac{\sqrt{p_1 p_2} - \delta}{p_1}\right) + \left(\frac{\sqrt{p_1 p_2} + \delta}{C}\right) \log\left(\frac{\sqrt{p_1 p_2} + \delta}{p_2}\right)$$

$$\implies \limsup_{t\to\infty} -\frac{1}{t} \log(P_e(D)) \leq \log\left(\frac{1}{C}\right) \text{ (letting } \delta \downarrow 0).$$

## Appendix C. Proof of Theorem 3

We will begin by proving the first assertion in the theorem statement which provides an upper bound on the probability of error for $t \leq \min\left\{\frac{d_1 + d_m}{2d_1}N, \frac{16Nd_1}{(d_1 - d_m)^2}\right\}$. Let $S_i(t)$ denote the number of distinct samples seen from community $C_i$ in $t$ samples. We have the following lemma:

**Lemma 17.** *The probability of error of the DSM algorithm is bounded as*

$$P_e(D) \leq \sum_{i=2}^{m} P(S_i(t) - S_1(t) > 0) + \frac{1}{2}P(S_i(t) = S_1(t)).$$

*Proof.* For any $i \in 2, 3, \ldots, m$, it is clear that when $S_i(t) - S_1(t) > 0$, DSM will erroneously output $i$ as the index of the community mode. Furthermore, since DSM breaks ties arbitrarily, with some positive probability (bounded by $1/2$) it makes the same error when $S_i(t) = S_1(t)$. Together with the union bound over all $i \in 2, 3, \ldots, m$, this gives the above result. □

Next, for each $k \in \{2, 3, \ldots, m\}$ let $Z_k$ be the random variable denoting the number of samples observed from communities $C_1$ and $C_k$.[9] We note that the expected value of $Z_k$ is given by

$$E[Z_k] = \frac{(d_1 + d_k)t}{N}. \tag{C.1}$$

---

[9]Note that $Z_k$ corresponds to the total number of samples from communities $C_1$ and $C_k$, not necessarily distinct.

Define events $E_{k1} = \{Z_k \in [(1 - \epsilon_k)E[Z_k], (1 + \epsilon_k)E[Z_k]]\}$ and $E_{k2} = E_{k1}^c$, with

$$\epsilon_k = \frac{\sqrt{\frac{9}{64}\beta_k^4 + \frac{3}{2}\beta_k^2} - \frac{3}{8}\beta_k^2}{2} \text{ where } \beta_k = \frac{d_1 - d_k}{d_1 + d_k}. \tag{C.2}$$

It is easy to verify that $\beta_k < 1$ and $\epsilon_k \leq \min\{\beta_k, 1/2\}$. Then, we have

$$P(S_k(t) - S_1(t) > 0) + \frac{1}{2}P(S_k(t) = S_1(t))$$

$$\leq P(S_k(t) - S_1(t) > 0|E_{k1})P(E_{k1}) + P(S_k(t) - S_1(t) > 0|E_{k2})P(E_{k2})$$

$$+ \frac{1}{2}P(S_k(t) = S_1(t)|E_{k1})P(E_{k1}) + \frac{1}{2}P(S_k(t) = S_1(t)|E_{k2})P(E_{k2})$$

$$\leq P(S_k(t) - S_1(t) \geq 0|E_{k1})P(E_{k1}) + P(S_k(t) - S_1(t) > 0|E_{k2})P(E_{k2}) + \frac{1}{2}P(S_k(t) = S_1(t)|E_{k2})P(E_{k2}). \tag{C.3}$$

Note that the LHS above appears for each $k \in \{2, 3, \dots, m\}$ in the upper bound on $P_e(D)$ in Lemma 17. We will bound the terms in the RHS separately, and then combine them together to get an overall upper bound on $P_e(D)$. To begin with, note that

$$E[S_i(t)|Z_k] = d_i\left[1 - \left(1 - \frac{1}{d_1 + d_k}\right)^{Z_k}\right], \text{ for } i \in \{1, k\}. \tag{C.4}$$

We consider the function $f(x_1, x_2, x_3, \dots, x_t) = S_k(t) - S_1(t)$ where $x_i$ is the identity of the individual sampled at the $i$-th instant. Note that for any $i \in \{1, 2, \dots, t\}$ and for all $x_1, x_2, x_3, \dots, x_t, x_i' \in \{1, 2, \dots, N\}$, we have $|f(x_1, x_2, \dots, x_i, \dots, x_t) - f(x_1, x_2, \dots, x_i', \dots, x_t)| \leq c_i \triangleq 2\mathbb{1}_{x_i \in C_1 \cup C_k}$. Then, conditioning on $Z_k$ and applying McDiarmid's inequality, we get

$$P(f - E[f|Z_k] \geq t'|Z_k) \leq P(|f - E[f|Z_k]| \geq t'|Z_k) \leq \exp\left(-\frac{2t'^2}{\sum_{i=1}^{t} c_i^2}\right) = \exp\left(-\frac{t'^2}{2Z_k}\right).$$

Plugging in $t' = -E[f|Z_k]$, and computing $E[f|Z_k]$ using Equation (C.4), we obtain

$$P(f \geq 0|Z_k) = P(S_k(t) - S_1(t) \geq 0|Z_k) \leq exp\left(-\frac{(d_1 - d_k)^2\left[1 - \left(1 - \frac{1}{d_1 + d_k}\right)^{Z_k}\right]^2}{2Z_k}\right). \tag{C.5}$$

We will start with deriving an upper bound on the first term in the RHS of equation (C.3) given by $P(S_k(t) - S_1(t) \geq 0|E_{k1})P(E_{k1})$. Conditioned on the event $E_{k1}$, we have $Z_k \in [(1 - \epsilon_k)E[Z_k], (1 + \epsilon_k)E[Z_k]]$. Furthermore, from the statement of the first part of the theorem statement and the definitions of $\epsilon_k, \beta_k$ from equation (C.2), we have the following sequence of assertions:

$$t \leq \frac{d_1 + d_k}{2d_1}N \Rightarrow \beta_k = \frac{d_1 - d_k}{d_1 + d_k} \leq \frac{N}{t} - 1 \Rightarrow \epsilon_k \leq \frac{N}{t} - 1 \Rightarrow Z_k \leq (1 + \epsilon_k)\frac{t(d_1 + d_k)}{N} \leq d_1 + d_k.$$

Using the above inequalities and the Taylor series expansion, we have

$$\left[1 - \left(1 - \frac{1}{d_1 + d_k}\right)^{Z_k}\right] \geq \left[\frac{Z_k}{d_1 + d_k} - \frac{Z_k^2}{2(d_1 + d_k)^2}\right] \geq \frac{Z_k}{2(d_1 + d_k)}. \tag{C.6}$$

27

Plugging the bound above in equation (C.5), and using $Z_k \geq (1 - \epsilon_k)E[Z_k] = (1 - \epsilon_k)(d_1 + d_k)t/N$, we have

$$P(S_k(t) - S_1(t) \geq 0|E_{k1}) \times P(E_{k1}) \leq P(S_k(t) - S_1(t) \geq 0|E_{k1}) \leq exp\left(-\frac{t(1 - \epsilon_k)(d_1 - d_k)^2}{8N(d_1 + d_k)}\right),$$
(C.7)

thus giving us an upper bound on the first term in the RHS of equation (C.3).

For bounding the sum of the second and third terms in the RHS of equation (C.3), we use the following lemma:

**Lemma 18.** *For any $k \in \{2, 3, \ldots, m\}$ so that $d_k \leq d_1$ and for any $l \geq 0$, we have*

$$P(S_k(t) - S_1(t) > 0|Z_k = l) + \frac{1}{2}P(S_k(t) = S_1(t)|Z_k = l) \leq \frac{1}{2}$$

*Proof.* Note that the theorem statement is equivalent to showing that, when $d_k \leq d_1$,

$$P(S_k(t) - S_1(t) > 0|Z_k = l) \leq P(S_k(t) - S_1(t) < 0|Z_k = l),$$

which says that, conditioned on the total number of samples from communities 1 and $k$ together being some fixed $l$, the likely event is that the community 1, whose size is at least that of community $k$, will have as many or more distinct individuals than community $k$. Given $d_k \leq d_1$, this is intuitive and while it can be argued formally, we skip the argument here for brevity. $\square$

Using Lemma 18, we get that the second and third terms in the RHS of equation (C.3) are bounded as

$$P(S_k(t) - S_1(t) > 0|E_{k2})P(E_{k2}) + \frac{1}{2}P(S_k(t) = S_1(t)|E_{k2})P(E_{k2}) \leq \frac{1}{2}P(E_{k2}).$$

Further, using Chernoff's inequality for $P(E_{k2})$ and $E[Z_k] = (d_1 + d_k)t/N$, we have

$$\frac{1}{2}P(E_{k2}) = \frac{1}{2}P(|Z_k - E[Z_k]| > \epsilon_k) \leq \exp\left(-\frac{\epsilon_k^2(d_1 + d_k)t}{3N}\right).$$
(C.8)

Finally, combining Lemma 17, equation (C.7), and equation (C.8), we get the following upper bound on $P_e(D)$.

$$P_e(D) \leq \sum_{k=2}^{m} \exp\left(-\frac{t(1 - \epsilon_k)(d_1 - d_k)^2}{8N(d_1 + d_k)}\right) + \exp\left(-\frac{\epsilon_k^2(d_1 + d_k)t}{3N}\right).$$

From the value of $\epsilon_k$ in equation (C.2), we have that the exponents in the two terms of the summation above are equal. Thus, we have

$$P_e(D) \leq \sum_{k=2}^{m} 2\exp\left(-\frac{t(1 - \epsilon_k)(d_1 - d_k)^2}{8N(d_1 + d_k)}\right) \leq \sum_{k=2}^{m} 2\exp\left(-\frac{t(d_1 - d_k)^2}{16N(d_1 + d_k)}\right) \leq \sum_{k=2}^{m} 2\exp\left(-\frac{t(d_1 - d_k)^2}{32Nd_1}\right),$$
(C.9)

where the first inequality is true because $\epsilon_k \leq 1/2$; and the second inequality follows since $d_k \leq d_1$ for all $k \in \{2, 3, \ldots, m\}$.

The next result comments on the shape of the function $f(x) = \exp(-\frac{t(d_1 - x)^2}{32Nd_1})$, which appears in equation (C.9) above.

**Lemma 19.** *The function $f(x) = \exp\left(-\frac{t(d_1-x)^2}{32Nd_1}\right)$ is concave for any $x \geq d_m$ and $t \leq \frac{16Nd_1}{(d_1-d_m)^2}$.*

*Proof.* We differentiate $f(x)$ twice to confirm that it is concave.

$$f''(x) = \frac{t}{16Nd_1}\exp\left(-\frac{t(d_1-x)^2}{32Nd_1}\right)\left(\frac{t}{16Nd_1}(d_1-x)^2 - 1\right)$$

Using the inequality $t \leq \frac{16Nd_1}{(d_1-d_m)^2}$, we have that

$$f''(x) \leq \frac{t}{16Nd_1}\exp\left(-\frac{t(d_1-x)^2}{32Nd_1}\right)\left(\frac{(d_1-x)^2}{(d_1-d_m)^2} - 1\right)$$

which implies $f''(x) \leq 0$ since $x \geq d_m$. $\qquad\square$

From (C.9) and using Lemma 19, we have from Jensen's inequality that for $t \leq \min\left\{\frac{d_1+d_m}{2d_1}N, \frac{16Nd_1}{(d_1-d_m)^2}\right\}$

$$P_e(D) \leq 2\sum_{k=2}^{m}\exp\left(-\frac{t(d_1-d_k)^2}{32Nd_1}\right) \leq 2(m-1)\exp\left(-\frac{t\left(d_1-\frac{\sum_{k=2}^{m}d_i}{m-1}\right)^2}{32Nd_1}\right),$$

which proves the first assertion in the theorem statement.

For the second assertion in the theorem statement, note that the algorithm will certainly not make an error if the number of distinct individuals seen from the $i$-th community, $S_i(t) \geq d_2 + 1$, where $d_2$ denotes the size of the second-largest community. Hence, the probability of error is bounded as $P_e(D) \leq P(S_1(t) \leq d_2)$. Further, note that if the event $\{S_1(t) \leq d_2\}$ occurs, then there exists a set of $d_1 - d_2$ individuals in $C_1$ which remain unsampled in the $t$ samples. Thus, we have

$$P_e(D) \leq P(S_1(t) \leq d_2) \leq \binom{d_1}{d_2}\left(1 - \frac{d_1 - d_2}{N}\right)^t.$$

## Appendix D. Proof of Theorem 4

This proof is similar in spirit to the proof of [23, Theorem 1]. Consider an instance $D = (d_1, d_2, \ldots, d_m)$. First, we note that since $(S_j(t))$, $1 \leq j \leq m$) is a sufficient statistic for $D$, it suffices to restrict attention to (consistent) algorithms whose output depends only on the vector $(S_j(t)$, $1 \leq j \leq m)$. Given this restriction, we track the temporal evolution of the vector $S(k) = (S_j(k)$, $1 \leq j \leq m)$, where $S_j(k)$ is the number of distinct individuals from community $j$ seen in the first $k$ oracle queries. This evolution can be modeled as an absorbing Markov chain over state space $\prod_{j=1}^{m}\{0, 1, \cdots d_i\}$, with $S(0) = (0, 0, \cdots, 0)$. Next, let us write down the transition probabilities $q_D(s, s')$ for each state pair $(s, s')$. Note that from state $s$, the chain can transition to the states $s + e_j$ for $1 \leq j \leq m$, where the vector $e_j$ has 1 in the $j$th position and 0 elsewhere, or remain in state $s$. Moreover, $q_D(s, s + e_j) = (d_j - s_j)/N$, and $q_D(s, s) = \frac{\sum_{j=1}^{m} s_j}{N}$. Recall that by assumption, community 1 is the largest community for the instance $D$. Let us consider an alternate instance $D' = (d'_1, d'_2, \ldots, d'_m)$ such that $d'_1 = d_2 - 1$, $d'_j = d_j \; \forall j \neq 1$, and

29

$N' = N - d_1 + d_2 - 1$. Note that the community mode under the alternate instance $D'$ is different from that under the original instance $D$. Thus, for state $s$ that is feasible under both $D$ and $D'$,

$$\log\left(\frac{q_{D'}(s, s)}{q_D(s, s)}\right) = \log\left(\frac{N}{N - d_1 + d_2 - 1}\right).$$

Similarly, for state pair $(s, s + e_j)$ that is feasible under both $D$ and $D'$,

$$\log\left(\frac{q_{D'}(s, s + e_j)}{q_D(s, s + e_j)}\right) = \log\left(\frac{N}{N - d_1 + d_2 - 1}\right), j \neq 1,$$

$$\log\left(\frac{q_{D'}(s, s + e_1)}{q_D(s, s + e_1)}\right) = \log\left(\frac{N(d_2 - 1 - s_1)}{(N - d_1 + d_2 - 1)(d_1 - s_1)}\right) = \log\left(\frac{N}{N - d_1 + d_2 - 1}\right) + \log\left(\frac{d_2 - 1 - s_1}{d_1 - s_1}\right).$$

Therefore, for any state pair $(s, s')$ such that $q_D(s, s'), q_{D'}(s, s') > 0$, we have

$$\log\left(\frac{q_{D'}(s, s')}{q_D(s, s')}\right) \leq \log\left(\frac{N}{N - d_1 + d_2 - 1}\right). \tag{D.1}$$

Next, let $\mathbb{P}_D, \mathbb{P}_{D'}$ denote the probability measures induced by the algorithm under consideration under the instances $D$ and $D'$, respectively. Then, given a state evolution sequence $(S(1), \cdots, S(t))$, the log-likelihood ratio is given by

$$\log\frac{\mathbb{P}_{D'}(S(1), \cdots, S(t))}{\mathbb{P}_D(S(1), \cdots, S(t))} = \sum_{s,s'} N(s, s', t) \log\left(\frac{q_{D'}(s, s')}{q_D(s, s')}\right),$$

where $N(s, s', t)$ represents the number of times the transition from state $s$ to state $s$ occurs over the course of $t$ queries. Combining with (D.1), we get

$$\log\frac{\mathbb{P}_{D'}(S(1), \cdots, S(t))}{\mathbb{P}_D(S(1), \cdots, S(t))} \leq t \log\left(\frac{N}{N - d_1 + d_2 - 1}\right),$$

which implies

$$D(\mathbb{P}_{D'}\|\mathbb{P}_D) = E_{D'}\left[\log\frac{\mathbb{P}_{D'}(S(1), \cdots, S(t))}{\mathbb{P}_D(S(1), \cdots, S(t))}\right] \leq t \log\left(\frac{N}{N - d_1 + d_2 - 1}\right), \tag{D.2}$$

where $D(\cdot\|\cdot)$ denotes the Kullback-Leibler divergence. On the other hand, since the algorithm produces an estimate $\hat{h}^*$ of the community mode based solely on $S(t)$, we have from the data-processing inequality (see [24]) that

$$D(\mathbb{P}_{D'}\|\mathbb{P}_D) \geq D(Ber(\mathbb{P}_{D'}(\hat{h}^* = 1))\|Ber(\mathbb{P}_D(\hat{h}^* = 1))), \tag{D.3}$$

where $Ber(x)$ denotes the Bernoulli distribution with parameter $x \in (0, 1)$. Recall that the community mode under $D$ is community 1, while it is community 2 under $D'$. Then from the definition of consistent algorithms, for every $\epsilon > 0$, $\exists t_0(\epsilon)$ such that for $t \geq t_0(\epsilon)$, $\mathbb{P}_{D'}(\hat{h}^* = 1) \leq \epsilon \leq \mathbb{P}_D(\hat{h}^* = 1)$. Thus, we have

$$D(Ber(\mathbb{P}_{D'}(\hat{h}^* = 1))\|Ber(\mathbb{P}_D(\hat{h}^* = 1))) \geq D(Ber(\epsilon)\|Ber(\mathbb{P}_D(\hat{h}^* = 1)))$$

$$\geq \epsilon \log\left(\frac{\epsilon}{\mathbb{P}_D(\hat{h}^* = 1)}\right) + (1 - \epsilon) \log\left(\frac{1 - \epsilon}{\mathbb{P}_D(\hat{h}^* \neq 1)}\right) \geq \epsilon \log(\epsilon) + (1 - \epsilon) \log\left(\frac{1 - \epsilon}{\mathbb{P}_D(\hat{h}^* \neq 1)}\right).$$

Using $\epsilon \to 0$ and $\mathbb{P}_D(\hat{h}^* \neq 1) = P_e(D)$, we have

$$D(Ber(\mathbb{P}_{D'}(\hat{h}^* = 1))\|Ber(\mathbb{P}_D(\hat{h}^* = 1))) \geq -\log(P_e(D)).$$

Finally, combining with (D.2) and (D.3), we have that

$$\liminf_{t \to \infty} \frac{\log(P_e(D))}{t} \geq -\log\left(\frac{N}{N - (d_1 - d_2 + 1)}\right).$$

## Appendix E. Proof of Theorem 5

Note that

$$P_e(D) \leq \sum_{r=1}^{b-1} P(C_1 \text{ gets eliminated in round } r).$$

Let $S_i(K)$ denote the number of (immediate pairwise) collisions recorded in $C_i$ after $K$ pairs of samples. Since at least one of the smallest $r$ communities is guaranteed to be present during round $r$,

$$P_e(D) \leq \sum_{r=1}^{b-1} \sum_{j=b+1-r}^{b} P(S_j(K_r) - S_1(K_r) \leq 0)$$

$$\leq \sum_{r=1}^{b-1} rP(S_{b+1-r}(K_r) - S_1(K_r) \leq 0). \tag{E.1}$$

Denoting, for $i \neq 1$, $f_i(K) := S_i(K) - S_1(K)$, we now derive an upper bound on $P(f_i(K) \leq 0)$. Applying Chernoff's inequality, for $\lambda \leq 0$,

$$P(f_i(K) \leq 0) \leq E\left[e^{\lambda f_i(K)}\right]$$

$$= \left[\frac{1}{d_1 d_i} + \left(1 - \frac{1}{d_1}\right)\left(1 - \frac{1}{d_i}\right) + e^\lambda \left(1 - \frac{1}{d_1}\right)\frac{1}{d_i} + e^{-\lambda}\left(1 - \frac{1}{d_i}\right)\frac{1}{d_1}\right]^K.$$

Setting $e^\lambda = \sqrt{\frac{d_i - 1}{d_1 - 1}}$,

$$P(f_i(K) \leq 0) \leq \left(1 - \frac{(\sqrt{d_1 - 1} - \sqrt{d_i - 1})^2}{d_1 d_i}\right)^K \leq \exp\left(-\frac{K(\sqrt{d_1 - 1} - \sqrt{d_i - 1})^2}{d_1 d_i}\right).$$

Since $d_1 > d_i$, $(\sqrt{d_1 - 1} - \sqrt{d_i - 1})^2 > \frac{((d_1-1)-(d_i-1))^2}{4(d_1-1)} > \frac{((d_1-1)-(d_i-1))^2}{4d_1} = \frac{(d_1-d_i)^2}{4d_1}$.

$$\Rightarrow P(f_i(K) \leq 0) \leq \exp\left(-\frac{K(d_1 - d_i)^2}{4d_1^2 d_i}\right).$$

Substituting the above into (E.1),

$$P_e(D) \leq \sum_{r=1}^{b-1} r \exp\left(-\frac{K_r(d_1 - d_{b+1-r})^2}{4d_1^2 d_{b+1-r}}\right).$$

31

Since $K_r = \left\lceil \frac{1}{\overline{\log}(b)} \frac{t/2-b}{b+1-r} \right\rceil$, where $\overline{\log}(b) = \frac{1}{2} + \sum_{i=2}^{b} \frac{1}{i}$ and $\Delta_i = \frac{1}{d_i} - \frac{1}{d_1}$,

$$P_e(D) \le \sum_{r=1}^{b-1} r \exp\left(-\frac{K_r d_{b+1-r} \Delta_{(b+1-r)}^2}{4}\right).$$

For $H^c(D) = \max_{i \in [2:b]} \frac{i \Delta_i^{-2}}{d_i}$,

$$K_r d_{b+1-r} \Delta_{(b+1-r)}^2 \ge \frac{(t/2-b)}{\overline{\log}(b) H^c(D)}$$

$$\Rightarrow P_e(D) \le \frac{b(b-1)}{2} \exp\left(-\frac{(t/2-b)}{4\overline{\log}(b) H^c(D)}\right).$$

## Appendix F.  Proof of Theorem 10

Let $P_e^i(D)$ denote the probability of the community mode being eliminated at the $i$th step; i.e, for $i \le b-1$, $P_e^i(D)$ denotes the probability of removing box 1 in phase $i$ of SR, and $P_e^b(D)$ denotes the probability of choosing the wrong community from box 1 after this box survived the $(b-1)$ SR phases. Then, we have

$$P_e(D) = \sum_{i=1}^{b-1} P_e^i(D) + P_e^b(D),$$

$$P_e^i(D) \le \binom{d_{11}}{c_{b-i+1}} \exp\left(-K_i \log\left(\frac{N_1}{N_1 - d_{11} + c_{b-i+1}}\right)\right) \quad (1 \le i \le b-1),$$

$$P_e^b(D) \le \binom{d_{11}}{c_1} \exp\left(-K_{b-1} \log\left(\frac{N_1}{N_1 - d_{11} + c_1}\right)\right),$$

where the second and third statements are based on a coupon collector argument, similar to the one employed in the proof of Theorem 6 for the separated community setting. The proof is now completed by substituting the values of $K_r$, and using the definition of $H^b(D)$.

## Appendix G.  Proof of Theorem 11

We show that ENDS-SR has the same decay rate as DS-SR. Recall that the comparison function used in ENDS-SR is

$$\frac{S_{ij} N_i}{E[S_i]},$$

where $S_{ij}$ is the number of distinct samples from community $i$ in box $j$, and $S_i$ is the number of distinct samples from box $i$. At the end of $r$ rounds,

$$E[S_i] = N_i \left(1 - \left(1 - \frac{1}{N_i}\right)^{K_r}\right).$$

32

Similar to the coupon collector argument in the proof of Theorem 10, we let $P_e^i(D)$ be the probability of the community mode being eliminated in the $i$th step. We have that

$$P_e(D) \le \sum_{i=1}^{b} P_e^i(D).$$

After $r \le b - 1$ rounds/phases, the comparison function for the largest community equals

$$\frac{S_{11}}{\left(1 - \left(1 - \frac{1}{N_1}\right)^{K_r}\right)}.$$

For some community $j$ in box $i$, the comparison function is

$$\frac{S_{ij}}{\left(1 - \left(1 - \frac{1}{N_i}\right)^{K_r}\right)} \le \frac{c_i}{\left(1 - \left(1 - \frac{1}{N_m}\right)^{K_r}\right)},$$

where $N_m = \max_i N_i$. Thus, if we have

$$S_{11} > \frac{c_{b-r+1}\left(1 - \left(1 - \frac{1}{N_1}\right)^{K_r}\right)}{\left(1 - \left(1 - \frac{1}{N_m}\right)^{K_r}\right)},$$

then the community mode cannot be eliminated in the $r$th round. For round $r = b$, we just note that

$$S_{11} > c_1$$

is sufficient for the community mode estimate to be correct. Applying the coupon collector argument on these events, by using the notation

$$f_i(K) := \frac{c_i\left(1 - \left(1 - \frac{1}{N_1}\right)^{K}\right)}{\left(1 - \left(1 - \frac{1}{N_m}\right)^{K}\right)},$$

we have

$$P_e(D) \le \sum_{i=1}^{b-1} \binom{d_{11}}{f_{b-i+1}(K_i)} \exp\left(-K_i \log\left(\frac{N_1}{N_1 - d_{11} + f_{b-i+1}(K_i)}\right)\right) + \binom{d_{11}}{c_1} \exp\left(-K_b \log\left(\frac{N_1}{N_1 - d_{11} + c_1}\right)\right).$$

We note that, as $t \to \infty$, $f_i(t) \to c_i$, which then implies the statement of the theorem.

## Appendix H. Proof of Theorem 13

We first state the following lemma (analogous to Lemma 9) for this setting (the proof is straightforward and omitted):

**Lemma 20.** *For any algorithm $\mathcal{A}$ and instance $D$, there must exist a box $a \in [2:b]$ such that $E_D[N_a(t)] \leq \frac{t}{(\log(N_1) - \log(N_1 - d_{11} + c_a))H_2^b(D)}$, where $N_a(t)$ denotes the number of times box $a$ is sampled in $t$ queries under $\mathcal{A}$.*

*Proof of Theorem 13.* Given an instance $D$, we construct an alternate instance $D^{[a]}$ by changing the size of the largest community in box $a$ (corresponding to the one specified by Lemma 20) from $c_a$ to $g_a' = c_a + N_a' - N_a$. [10] Note that the size of box $a$ changes from $N_a$ to $N_a' = N_a + g_a' - c_a$. Furthermore, we can see that the community mode under instance $D^{[a]}$ is different from the one under the original instance $D$, since

$$g_a' = c_a + N_a' - N_a \geq c_a + \frac{N_1(N_a - c_a + d_{11})}{(N_1 - d_{11} + c_a)} - N_a > c_a + (N_a - c_a + d_{11}) - N_a = d_{11}.$$

Following steps similar to the proof of Theorem 8, we get

$$D(\mathbb{P}_D, \mathbb{P}_{D^{[a]}}) \leq E_D[N_a(t)] \log\left(\frac{N_a'}{N_a}\right).$$

From the definition of $\Gamma$, it follows that $\frac{N_a'}{N_a} = \left(\frac{N_1}{N_1 - d_{11} + c_a}\right)^{\Gamma}$. Thus, invoking Lemma 20, we have

$$D(\mathbb{P}_D, \mathbb{P}_{D^{[a]}}) \leq \frac{t\Gamma}{H_2^b(D)}.$$

Finally, similar to the proof of Theorem 8, we use Lemma 21 to get

$$\max\left(P_e(D), P_e(D^{[a]})\right) \geq \frac{1}{4} \exp\left(-\frac{t\Gamma}{H_2^b(D)}\right)$$

which matches the statement of the theorem.

Finally, we show that

$$H_2^b(D^{[a]}) \leq H_2^b(D) \Leftrightarrow \sum_{i \neq a} \frac{1}{\log(N_a') - \log(N_a' - g_a' + c_i')} \leq \sum_{i \neq 1} \frac{1}{\log(N_1) - \log(N_1 - d_{11} + c_i)}$$

We do this in two steps:

1. Firstly, for each $i \notin \{1, a\}$, we show that the term corresponding to box $i$ in the sum on the left is smaller than the corresponding term in the sum on the right, i.e.,

$$\frac{1}{\log(N_a') - \log(N_a' - g_a' + c_i')} \leq \frac{1}{\log(N_1) - \log(N_1 - d_{11} + c_i)}$$

or equivalently, $\quad \frac{N_1}{N_1 - d_{11} + c_i} \leq \frac{N_a'}{N_a' - g_a' + c_i'}.$

This follows from the following sequence of inequalities.

$$\frac{N_1}{(N_1 - d_{11} + c_i)} = \frac{N_1(N_a - c_a + c_i)}{(N_1 - d_{11} + c_i)(N_a - c_a + c_i)} \leq \frac{N_a'}{N_a - c_a + c_i} = \frac{N_a'}{N_a' - g_a' + c_i'}$$

where the last step follows since $N_a' = N_a + g_a' - c_a$ and $c_i' = c_i$ for $i \notin \{1, a\}$.

---

[10] We use $g_a'$ and not $c_a'$ to denote the new size of this community because in the alternate instance $D^{[a]}$, this community is the largest community, and is thus no longer the *competing* community in box $a$.

2. Secondly, we show that the term corresponding to box 1 in the sum on the left is smaller than the term corresponding to box $a$ in the sum on the right, i.e,

$$\frac{1}{\log(N'_a) - \log(N'_a - g'_a + c'_1)} \leq \frac{1}{\log(N_1) - \log(N_1 - d_{11} + c_a)}$$

or equivalently,     $\frac{N_1}{(N_1 - d_{11} + c_a)} \leq \frac{N'_a}{N'_a - g'_a + d_{11}}.$

This follows from the following sequence of inequalities.

$$\frac{N_1}{N_1 - d_{11} + c_a} = \frac{N_1(N_a - c_a + d_{11})}{(N_a - c_a + d_{11})(N_1 - d_{11} + c_a)} \leq \frac{N'_a}{N_a - c_a + d_{11}} = \frac{N'_a}{N'_a - g'_a + d_{11}},$$

where the last step is true because $N'_a - g'_a = N_a - c_a$.

This completes the proof.

$\square$

## Appendix I.  Other Lemmas

**Lemma 21.** *Let $\rho_0$ and $\rho_1$ be two probability distributions supported on some set $\chi$, with $\rho_1$ absolutely continuous with respect to $\rho_0$. Then for any measurable function $\phi : \chi \to \{0, 1\}$,*

$$P_{X \sim \rho_0}(\phi(X) = 1) + P_{X \sim \rho_1}(\phi(X) = 0) \geq \frac{1}{2} \exp(-D(\rho_0 \| \rho_1))$$

*Proof.*  This is [17, Lemma 20].     $\square$

**Lemma 22.** $\frac{H(D)}{2} \leq H_2(D) \leq \overline{\log}(b)H(D).$

*Proof.*  For the inequality on the left, we note that

$$H_2(D) = \sum_{i=2}^{b} \frac{1}{\log(d_1) - \log(d_i)} \geq \sum_{i=2}^{j} \frac{1}{\log(d_1) - \log(d_i)} \geq \frac{j-1}{\log(d_1) - \log(d_j)} \forall j \in [2 : b]$$

Since this is true for all $j \in [2 : b]$, taking the max of these values and using that $j - 1 \geq \frac{j}{2}, j \geq 2$ we have

$$H_2(D) \geq \max_{j \neq 1} \frac{j/2}{\log(d_1) - \log(d_j)} = \frac{H(D)}{2}$$

For the inequality on the right, we multiply and divide each term in the summation of $H_2(D)$ by $i$:

$$H_2(D) = \sum_{i=2}^{b} \frac{i}{i(\log(d_1) - \log(d_i))} \leq \sum_{i=2}^{b} \frac{H(D)}{i} \leq \overline{\log}(b)H(D)$$

This completes the proof of both inequalities in the statement of the lemma.     $\square$

# References

[1] M. Finkelstein, H. G. Tucker, J. A. Veeh, Confidence intervals for the number of unseen types, Statistics & Probability Letters 37 (4) (1998) 423–430.

[2] C. Budianu, S. Ben-David, L. Tong, Estimation of the number of operating sensors in large-scale sensor networks with mobile access, IEEE Transactions on Signal Processing 54 (5) (2006) 1703–1715.

[3] M. Bressan, E. Peserico, L. Pretto, Simple set cardinality estimation through random sampling, arXiv preprint arXiv:1512.07901 (2015).

[4] X. Chen, W. Huang, W. Chen, J. C. Lui, Community exploration: from offline optimization to online learning, in: Proceedings of the 32nd International Conference on Neural Information Processing Systems, 2018.

[5] S. Bubeck, D. Ernst, A. Garivier, Optimal discovery with probabilistic expert advice: finite time analysis and macroscopic optimality, Journal of Machine Learning Research 14 (2013) 601–623.

[6] T. Lattimore, C. Szepesvári, Bandit algorithms, Cambridge University Press, 2020.

[7] J.-Y. Audibert, S. Bubeck, Best Arm Identification in Multi-Armed Bandits, in: COLT - 23th Conference on Learning Theory - 2010, Haifa, Israel, 2010, p. 13 p.
URL https://hal-enpc.archives-ouvertes.fr/hal-00654404

[8] C. Caferov, B. Kaya, R. O'Donnell, A. C. Say, Optimal bounds for estimating entropy with pmf queries, in: International Symposium on Mathematical Foundations of Computer Science, Springer, 2015, pp. 187–198.

[9] J. Acharya, A. Orlitsky, A. T. Suresh, H. Tyagi, Estimating rényi entropy of discrete distributions, IEEE Transactions on Information Theory 63 (1) (2016) 38–56.

[10] Y. Hao, A. Orlitsky, Data amplification: Instance-optimal property estimation, arXiv preprint arXiv:1903.01432 (2019).

[11] Y. Wu, P. Yang, Sample complexity of the distinct elements problem, Mathematical Statistics and Learning 1 (1) (2018) 37–72.

[12] Y. Hao, A. Orlitsky, Unified sample-optimal property estimation in near-linear time, in: Advances in Neural Information Processing Systems, 2019, pp. 11104–11114.

[13] H. Chernoff, Estimation of the mode, Annals of the Institute of Statistical Mathematics 16 (1) (1964) 31–41.

[14] E. Parzen, On estimation of a probability density function and mode, The Annals of Mathematical Statistics 33 (3) (1962) 1065–1076.

[15] D. Shah, T. Choudhury, N. Karamchandani, A. Gopalan, Sequential mode estimation with oracle queries, Proceedings of the AAAI Conference on Artificial Intelligence 34 (04) (2020) 5644–5651. doi:10.1609/aaai.v34i04.6018.
URL https://ojs.aaai.org/index.php/AAAI/article/view/6018

[16] Z. Karnin, T. Koren, O. Somekh, Almost optimal exploration in multi-armed bandits, in: S. Dasgupta, D. McAllester (Eds.), Proceedings of the 30th International Conference on Machine Learning, Vol. 28 of Proceedings of Machine Learning Research, PMLR, Atlanta, Georgia, USA, 2013, pp. 1238–1246.
URL https://proceedings.mlr.press/v28/karnin13.html

[17] E. Kaufmann, O. Cappé, A. Garivier, On the Complexity of Best Arm Identification in Multi-Armed Bandit Models, Journal of Machine Learning Research 17 (2016) 1–42.

[18] A. Carpentier, A. Locatelli, Tight (lower) bounds for the fixed budget best arm identification bandit problem, in: V. Feldman, A. Rakhlin, O. Shamir (Eds.), 29th Annual Conference on Learning Theory, Vol. 49 of Proceedings of Machine Learning Research, PMLR, Columbia University, New York, New York, USA, 2016, pp. 590–604.
URL https://proceedings.mlr.press/v49/carpentier16.html

[19] Z. Karnin, T. Koren, O. Somekh, Almost optimal exploration in multi-armed bandits, in: International Conference on Machine Learning, 2013, pp. 1238–1246.

[20] Properati, Real estate listings - brazil, https://data.world/properati/real-estate-listings-brazil, accessed: 2021-05-24 (2016).

[21] M. Cox, Inside airbnb - new york city, http://insideairbnb.com/get-the-data.html, accessed: 2021-05-24 (2021).

[22] M. Jolly, Trending youtube video statistics, https://www.kaggle.com/datasnaek/youtube-new, accessed: 2021-05-24 (2019).

[23] V. Moulos, Optimal best markovian arm identification with fixed confidence, in: H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, R. Garnett (Eds.), Advances in Neural Information Processing Systems, 2019.

[24] T. M. Cover, J. A. Thomas, Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing), Wiley-Interscience, USA, 2006.