

Hadamard Response: Estimating Distributions Privately, Efficiently, and with Little Communication

Jayadev Acharya*
Cornell University
acharya@cornell.edu

Ziteng Sun*
Cornell University
zs335@cornell.edu

Huanyu Zhang*
Cornell University
hz388@cornell.edu

June 27, 2018

Abstract

We study the problem of estimating k -ary distributions under ε -local differential privacy. n samples are distributed across users who send privatized versions of their sample to a central server. All previously known sample optimal algorithms require linear (in k) communication from each user in the high privacy regime ($\varepsilon = O(1)$), and run in time that grows as $n \cdot k$, which can be prohibitive for large domain size k .

We propose *Hadamard Response (HR)*, a local privatization scheme that requires no shared randomness and is symmetric with respect to the users. Our scheme has order optimal sample complexity for all ε , a communication of at most $\log k + 2$ bits per user, and nearly linear running time of $\tilde{O}(n + k)$.

Our encoding and decoding are based on Hadamard matrices, and are simple to implement. The statistical performance relies on the coding theoretic aspects of Hadamard matrices, ie, the large Hamming distance between the rows. An efficient implementation of the algorithm using the Fast Walsh-Hadamard transform gives the computational gains.

We compare our approach with Randomized Response (RR), RAPPOR, and subset-selection mechanisms (SS), both theoretically, and experimentally. For $k = 10000$, our algorithm runs about 100x faster than SS, and RAPPOR.

1 Introduction

Estimating the underlying probability distribution from data samples is a quintessential statistical problem. Given samples from an unknown distribution p , the goal is to obtain an estimate \hat{p} of p . The problem has a rich, and vast literature (see e.g. [6, 39, 17, 18], and many others), with the primary goal of statistical efficiency, namely minimizing the sample complexity for estimation, which is the first resource we consider.

1. Utility. What is the **sample complexity** of estimation?

In many applications, data contains sensitive information, and preserving the privacy of individuals is paramount. Without proper precautions, sensitive information can be inferred as evidenced by well publicized data leaks over the past decade, including de-anonymization of public health records in Massachusetts [41], de-anonymization of Netflix users [37] and de-anonymization of individuals participating in the genome wide association study [29]. On the policy side, the

*This research is supported by NSF-CCF-CRII 1657471, and a grant from Cornell University.

EU General Data Protection Regulation are now in effect, putting strict regulations on the data collection methods across the EU (visit <http://www.eugdpr.org>).

Private data release and computation on data has been studied in several fields, including statistics, machine learning, database theory, algorithm design, and cryptography (See e.g., [45, 14, 22, 46, 23, 42, 13]). *Differential Privacy (DP)* [24] has emerged as one of the most popular notions of privacy (see [24, 46, 26, 9, 36, 32], references therein, and the recent book [25]). DP has been adopted by several companies including Google, and Apple [21, 27].

A particularly popular privacy setting is *local differential privacy (LDP)* [45, 23], where users do not trust the data collector, and privatize their data before releasing. We study distribution estimation under LDP. Distribution estimation with privacy is an important problem. For example, understanding the drug usage habits of the entire population (the distribution) is crucial for policy design. Understanding the internet traffic distribution is important for ad-placement. In both these applications, preserving individual privacy is essential.

2. Privacy. How much information about a user is leaked by the scheme?

There are inherent trade-offs between utility and privacy. Sample privacy trade-offs have been recently studied for various problems, including distribution estimation [23, 31, 47, 43, 20, 35].

However, two crucial resources have not been considered in private distribution estimation, computation, and communication. In applications where the underlying dimensionality is high, or the number of samples is large, it is imperative to have computationally efficient algorithms. Internet companies collect information about user’s browsing history over a large number of users and websites, and large departmental stores collect purchase statistics over a large number of users and products. In these problems, algorithms with high computational overhead are prohibitive, even if they have optimal sample complexity. There has been recent interest in computationally efficient distribution estimation in the non-private setting (see e.g., [15, 1, 33, 12, 16, 40, 3]).

3. Computational Complexity. What is the running time of the algorithm?

In distributed applications, communication (both with and without privacy) is critical. For example, a large fraction of internet traffic is on hand-held devices with limited uplink capacity due to limited battery power, limited uplink bandwidth, or expensive data rates. Similarly, in large scale distributed machine learning problems, communication from processors to the server is the bottleneck since local computations are fast. Communication limited distributed distribution estimation has been studied in the non-private setting (e.g., [48, 4, 19, 2, 28]).

In the context of private estimation tasks, the problem of finding the heavy hitters, and learning properties under local differential privacy under the assumption of public randomness, where the server can send communication to the clients to reduce communication from user end has received much attention recently [8, 7, 5, 30, 44, 11]. However, these algorithms require shared randomness, as well as asymmetric schemes, where each user can use a different privatization mechanism. [7] uses a Hadamard transform, but they use it to form orthogonal basis and reduce storage, which is different from us.

4. Communication Complexity. How many bits are communicated?

In this work, we consider discrete distribution estimation under the aforementioned four resources. We provide the first algorithm that is simultaneously sample order optimal for any privacy value, has logarithmic communication per symbol, and runs in linear time in the input and output size.

1.1 Organization.

In Section 2 we describe the problem set-up, in Section 2.1 and 2.2 we describe prior privatization

schemes, and our results. In Section 3, we provide a family of ε -LDP privatization schemes. In Section 4, we specialize and design schemes that are optimal in the most interesting regime of high privacy. Finally in Section 5 we will describe how to extend these schemes to general ε .

2 Preliminaries

Local Differential Privacy (LDP). Suppose x is a private information that takes values in a set \mathcal{X} with k elements (wlog let $\mathcal{X} = [k] := \{0, 1, \dots, k-1\}$). A privatization mechanism is a randomized mapping Q from $[k]$ to an output set \mathcal{Z} (which can be arbitrary), that maps $x \in \mathcal{X}$ to $z \in \mathcal{Z}$ with probability $Q(z|x)$. The output z of this mapping, called the privatized sample, is then released. Q is ε -locally differentially private (ε -LDP) [23] if for all $x, x' \in \mathcal{X}$,

$$\sup_{z \in \mathcal{Z}} \frac{Q(z|x)}{Q(z|x')} \leq e^\varepsilon. \quad (1)$$

Small values of ε are more stringent and is the high privacy regime, and large values of ε is the low privacy regime. When \mathcal{X} and \mathcal{Z} are both discrete, the mechanism Q is described by a stochastic matrix of size $|\mathcal{X}| \times |\mathcal{Z}|$ whose (x, z) th entry is $Q(z|x)$. Q is ε -LDP if the ratio of *any two entries* in a column of this matrix is at most e^ε .

Randomness, and Symmetry. A scheme that requires shared/public randomness requires the generation of shared randomness at the server, which needs to be communicated to the users. Symmetric schemes are those where each user uses the same privatization scheme [38]. In this paper, we consider schemes that are symmetric and require no shared randomness. Other such schemes include RAPPOR, Randomized Response, and subset selection methods, described later. We note that the literature on heavy hitter estimation has mostly considered schemes with shared randomness [8, 7, 11], and it will be interesting to see if our methods can provide improved algorithms for the heavy hitter problem.

LDP distribution estimation. Let $\Delta_k = \left\{ p(0), \dots, p(k-1) : p(x) \geq 0, \sum_{x=0}^{k-1} p(x) = 1 \right\}$ be the set of all distributions over $[k]$. Let X_1, \dots, X_n be independent samples drawn from an *unknown* $p \in \Delta_k$, where X_i is the private (sensitive) data with the i th user. Each user maps X_i through an ε -LDP Q , to obtain Z_i . The task at the server, upon observing the privatized samples Z_1, \dots, Z_n , is to output $\hat{p} : \mathcal{Z}^n \rightarrow \Delta_k$, an estimate of p . Let $d : \Delta_k \times \Delta_k \rightarrow \mathbb{R}_+$ be a distance measure between distributions in Δ_k . Private distribution estimation task is the following:

Given $\alpha > 0, \varepsilon > 0, d : \Delta_k \times \Delta_k \rightarrow \mathbb{R}$, design an ε -LDP Q , and a corresponding estimation \hat{p} , such that $\forall p \in \Delta_k$, with probability at least 0.9, $d(\hat{p}, p) < \alpha$.

The *sample complexity* is the least n for which such an ε -LDP scheme Q , and a corresponding \hat{p} exists. The *communication complexity* is the number of bits to send Z_i to the server. The *computational complexity* is the total time to estimate \hat{p} from Z_1, \dots, Z_n at the server and to privatize X_i using Q at the users.

We will use ℓ_1 , and ℓ_2 distance in this paper. For $r \geq 0$, the ℓ_r distance between $p, q \in \Delta_k$ is $\ell_r(p, q) := (\sum_x |p(x) - q(x)|^r)^{1/r}$. In non-private setting, the sample complexity of distribution estimation under these distances is known even including precise constants [10, 34].

ε	k -RR	RAPPOR	k -SS	ε -HR
$(0, 1)$	$\frac{k^3}{\varepsilon^2 \alpha^2}$	$\frac{k^2}{\varepsilon^2 \alpha^2}$	$\frac{k^2}{\varepsilon^2 \alpha^2}$	$\frac{k^2}{\varepsilon^2 \alpha^2}$
$(1, \log k)$	$\frac{k^3}{e^{2\varepsilon} \alpha^2}$	$\frac{k^2}{e^{\varepsilon/2} \alpha^2}$	$\frac{k^2}{e^\varepsilon \alpha^2}$	$\frac{k^2}{e^\varepsilon \alpha^2}$
$(\log k, 2 \log k)$	$\frac{k}{\alpha^2}$	$\frac{k^2}{e^{\varepsilon/2} \alpha^2}$	$\frac{k}{\alpha^2}$	$\frac{k}{\alpha^2}$
$(2 \log k, +\infty)$	$\frac{k}{\alpha^2}$	$\frac{k}{\alpha^2}$	$\frac{k}{\alpha^2}$	$\frac{k}{\alpha^2}$

Table 1: Sample complexity, up to constant factors, under ℓ_1 distance for the different methods. The sample complexity under ℓ_2 distance is exactly a factor k smaller in each cell above.

ε	k -RR	RAPPOR	k -SS	ε -HR
$(0, 1)$	$\log k$	k	k	$\log k$
$(1, \log k)$	$\log k$	$\frac{k}{e^{\varepsilon/2}}$	$\frac{k}{e^\varepsilon}$	$\log k$
$(\log k, 2 \log k)$	$\log k$	$\frac{k}{e^{\varepsilon/2}}$	$\log k$	$\log k$
$(2 \log k, +\infty)$	$\log k$	$\log k$	$\log k$	$\log k$

Table 2: Communication requirements for distribution estimation techniques.

2.1 The privatization mechanisms

We will now briefly describe RR, RAPPOR, the most popular ε -LDP schemes using no interaction and public randomness. We will also mention SS, and our proposed HR. For a detailed description of RAPPOR and SS, please refer to Section C.

k -Randomized Response (RR). The k -RR mechanism [45, 31] is an ε -LDP Q_{RR} with $\mathcal{Z} = \mathcal{X} = [k]$, such that

$$Q_{\text{RR}}(z|x) := \begin{cases} \frac{e^\varepsilon}{e^\varepsilon + k - 1} & \text{if } z = x, \\ \frac{1}{e^\varepsilon + k - 1} & \text{otherwise.} \end{cases} \quad (2)$$

k -RAPPOR. The randomized aggregatable privacy-preserving ordinal response (RAPPOR) is an ε -LDP mechanism which was proposed in [23, 27]. Its simplest implementation k -RAPPOR maps $\mathcal{X} = [k]$ to $\mathcal{Z} = \{0, 1\}^k$. It first does a one hot encoding to the input $x \in [k]$ to obtain $\mathbf{y} \in \{0, 1\}^k$, such that $\mathbf{y}_j = 1$ for $j = x$, and $\mathbf{y}_j = 0$ for $j \neq x$. The privatized output of k -RAPPOR is a k -bit vector obtained by independently flipping each bit of \mathbf{y} with probability $\frac{1}{e^{\varepsilon/2} + 1}$.

Subset Selection techniques. [43, 47] propose an ε -LDP scheme that maps $x \in [k]$ to subsets of $[k]$ of size $\lceil k/(e^\varepsilon + 1) \rceil$. The scheme is described in detail in Section C.

Hadamard Response. We propose Hadamard Response (HR), an ε -LDP scheme with $\mathcal{Z} = [K]$, for some $k \leq K \leq 4k$. The algorithm is described in Section 4 for high privacy, and in Section 5 for general privacy.

2.2 Previous Results

To estimate distributions in Δ_k to ℓ_1 distance α under ε -LDP, the sample, communication and time requirements of the various schemes are given in Table 1, 2 and 3 respectively..

k -RR	k -RAPPOR	Subset selection	ε -HR
$n + k$	$n + k + \frac{nk}{e^{\varepsilon/2}}$	$n + k + \frac{nk}{e^{\varepsilon}}$	$n + k$

Table 3: Time bounds for distribution estimation. The running times are described in Section C. These are upper bounds up to logarithmic factors.

The sample complexity is given in **Table 1**. The entries in *green* boxes are sample-order optimal, namely there is a matching lower bound [47]. Note that RR is sample-optimal in the low privacy regime (last two rows), and is *highly sub-optimal* in the high privacy regime ($\varepsilon = O(1)$). RAPPOR is optimal for high-privacy, but sub-optimal for medium privacy. SS, and our proposed HR are sample-order-optimal for all ε . The sample complexity arguments for RR, RAPPOR, and SS can be found in [31, 47].

Table 2 describes the communication requirements of various schemes. However, it is not clear how to measure the communication requirements, since for a given privatization scheme, there might be communication protocols requiring fewer bits than others. For example, RAPPOR is described as giving k bits as its output, but perhaps these k bits can be compressed further requiring much smaller communication. We get around such concerns by observing that, once the input distribution p and the privatization mechanism Q is fixed, the output distribution of the privatized sample Z is fixed. By Shannon’s source coding theorem, to *faithfully* send Z to the server requires at least $H(Z)$ bits of communication. The entries in the table are derived by considering the input distribution to be near uniform, and evaluating the entropy of the output of the mechanisms. For RR, $\log k$ bits of communication follows from $\mathcal{Z} = [k]$. Note that in this paper all logarithms are in base 2. The communication requirements for RAPPOR, and SS are derived in Section C (Theorems 9, and Theorem 10 respectively).

Table 3 describes the total running time lower bounds for faithfully implementing the known schemes. The argument is that at the server, the computation complexity is at least the number of bits that need to be read, which is the amount of communication from the users. If there are n users, then $n \cdot H(Z)$ serves as our time complexity bound, and these form the entries in the table.

2.3 Motivation and Our Results

Our work is motivated by the first three columns of the tables, which captures the apparent sample-communication-computation trade-offs present in the existing schemes. We elaborate this point in the most interesting regime of high privacy. For simplicity, fix $\varepsilon = 1$, and $\alpha = 0.1$ (chosen arbitrarily!), and treat them as fixed constants in this paragraph. In this setting, from **Table 1**, note that the optimal sample complexity is $\Theta(k^2)$, achieved by RAPPOR, and SS, while RR has a sub-optimal sample complexity of $\Theta(k^3)$. Now consider the communication requirements. $\mathcal{Z} = [k]$ for RR, requiring only $\log k$ bits. A straight-forward computation shows that any input distribution to the RAPPOR mechanism induces an output distribution over $\{0, 1\}^k$ with entropy at least $\Omega(k)$, thus requiring $\Omega(k)$ bits to *faithfully* send the privatized samples to the server. SS also requires $\Omega(k)$ bits in this regime. These are formally shown in **Theorem 9** and **Theorem 10**. As for the running time at the server end, a bound of $\Omega(k^3)$ for all these three methods follows from the total communication to the server ($\#samples \times \#bits$ per sample), which is a factor k larger than the $\Theta(k^2)$ optimal sample complexity bound.

Our main result is the following, which is formally stated in **Theorem 2**, and **Theorem 7**.

Theorem 1. *We propose a simple algorithm for ε -LDP distribution estimation that for all param-*

eter regimes, is sample optimal, runs in near-linear time in the number of samples, and has only a logarithmic communication complexity in the domain size, for both the ℓ_1 , and ℓ_2 distance.

Going back to the high privacy regime, considered before, this shows that our scheme has a running time of $\tilde{O}(k^2)$, which is nearly linear in the optimal sample complexity under ℓ_1 distance.

3 A family of ε -LDP schemes

We first propose a general family of LDP schemes, and then carefully choose schemes from this family that are sample-optimal, communication and computationally efficient for distribution estimation.

The scheme involves the following steps:

1. Choose an integer K , and let the output alphabet be $\mathcal{Z} = [K]$.
2. Choose a positive integer $s \leq K$.
3. For each $x \in \mathcal{X} = [k]$, pick $C_x \subseteq [K]$ with $|C_x| = s$.
4. The privatization scheme from $[k]$ to $[K]$ is then given by:

$$Q(z|x) := \begin{cases} \frac{e^\varepsilon}{se^\varepsilon + K - s} & \text{if } z \in C_x, \\ \frac{1}{se^\varepsilon + K - s} & \text{if } z \in \mathcal{Z} \setminus C_x. \end{cases} \quad (3)$$

This scheme satisfies (1), and is ε -LDP. This privatization scheme chooses a set C_x for each x and assigns the elements in C_x a higher probability than those not in C_x . We also note that RR is a special case of this construction when $K = k$, $s = 1$, and $C_x = \{x\}$. We know from the last section that RR is sub-optimal in the high privacy regime. Our general inspiration comes from coding theory, and we select s , and C_x carefully in order to send more information across Q than RR.

In Section 4 we give an optimal scheme in the high privacy regime, and extend it to the general case in Section 5

4 Optimal scheme for high privacy regime

Privatization scheme. If for two x , and x' , $C_x = C_{x'}$, then we cannot tell them apart. Therefore, the hope is that the farther apart C_x and $C_{x'}$ are, the easier it is to tell them apart. With this in mind, we specify a particular choice of parameters for our scheme, which turns out to be sample-optimal in the high privacy regime. In particular, our privatization scheme will satisfy the following:

An optimal privatization for high privacy

Choose K , and C_x 's such that (We will show in Section 4.1 how to satisfy these conditions.):

- C1.** K is between k and $2k$, and $s = K/2$, namely for all $x \in [k]$, $|C_x| = \frac{K}{2}$.
- C2.** For any $x, x' \in [k]$, and $x \neq x'$, $|\Delta(C_x, C_{x'})| = |(C_x \setminus C_{x'}) \cup (C_{x'} \setminus C_x)| = \frac{K}{2}$.

Use (3) for privatization.

Performance. We will show that for $\varepsilon = O(1)$, this privatization is sample-order-optimal, namely there is a corresponding estimator $\hat{p} : [K]^n \rightarrow \Delta_k$ that is sample-optimal. Before describing the estimation procedure, we provide the statistical guarantees.

Theorem 2. For any privatization scheme satisfying **C1**, **C2**, there is a corresponding estimation scheme $\hat{p} : [K]^n \rightarrow \Delta_k$, such that

$$\mathbb{E} \left[\ell_2^2(\hat{p}, p) \right] \leq \frac{4k(e^\varepsilon + 1)^2}{n(e^\varepsilon - 1)^2}, \text{ and } \mathbb{E} [\ell_1(\hat{p}, p)] \leq \sqrt{\frac{4k^2(e^\varepsilon + 1)^2}{n(e^\varepsilon - 1)^2}}. \quad (4)$$

The sample optimality, and small communication for high privacy is an immediate corollary.

Corollary 3. When $\varepsilon = O(1)$, the sample complexity of this scheme for estimation to ℓ_1 distance α is $O(k^2/\varepsilon^2\alpha^2)$ samples, and for ℓ_2^2 distance is $O(k/\varepsilon^2\alpha^2)$. Further, the communication from each user is at most $\log(k) + 1$ bits. This is sample-optimal for $\varepsilon = O(1)$ for both ℓ_1 (Table 1) and ℓ_2^2 (see [47]).

Proof. Applying Markov's inequality in Theorem 2, and substituting $e^\varepsilon + 1 = \Theta(1)$, and $e^\varepsilon - 1 = \Theta(\varepsilon)$ when $\varepsilon = O(1)$ gives the sample complexity bounds. The communication bounds are from $\log K \leq \log(k) + 1$. \square

Estimation. Suppose $Q_{K,\varepsilon}$ is an ε -LDP scheme satisfying **C1**, and **C2**. For an input distribution p over $[k]$, let $p(C_x)$ be the probability that the privatized sample $Z \in C_x$. Using $|C_x| = K/2$, and **C2**, it follows that $|C_x \setminus C_{x'}| = K/4$, and $|C_x \cap C_{x'}| = K/4$. Therefore,

$$\begin{aligned} p(C_x) &= p(x) \left(\sum_{z \in C_x} Q_{K,\varepsilon}(z|x) \right) + \sum_{x' \neq x} p(x') \left(\sum_{z \in C_x \setminus C_{x'}} Q_{K,\varepsilon}(z|x') + \sum_{z \in C_x \cap C_{x'}} Q_{K,\varepsilon}(z|x') \right) \\ &= p(x) |C_x| \frac{e^\varepsilon}{(se^\varepsilon + K - s)} + \sum_{x' \neq x} p(x') \left(\frac{|C_x \setminus C_{x'}| \cdot 1}{se^\varepsilon + K - s} + \frac{|C_x \cap C_{x'}| \cdot e^\varepsilon}{se^\varepsilon + K - s} \right) \end{aligned} \quad (5)$$

$$= \frac{1}{2} + \frac{e^\varepsilon - 1}{2(e^\varepsilon + 1)} p(x), \quad (6)$$

where (5) follows from (3), and (6) by plugging $s = K/2$, and from **C2**. We can rewrite this as

$$p(x) = \frac{2(e^\varepsilon + 1)}{e^\varepsilon - 1} \left(p(C_x) - \frac{1}{2} \right). \quad (7)$$

This forms the basis of our estimation. From the privatized samples, we estimate of $p(C_x)$, and from that we estimate p . The entire scheme is given below.

An optimal distribution estimation scheme for high privacy

Input: k, ε , privatized samples Z_1, \dots, Z_n

1. For each $x \in [k]$, estimate $p(C_x)$ with its empirical probability:

$$\widehat{p(C_x)} := \sum_{j=1}^n \frac{\mathbb{I}\{Z_j \in C_x\}}{n}. \quad (8)$$

2. Estimate \hat{p} as:

$$\hat{p}(x) := \frac{2(e^\varepsilon + 1)}{e^\varepsilon - 1} \left(\widehat{p(C_x)} - \frac{1}{2} \right). \quad (9)$$

Proof of Theorem 2.¹ Let $p(C), \widehat{p(C)}$, be the vector of probabilities of $p(C_x)$'s and $\widehat{p(C_x)}$'s respectively. From (7) and (9),

$$\mathbb{E} \left[\ell_2^2(\widehat{p}, p) \right] = \frac{4(e^\varepsilon + 1)^2}{(e^\varepsilon - 1)^2} \mathbb{E} \left[\ell_2^2(\widehat{p(C)}, p(C)) \right].$$

From (8), $\mathbb{E} \left[\widehat{p(C_x)} \right] = \mathbb{E} [\mathbb{I} \{Z_j \in C_x\}] = p(C_x)$. Therefore,

$$\mathbb{E} \left[\ell_2^2(\widehat{p(C)}, p(C)) \right] = \mathbb{E} \left[\sum_{x \in [k]} (\widehat{p(C_x)} - p(C_x))^2 \right] = \sum_{x \in [k]} \mathbb{E} \left[(\widehat{p(C_x)} - p(C_x))^2 \right] = \sum_{x \in [k]} \text{Var}(\widehat{p(C_x)}).$$

By the independence of Z_i 's, $\widehat{p(C_x)}$ is the average of n independent Bernoulli random variables each with expectation $p(C_x)$. Hence,

$$\sum_{x \in [k]} \text{Var}(\widehat{p(C_x)}) = \sum_{x \in [k]} \frac{1}{n} \cdot p(C_x)(1 - p(C_x)) \leq \frac{1}{n} \sum_{x \in [k]} p(C_x) \leq \frac{k}{n}.$$

Plugging this bound in the previous expression gives the bound on ℓ_2^2 distance of the theorem.

$$\mathbb{E} \left[\ell_2^2(\widehat{p}, p) \right] \leq \frac{4k(e^\varepsilon + 1)^2}{n(e^\varepsilon - 1)^2}. \quad (10)$$

Using $k \cdot \ell_2^2(\widehat{p}, p) \geq \ell_1(\widehat{p}, p)^2$ with (10) gives the desired bound on $\mathbb{E} [\ell_1(\widehat{p}, p)]$. \square

4.1 Computational complexity and Hadamard matrices.

We showed the sample, and communication complexity guarantees. However, two questions are still unanswered:

- How to choose K , and design C_x 's that satisfy **C1**, **C2**?
- What is the time complexity of privatization and estimation?

We now address these questions. We start with the computational requirements of the proposed scheme, assuming **C1**, **C2**.

Computation at users. Given C_x 's, each user needs to implement (3). This requires uniform sampling from C_x 's, as well as from $[K] \setminus C_x$. We will design schemes to do this in time $O(\log K)$.

Computation at the server. The server needs to implement (8) and (9). Note that (9) can be implemented in time $O(k)$ after implementing (8). However, a straightforward implementation of (8) requires $n \cdot k$ time, since for each x we iterate over all the samples, giving running time of $O(n \cdot k)$. In particular, in the high privacy regime (say with $\varepsilon = 1$, and $\alpha = 0.1$) the sample complexity is $O(k^2)$ but the time requirement will be $O(k^3)$. We now show how to design a privatization to satisfy **C1**, **C2**, and for which we can implement (8) in time only $\tilde{O}(n + k)$.

Hadamard Response (HR) for high privacy. Suppose K is a power of two, and let $H_K \in \{\pm 1\}^{K \times K}$ be the Hadamard matrix of size $K \times K$ designed by the well known Sylvester's construction as follows. Let $H_1 = [1]$, and for $m = 2^j$, for $j \geq 1$, then

$$H_m := \begin{bmatrix} H_{m/2} & H_{m/2} \\ H_{m/2} & -H_{m/2} \end{bmatrix}.$$

¹A technicality here is that $\widehat{p}(x)$'s can be negative, but we can project \widehat{p} onto the simplex with the same order performance. We therefore only analyze the performance of \widehat{p} described in (9).

Some standard properties of Hadamard matrices that we use are the following:

- (i) The number of +1's in each row except the first is $K/2$,
- (ii) Any two rows agree (and disagree) on exactly $K/2$ locations,
- (iii) Vector multiplication with H_K is possible in time $O(K \log K)$ with Fast Walsh Hadamard transform,
- (iv) We can uniformly sample from the +1's (and the -1's) in any row in time $O(\log K)$.

We now describe the parameters for the privacy mechanism:

1. Choice of K : Let $K = 2^{\lceil \log_2(k+1) \rceil} \geq k + 1$, the smallest power of 2 larger than k . To satisfy **C1**, we will choose $s = K/2$.

2. Choice of C_x 's: Map the symbols $[k] = \{0, \dots, k-1\}$ to rows of H_K as follows: map 0 to the second row, 1 to the third row, and so on. In other words, x is mapped to row $x + 1$. Given any x , we choose $C_x \subset [K]$ to be the column indices with a '+1' in the $(x + 1)$ th row of H_K . By Property (i) and (ii) of H_K , both **C1**, and **C2** are satisfied. This implies a privatization scheme with optimal sample and communication complexity in the high privacy regime.

Fast computation with HR. By Property (iv), we can efficiently implement the privatization scheme at the users. We will now provide an efficient implementation of (8). Let $q = (q(0), \dots, q(K-1))$ be the vector of the empirical distribution of Z_1, \dots, Z_n over $[K] = \{0, \dots, K-1\}$, namely

$$q(z) = \sum_{i=1}^n \frac{\mathbb{I}\{Z_i = z\}}{n}.$$

We can compute q in linear time with a single pass over Z_1, \dots, Z_n . Consider the matrix vector product $\mathbf{c} = H_K \cdot q$. For $x \in [k]$, the $(x + 1)$ th entry of $H_K \cdot q$ is $\sum_{z=0}^{K-1} H_K(x + 1, z) \cdot q(z)$. Now note that the +1's in the $(x + 1)$ th column correspond to C_x by construction, therefore

$$\sum_{z=0}^{K-1} H_K(x + 1, z) \cdot q(z) = \sum_{z \in C_x} q(z) - \sum_{z \in [K] \setminus C_x} q(z) = 2\widehat{p(C_x)} - 1 \quad (11)$$

$$= \left(\frac{e^\varepsilon - 1}{e^\varepsilon + 1} \right) \hat{p}(x), \quad (12)$$

where (11) follows from observing that $\sum_{z \in C_x} q(z) = \widehat{p(C_x)}$ from (8), and (12) follows from (9). Therefore the estimator \hat{p} is simply entries of a Hadamard vector product, appropriately normalized. By property (iii), this can be done in time $O(K \log K) = O(k \log k)$. This computational advantage is captured in the following theorem:

Theorem 4. *HR is an ε -LDP mechanism satisfying Theorem 2 that has a running time $\tilde{O}(n + k)$.*

5 General privacy regimes

Recall that RR is optimal for the low-privacy regime, which corresponds to $s = 1$. In the high privacy regime, we used $s = O(k)$. For general ε 's, we propose schemes within the framework of Section 3 that interpolate between HR and RR, while achieving the optimal sample complexity for every ε . In general, our choice of s will be close to $\max\{k/e^\varepsilon, 1\}$ as we will see below.

Privatization scheme in general privacy regime. We describe how we choose K , s , and C_x 's, in the scheme from Section 3.

We will consider block-structured matrices that interpolate between Hadamard matrices in high privacy regime and the identity matrix (corresponding to RR) in the low privacy regime.

Definition 5. Let B and b be powers of two, and let $K = B \cdot b$. A (B, b) -“reduced” Hadamard matrix is a $K \times K$ matrix with entries in $\{-1, +1\}$ defined as:

$$H_K^b := \begin{bmatrix} H_b & P_b & \dots & P_b \\ P_b & H_b & \dots & P_b \\ \vdots & \vdots & \ddots & \vdots \\ P_b & P_b & \dots & H_b \end{bmatrix},$$

where H_b is the $b \times b$ Hadamard matrix, and P_b is the $b \times b$ matrix with all entries ‘ -1 ’. Note that there are B occurrences of H_b along the diagonal.

1. Choice of K : Let B be the largest power of 2 less than $\min\{e^\varepsilon, 2k\}$, and b is the smallest power of 2 larger than $\frac{k}{B} + 1$, i.e.,

$$B := 2^{\lceil \log_2 \min\{e^\varepsilon, 2k\} \rceil - 1}, \quad b := 2^{\lceil \log_2(\frac{k}{B} + 1) \rceil}.$$

Let $K = B \cdot b$. A simple computation shows that $K \leq 4k$, implying that the communication from the users is at most $\log k + 2$ bits.

2. We will choose $s = b/2$.
3. Choice of C_x : From Property (i) of Hadamard matrices in the last section, the number of $+1$'s in the rows of H_K^b corresponding to the first rows of the embedded H_b 's is b , and for all other rows it is $b/2$. Similar to the high privacy regime, we map each x to a distinct row r_x of H_K^b with $b/2$ entries as $+1$'s. A simple way to do as before would be to map $0 \in [k]$ to the second row of H_K^b , and $x \in [k]$ is assigned to row $r_{x-1} + 1$, if it is not the first row of an embedded H_b , otherwise we assign it to $r_{x-1} + 2$. As before, let C_x be the column's with a $+1$ in the r_x th row of H_K^b .
4. The privatization mechanism then applies (3), which can be done in time $O(\log k)$ at each user by Property (iv) of Hadamard matrices.

Estimation scheme in general privacy regime. In the high privacy regime, we related $p(C_x)$ to $p(x)$ in (7). We will do the same here, however, because of the block-structure, the inputs that map to different blocks behave differently. Let $S_i \subseteq [K]$ be the columns of the i th embedded H_b block. Similar to $p(C_x)$, let $p(S_i)$ to be the probability that the output $z \in S_i$, when the input distribution is p . In other words,

$$S_i := \left\{ z \mid \lfloor \frac{z}{b} \rfloor = i \right\}, \quad p(S_i) := \sum_x p(x) \left(\sum_{z \in S_i} Q(z|x) \right).$$

Similar to (7), the following lemma relates $p(C_x)$, $p(S_i)$, and $p(x)$. It is proved in Section A.

Lemma 6. For the input distribution p , and $x \in [k]$ such that r_x is in the i th embedded H_b ,

$$p(C_x) - \frac{1}{2}p(S_i) = \frac{e^\varepsilon - 1}{2(2B - 1 + e^\varepsilon)}p(x). \quad (13)$$

With this lemma, our estimation algorithm is the following:

Distribution estimation for general privacy

Input: k, ε , privatized samples Z_1, \dots, Z_n

1. For each $x \in [k]$, estimate $p(C_x)$ with its empirical probability:

$$\widehat{p(C_x)} := \sum_{j=1}^n \frac{\mathbb{I}\{Z_j \in C_x\}}{n}. \quad (14)$$
2. For each $i \in B$, estimate $p(S_i)$ with its empirical probability:

$$\widehat{p(S_i)} := \sum_{j=1}^n \frac{\mathbb{I}\{Z_j \in S_i\}}{n}. \quad (15)$$
3. The estimator \hat{p} is then given by

$$\hat{p}(x) = \frac{2(2B - 1 + e^\varepsilon)}{e^\varepsilon - 1} \cdot \left(\widehat{p(C_x)} - \frac{1}{2} \widehat{p(S_i)} \right). \quad (16)$$

5.1 Performance

Our main performance bound on this scheme is given below. The analysis is similar to the high privacy regime and is given in Section B.

Theorem 7. *For all values of ε , and k , and the privatization scheme above, there is an estimate \hat{p} such that*

$$\mathbb{E} \left[\ell_2^2(\hat{p}, p) \right] \leq \frac{36(k + (e^\varepsilon - 1)b)e^\varepsilon}{n(e^\varepsilon - 1)^2}, \quad \mathbb{E} [\ell_1(\hat{p}, p)] \leq \sqrt{36 \frac{k(k + (e^\varepsilon - 1)b)e^\varepsilon}{n(e^\varepsilon - 1)^2}}.$$

The running time is $\tilde{O}(n + k)$, and communication is at most $\log k + 2$ bits.

Corollary 8. *Plugging the values of b in different regimes we obtain*

$$\mathbb{E} [\ell_1(\hat{p}, p)] \leq \begin{cases} O\left(\sqrt{\frac{k^2}{n\varepsilon^2}}\right), & \text{if } \varepsilon < 1, \\ O\left(\sqrt{\frac{k^2}{ne^\varepsilon}}\right), & \text{if } 1 < \varepsilon < \log k, \\ O\left(\sqrt{\frac{k}{n}}\right), & \text{if } \varepsilon > \log k. \end{cases}$$

Applying Markov's inequality, we obtain all the sample complexity bounds for HR described in the last column of Table 1.

6 Experiments.

We experimentally compare our algorithm with RR, RAPPOR and SS. We set $k \in \{100, 1000, 5000, 10000\}$, $n \in \{50000, 100000, 150000, \dots, 1000000\}$, and $\varepsilon \in \{0.1, 0.5, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$. We consider geometric distributions $Geo(\lambda)$, where $p(i) \propto (1 - \lambda)^i \lambda$, Zipf distributions $Zipf(k, t)$ where $p(i) \propto (i + 1)^{-t}$, two-step distributions, and uniform distributions. For every setting of (k, p, n, ε) , and

for each scheme, we simulate 30 runs, and compute the averaged ℓ_1 error, and averaged decoding time at the server. Our code and data for the experiments can be found at https://github.com/jlyx417353617/hadamard_response.

In a nutshell, we observe that in each regime, the statistical performance of HR is comparable to the best possible. Moreover, the decoding time of HR is similar to that of RR. In comparison to RAPPOR and SS, our running times can be orders of magnitude smaller, particularly for large k , and small ε . We remark that we implement RAPPOR, and SS such that their running time is almost linear in the time needed to read the already compressed communication from the users.

We describe some of our experimental results here. Figure 1 plots the ℓ_1 error for estimating geometric distribution for $k = 1000$. Note that for $\varepsilon = 0.5$, and $\varepsilon = 7$, our performance matches with the best schemes. In all the plots SS has the best statistical performance, however that can come at the cost of higher communication, and computation. Figures 2 captures similar statistical performance results for the uniform distribution for $k = 1000$. For larger k such as $k = 10000$, the performance is shown in figure 3.

The running time of our algorithm is theoretically a factor $k/\log k$ smaller than RAPPOR and subset selection. This is evident from the plots which show that for large k the running times of RAPPOR and SS are orders of magnitude more than HR, and RR.

Figure 4 shows the decoding time for the algorithms when $k = 100, 1000, 5000, 10000$ and $\varepsilon = 1$. It can be seen that our algorithm is orders of magnitudes faster in comparison to k -RAPPOR and k -SS. The gap in computation gets larger when k is larger, which is consistent with our theoretical analysis. For example, for $k = 10000$, our algorithm runs 100x faster than SS, and RAPPOR.

To compare the decoding time more fairly, we use a fast implementation for RAPPOR and SS in the middle and low privacy regime. We encode the k bit vector into a list of locations where there is a ‘+1’. The decoder uses this list to compute the histogram. The time requirement is $O(\frac{nk}{1+e^\varepsilon})$ and $O(\frac{nk}{1+e^{\varepsilon/2}})$ for SS and k -RAPPOR respectively in expectation (a naive implementation takes $O(nk)$ time). This is the best anyone can do to faithfully implement the algorithms.

Figure 5 shows the decoding time in middle privacy regime. We use fast implementation for RAPPOR and SS here. We can see our proposed algorithm is still saving a lot of time comparing to k -RAPPOR and SS. For low privacy regime, essentially everything breaks down to Randomized Response, so we won’t show the plots for the time comparison in this regime.

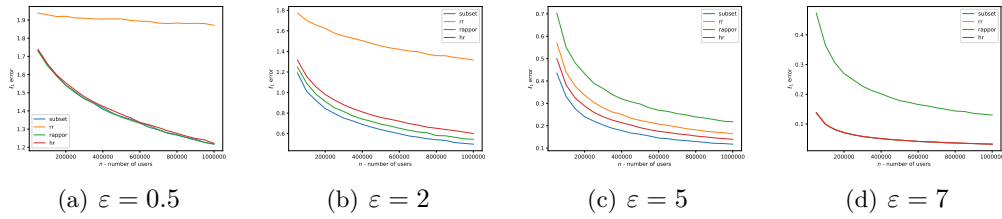


Figure 1: ℓ_1 -error comparison between four algorithms $k = 1000$ and $p \sim Geo(0.8)$

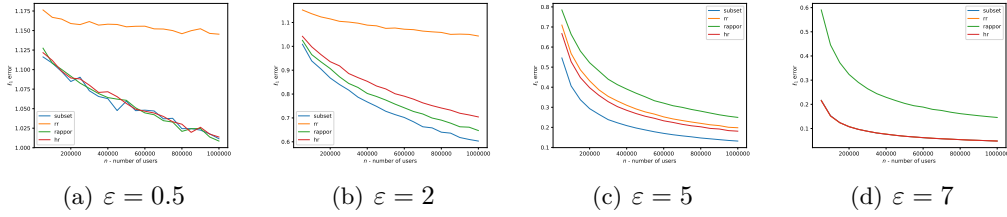


Figure 2: ℓ_1 -error comparison between four algorithms $k = 1000$ and $p \sim U[k]$

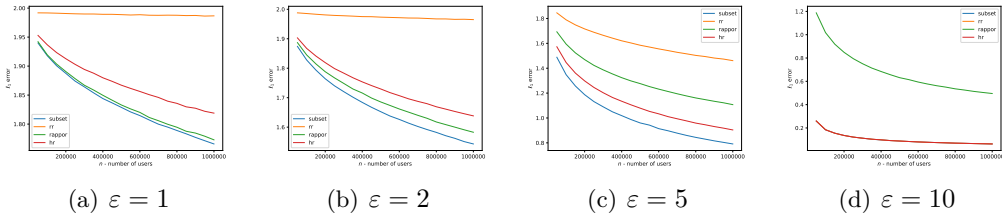


Figure 3: ℓ_1 -error comparison between four algorithms $k = 10000$ and $p \sim Geo(0.8)$

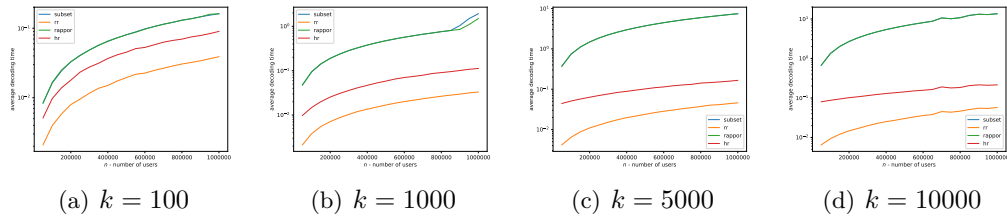


Figure 4: Decoding time comparison between four algorithms for $\varepsilon = 1$ and $p \sim Geo(0.8)$ and different values of k

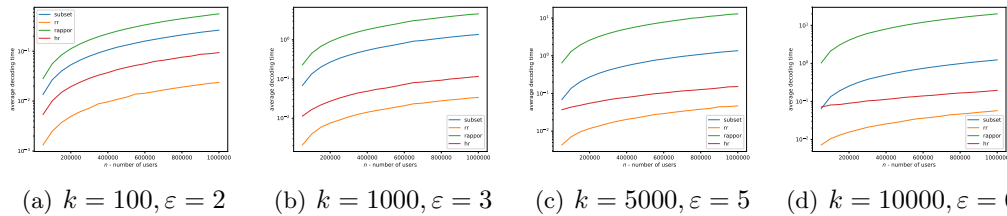


Figure 5: Decoding time comparison between four algorithms in middle privacy regime and $p \sim Geo(0.8)$. Note that the decoding times are in logarithmic scale.

References

- [1] F. Abramovich, Y. Benjamini, D. L. Donoho, and I. M. Johnstone. Adapting to unknown sparsity by controlling the false discovery rate. *The Annals of Statistics*, 34(2):584–653, 2006.
- [2] J. Acharya, C. L. Canonne, and H. Tyagi. Distributed simulation and distributed inference. *CoRR*, abs/1804.06952, 2018.
- [3] J. Acharya, I. Diakonikolas, J. Li, and L. Schmidt. Sample-optimal density estimation in nearly-linear time. In *Proceedings of the 28th Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '17, pages 1278–1289, Philadelphia, PA, USA, 2017. SIAM.
- [4] M. F. Balcan, A. Blum, S. Fine, and Y. Mansour. Distributed learning, communication complexity and privacy. In *Conference on Learning Theory*, pages 26–1, 2012.
- [5] S. Banerjee, N. Hegde, and L. Massoulié. The price of privacy in untrusted recommendation engines. In *Communication, control, and computing (Allerton), 2012 50th annual Allerton conference on*, pages 920–927. IEEE, 2012.
- [6] R. Barlow, D. Bartholomew, J. Bremner, and H. Brunk. *Statistical Inference under Order Restrictions*. Wiley, New York, 1972.
- [7] R. Bassily, K. Nissim, U. Stemmer, and A. G. Thakurta. Practical locally private heavy hitters. In *Advances in Neural Information Processing Systems*, pages 2285–2293, 2017.
- [8] R. Bassily and A. Smith. Local, private, efficient protocols for succinct histograms. In *STOC*, pages 127–135. ACM, 2015.
- [9] A. Blum, K. Ligett, and A. Roth. A learning theory approach to noninteractive database privacy. *Journal of the ACM (JACM)*, 60(2):12, 2013.
- [10] D. Braess and T. Sauer. Bernstein polynomials and learning theory. *Journal of Approximation Theory*, 128(2):187–206, 2004.
- [11] M. Bun, J. Nelson, and U. Stemmer. Heavy hitters and the structure of local privacy. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, pages 435–447. ACM, 2018.
- [12] S. O. Chan, I. Diakonikolas, R. A. Servedio, and X. Sun. Efficient density estimation via piecewise polynomial approximation. In *STOC*, 2014.
- [13] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12:1069–1109, 2011.
- [14] T. Dalenius. Towards a methodology for statistical disclosure control. *Statistisk Tidskrift*, 15:429–444, 1977.
- [15] S. Dasgupta. Learning mixtures of gaussians. In *Annual Symposium on Foundations of Computer Science (FOCS)*, 1999.
- [16] C. Daskalakis and G. Kamath. Faster and sample near-optimal algorithms for proper learning mixtures of gaussians. In *COLT*, 2014.

- [17] L. Devroye and L. Györfi. *Nonparametric Density Estimation: The L_1 View*. John Wiley & Sons, 1985.
- [18] L. Devroye and G. Lugosi. *Combinatorial methods in density estimation*. Springer, 2001.
- [19] I. Diakonikolas, E. Grigorescu, J. Li, A. Natarajan, K. Onak, and L. Schmidt. Communication-efficient distributed learning of discrete distributions. In *NIPS*, pages 6394–6404. Curran Associates, Inc., 2017.
- [20] I. Diakonikolas, M. Hardt, and L. Schmidt. Differentially private learning of structured discrete distributions. In *NIPS*, pages 2566–2574, 2015.
- [21] Differential Privacy Team, Apple. Learning with privacy at scale, December 2017.
- [22] I. Dinur and K. Nissim. Revealing information while preserving privacy. In *PODS*, pages 202–210, New York, NY, USA, 2003. ACM.
- [23] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy and statistical minimax rates. In *FOCS*, pages 429–438. IEEE, 2013.
- [24] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, pages 265–284. Springer, 2006.
- [25] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [26] C. Dwork, G. N. Rothblum, and S. Vadhan. Boosting and differential privacy. In *FOCS*, pages 51–60, 2010.
- [27] Ú. Erlingsson, V. Pihur, and A. Korolova. Rappor: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 ACM SIGSAC conference on computer and communications security*, pages 1054–1067. ACM, 2014.
- [28] Y. Han, P. Mukherjee, A. Özgür, and T. Weissman. Distributed statistical estimation of high-dimensional and nonparametric distributions with communication constraints. In *ISIT*, 2018.
- [29] N. Homer, S. Szelingner, M. Redman, D. Duggan, W. Tembe, J. Muehling, J. V. Pearson, D. A. Stephan, S. F. Nelson, and D. W. Craig. Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays. *PLoS Genetics*, 4(8):1–9, 2008.
- [30] J. Hsu, S. Khanna, and A. Roth. Distributed private heavy hitters. In *International Colloquium on Automata, Languages, and Programming*, pages 461–472. Springer, 2012.
- [31] P. Kairouz, K. Bonawitz, and D. Ramage. Discrete distribution estimation under local privacy. *arXiv preprint arXiv:1602.07387*, 2016.
- [32] P. Kairouz, S. Oh, and P. Viswanath. The composition theorem for differential privacy. *IEEE Transactions on Information Theory*, 63(6):4037–4049, 2017.
- [33] A. T. Kalai, A. Moitra, and G. Valiant. Efficiently learning mixtures of two gaussians. In *STOC*, 2010.

- [34] S. Kamath, A. Orlitsky, V. Pichapathi, and A. T. Suresh. On learning distributions from their samples. *In preparation*, 2015.
- [35] J. Lei. Differentially private m-estimators. In *Advances in Neural Information Processing Systems*, pages 361–369, 2011.
- [36] F. McSherry and K. Talwar. Mechanism design via differential privacy. In *FOCS*, pages 94–103. IEEE, 2007.
- [37] A. Narayanan and V. Shmatikov. Robust de-anonymization of large sparse datasets. In *Security and Privacy.IEEE Symposium on*, pages 111–125, 2008.
- [38] O. Sheffet. Locally private hypothesis testing. *CoRR*, abs/1802.03441, 2018.
- [39] B. W. Silverman. *Density Estimation*. Chapman and Hall, London, 1986.
- [40] A. T. Suresh, A. Orlitsky, J. Acharya, and A. Jafarpour. Near-optimal-sample estimators for spherical gaussian mixtures. In *NIPS*, pages 1395–1403. Curran Associates, Inc., 2014.
- [41] L. Sweeney. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05):557–570, 2002.
- [42] M. J. Wainwright, M. I. Jordan, and J. C. Duchi. Privacy aware learning. In *Advances in Neural Information Processing Systems*, pages 1430–1438, 2012.
- [43] S. Wang, L. Huang, P. Wang, Y. Nie, H. Xu, W. Yang, X. Li, and C. Qiao. Mutual information optimally local private discrete distribution estimation. *CoRR*, abs/1607.08025, 2016.
- [44] T. Wang and J. Blocki. Locally differentially private protocols for frequency estimation. In *Proceedings of the 26th USENIX Security Symposium*, 2017.
- [45] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965.
- [46] L. Wasserman and S. Zhou. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 105(489):375–389, 2010.
- [47] M. Ye and A. Barg. Optimal schemes for discrete distribution estimation under locally differential privacy. *CoRR*, abs/1702.00610, 2017.
- [48] Y. Zhang, M. J. Wainwright, and J. C. Duchi. Communication-efficient algorithms for statistical optimization. In *NIPS*, pages 1502–1510. 2012.

A Proof of Lemma 6

Let $T_i = \{x \in [k] \mid r_x \text{ is in the } i\text{th } H_b \text{ block}\}$ be the set of symbols such that r_x is in the i th H_b block. From the description of r_x , we obtain

$$T_i := \left\{ x \mid \lfloor \frac{x}{b-1} \rfloor = i \right\}, \quad p(T_i) := \sum_{x \in T_i} p(x) \quad \text{and} \quad \sum_i p(T_i) = 1. \quad (17)$$

We will prove that

$$p(C_x) = \frac{1}{2B-1+e^\varepsilon} + \frac{e^\varepsilon-1}{2(2B-1+e^\varepsilon)}p(x) + \frac{e^\varepsilon-1}{2(2B-1+e^\varepsilon)}p(T_i), \quad \text{and} \quad (18)$$

$$p(S_i) = \frac{e^\varepsilon-1}{2B-1+e^\varepsilon}p(T_i) + \frac{2}{2B-1+e^\varepsilon} \quad (19)$$

Then note that (18)– $\frac{1}{2}$ (19) gives Lemma 6.

Proof of (18). Recall that for any x ,

$$p(C_x) = \sum_{x'} p(x')Q(Z \in C_x | X = x'). \quad (20)$$

For any $x, x' \in [k]$, by (3) and $s = b/2$,

$$Q(Z \in C_x | X = x') = \frac{2e^\varepsilon}{be^\varepsilon + 2K - b} \times |C_x \cap C_{x'}| + \frac{2}{be^\varepsilon + 2K - b} \times |C_x \setminus C_{x'}|. \quad (21)$$

There are three cases:

- $x' = x$. In this case, $|C_x \cap C_{x'}| = s = b/2$, and $|C_x \setminus C_{x'}| = \emptyset$.
- $x' \in T_{\lfloor \frac{x}{b-1} \rfloor} \setminus \{x\}$: When this happens, then by the Property (ii) of Hadamard matrices, $|C_x \cap C_{x'}| = s/2 = b/4$, and $|C_x \setminus C_{x'}| = s/2 = b/4$.
- $x' \notin T_{\lfloor \frac{x}{b-1} \rfloor}$: The symbols $x' \notin T_{\lfloor \frac{x}{b-1} \rfloor}$ satisfy $|C_x \cap C_{x'}| = \emptyset$, and $|C_x \setminus C_{x'}| = b/2$.

Plugging these in (21), and using (20) with $K = Bb$, we obtain

$$\begin{aligned} p(C_x) &= p(x) \frac{e^\varepsilon}{2B-1+e^\varepsilon} + p\left(T_{\lfloor \frac{x}{b-1} \rfloor} \setminus \{x\}\right) \frac{e^\varepsilon+1}{2(2B-1+e^\varepsilon)} + p\left([k] \setminus T_{\lfloor \frac{x}{b-1} \rfloor}\right) \frac{1}{2B-1+e^\varepsilon} \\ &= \frac{e^\varepsilon \cdot p(x)}{2B-1+e^\varepsilon} + \frac{e^\varepsilon+1}{2(2B-1+e^\varepsilon)} \left(p\left(T_{\lfloor \frac{x}{b-1} \rfloor}\right) - p(x) \right) + \frac{1}{2B-1+e^\varepsilon} \left(1 - p\left(T_{\lfloor \frac{x}{b-1} \rfloor}\right) \right) \\ &= \frac{1}{2B-1+e^\varepsilon} + \frac{e^\varepsilon-1}{2(2B-1+e^\varepsilon)}p(x) + \frac{e^\varepsilon-1}{2(2B-1+e^\varepsilon)}p\left(T_{\lfloor \frac{x}{b-1} \rfloor}\right). \end{aligned}$$

Proof of (19). Recall that $p(S_i)$ is the probability that the output is in S_i , when input distribution is p . Note that for $x \in T_i$, $C_x \subset S_i$, hence $|S_i \cap C_x| = |C_x| = b/2$ and $C_x \setminus S_i = \emptyset$. For $x \notin T_i$, $C_x \cap S_i = \emptyset$ and $|C_x \setminus S_i| = |C_x| = b/2$. Again using (21) and replacing C_x with S_i , we obtain

$$p(S_i) = \frac{e^\varepsilon+1}{2B-1+e^\varepsilon} p(T_i) + \frac{2}{2B-1+e^\varepsilon} (1 - p(T_i)).$$

Rearranging the terms gives (19).

B Proof of sample complexity bounds (Theorem 7)

We will prove Theorem 7 from Lemma 6. The proof follows the general approach used in Section 4 for the high privacy regime. Subtracting (13) of Lemma 6 from (16), we obtain

$$\hat{p}(x) - p(x) = \frac{2(2B - 1 + e^\varepsilon)}{e^\varepsilon - 1} \cdot \left(\left(\widehat{p(C_x)} - p(C_x) \right) - \frac{1}{2} \left(\widehat{p(S_i)} - p(S_i) \right) \right), \quad (22)$$

where recall that S_i is the output columns corresponding to the r_x th block. Squaring both sides, and observing that for any reals $(a - b)^2 \leq 2(a^2 + b^2)$, we obtain

$$(\hat{p}(x) - p(x))^2 \leq \frac{8(2B - 1 + e^\varepsilon)^2}{(e^\varepsilon - 1)^2} \cdot \left(\left(\widehat{p(C_x)} - p(C_x) \right)^2 + \frac{1}{4} \left(\widehat{p(S_i)} - p(S_i) \right)^2 \right). \quad (23)$$

Now recall that $\widehat{p(C_x)}$ is average of independent Bernoulli's with mean $p(C_x)$, and $\widehat{p(S_i)}$ is the average of independent Bernoulli's with mean $p(S_i)$. Therefore,

$$\mathbb{E} \left[\left(\widehat{p(C_x)} - p(C_x) \right)^2 \right] = \frac{1}{n} p(C_x)(1 - p(C_x)) < \frac{1}{n} p(C_x), \quad (24)$$

and

$$\mathbb{E} \left[\left(\widehat{p(S_i)} - p(S_i) \right)^2 \right] = \frac{1}{n} p(S_i)(1 - p(S_i)) < \frac{1}{n} p(S_i). \quad (25)$$

Summing over x in (18) and using (17), we obtain

$$\begin{aligned} \sum_x \mathbb{E} \left[\left(\widehat{p(C_x)} - p(C_x) \right)^2 \right] &< \frac{1}{n} \sum_x \left(\frac{1}{2B - 1 + e^\varepsilon} + \frac{e^\varepsilon - 1}{2(2B - 1 + e^\varepsilon)} p(x) + \frac{e^\varepsilon - 1}{2(2B - 1 + e^\varepsilon)} p(T_i) \right) \\ &\leq \frac{1}{n} \left(\frac{k}{2B - 1 + e^\varepsilon} + \frac{e^\varepsilon - 1}{2(2B - 1 + e^\varepsilon)} + \frac{b(e^\varepsilon - 1)}{2(2B - 1 + e^\varepsilon)} \right). \end{aligned} \quad (26)$$

where the last inequality follows since each T_i is of size at most b , and $\sum_i p(T_i) = 1$, implying that $\sum_x p(T_i) \leq b$. Similarly, summing over x in (19),

$$\sum_x p(S_i) \leq \frac{b(e^\varepsilon - 1)}{2B - 1 + e^\varepsilon} + \frac{2k}{2B - 1 + e^\varepsilon}. \quad (27)$$

Summing over x in (23) and taking the expectations, and plugging the bounds above with the observation that $B < e^\varepsilon$ by design, we obtain the bound on $\mathbb{E} [\ell_2(\hat{p}, p)^2]$, proving the theorem.

$$\begin{aligned} \mathbb{E} [\ell_2^2(\hat{p}, p)] &\leq \frac{1}{n} \frac{8(2B - 1 + e^\varepsilon)^2}{(e^\varepsilon - 1)^2} \left(\left(\frac{2k + (b + 1)(e^\varepsilon - 1)}{2(2B - 1 + e^\varepsilon)} \right) + \frac{1}{4} \left(\frac{4k + 2b(e^\varepsilon - 1)}{2(2B - 1 + e^\varepsilon)} \right) \right) \\ &= \frac{1}{n} \frac{4(2B - 1 + e^\varepsilon)}{(e^\varepsilon - 1)^2} \left(3k + \left(\frac{3}{2}b + 1 \right) (e^\varepsilon - 1) \right). \\ &\leq \frac{36e^\varepsilon(k + b(e^\varepsilon - 1))}{n(e^\varepsilon - 1)^2}. \end{aligned} \quad (28)$$

C Description and performance of RAPPOR and SS

C.1 k -RAPPOR.

Recall from Section 2.2 the privatization mechanism of RAPPOR. For input $x \in [k]$, $\mathbf{y} \in \{0, 1\}^k$ is such that $\mathbf{y}_j = 1$ for $j = x$, and $\mathbf{y}_j = 0$ for $j \neq x$. The privatized output of RAPPOR is a k bit vector \mathbf{z} such that

$$Q(\mathbf{z}_j = \mathbf{y}_j) = \frac{e^{\varepsilon/2}}{e^{\varepsilon/2} + 1}, \text{ and } Q(\mathbf{z}_j = 1 - \mathbf{y}_j) = \frac{1}{e^{\varepsilon/2} + 1}.$$

[31] analyze the sample complexity of RAPPOR (See Table 1). We will consider the communication requirements now in Theorem 9.

Communication. The output of k -RAPPOR mechanism is described above with k bits. We now consider the communication requirements for any algorithm that faithfully sends the output of RAPPOR privatization to the server. By Shannon's coding theorem, any algorithm to do this requires at least $H(Z|p)$ bits even if it knows the distribution p .

Theorem 9. *The entropy Z of the output of RAPPOR for any input distribution satisfies*

$$H(Z) \geq \begin{cases} \Omega(k) & \text{when } \varepsilon < 1, \\ \Omega\left(\frac{k}{e^{\varepsilon/2}}\right) & \text{when } 1 < \varepsilon < 2 \log k, \end{cases}$$

and for the uniform input distribution $H(Z) \geq \log k$ when $\varepsilon > 2 \log k$.

Proof. For any input x , the outputs \mathbf{z}_j for $j \neq x$ are all i.i.d. $B\left(\frac{1}{1+e^{\varepsilon/2}}\right)$ random variables, where $B(r)$ is a Bernoulli random variable with bias r . Therefore the entropy of the output is at least $(k-1) \cdot h(1/(1+e^{\varepsilon/2}))$, where $h(r) := -r \log r - (1-r) \log(1-r)$ is the entropy of a $B(r)$ random variable. Note that $h(r) > -r \log r$. Therefore, $(k-1)h(1/(1+e^{\varepsilon/2})) > (k-1) \frac{\log(1+e^{\varepsilon/2})}{1+e^{\varepsilon/2}}$. For $\varepsilon < 1$ this bound reduces to $\Omega(k)$, and for any $\varepsilon < 2 \log k$ ignoring the logarithmic term gives the theorem. For the uniform input distribution, when $\varepsilon > 2 \log k$, we note that the output is nearly uniform on the basis vectors, giving the bound. \square

C.2 Subset Selection Approaches.

The papers [43, 47] propose sample optimal privacy mechanisms for all ranges of ε . The mechanism is as follows. The output is again k bits, and suppose $d = \lceil k/(e^\varepsilon + 1) \rceil$. The output is $\mathcal{Z} = \mathcal{Z}_{k,d}$, where $\mathcal{Z}_{k,d}$ is the set of all binary strings with Hamming weight d . For an $i \in [k]$, let $\mathcal{Z}_{k,d}^i$ be the elements in $\mathcal{Z}_{k,d}$ with 1 in the i th location. Then note that $|\mathcal{Z}_{k,d}| = \binom{k}{d}$, and $|\mathcal{Z}_{k,d}^i| = \binom{k-1}{d-1}$. Then, for $Z_1 \dots Z_k \in \mathcal{Z}_{k,d}$,

$$Q(Z_1 \dots Z_k | i) = \begin{cases} \frac{e^\varepsilon}{\binom{k-1}{d-1} e^{\varepsilon} + \binom{k-1}{d}}, & \text{for } Z_1 \dots Z_k \in \mathcal{Z}_{k,d}^i, \\ \frac{1}{\binom{k-1}{d-1} e^{\varepsilon} + \binom{k-1}{d}}, & \text{for } Z_1 \dots Z_k \in \mathcal{Z}_{k,d} \setminus \mathcal{Z}_{k,d}^i. \end{cases}$$

Communication. We can characterize the communication complexity by computing the entropy of the output distributions of the mechanism. The computations are similar to the last section, we simply state the entropy bounds that imply the communication bounds.

Suppose the underlying distribution is uniform. In this case, the output distribution is uniform among all possible strings of weight equal to d . Therefore the entropy of the output string is

identical to $\log \binom{k}{d}$, which is the optimal communication complexity per user. Note that for small ε this communication is strictly undesirable!

Theorem 10. *The entropy Z of the output of SS for any input distribution satisfies*

$$H(Z) \geq \begin{cases} \Omega(k) & \text{when } \varepsilon < 1, \\ \Omega\left(\frac{k}{e^\varepsilon}\right) & \text{when } 1 < \varepsilon < \log k, \end{cases}$$

and for the uniform input distribution $H(Z) \geq \log k$ when $\varepsilon > \log k$.

While the final mechanism is sample order optimal for all values of ε , the communication cost for each user depends critically on the value of ε . Table 2 characterizes the communication cost for all three well known mechanisms and our proposed method.