

The Optimality of Correlated Sampling

Mohammad Bavarian* Badih Ghazi† Elad Haramaty‡
 Pritish Kamath§ Ronald L. Rivest¶ Madhu Sudan||

July 31, 2020

Abstract

In the *correlated sampling* problem, two players are given probability measures P and Q respectively, over the same measurable space and access to shared randomness. Without any interaction, the two players are each required to output an element sampled according to their respective measures, while trying to minimize the probability that their outputs disagree. A well-known strategy due to Kleinberg & Tardos and Holenstein, with a close variant (for a similar problem) due to Broder, solves this task with disagreement probability at most $2\delta/(1 + \delta)$, where δ is the total variation distance between P and Q . This strategy has been used in several different contexts including sketching algorithms, approximation algorithms based on rounding linear programming relaxations, the study of parallel repetition and cryptography.

In this paper, we give a surprisingly simple proof that this strategy is in fact tight. Specifically, for every $\delta \in (0, 1)$, we show that any correlated sampling strategy should have disagreement probability at least $2\delta/(1 + \delta)$. This partially answers a recent question of Rivest.

Our proof is based on studying a new problem that we call *constrained agreement*. Here, the two players are given subsets $A \subseteq [n]$ and $B \subseteq [n]$ respectively and their goal is to output an element $i \in A$ and $j \in B$ respectively while minimizing the probability that $i \neq j$. We prove tight bounds for this question, which in turn imply tight bounds for correlated sampling. Though we settle basic questions about the two problems, our formulation leads to more fine-grained questions that remain open.

*Department of Mathematics and Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge MA 02139. Supported in part by NSF Award CCF-1420692. bavarian@mit.edu.

†Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge MA 02139. Supported in part by NSF CCF-1420956, NSF CCF-1420692 and CCF-1217423. badih@mit.edu.

‡Harvard John A. Paulson School of Engineering and Applied Sciences. Part of this work supported by NSF Award CCF-1565641. seladh@gmail.com.

§Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge MA 02139. Supported in part by NSF CCF-1420956 and NSF CCF-1420692. prish@mit.edu.

¶Institute Professor, MIT. This work supported by the Center for Science of Information (CSoI), an NSF Science and Technology Center, under grant agreement CCF-0939370.

||Harvard John A. Paulson School of Engineering and Applied Sciences. Part of this work supported by NSF Award CCF-1565641 and a Simons Investigator Award. madhu@cs.harvard.edu.

1 Introduction

In this work, we study *correlated sampling*, a very basic task, variants of which have been considered in the context of sketching algorithms [Bro97], approximation algorithms based on rounding linear programming relaxations [KT02, Cha02], the study of parallel repetition [Hol07, Rao11, BHH⁺08] and very recently cryptography [Riv16].

The *correlated sampling problem* is defined as follows: Alice and Bob are given probability measures P and Q respectively over a measurable space Ω . They also have access to shared randomness, modeled as a suitably chosen probability space \mathcal{R} . Without any interaction, Alice is required to output an element a distributed according to P and Bob is required to output an element b distributed according to Q . Their goal is to minimize the disagreement probability $\Pr[a \neq b]$, which we will compare against the total variation distance $d_{\text{TV}}(P, Q)$ defined as

$$d_{\text{TV}}(P, Q) = \frac{1}{2} \cdot \int |P(\omega) - Q(\omega)| d\omega.$$

More formally, a *correlated sampling strategy* is defined as follows, where we use Δ_Ω to denote the set of all probability measures over Ω .

Definition 1.1. A correlated sampling strategy for a measurable space Ω with error $\varepsilon : [0, 1] \rightarrow [0, 1]$ is specified by a probability space \mathcal{R} and a pair of functions $(f : \Delta_\Omega \times \mathcal{R} \rightarrow \Omega, g : \Delta_\Omega \times \mathcal{R} \rightarrow \Omega)$, measurable in their second argument, such that for all $P, Q \in \Delta_\Omega$ with $d_{\text{TV}}(P, Q) \leq \delta$, it holds that,

$$\begin{aligned} \text{[Correctness]} \quad & \{f(P, r)\}_{r \sim \mathcal{R}} = P \text{ and } \{g(Q, r)\}_{r \sim \mathcal{R}} = Q, \\ \text{[Error Guarantee]} \quad & \Pr_{r \sim \mathcal{R}} [f(P, r) \neq g(Q, r)] \leq \varepsilon(\delta). \end{aligned}$$

In the above, we used $\{f(P, r)\}_{r \sim \mathcal{R}}$ to denote the pushforward measure of $f(P, \cdot)$ under the probability measure in \mathcal{R} . For simplicity, we will often not mention \mathcal{R} explicitly when talking about correlated sampling strategies. While we define correlated sampling for general measurable spaces, we will mostly deal with Ω that is a finite discrete space.

A correlated sampling strategy is notably different from the notion of a *coupling* (cf. [Tho00] for an introduction), where we are required to have a *single* coupling function $h : \Delta_\Omega \times \Delta_\Omega \times \mathcal{R} \rightarrow \Omega \times \Omega$ such that for any probability measures P and Q it holds that $\{h(P, Q, r)_1\}_{r \sim \mathcal{R}} = P$ and $\{h(P, Q, r)_2\}_{r \sim \mathcal{R}} = Q$. In other words, a coupling function has the knowledge of both measures P and Q , whereas a correlated sampling strategy operates locally on the knowledge of P and Q . It is well known that for any coupling function h , it holds that $\Pr_{r \sim \mathcal{R}} [h(P, Q, r)_1 \neq h(P, Q, r)_2] \geq d_{\text{TV}}(P, Q)$ and that this bound is achievable. Since a correlated sampling strategy is a special case of coupling, it follows that $\varepsilon(\delta) \geq \delta$, yet a priori, it is unclear whether any non-trivial correlated sampling strategy can even exist, since the error ε is not allowed to depend on the size of Ω .

Somewhat surprisingly, there exists a simple strategy whose error can be bounded by roughly twice the total variation distance (and in particular does not degrade with the size of Ω). Variants of this strategy have been rediscovered multiple times in the literature yielding the following theorem.

Theorem 1.2 ([Bro97, KT02, Hol07]). For all finite discrete spaces Ω , there exists a correlated sampling strategy with error $\varepsilon : [0, 1] \rightarrow [0, 1]$ such that,

$$\forall \delta \in [0, 1], \quad \varepsilon(\delta) \leq \frac{2 \cdot \delta}{1 + \delta}. \tag{1}$$

Strictly speaking, the work of Broder [Bro97] does not consider the general correlated sampling problem. Rather it gives a strategy (the “MinHash strategy”) which can be interpreted as a correlated sampling strategy for the special case where P and Q are *flat* distributions, i.e., they are uniform over some subsets of Ω . In particular, if $P = \mathcal{U}(A)$ and $Q = \mathcal{U}(B)$ are distributions that are uniform over sets $A, B \subseteq \Omega$ respectively, then the MinHash strategy gives an error probability of $1 - \frac{|A \cap B|}{|A \cup B|}$, also known as the *Jaccard distance* between A and B . In the special case when $|A| = |B|$, this is equivalent to the bound above.

The technique can be generalized to other (non-flat) distributions to get the bound in [Theorem 1.2](#), thereby yielding a strategy due to Kleinberg & Tardos and Holenstein.¹ Several variants of this (sometimes referred to as “consistent sampling” protocols) have been used in several applied works, e.g., [Man94, GP06, MMT10, HMT14].

Given [Theorem 1.2](#), a natural and basic question is whether the bound on the error can be improved; the only lower bound we are aware of is the one that arises from coupling, namely $\varepsilon(\delta) \geq \delta$. This question was very recently raised by Rivest [Riv16] in the context of a new encryption scheme and was one of the motivations for this work. We give a surprisingly simple proof that the bound in [Theorem 1.2](#) is actually tight!

Theorem 1.3 (Main Result). *For all $\delta, \gamma \in (0, 1)$, there exists a (sufficiently large) finite discrete space Ω for which any correlated sampling strategy with error $\varepsilon : [0, 1] \rightarrow [0, 1]$ satisfies*

$$\varepsilon(\delta) \geq \frac{2 \cdot \delta}{1 + \delta} - \gamma. \tag{2}$$

Organization of the paper. In [Section 2](#), we prove [Theorem 1.3](#). In [Section 3](#), we consider the setting where Ω is of a fixed finite size, which was the question originally posed by Rivest [Riv16]. In this regime, there turns out to be a surprising strategy that gets better error than [Theorem 1.2](#) in a very special case. However, it was conjectured in [Riv16] that in fact a statement like [Theorem 1.3](#) holds in every other case and we make progress on this conjecture by proving it in one such case. We conclude with some more observations and open questions in [Section 4](#). Finally, for completeness, we describe in [Appendix A](#), the correlated sampling strategies of Broder, Kleinberg & Tardos and Holenstein underlying [Theorem 1.2](#).

Acknowledgements. We would like to thank anonymous reviewers for their feedback that has helped improve the presentation of this paper.

2 Lower Bound on Correlated Sampling

In order to prove [Theorem 1.3](#), we first introduce the following *constrained agreement* problem which is relaxation of the correlated sampling problem. Alice and Bob are given sets $A \subseteq \Omega$ and $B \subseteq \Omega$ respectively, where the pair (A, B) is sampled from some (known) distribution \mathcal{D} . Alice and Bob are required to output elements $a \in A$ and $b \in B$ respectively, such that the disagreement probability $\Pr_{(A,B) \sim \mathcal{D}}[a \neq b]$ is minimized.

¹strictly speaking, if P and Q are flat over different sized subsets, the above bound is weaker than that obtained from a direct application of the MinHash strategy! See [Section 4](#) for a discussion.

This can be viewed as a relaxation of the correlated sampling problem by first considering the case of *flat* distributions in [Definition 1.1](#) and relaxing the restrictions of $\{f(P, r)\}_{r \sim \mathcal{R}} = P$ and $\{g(Q, r)\}_{r \sim \mathcal{R}} = Q$ to only requiring that $f(P, r) \in \text{supp}(P)$ and $g(Q, r) \in \text{supp}(Q)$ for all $r \in \mathcal{R}$. This makes it a constraint satisfaction problem and we consider a distributional version of the same.

In the following definition, we use 2^Ω to denote the powerset of Ω .

Definition 2.1. A constrained agreement strategy for a finite discrete space Ω and a probability measure \mathcal{D} over $2^\Omega \times 2^\Omega$ is specified by a pair of functions $(f : 2^\Omega \rightarrow \Omega, g : 2^\Omega \rightarrow \Omega)$ with error $\text{err}_{\mathcal{D}}(f, g) :=$ smallest $\varepsilon \in [0, 1]$, such that the following hold

$$\text{[Correctness]} \quad \forall A, B \subseteq \Omega : f(A) \in A \text{ and } g(B) \in B,$$

$$\text{[Error guarantee]} \quad \Pr_{(A, B) \sim \mathcal{D}} [f(A) \neq g(B)] \leq \varepsilon.$$

Note that since the constrained agreement problem is defined with respect to a (known) probability measure \mathcal{D} on pairs of sets, we can require without loss of generality, that the strategies (f, g) be deterministic (since any randomized strategy can be derandomized with no degradation in the error).

In order to prove [Theorem 1.3](#), we characterize the optimal constrained agreement strategy in the special case when $\mathcal{D} = \mathcal{D}_p$ where every element $\omega \in \Omega$ is independently included in each of A and B with probability p .

Lemma 2.2. For all $p \in [0, 1]$, any constrained agreement strategy (f, g) for a finite discrete space Ω and distribution $\mathcal{D} = \mathcal{D}_p$ over $2^\Omega \times 2^\Omega$, any has error $\text{err}_{\mathcal{D}_p}(f, g) \geq \frac{2(1-p)}{2-p}$.

Proof. For ease of notation, let $\Omega = [n]$. Let (f, g) be a constrained agreement strategy. We will construct functions f^* and g^* such that $\text{err}_{\mathcal{D}_p}(f, g) \geq \text{err}_{\mathcal{D}_p}(f^*, g^*) \geq \frac{2(1-p)}{2-p}$.

For every $i \in [n]$, let $\beta_i := \Pr_B[g(B) = i]$. Without loss of generality (by suitably permuting $[n]$), we can assume that $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$. Since A and B are independently sampled in \mathcal{D}_p , it follows that when Bob's strategy is fixed to g , the strategy of Alice that results in the largest agreement probability is simply $f^*(A) := \arg\max_{i \in A} \beta_i = \min\{i : i \in A\}$ for all $A \subseteq [n]$.

So far we have $\text{err}_{\mathcal{D}_p}(f, g) \geq \text{err}_{\mathcal{D}_p}(f^*, g)$. We can repeat the same process again. For every $i \in [n]$, define $\alpha_i := \Pr_A[f^*(A) = i]$. Due to the specific choice of f^* , it holds that $\alpha_i = (1-p)^{i-1}p$ and hence $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$. Thus, when Alice's strategy is fixed to f^* , the strategy of Bob that results in the largest agreement probability is given by $g^*(B) = \arg\max_{i \in B} \alpha_i = \min\{i : i \in B\}$ for all $B \subseteq [n]$.

Thus, we get $\text{err}_{\mathcal{D}_p}(f, g) \geq \text{err}_{\mathcal{D}_p}(f^*, g) \geq \text{err}_{\mathcal{D}_p}(f^*, g^*)$ where

$$\begin{aligned} \text{err}_{\mathcal{D}_p}(f^*, g^*) &:= 1 - \Pr_{(A, B) \sim \mathcal{D}_p} [f^*(A) = g^*(B)] \\ &= 1 - \sum_{i=1}^n \Pr_A[f^*(A) = i] \cdot \Pr_B[g^*(B) = i] \\ &= 1 - \sum_{i=1}^n (1-p)^{2(i-1)} \cdot p^2 \\ &\geq 1 - \frac{p}{2-p} = \frac{2(1-p)}{2-p}, \end{aligned}$$

Thus, we conclude that $\text{err}_{\mathcal{D}_p}(f, g) \geq \frac{2(1-p)}{2-p}$. □

Before turning to the proof of [Theorem 1.3](#), we note a simple fact.

Fact 2.3. For flat distributions $P = \mathcal{U}(A)$ and $Q = \mathcal{U}(B)$ with $A, B \subseteq \Omega$, it holds that,

$$d_{\text{TV}}(P, Q) = 1 - \frac{|A \cap B|}{\max\{|A|, |B|\}}$$

Proof of [Theorem 1.3](#). Fix $\delta, \gamma \in (0, 1)$. Assume, for the sake of contradiction, that for all finite discrete spaces Ω there is a correlated sampling strategy (f^*, g^*) with error $\varepsilon(\delta) < \frac{2\cdot\delta}{1+\delta} - \gamma$. Let $\delta' \in (0, \delta)$ be such that

$$\frac{2\cdot\delta}{1+\delta} - \gamma < \frac{2\cdot\delta'}{1+\delta'} < \frac{2\cdot\delta}{1+\delta}. \quad (3)$$

Consider the distribution \mathcal{D}_p over pairs (A, B) of subsets $A, B \subseteq [n]$ where each $i \in [n]$ is independently included in each of A and B with probability $p := 1 - \delta'$. We then have that $\mathbb{E}[|A|] = \mathbb{E}[|B|] = p \cdot n$, and $\mathbb{E}[|A \cap B|] = p^2 \cdot n$. Moreover, by the Chernoff bound, we have that

$$\Pr_A[|A| - p \cdot n > p \cdot n^{0.8}] \leq 2 \cdot e^{-p \cdot n^{0.6}/3},$$

$$\Pr_B[|B| - p \cdot n > p \cdot n^{0.8}] \leq 2 \cdot e^{-p \cdot n^{0.6}/3},$$

and

$$\Pr_{A,B}[|A \cap B| - p^2 \cdot n > p^2 \cdot n^{0.8}] \leq 2 \cdot e^{-p^2 \cdot n^{0.6}/3}.$$

Hence, by the union bound and using $p^2 \leq p$, we get that with probability at least $1 - 6 \cdot e^{-p^2 \cdot n^{0.6}/3}$, we have that $||A| - p \cdot n| \leq pn^{0.8}$, $||B| - p \cdot n| \leq pn^{0.8}$ and $||A \cap B| - p^2 \cdot n| \leq p^2 n^{0.8}$. Thus, for the distributions $P = \mathcal{U}(A)$ and $Q = \mathcal{U}(B)$, it holds with probability at least $1 - 6 \cdot e^{-p^2 \cdot n^{0.6}/3}$ that

$$\begin{aligned} d_{\text{TV}}(P, Q) &= 1 - \frac{|A \cap B|}{\max\{|A|, |B|\}} \\ &\leq 1 - p + o_n(1) \\ &= \delta' + o_n(1) \\ &< \delta \quad \text{for sufficiently large } n. \end{aligned}$$

The assumed strategy (f^*, g^*) is such that $\Pr_{r \sim \mathcal{R}}[f(P, r) \neq g(Q, r)] \leq \frac{2\delta}{(1+\delta)} - \gamma$ when $d_{\text{TV}}(P, Q) \leq \delta$ and at most 1 otherwise. In our random choice of the distribution pair (P, Q) , the probability of $d_{\text{TV}}(P, Q) > \delta$ is at most $o_n(1)$. Thus, $\Pr_{(P,Q), r \sim \mathcal{R}}[f(P, r) \neq g(Q, r)] \leq \frac{2\cdot\delta}{1+\delta} - \gamma + o_n(1)$ when applied on the randomly sampled (P, Q) . In particular, by averaging, there exists a deterministic constrained agreement strategy with no worse disagreement probability. That is,

$$\exists (f, g), \quad \text{err}_{\mathcal{D}_p}(f, g) \leq \frac{2\cdot\delta}{1+\delta} - \gamma + o_n(1). \quad (4)$$

But from [Lemma 2.2](#) we have that,

$$\forall (f, g), \quad \text{err}_{\mathcal{D}_p}(f, g) \geq \frac{2(1-p)}{2-p} = \frac{2\cdot\delta'}{1+\delta'} \quad (5)$$

Putting [Equations \(4\)](#) and [\(5\)](#) together contradicts [Equation \(3\)](#) for sufficiently large n . \square

3 Correlated Sampling over a Fixed Finite Sized Universe

While we seem to have proved the optimality of the correlated sampling strategy, [Theorem 1.3](#) requires the universe to be of sufficiently large size. In particular, it does not say that the strategy underlying [Theorem 1.2](#) is optimal for a fixed finite sized universe. The quest for understanding optimality in this setting was motivated by the new encryption scheme proposed by Rivest [[Riv16](#)]. But as we will see shortly, this quest is not entirely straightforward!

In order to elaborate on this, it will be useful to formally define restricted versions of the correlated sampling strategy which are required to work only when the input pair (P, Q) is promised to lie in a given relation $\mathcal{G} \subseteq \Delta_\Omega \times \Delta_\Omega$.

Definition 3.1. For a finite discrete space Ω and a relation $\mathcal{G} \subseteq \Delta_\Omega \times \Delta_\Omega$, a \mathcal{G} -restricted correlated sampling strategy with error ε is specified by a probability space \mathcal{R} , a pair of functions $(f : \Delta_\Omega \times \mathcal{R} \rightarrow \Omega, g : \Delta_\Omega \times \mathcal{R} \rightarrow \Omega)$ if the following hold for all distribution pairs $(P, Q) \in \mathcal{G}$,

$$\text{[Correctness]} \quad \{f(P, r)\}_{r \sim \mathcal{R}} = P \text{ and } \{g(Q, r)\}_{r \sim \mathcal{R}} = Q,$$

$$\text{[Error Guarantee]} \quad \Pr_{r \sim \mathcal{R}} [f(P, r) \neq g(Q, r)] \leq \varepsilon.$$

For the rest of this section, we will consider a special kind of \mathcal{G} -restriction corresponding to Alice and Bob having “flat distributions”.

Definition 3.2. For the discrete space $\Omega = [n]$, the relation $\mathcal{G}_{a,b,\ell}^n \subseteq \Delta_{[n]} \times \Delta_{[n]}$ is defined to consist of “flat” distribution pairs (P, Q) corresponding to sets $A, B \subseteq [n]$ such that $P = \mathcal{U}(A)$, $Q = \mathcal{U}(B)$ and $|A| = a$, $|B| = b$, $|A \cap B| = \ell$. (For the relation to be non-empty, it is required that $\ell \leq \min\{a, b\}$ and $a + b - \ell \leq n$.)

Recall from [Fact 2.3](#), that for all $(P, Q) \in \mathcal{G}_{a,b,\ell}^n$ with $P = \mathcal{U}(A)$ and $Q = \mathcal{U}(B)$, is given by

$$d_{\text{TV}}(P, Q) = 1 - \frac{|A \cap B|}{\max\{|A|, |B|\}} = 1 - \frac{\ell}{\max\{a, b\}}.$$

Moreover, the MinHash strategy applied on input pairs $(P, Q) \in \mathcal{G}_{a,b,\ell}^n$ has a disagreement probability

$$1 - \frac{|A \cap B|}{|A \cup B|} = 1 - \frac{\ell}{a + b - \ell}.$$

One might suspect that this is optimal for all values of n, a, b and ℓ . But rather surprisingly, in the very special case where $|A \cap B| = 1$ and $|A \cup B| = n$, Rivest [[Riv16](#)] gave a strategy with smaller error probability than the above! While we don’t know of any applications for this strategy itself, its purpose here is to illustrate that there can be strategies which do better than the MinHash strategy in some special cases.

Theorem 3.3 ([\[Riv16\]](#)). For all $a, b \in \mathbb{Z}_{\geq 1}$ there exists a $\mathcal{G}_{a,b,1}^{a+b-1}$ -restricted correlated sampling strategy with error at most $1 - 1/\max\{a, b\}$.

For completeness, we describe this strategy in [Section 3.1](#). Note that for $(P, Q) \in \mathcal{G}_{a,b,1}^{a+b-1}$,

$$d_{\text{TV}}(P, Q) = 1 - \frac{1}{\max\{a, b\}} < 1 - \frac{1}{a + b - 1}.$$

This naturally leads to the question: *Is there a better correlated sampling strategy for larger intersection sizes?* In fact MinHash strategy was conjectured to be optimal in every other case (i.e. $\ell > 1$) by Rivest [[Riv16](#)] and this is necessary for proving the security of his proposed encryption scheme.

Conjecture 3.4 ([Riv16]). For every collection of positive integers $n \geq a, b \geq \ell \geq 2$ such that $n \geq a + b - \ell$, any $\mathcal{G}_{a,b,\ell}^n$ -restricted correlated sampling strategy makes error at least $1 - \ell/(a + b - \ell)$.

As partial progress towards this conjecture, we prove that in the other extreme (as compared to [Theorem 3.3](#)), the above conjecture does hold. In particular, we show the following theorem in [Section 3.2](#).

Theorem 3.5. For all $a = b \geq 1, \ell = a - 1$ and $n \geq a + 1$, any $\mathcal{G}_{a,b,\ell}^n$ -restricted correlated sampling strategy makes error at least $1 - \ell/(a + b - \ell)$.

3.1 Correlated Sampling Strategy of Rivest [Riv16]

In order to prove [Theorem 3.3](#), we recall the well-known Hall's Theorem.

Lemma 3.6 (Hall; cf. [vLW01]). Fix any bipartite graph G on vertex sets L and R (with $|L| \leq |R|$). There exists a matching that entirely covers L if and only if for every subset $S \subseteq L$, we have that $|S| \leq |N_G(S)|$, where $N_G(S)$ denotes the set of neighbors in G of vertices in S .

Proof of Theorem 3.3. First, let's consider the case where $a = b$. Let $\binom{[n]}{a}$ denote the set of all subsets $A \subseteq [n]$ with $|A| = a$. Consider the bipartite graph G on vertices $\binom{[n]}{a} \times \binom{[n]}{a}$, with an edge between vertices A and B if $|A \cap B| = 1$. It is easy to see that G is a -regular (since $n = 2a - 1$). Iteratively using [Lemma 3.6](#), we get that the edges of G can be written as a disjoint union of a matchings. Let's denote these as M_1, M_2, \dots, M_a .

The $\mathcal{G}_{a,a,1}^{2a-1}$ -restricted correlated sampling strategy of Alice and Bob is as follows: Use the shared randomness to sample a random index $r \in [a]$ and consider the matching M_r . If (A, B') is the edge present in M_r , then Alice outputs the unique element in $A \cap B'$. Similarly, if (A', B) is the edge present in M_r , then Bob outputs the unique element in $A' \cap B$. This strategy is summarized in [Algorithm 1](#).

Algorithm 1: Rivest's strategy [Riv16]

Alice's input: $A \subseteq [n]$

Bob's input: $B \subseteq [n]$

\mathcal{G} -restriction: $|A| = |B| = a, |A \cap B| = 1$ and $A \cup B = [n]$ i.e. $n = a + b - 1$

Pre-processing: Let G be the bipartite graph on vertices $\binom{[n]}{a} \times \binom{[n]}{a}$, with an edge between vertices A and B if $|A \cap B| = 1$. Decompose the edges of G into a disjoint matchings M_1, \dots, M_a .

Shared randomness: Index $r \sim \mathcal{U}([a])$

Strategy:

- Let (A, B') and (A', B) be edges present in M_r .
- $f(A, r) :=$ unique element in $A \cap B'$.
- $g(B, r) :=$ unique element in $A' \cap B$.

It is easy to see that both Alice's and Bob's outputs are uniformly distributed in A and B respectively. Moreover, the probability that they output the same element, is exactly $1/a$, which is

the probability of choosing the unique matching M_r which contains the edge (A, B) (i.e. enforcing $A' = A$ and $B' = B$).

The strategy in the general case of $a \neq b$ is obtained by a simple reduction to the case above. Suppose w.l.o.g. that $a > b$. Alice and Bob get sets $A \subseteq [n]$ and $B \subseteq [n]$ such that $|A| = a$, $|B| = b$ and $|A \cap B| = 1$ and $A \cup B = [n]$. We extend the universe by adding $(a - b)$ dummy elements to get a universe of size $(2a - 1)$ (note, $n = a + b - 1$). Moreover, whenever Bob gets set B , he extends it to B' by adding all the dummy elements to B and thus $|B'| = a$ while having $|A \cap B'| = 1$ and $|A \cup B'| = 2a - 1$. Now, Alice and Bob can use the $\mathcal{G}_{a,a,1}^{2a-1}$ -restricted correlated sampling strategy from above on the input pair (A, B') . This achieves an error of $1 - 1/a = 1 - 1/\max\{a, b\}$. However, Bob's output is uniformly distributed over B' and not B . To fix this, Bob can simply output a uniformly random element of B whenever the above strategy requires him to return an element of $B' \setminus B$. It is easy to see that this doesn't change the error probability. \square

3.2 Proof of Theorem 3.5

Proof of Theorem 3.5. Let $A, B \subseteq [n]$ be such that $a = |A| = |B| = |A \cap B| + 1$ and let $P = \mathcal{U}(A)$ and $Q = \mathcal{U}(B)$. For simplicity, we can assume without loss of generality that $A \cup B = [n]$. Thus, $n = a + 1$ and $\ell = a - 1$. Assume for the sake of contradiction that there is a $\mathcal{G}_{a,a,a-1}^{a+1}$ -correlated sampling strategy with disagreement probability $< 1 - \ell/(2a - \ell) = 2/n$. Let \mathcal{D} be the uniform distribution over pairs (A, B) of subsets of $[n]$ satisfying $A \cup B = [n]$ and $|A| = |B| = |A \cap B| + 1$. Note that \mathcal{D} is not a product distribution over (A, B) , unlike in Lemma 2.2, which is what makes it challenging to analyze. By an averaging argument, there is a deterministic strategy pair (f, g) such that,

$$\Pr_{(A,B) \sim \mathcal{D}} [f(A) \neq g(B)] < \frac{2}{n}. \quad (6)$$

Let $i := \operatorname{argmax}_{\ell \in [n]} \left| \left\{ A \in \binom{[n]}{n-1} : f(A) = \ell \right\} \right|$ be the element that is most frequently output by Alice's strategy f , and denote its number of occurrences by $k := \left| \left\{ A \in \binom{[n]}{n-1} : f(A) = i \right\} \right|$. We consider three different cases depending on the value of k :

- (i) If $k \leq n - 3$, then consider any $B \subseteq [n]$ with $|B| = n - 1$. For any value of $f(B) \in B$, the conditional error probability $\Pr[f(A) \neq g(B) \mid B]$ is at least $2/(n - 1)$. Averaging over all such B , we get a contradiction to Equation (6).
- (ii) If $k = n - 2$, let $A_1 \neq A_2$ be the two subsets of $[n]$ with $|A_1| = |A_2| = n - 1$ such that $f(A_1) \neq i$ and $f(A_2) \neq i$. For all $B \subseteq [n]$ with $|B| = n - 1$ such that $B \neq A_1$ and $B \neq A_2$, the conditional error probability $\Pr[f(A) \neq g(B) \mid B]$ is at least $2/(n - 1)$. Note that there are $n - 2$ such B 's, and that either A_1 or A_2 is the set $[n] \setminus \{i\}$. If $B = [n] \setminus \{i\}$, then the conditional disagreement probability $\Pr[f(A) \neq g(B) \mid B]$ is at least $(n - 2)/(n - 1)$. Averaging over all B , we get that

$$\begin{aligned} \Pr_{(A,B) \sim \mathcal{D}} [f(A) \neq g(B)] &\geq \binom{2}{n-1} \cdot \binom{n-2}{n} + \binom{n-2}{n-1} \cdot \binom{1}{n} \\ &\geq \frac{2}{n}, \end{aligned}$$

where the last inequality holds for all $n \geq 2$. This contradicts Equation (6).

(iii) If $k = n - 1$, then the only subset A_1 of $[n]$ with $|A_1| = n - 1$ and such that $f(A_1) \neq i$ is $A_1 = [n] \setminus \{i\}$. For all $B \neq A_1$, the conditional error probability $\Pr[f(A) \neq g(B) | B]$ is at least $1/(n - 1)$. On the other hand, if $B = A_1$, then the conditional error probability is equal to 1. Averaging over all B , we get that

$$\Pr_{(A,B) \sim \mathcal{D}}[f(A) \neq g(B)] \geq \left(\frac{1}{n-1}\right) \cdot \left(\frac{n-1}{n}\right) + 1 \cdot \left(\frac{1}{n}\right) = \frac{2}{n},$$

which contradicts Equation (6). □

Remark. In [KT02], the correlated sampling strategy is used to give a randomized rounding procedure for a linear program. The factor 2 loss in the correlated sampling strategy translates into an integrality gap of at most 2. In fact, they also prove that the integrality gap is roughly tight. As pointed out by an anonymous reviewer, their proof essentially establishes Theorem 3.5 under the assumption that $f = g$.

4 Other Observations and Open Questions

In the context of Conjecture 3.4, even in the setting where the set sizes are allowed to vary slightly, our knowledge is somewhat incomplete. Lemma 2.2 shows optimality of the MinHash strategy when $(A, B) \sim \mathcal{D}_p$. In this case, A and B are independent and p -biased each, so $|A| \approx p \cdot n$, $|B| \approx p \cdot n$ and $|A \cap B| \approx p^2 \cdot n$. A simple reduction to Lemma 2.2 also implies the optimality of the MinHash strategy in the case where A and B are “positively correlated”. Specifically for parameters $\alpha > p$, consider the following distribution $\mathcal{D}_{p,\alpha}$ on pairs (A, B) of subsets of $[n]$, where we first sample $S \subseteq [n]$ by independently including each element of $[n]$ with probability p/α , and then independently including every $i \in S$ in each of A and B with probability α . In this case, $|A| \approx p \cdot n$, $|B| \approx p \cdot n$ and $|A \cap B| \approx \alpha p \cdot n > p^2 \cdot n$. Even if we reveal S to both Alice and Bob, Lemma 2.2 implies a lower bound of $2 \cdot \delta / (1 + \delta)$ on the error probability (for large enough n). It is unclear if the optimality holds even in the case where A and B are “negatively-correlated”, i.e., when $|A| \approx p \cdot n$, $|B| \approx p \cdot n$ and $|A \cap B| \ll p^2 \cdot n$.

Finally as alluded to in the introduction, in the setting where P and Q are flat distributions on different sized subsets of the universe, there is a strategy with lower error than provided in Theorem 1.2. In particular, for $P = \mathcal{U}(A)$ and $Q = \mathcal{U}(B)$ where $|A| \neq |B|$, the MinHash strategy gives an error probability of $1 - \frac{|A \cap B|}{|A \cup B|}$ (which is the Jaccard distance between A and B). However, naïvely using the strategy of Kleinberg-Tardos/Holenstein would give an error probability of $1 - \frac{|A \cap B|}{|A \cup B| + ||A| - |B||}$ which is higher when $|A| \neq |B|$. This implies that the strategy of Kleinberg-Tardos/Holenstein is not “always optimal”. Thus, it will be interesting to identify the right measure that captures the minimum error of a general \mathcal{G} -restricted correlated sampling strategy.

A Correlated Sampling Strategies of [Bro97, KT02, Hol07]

For completeness, we describe the correlated sampling strategies of Broder and of Kleinberg & Tardos and Holenstein thereby proving Theorem 1.2.

Broder’s Min Hash Strategy. Consider the case of *flat distributions*, where the distributions P and Q are promised to be of the following form: there exist $A, B \subseteq [n]$ such that $P = \mathcal{U}(A)$ and $Q = \mathcal{U}(B)$. In this case, it is easy to show that the strategy given in [Algorithm 2](#) achieves an error probability of $1 - \frac{|A \cap B|}{|A \cup B|}$. Since π is a random permutation, $f(A, \pi)$ is uniformly distributed over A and $g(B, \pi)$ is uniformly distributed over B . Let i_0 be the smallest index such that $\pi(i_0) \in A \cup B$. The probability that $\pi(i_0) \in A \cap B$ is exactly $\frac{|A \cap B|}{|A \cup B|}$, and this happens precisely when $f(A, \pi) = g(B, \pi)$. Hence, we get the claimed error probability.

Algorithm 2: MinHash strategy [Bro97]

Alice’s input: $A \subseteq [n]$

Bob’s input: $B \subseteq [n]$

Shared randomness: a random permutation $\pi : [n] \rightarrow [n]$

Strategy:

- $f(A, \pi) = \pi(i_A)$, where i_A is the smallest index such that $\pi(i_A) \in A$.
- $g(B, \pi) = \pi(i_B)$, where i_B is the smallest index such that $\pi(i_B) \in B$.

The correlated sampling strategy of [\[KT02, Hol07\]](#) follows a similar approach.

Proof of Theorem 1.2. Given a finite discrete space Ω and probability measures P and Q over Ω , define $A := \{(\omega, p) \in \Omega \times [0, 1] : p < P(\omega)\}$ and $B := \{(\omega, q) \in \Omega \times [0, 1] : q < Q(\omega)\}$. Also for all $\omega \in \Omega$, define $A_\omega := A \cap (\{\omega\} \times [0, 1])$ and $B_\omega := B \cap (\{\omega\} \times [0, 1])$.

The strategy of [\[KT02, Hol07\]](#) can be intuitively understood as follows: Alice and Bob use the MinHash strategy on inputs A and B over the universe $\Omega \times [0, 1]$, to obtain elements (ω_A, p_A) and (ω_B, p_B) respectively, and simply output ω_A and ω_B respectively. However, this by itself is not well defined since $\Omega \times [0, 1]$ is not a finite set. Nevertheless, the MinHash strategy can be modified to instead have a (countably) infinite sequence of points sampled i.i.d. from the uniform measure over $\Omega \times [0, 1]$, instead of a permutation π . This strategy is summarized in [Algorithm 3](#).

Algorithm 3: Kleinberg-Tardos’ / Holenstein’s strategy [KT02, Hol07]

Alice’s input: $P \in \Delta_\Omega$; let $A := \{(\omega, p) \in \Omega \times [0, 1] : p < P(\omega)\}$

Bob’s input: $Q \in \Delta_\Omega$; let $B := \{(\omega, q) \in \Omega \times [0, 1] : q < Q(\omega)\}$

Shared randomness: An infinite sequence $\pi = ((\omega_1, r_1), (\omega_2, r_2), \dots)$ where each (ω_i, r_i) is i.i.d. sampled uniformly from $\Omega \times [0, 1]$.

Strategy:

- $f(P, \pi) := \omega_{i_A}$, where i_A is the smallest index such that $(\omega_{i_A}, r_{i_A}) \in A$
- $g(Q, \pi) := \omega_{i_B}$, where i_B is the smallest index such that $(\omega_{i_B}, r_{i_B}) \in B$

Let μ be the uniform measure over $\Omega \times [0, 1]$. Observe that $\mu(A) = \mu(B) = 1/|\Omega|$ and for all $\omega \in \Omega$, we have $\mu(A_\omega) = P(\omega)/|\Omega|$ and $\mu(B_\omega) = Q(\omega)/|\Omega|$. Similar to the analysis of the MinHash strategy,

for Alice's chosen index i_A , we have (ω_{i_A}, r_{i_A}) is uniform over A . Thus, $\Pr[f(P, \pi) = \omega]$ is precisely $\mu(A_\omega)/\mu(A) = P(\omega)$. Thus, $f(P, \pi)$ is distributed according to P and similarly, $g(B, \pi)$ is distributed according to Q . Finally, $\Pr[f(P, \pi) = g(Q, \pi)] \geq \Pr[i_A = i_B]$. To bound this probability, note that $\mu(A \cap B) = (1 - \delta)/|\Omega|$ and $\mu(A \cup B) = (1 + \delta)/|\Omega|$.

$$\Pr[f(P, \pi) = g(Q, \pi)] \geq \Pr[i_A = i_B] = \frac{\mu(A \cap B)}{\mu(A \cup B)} = \frac{1 - \delta}{1 + \delta} = 1 - \frac{2\delta}{1 + \delta}$$

We ignore the possibility that no index i_A exists satisfying $(\omega_{i_A}, r_{i_A}) \in A$ (similarly for B) since this happens with probability 0. \square

References

- [BHH⁺08] Boaz Barak, Moritz Hardt, Ishay Haviv, Anup Rao, Oded Regev, and David Steurer. Rounding parallel repetitions of unique games. In *Proceedings of the 49th annual IEEE Symposium on Foundations of Computer Science FOCS*, pages 374–383. IEEE, 2008. 2
- [Bro97] Andrei Z Broder. On the resemblance and containment of documents. In *Compression and Complexity of Sequences 1997. Proceedings*, pages 21–29. IEEE, 1997. 2, 3, 9, 10
- [Cha02] Moses S Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the 34th annual ACM Symposium on Theory of computing (STOC)*, pages 380–388. ACM, 2002. 2
- [GP06] Sreenivas Gollapudi and Rina Panigrahy. A dictionary for approximate string search and longest prefix search. In *Proceedings of the 15th ACM international conference on Information and knowledge management*, pages 768–775. ACM, 2006. 3
- [HMT14] Bernhard Haeupler, Mark Manasse, and Kunal Talwar. Consistent weighted sampling made fast, small, and easy. *arXiv preprint arXiv:1410.4266*, 2014. 3
- [Hol07] Thomas Holenstein. Parallel repetition: simplifications and the no-signaling case. In *Proceedings of the 39th annual ACM Symposium on Theory of computing (STOC)*, pages 411–419. ACM, 2007. 2, 9, 10
- [KT02] Jon Kleinberg and Eva Tardos. Approximation algorithms for classification problems with pairwise relationships: Metric labeling and markov random fields. *Journal of the ACM (JACM)*, 49(5):616–639, 2002. 2, 9, 10
- [Man94] Udi Manber. Finding similar files in a large file system. In *Usenix Winter*, volume 94, pages 1–10, 1994. 3
- [MMT10] Mark Manasse, Frank McSherry, and Kunal Talwar. Consistent weighted sampling. *Technical Report MSR-TR 2010-73*, 2010. 3
- [Rao11] Anup Rao. Parallel repetition in projection games and a concentration bound. *SIAM Journal on Computing*, 40(6):1871–1891, 2011. 2

- [Riv16] Ronald L. Rivest. Symmetric encryption via keyrings and ecc. <https://people.csail.mit.edu/rivest/pubs/Riv16u.pdf>, 2016. 2, 3, 6, 7
- [Tho00] H. Thorisson. *Coupling, Stationarity, and Regeneration*. Probability and Its Applications. Springer New York, 2000. 2
- [vLW01] Jacobus Hendricus van Lint and Richard Michael Wilson. *A course in combinatorics*. Cambridge university press, 2001. 7