

# Improving Frequency Estimation under Local Differential Privacy

Milan Lopushaa-Zwakenberg\*, Zitao Li†, Boris Škorić\* and Ninghui Li†

\*Department of Mathematics and Computer Science

Eindhoven University of Technology

m.a.lopuhaa@tue.nl, b.skoric@tue.nl

†Department of Computer Sciences

Purdue University

li2490@purdue.edu, ninghui@cs.purdue.edu

**Abstract**—Local Differential Privacy protocols are stochastic protocols used in data aggregation when individual users do not trust the data aggregator with their private data. In such protocols there is a fundamental tradeoff between user privacy and aggregator utility. In the setting of frequency estimation, established bounds on this tradeoff are either nonquantitative, or far from what is known to be attainable. In this paper, we use information-theoretical methods to significantly improve established bounds. We also show that the new bounds are attainable for binary inputs. Furthermore, our methods lead to improved frequency estimators, which we experimentally show to outperform state-of-the-art methods.

**Index Terms**—Local Differential Privacy, frequency estimation, accuracy bound, privacy-utility tradeoff.

## 1. Introduction

In a context where a data aggregator collects potentially sensitive data, there is an inherent tension between the aggregator’s desire to obtain accurate population statistics and the individuals’ desire to protect their private data. One approach to protect privacy is offered by *Local Differential Privacy* (LDP) protocols [18]. Under this framework, each user randomises their private data before sending it to the aggregator. This hides the users’ true data, while for a large population size the randomness of the users cancels out, allowing the aggregator to obtain accurate estimates of the population statistics. Because of these properties, LDP mechanisms are widely used in industry by companies such as Apple [22], Google [15], and Microsoft [10].

One of the main settings in which LDP protocols are used is that of *frequency estimation* [27, 15, 26]. In this setting, every user has a private data item from a (finite) set  $\mathcal{A}$ , and the aggregator’s goal is to determine the frequencies of the elements of  $\mathcal{A}$  among the user population. For example, Chrome uses the LDP protocol RAPPOR to estimate the relative frequencies of homepages and used search engines among its userbase [15]. In this paper, we focus on frequency estimation, as it is both widely studied in the literature, and frequently applied in industry.

In the LDP setting there is a tradeoff between user privacy and frequency estimation accuracy. Intuitively, the more ‘random’ the privacy protocol is, the better it will hide an individual’s private data, but the more noisy the

aggregator’s estimations will be. Therefore, it is natural to ask the following question:

**Question 1.1.** *How can one characterise the relation between user privacy and aggregator utility in the LDP setting?*

Of course, to answer this question, we need to choose suitable privacy and utility metrics. As a privacy metric, we focus on  $\varepsilon$ -LDP, the *de facto* standard privacy metric for stochastic privacy protocols. The privacy parameter  $\varepsilon \in \mathbb{R}_{\geq 0}$  provides a measure of the worst-case leakage of the protocol; the advantage of this utility metric is that in practice it is often easily computed for specific protocols, and it provides privacy guarantees that hold in all situations. As a utility metric, we consider the mean squared error (MSE) for the frequency estimation, which is an often-used metric in the literature [14, 26, 25].

There are two ways to approach Question 1.1. The first method is to study the tradeoff between privacy and utility for specific protocols, either analytically or experimentally; the second method is to prove theoretic bounds on this tradeoff that hold for any protocol. Unfortunately, there is still a large gap between what can currently be achieved by these two approaches. To be more precise, given  $a = \#\mathcal{A}$  categories,  $n \gg 0$  users and LDP parameter  $\varepsilon \ll 1$ , the currently best achievable MSE is  $\frac{4a}{n\varepsilon^2}$  [14, 26], while the theoretical lower bound on the MSE is either  $\frac{a}{64n\varepsilon^2}$  or  $\Omega(\frac{a}{n\varepsilon^2})$ , depending on the precise utility metric (see Section 1.1). While this shows that asymptotically we can attain the optimal behaviour in  $a, n, \varepsilon$  up to a constant, in practice one needs to know this constant if one is to determine a satisfactory level of privacy based on  $a, n$ , and utility demands. Therefore, it is important to bridge the gap between practice and theoretical lower bounds.

One of the factors that plays a role in the existence of this gap is that a privacy protocol typically consists of two algorithms  $(\mathcal{Q}, \Phi)$ , where  $\mathcal{Q}$  is employed by the user to randomise their data, and  $\Phi$  is used by the aggregator to obtain frequency estimations from the randomised data. However, it is unknown what the optimal  $\Phi$  is given  $\mathcal{Q}$  [25]. Typically, one first uses an unbiased, linear estimator  $\Phi'$  called a frequency oracle [27, 15, 26], and then post-processes the results to make them more in line with what the aggregator expects a frequency distribution to look like [16, 25]. However, theoretical analyses of the utility of postprocessing methods are lacking. Furthermore, the

estimator  $\Phi'$  typically discards at least some of the information present in the randomised data, which leaves open the possibility that postprocessing the outcome of  $\Phi'$  does not lead to the optimal  $\Phi$ .

## 1.1. Our contributions

This paper presents two contributions towards answering Question 1.1, by working towards closing the gap in the privacy-utility tradeoff between existing protocols and theoretical lower bounds. To outline these contributions we first need to introduce some notation. In this paper, we consider two different, but related estimation problems: (i) the aggregator wants to estimate the probability distribution  $P$  from which the users' private data is drawn, and (ii) the aggregator wants to estimate the actual frequencies  $F$  of the private data. The distribution  $P$  is unknown to the aggregator, and as such we can consider it to be a random variable itself; its distribution reflects the aggregator's prior knowledge. It turns out that for a given protocol  $\mathcal{Q}$ , the optimal estimators  $\Pi$  for  $P$  and  $\Phi$  for  $F$  can be stated explicitly in terms of this prior distribution:

**Theorem 1.2** (Informal version of Theorem 3.1). *Let  $\mathcal{Q}$  be a privacy protocol, and for each  $i \leq n$ , let  $Y_i := \mathcal{Q}(X_i)$  be the randomisation of user  $i$ 's private data. Then given  $\vec{Y} = \vec{y}$ , the optimal estimators  $\Pi_{\text{opt}}$  for  $P$  and  $\Phi_{\text{opt}}$  for  $F$  are given by*

$$\Pi_{\text{opt}}(\vec{y}) = \mathbb{E}[P | \vec{Y} = \vec{y}], \quad (1)$$

$$\Phi_{\text{opt}}(\vec{y}) = \mathbb{E}[F | \vec{Y} = \vec{y}]. \quad (2)$$

This theorem 'solves' the problem of postprocessing. Unfortunately, computing  $\Pi_{\text{opt}}$  and  $\Phi_{\text{opt}}$  can be time-consuming in practice: if  $b$  is the size of the output space of  $\mathcal{Q}$ , then the time complexity is  $\mathcal{O}(n^{a(b-1)})$  under moderate assumptions on  $\mu$ , which grows unfeasibly large for large  $n$ . However, the following theorem shows that we can approximate  $\Pi_{\text{opt}}$  and  $\Phi_{\text{opt}}$  by estimators which are much easier to compute.

**Theorem 1.3** (Informal version of Theorem 3.3). *Let  $\Pi_{\text{MLE}}(\vec{y})$  be the maximum likelihood estimator of  $P$  given  $\vec{y}$ . Then  $\Pi_{\text{MLE}} \rightarrow \Pi_{\text{opt}}$  in probability as  $n \rightarrow \infty$ . Furthermore, one finds  $\Pi_{\text{MLE}}(\vec{y})$  by solving a convex optimisation problem of dimension  $(a-1)$  of which the complexity does not depend on  $n$ .*

Furthermore, we perform experiments that show that  $\Pi_{\text{MLE}}$  outperforms state-of-the-art postprocessing methods, demonstrating the validity of our approach. Since  $F \approx P$  for large  $n$ , we can also use  $\Pi_{\text{MLE}}$  to approximate  $\Phi_{\text{opt}}$ . This reduces the problem of finding the optimal  $(\mathcal{Q}, \Pi)$  or  $(\mathcal{Q}, \Phi)$  to finding the optimal  $\mathcal{Q}$ .

On the other hand, we use the information-theoretic methods from [19] to find new lower bounds:

**Theorem 1.4** (Informal version of Theorem 4.5). *For  $n \gg 0$  and  $\varepsilon \ll 1$ , one has, for any  $\mathcal{Q}$ ,  $\Pi$  and  $\Phi$ , and any distribution for  $P$ :*

$$\text{MSE}^{\text{distr}}(\mathcal{Q}, \Pi), \text{MSE}^{\text{freq}}(\mathcal{Q}, \Phi) \geq \frac{a}{n\varepsilon^2},$$

where  $\text{MSE}^{\text{distr}}(\mathcal{Q}, \Pi)$  is the expected mean squared error for distribution estimation, over the distributions of  $P$  and

TABLE 1: New results compared to known lower bounds and attainable MSE for  $\varepsilon \ll 1$ .

Distribution estimation		
Known lower bound [13]	$\frac{a}{64n\varepsilon^2}$	worst case
Attainable [1, 13]	$\frac{4a}{n\varepsilon^2}$	
New lower bound	$\frac{a}{n\varepsilon^2}$	any case, for $n \rightarrow \infty$
Frequency estimation		
Known lower bound [4]	$\Omega(\frac{a}{n\varepsilon^2})$	worst case
Attainable [26]	$\frac{4a}{n\varepsilon^2}$	
New lower bound	$\frac{a}{n\varepsilon^2}$	any case, for $n \rightarrow \infty$

the stochastic function  $\mathcal{Q}$ , and  $\text{MSE}^{\text{freq}}(\mathcal{Q}, \Phi)$  is the mean squared error for frequency estimation.

This result significantly improves known lower bounds, as can be seen from Table 1. Not only does this provide us with a constant for  $\text{MSE}^{\text{freq}}(\mathcal{Q}, \Phi)$  where first none was known, and significantly improves the constant for  $\text{MSE}^{\text{distr}}(\mathcal{Q}, \Pi)$ , it also holds for any prior distribution, rather than just the worst-case one. The downside, however, is that it only holds for asymptotically large  $n$ . However, one of the settings in which LDP is typically employed is one where  $n \gg a$ . It should be noted that in this introduction we write the results in terms of the limit case  $\varepsilon \rightarrow 0$  for simplicity, but our results also improve known bounds for any  $\varepsilon$  (see Section 4).

While earlier works also use information-theoretic methods to derive lower bounds, our methods rely on and expand upon the description of the asymptotic behaviour of the mutual information  $I(P; \vec{Y})$  in [19]. Since this description gives a limit, rather than a lower bound, this allows us to be more precise than earlier work.

The structure of this paper is as follows. In Section 2 we introduce the mathematical setting for this paper. In Section 3, we discuss how to find, and approximate computationally, the optimal  $\Pi$  given  $\mathcal{Q}$ . In Section 4, we prove the new lower bounds for the MSE. In Section 5, we evaluate the methods from Section 3 experimentally.

## 1.2. Related work

A lower bound for distribution estimation is given in Proposition 6 of [14]; from the proof in the arXiv version [13] we find that we can express this lower bound as  $\frac{a}{64n\varepsilon^2}$  for  $\varepsilon \ll 1$ . Their strategy is to first relate frequency estimation to a binary hypothesis testing problem, and then use information theory to bound the accuracy of this testing problem. They also show that a frequency MSE of  $\frac{4a}{n\varepsilon^2}$  can be attained by providing a specific protocol, namely a version of the privacy protocol RAPPOR [15]. Another privacy protocol with the same frequency MSE, inspired by coding theory, is described in [1].

The work in [4] looks at the 'dual' problem, i.e. given  $\varepsilon$ ,  $a$ , and an acceptable frequency MSE<sup>1</sup>  $\kappa$ , what is the minimal  $n$  for which a protocol exists that satisfies these criteria? In Theorem 1.7, they find that  $n \geq \Omega(\frac{a}{\kappa\varepsilon^2})$  (even when allowing for protocols that only offer  $(\varepsilon, \delta)$ -LDP, a weaker privacy notion), which we can restate as  $\kappa \geq \Omega(\frac{a}{n\varepsilon^2})$ . Their result is based on techniques in [2] to bound

<sup>1</sup>In the notation of [4] we would have  $\kappa = \alpha^2$ .

the mutual information between the input and output of  $Q$  in terms of the privacy parameters  $(\varepsilon, \delta)$ .

In [26], it is proven that RAPPOR attains a frequency MSE of  $\frac{4a}{n\varepsilon^2}$ . It is also shown there that this can be improved upon by tweaking its parameters, leading to a protocol called Optimised Unary Encoding (UE). This does not affect its asymptotic behaviour for  $\varepsilon \ll 1$ .

Instead of looking at the MSE, which is essentially the  $\ell_2$ -distance between the true value and its estimation, one can also consider other error metrics, such as  $\ell_1$  or  $\ell_\infty$ . An overview of the behaviour of these error metrics is found in [6].

There is much literature on studying LDP via information theory, and using this to derive properties on privacy or utility [2, 9, 14, 17]. We mostly rely on the techniques of [19], where information-theoretic properties are stated as (asymptotic) equalities, rather than inequalities; this allows us to obtain tighter bounds than earlier work.

There exists a body of work on post-processing LDP results [7, 16, 24, 25]. Heuristically, these approaches are based on the idea that one can improve an estimation by taking advantage of the knowledge one has of what a distribution should look like. In [16], the estimation is adjusted to adhere more to a power law-type distribution; in [25], the estimation is adjusted to account for the fact that the tallies should be nonnegative and add up to 1. The validity of these approaches is supported by empirical evidence, but a theoretical analysis is lacking. In this paper, we formalise the intuition that frequency estimation can be improved by taking prior knowledge into account, leading to the optimal  $\Pi$  and  $\Phi$  given  $Q$ .

## 2. Preliminaries

In this section, we introduce our setting and various concepts that play a role in this paper.

### 2.1. Setting

We consider the setting of *Local Differential Privacy* (LDP). There are  $n$  users, and user  $i$  has a private data item  $X_i$  from a finite set  $\mathcal{A}$ . Let  $\mathcal{P}_{\mathcal{A}}$  be the space of probability distributions on  $\mathcal{A}$ ; then we assume that there is a distribution  $P \in \mathcal{P}_{\mathcal{A}}$  such that each  $X_i$  is drawn independently from  $P$ . The distribution  $P$  itself is unknown to the aggregator, so we consider it to be a random variable itself, taken from a continuous prior distribution  $\mu$  on  $\mathcal{P}_{\mathcal{A}}$ . The distribution  $\mu$  is known to the aggregator, and reflects their prior knowledge. The aggregator publishes a *privacy protocol*, i.e. a random function  $Q: \mathcal{A} \rightarrow \mathcal{B}$ . User  $i$  calculates  $Y_i := Q(X_i)$  and sends it to the aggregator. The aggregator is interested in one of two quantities unavailable to them:

- 1) The aggregator wants to know  $P$  as accurately as possible. This occurs, for instance, when the aggregator is a scientist whose userbase is a sample of a greater population [27]. The aggregator is not concerned with this specific userbase, but rather with the characteristics of the general population, which are modelled by  $P$ .
- 2) For  $a = \#\mathcal{A}$ , define the frequency vector  $F \in \mathbb{R}_{\geq 0}^a$  by  $F_x = \frac{\#\{i: X_i=x\}}{n}$  for all  $x \in \mathcal{A}$ . Then the aggregator wants to know  $F$  as accurately as possible. This

$\mathcal{A}$	input space
$a$	$\#\mathcal{A}$
$n$	number of users
$\vec{X}$	vector of private data
$F$	frequencies of private data
$T$	tallies of private data
$\mathcal{P}_{\mathcal{A}}$	space of prop. distr. on $\mathcal{A}$
$P$	prob. distr. of private data
$\mu$	probability measure of $P$
$Q$	privacy protocol
$Q$	matrix associated to $Q$
$\mathcal{B}$	output space associated to $Q$
$b$	$\#\mathcal{B}$
$\varepsilon$	LDP-parameter
$\vec{Y}$	vector of perturbed data
$S$	tallies of perturbed data
$\Pi$	estimator of $P$ given $\vec{Y}$
$\Phi$	estimator of $F$ given $\vec{Y}$

TABLE 2: Notation employed in this paper.

occurs, for instance, when the users are customers of a service, and the service provider wants to know statistics about its customer base [15].

Note that for large  $n$ , one has  $F \approx P$ , so these goals are closely aligned in practice. Nevertheless, we make a distinction between these two cases, as the means of obtaining their lower bounds are different. To estimate  $P$ , the aggregator employs a distribution estimator  $\Pi$ , which takes as input the perturbed data  $\vec{Y}$ , and outputs an estimation  $\hat{P} = \Pi(\vec{Y}) \in \mathbb{R}^a$  of  $P$ . To estimate  $F$ , the aggregator likewise computes  $\hat{F} = \Phi(\vec{Y}) \in \mathbb{R}^a$ . The setting is depicted in Figure 1, and our notation is listed in Table 2; some notation will be defined later in this section. In particular, for  $\alpha \in \mathcal{A}$ , and  $\beta \in \mathcal{B}$ , we will make use of the notations

$$T_\alpha = \#\{i : X_i = \alpha\} = nF_\alpha, \quad (3)$$

$$S_\beta = \#\{i : Y_i = \beta\}, \quad (4)$$

$$S_{\beta|\alpha} = \#\{i : X_i = \alpha, Y_i = \beta\}. \quad (5)$$

We write  $T$  and  $S$  for the vectors in  $\mathbb{Z}_{\geq 0}^a$  and  $\mathbb{Z}_{\geq 0}^b$ , where  $b = \#\mathcal{B}$ . For a particular  $\vec{x} \in \mathcal{A}^n$ ,  $\vec{y} \in \mathcal{B}^n$ , we write  $t_\alpha$ ,  $s_\beta$ ,  $s_{\beta|\alpha}$  for the associated tallies; we write  $t(\vec{x})$ , etc., if we need to disambiguate between different input data.

### 2.2. Privacy and utility metrics

Let  $\mathcal{A}$  be a finite set. A *privacy protocol* for  $\mathcal{A}$  is a random function  $Q: \mathcal{A} \rightarrow \mathcal{B}$ , where  $\mathcal{B}$  is another finite set. Upon identifying  $\mathcal{A} = \{1, \dots, a\}$  and  $\mathcal{B} = \{1, \dots, b\}$ , we can represent  $Q$  by a matrix  $Q \in \mathbb{R}^{b \times a}$ , with  $Q_{y|x} = \mathbb{P}(Q(x) = y)$ . The *de facto* standard way to measure the privacy of a protocol is via  $\varepsilon$ -Local Differential Privacy (LDP):

**Definition 2.1.** ( $\varepsilon$ -LDP [18]) *Let  $Q$  be a privacy protocol for  $\mathcal{A}$ , and let  $\varepsilon \in \mathbb{R}_{\geq 0}$ . We say that  $Q$  satisfies  $\varepsilon$ -LDP if for all  $x, x' \in \mathcal{A}$  and all  $y \in \mathcal{B}$  one has*

$$Q_{y|x'} \leq e^\varepsilon Q_{y|x}. \quad (6)$$

Intuitively, this means that for small  $\varepsilon$ , given the output  $y$ , it is difficult to decide whether the input was  $x$  or  $x'$ . The smaller  $\varepsilon$ , the more privacy the protocol offers. An advantage of this metric is that it is a worst-case approach

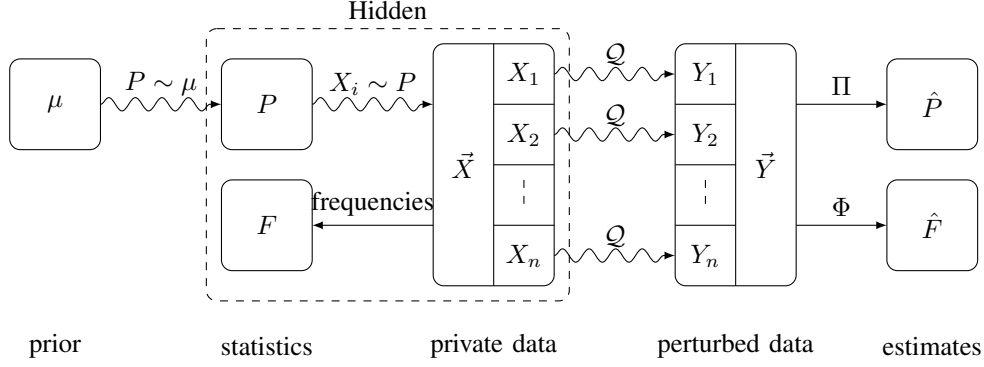


Figure 1: The LDP setting as used in this paper.

to privacy, ensuring strong privacy guarantees that hold in all situations.

On the side of utility, we measure the accuracy of the aggregator's estimation by taking the expected value (over  $P \sim \mu$ ) of the mean squared error:

**Definition 2.2.** We define the mean squared error of  $(Q, \Pi)$  and  $(Q, \Phi)$  to be

$$\text{MSE}_\mu^{\text{distr}}(Q, \Pi) := \mathbb{E}_{P, \vec{Y}} \|P - \Pi(\vec{Y})\|_2^2, \quad (7)$$

$$\text{MSE}_\mu^{\text{freq}}(Q, \Phi) := \mathbb{E}_{F, \vec{Y}} \|F - \Phi(\vec{Y})\|_2^2. \quad (8)$$

Note that these metrics depend on  $a$  and  $n$ . They also depend on  $\mu$ , as this affects the distributions of both  $P$  and  $F$ . Using this metric means that the best estimator is the one that gives on average the lowest squared error, when averaging over all possible input distributions, and all possible outputs. If one is more interested in worst-case performance, one can consider the following worst-case metrics [13, 4]:

**Definition 2.3.** We define the worst-case MSE of  $(Q, \Pi)$  and  $(Q, \Phi)$  to be

$$\text{WMSE}^{\text{distr}}(Q, \Pi) := \sup_p \mathbb{E}_{\vec{Y} | P=p} \left[ \|p - \Pi_n(\vec{Y})\|_2^2 \right], \quad (9)$$

$$\text{WMSE}^{\text{freq}}(Q, \Phi_n) := \max_f \mathbb{E}_{\vec{Y} | F=f} \left[ \|f - \Phi_n(\vec{Y})\|_2^2 \right]. \quad (10)$$

These are the metrics used by [4, 13] in Table 1. Note that

$$\text{WMSE}^{\text{distr}}(Q, \Pi) = \sup_\mu \text{MSE}_\mu^{\text{distr}}(Q, \Pi), \quad (11)$$

$$\text{WMSE}^{\text{freq}}(Q, \Phi) \geq \sup_\mu \text{MSE}_\mu^{\text{freq}}(Q, \Phi). \quad (12)$$

We prefer to use the MSE measures instead of the worst-case measures for two reasons: First, they align more closely with an aggregator's needs, since the aggregator will be interested in a protocol's behaviour on a typical database. Second, while worst case guarantees are helpful, the WMSE metrics do not fully give worst case guarantees, as they still involve an expected value over  $Q$ . In fact, since the LDP condition implies that any output has positive probability given any input, it is impossible to give a worst case guarantee on the estimation error.

Therefore it makes sense to consider average case metrics. Third, any lower bound on MSE will imply a lower bound on WMSE, making it more fruitful to look for lower bounds on MSE.

### 2.3. Probability and information theory

For a finite set  $\mathcal{A}$  of size  $a$ , we write  $\mathcal{P}_\mathcal{A}$  for the space of probability distributions on  $\mathcal{A}$ , i.e.

$$\mathcal{P}_\mathcal{A} = \left\{ p \in \mathbb{R}_{\geq 0}^a : \sum_{x \in \mathcal{A}} p_x = 1 \right\}. \quad (13)$$

For a discrete random variable  $X$  on  $\mathcal{A}$ , we write  $H(X)$  for its Shannon entropy, and for a continuous random variable  $Y$  on  $\mathbb{R}^d$ , we write  $h(Y)$  for its differential entropy. If  $P$  is a continuous random variable on  $\mathcal{P}_\mathcal{A}$ , we define its differential entropy as follows: Choose an identification  $\mathcal{A} = \{1, \dots, a\}$ , and define  $P' = (P_1, \dots, P_{a-1})$ , which is a continuous random variable on  $\mathbb{R}^{a-1}$ . We then define  $h(P) := h(P')$ ; this does not depend on the choice of enumeration of  $\mathcal{A}$ . As such, we fix such an enumeration for the rest of this paper, unless stated otherwise. Related information-theoretic measures, such as  $h(P|\vec{Y})$  and  $I(P; \vec{Y})$ , are defined similarly. For more details on information theory we refer to [8].

### 2.4. Assumptions on privacy protocols

Throughout this paper we make two technical but harmless assumptions on the privacy protocol  $Q$ :

- 1) We assume that the matrix  $Q$  has rank  $a$ , i.e. it is injective as a linear map. We do this because privacy protocols of lower rank are unable to distinguish all possible input distributions, which makes them unsuitable for frequency estimation [19].
- 2) We assume that for all  $x \in \mathcal{A}, y \in \mathcal{B}$  one has  $Q_{y|x} > 0$ . If there is an  $y$  such that  $Q_{y|x} = 0$  for all  $x$ , then we can remove  $y$  from  $\mathcal{B}$  without changing the protocol. If there is an  $y$  such that  $Q_{y|x} = 0$  but  $Q_{y|x'} > 0$ , then  $Q$  does not satisfy  $\varepsilon$ -LDP for any  $\varepsilon$ . Since we are only interested in privacy protocols that offer privacy, disregarding this case is harmless.

The main reason for these assumptions is that they simplify the expressions in Theorem 4.2 compared to [19], while still being general enough to describe all protocols that are used in practice.

### 3. Optimal estimators

We obtain our first upper bound for the MSE by giving formulas for the optimal estimators, given  $Q, \mu$ . Essentially, this is a well-known fact about minimum mean square error (MMSE) estimators. Recall that we have identified  $\mathcal{A} = \{1, \dots, a\}$ .

**Theorem 3.1.** *Let  $Q, \mu$  be given. The  $\Pi_{\text{opt}}$  and  $\Phi_{\text{opt}}$  minimising (7) and (8), respectively, are given by*

$$\Pi_{\text{opt}}(\vec{y}) = \mathbb{E}_{P|\vec{Y}=\vec{y}}[P], \quad (14)$$

$$\Phi_{\text{opt}}(\vec{y}) = \mathbb{E}_{F|\vec{Y}=\vec{y}}[F]. \quad (15)$$

For these estimators we have

$$\text{MSE}_{\mu}^{\text{distr}}(Q, \Pi_{\text{opt}}) = \sum_{x=1}^a \mathbb{E}_{\vec{y}} \text{Var}(P_x | \vec{Y} = \vec{y}), \quad (16)$$

$$\text{MSE}_{\mu}^{\text{freq}}(Q, \Phi_{\text{opt}}) = \sum_{x=1}^a \mathbb{E}_{\vec{y}} \text{Var}(F_x | \vec{Y} = \vec{y}). \quad (17)$$

*Proof.* For any estimator  $\Pi$  of  $P$  one has

$$\text{MSE}_{\mu}^{\text{distr}}(Q, \Pi) = \mathbb{E}_{\vec{y}} \mathbb{E}_{P|\vec{Y}=\vec{y}} [\|P - \Pi(\vec{y})\|_2^2] \quad (18)$$

$$= \mathbb{E}_{\vec{y}} \sum_{x=1}^a \mathbb{E}_{P|\vec{Y}=\vec{y}} [(P_x - \Pi(\vec{y})_x)^2] \quad (19)$$

$$\geq \mathbb{E}_{\vec{y}} \sum_{x=1}^a \text{Var}(P_x | \vec{Y} = \vec{y}). \quad (20)$$

One has equality in (20) if and only if  $\Pi(\vec{y})_x = \mathbb{E}_P[P_x | \vec{Y} = \vec{y}]$  for all  $x$  and  $\vec{y}$ . The proof for  $\Phi$  is similar.  $\square$

Theorem 3.1 formalises the intuition in [16, 25] that one gets better estimators by incorporating prior knowledge of the distribution. While Theorem 3.1 gives us a direct formula for the optimal estimators for a given privacy protocol, the disadvantage is that these can be computationally difficult to evaluate. With regards to  $\Pi_{\text{opt}}$ , we find that for any  $x$  we have

$$[\Pi_{\text{opt}}(\vec{y})]_x = \frac{1}{\mathbb{P}(\vec{Y} = \vec{y})} \int_{p \in \mathcal{P}_{\mathcal{A}}} \mu(p) p_x \prod_{\beta \in \mathcal{B}} (Q \cdot p)_{\beta}^{s_{\beta}} dp, \quad (21)$$

where  $Q \cdot p$  is matrix multiplication,  $s_{\beta}$  is as in (4), and  $\mu$  is the probability density function for  $p$ . For large  $a$ , the integral over  $\mathcal{P}_{\mathcal{A}}$  can be computationally involved. One approach is to do a Monte Carlo approximation of the integral. Typically,  $\mu$  will be a Dirichlet distribution, for which several efficient Monte Carlo methods exist [5, 3]. However, for large  $n$ , the  $s_{\beta}$  will be large as well. This leads to a spiky distribution, which needs more samples to approximate accurately. If the aggregator is interested in estimating  $P$ , we can circumvent this by giving an expression for  $\Pi_{\text{opt}}(\vec{y})$ , as well as for its MSE that does not involve any integration. Its complexity is stated in the following Proposition:

**Proposition 3.2.** *Let  $\Pi_{\text{opt}}$  be as in Theorem 3.1. Suppose  $\mu$  is a Dirichlet distribution, and let  $\vec{y} \in \mathcal{B}^n$ . Then the estimator  $\Pi_{\text{opt}}(\vec{y})$  and  $\text{MSE}_{\mu}^{\text{distr}}(Q, \Pi_{\text{opt}})$  can be calculated in time complexity  $\mathcal{O}(n^{b(a-1)})$  and  $\mathcal{O}(n^{ab-1})$ , respectively.*

This Proposition is proven in appendix A. Unfortunately, since  $n$  is typically very large in LDP settings, this can get computationally prohibitive. It therefore becomes useful to look for ways to approximate  $\Pi_{\text{opt}}$ . Below we discuss such an approximation method.

### 3.1. Approximation by MLE

For large  $n$ , the calculation in Proposition 3.2 will still be computationally involved. However, the following Theorem shows that we can efficiently and accurately approximate  $\Pi_{\text{opt}}$  by the maximum likelihood estimator:

**Theorem 3.3.** *For  $\vec{y} \in \mathcal{B}^n$ , let  $\Pi_{\text{MLE}}(\vec{y})$  be the maximum likelihood estimator of  $P$  given  $\vec{y}$ . Then  $\Pi_{\text{MLE}} \rightarrow \Pi_{\text{opt}}$  in probability as  $n \rightarrow \infty$ . Furthermore, one finds  $\Pi_{\text{MLE}}$  by solving an  $(a-1)$ -dimensional convex optimisation problem whose complexity does not depend on  $n$ .*

*Proof.* The convergence in probability is proven in [20]. Furthermore, since  $\mathbb{P}(\vec{Y} = \vec{y} | P = p) = \prod_{\beta} (Q \cdot p)_{\beta}^{s_{\beta}}$ , we find  $\Pi_{\text{MLE}}(\vec{y})$  by solving the optimisation problem

$$\min_{p \in \mathcal{P}_{\mathcal{A}}} \left\{ - \sum_{\beta \in \mathcal{B}} s_{\beta} \log((Q \cdot p)_{\beta}) \right\}. \quad (22)$$

The only effect of  $n$  is in the scaling of the objective function, but this does not influence the difficulty of minimisation.  $\square$

Since the objective function in (22) is smooth and convex in  $p$ , this can be solved quickly numerically. Using the MLE rather than the posterior also has the advantage that it is independent of the choice of  $\mu$ , making it a good choice of estimator when the prior is unknown.

Unfortunately, finding an expression for  $\Phi_{\text{MLE}}$  is more complicated. However, for large  $n$  we have  $F \approx P$ , so we can use the optimisation problem for  $P$  to get an approximate value for  $F$ . Note that our approach to approximating the MLE is different from that of [25], as there the MLE is approximated based on a noninjective transformation of the obfuscated tallies  $S$ , rather than on  $S$  itself. This transformation is protocol-specific, and hence this method does not easily extend to general protocols. The advantage of (22) is that it can be used to improve the estimation of any protocol.

In Section 5 we numerically evaluate how well  $\Pi_{\text{MLE}}$  and works as an estimator for both  $P$  and  $F$ .

### 4. Lower bounds on MSE

In this section, we prove our new lower bounds on  $\text{MSE}^{\text{distr}}(Q, \Pi)$  and  $\text{MSE}^{\text{freq}}(Q, \Phi)$  in terms of the LDP parameter  $\varepsilon$ . Our approach consists of multiple steps, but the general outline is as follows:

- 1) In Theorem 3.1 it is proven that one can state  $\text{MSE}^{\text{distr}}(Q, \Pi)$  and  $\text{MSE}^{\text{freq}}(Q, \Phi)$  as variances of  $P$  and  $F$  given  $\vec{Y}$ . We can bound these variances in terms of the entropies  $\mathbb{h}(P|\vec{Y})$  and  $\mathbb{H}(T|\vec{Y})$  (Theorem 4.1).
- 2) Using results about the asymptotic information-theoretic behaviour of LDP in [19], we express the

limit (as  $n \rightarrow \infty$ ) of these entropies in terms of properties of the matrix  $Q \in \mathbb{R}^{b \times a}$  (Theorem 4.2).

- 3) Finally, we bound these linear-algebraic constructs in terms of  $\varepsilon$ , giving the desired result (Theorem 4.5).

We also show that the MLE bound  $\frac{a}{n\varepsilon^2}$  is the optimal bound of this form, in the sense that it can be attained for  $a = 2$  for both distribution and frequency estimation (Corollary 4.7). Note that from a theoretical perspective, the intermediate steps can be of interest as well, as they offer lower bounds on the MSE from different viewpoints: information theory, linear algebra, and the LDP parameter  $\varepsilon$ .

#### 4.1. Reduction to information theory

The disadvantage of the optimal estimators of the previous section is that their MSEs can be hard to quantify. In this section, we give lower bounds for these, based on the information-theoretical quantities  $h(P|\vec{Y})$  and  $h(T|\vec{Y})$ . Intuitively, the lower the uncertainty about the value of  $P$  or  $T$ , the lower the average error on the estimation should be. While Fano's inequality expresses the same sentiment, we need the following version, which more closely aligns with the MSE rather than a binary fail/success estimator.

**Theorem 4.1.** *For any privacy protocol  $\mathcal{Q}$  one has*

$$\text{MSE}_\mu^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}}) \geq \frac{a}{2\pi e} e^{\frac{2}{a-1} h(P|\vec{Y})}, \quad (23)$$

$$\liminf_{n \rightarrow \infty} \frac{\text{MSE}_\mu^{\text{freq}}(\mathcal{Q}, \Phi_{\text{opt}})}{\frac{a}{2\pi en^2} e^{\frac{2}{a-1} H(F|\vec{Y})}} \geq 1. \quad (24)$$

*Proof.* We start with  $\Pi_{\text{opt}}$ . By (20) we have

$$\text{MSE}_\mu^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}}) = \mathbb{E}_{\vec{y}} \sum_{x=1}^a \text{Var}(P_x | \vec{Y} = \vec{y}) \quad (25)$$

$$= \frac{1}{a-1} \mathbb{E}_{\vec{y}} \sum_{\substack{x \leq a, \\ x' \neq x}} \text{Var}(P_{x'} | \vec{Y} = \vec{y}) \quad (26)$$

$$\geq \mathbb{E}_{\vec{y}} \sum_{x=1}^a \prod_{x' \neq x} \text{Var}(P_{x'} | \vec{Y} = \vec{y})^{\frac{1}{a-1}}. \quad (27)$$

Here the last inequality is the arithmetic-geometric mean inequality. Now consider the  $a$ -th summand; we again write  $P' = (P_1, \dots, P_{a-1})$ . The  $\{\text{Var}(P_{x'} | \vec{Y} = \vec{y})\}_{x' < a}$  are the diagonal coefficients of the positive definite matrix  $\text{Cov}(P' | \vec{Y} = \vec{y})$ . By Hadamard's inequality we have

$$\prod_{x' < a} \text{Var}(P_{x'} | \vec{Y} = \vec{y}) \geq \det \text{Cov}(P' | \vec{Y} = \vec{y}). \quad (28)$$

Furthermore, Theorem 9.6.5 of [8] shows that

$$\det \text{Cov}(P' | \vec{Y} = \vec{y}) \geq (2\pi e)^{1-a} e^{2h(P'|\vec{Y}=\vec{y})}. \quad (29)$$

hence

$$\prod_{x' < a} \text{Var}(P_{x'} | \vec{Y} = \vec{y})^{\frac{1}{a-1}} \geq \frac{1}{2\pi e} e^{\frac{2}{a-1} h(P'|\vec{Y}=\vec{y})}. \quad (30)$$

The discussion in section 2.3 shows us that  $h(P'|\vec{Y} = \vec{y}) = h(P|\vec{Y} = \vec{y})$ , and that this does not depend on the

enumeration of  $\mathcal{A}$ . It follows that (30) in fact holds for every summand in (27), hence we derive

$$\text{MSE}_\mu^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}}) \geq \frac{a}{2\pi e} \mathbb{E}_{\vec{y}} e^{\frac{2}{a-1} h(P|\vec{Y}=\vec{y})}. \quad (31)$$

Theorem 4.1 is now proven by the convexity of the exponential function, which tells us that

$$\mathbb{E}_{\vec{y}} e^{\frac{2}{a-1} h(P|\vec{Y}=\vec{y})} \geq e^{\frac{2}{a-1} \mathbb{E}_{\vec{y}} h(P|\vec{Y}=\vec{y})} = e^{\frac{2}{a-1} h(P|\vec{Y})}. \quad (32)$$

As for  $\Phi_{\text{opt}}$ , we let  $T = nF$  be as in (3). Then  $H(F|\vec{Y} = \vec{y}) = H(T|\vec{Y} = \vec{y})$  and  $\text{Var}(F_x | \vec{Y} = \vec{y}) = \frac{1}{n^2} \text{Var}(T_x | \vec{Y} = \vec{y})$ , for every  $x$  and  $\vec{y}$ . As  $n$  increases the covariance of  $T$  increases, and  $T$  approaches the discretisation (centered around the points of  $\mathbb{Z}^a$ ) of a continuous random variable  $\tilde{T}$ . It follows that  $H(T|\vec{Y} = \vec{y}) \approx h(\tilde{T}|\vec{Y} = \vec{y})$  for any  $\vec{y}$ . Analogous to the discussion on  $P$  above, we can show that

$$\mathbb{E}_{\vec{y}} \sum_{x=1}^a \text{Var}(T_x | \vec{Y} = \vec{y}) \approx \mathbb{E}_{\vec{y}} \sum_{x=1}^a \text{Var}(\tilde{T}_x | \vec{Y} = \vec{y}) \quad (33)$$

$$\geq \frac{a}{2\pi e} e^{\frac{2}{a-1} h(\tilde{T}|\vec{Y})} \quad (34)$$

$$\approx \frac{a}{2\pi e} e^{\frac{2}{a-1} H(T|\vec{Y})}, \quad (35)$$

with the approximations approaching equality as  $n \rightarrow \infty$ . It follows that  $\text{MSE}_\mu^{\text{freq}}(\mathcal{Q}, \Phi_{\text{opt}}) \gtrsim \frac{a}{2\pi en^2} e^{\frac{2}{a-1} H(F|\vec{Y})}$ .  $\square$

Unfortunately the result for frequencies only holds as  $n \rightarrow \infty$ , since it relies on approximating the discrete random variable  $F$  by a continuous one. For the interests of this paper, however, this is not too much of an inconvenience, since the results from [19] that we wish to apply require  $n \rightarrow \infty$  in the first place.

#### 4.2. Accuracy bounds from linear algebra

In the previous section, we gave a lower bound for the MSE in terms of the information-theoretic quantities  $h(P|\vec{Y})$  and  $H(F|\vec{Y})$ . We can obtain new lower bounds for the limit case by studying the behaviour of  $h(P|\vec{Y})$  and  $H(F|\vec{Y})$  as  $n \rightarrow \infty$ ; this expands on work in [19]. The resulting lower bound is weaker not in the sense that the bound is lower, but rather that it only applies to the limit case  $n \rightarrow \infty$ , rather than all  $n$ . However, the advantage is that this limit case can be formulated purely in terms of linear algebra, and as it does not depend on  $n$ , it is computationally more feasible for large amounts of users, which is typical in the LDP setting.

Before we can state the result, we first need a bit more notation. We fix an identification  $\mathcal{B} = \{1, \dots, b\}$ . For  $x \in \mathcal{A}$ , let  $w_x$  be the column vector  $(Q_{1|x}, \dots, Q_{b-1|x})^T \in \mathbb{R}^{b-1}$ . For  $x \in \mathcal{A}$  and  $p \in \mathcal{P}_A$ , the latter regarded as an  $a$ -dimensional column vector, we define matrices  $D_p$ ,  $E_x$  and  $G_p$  by

$$D_p := Q^T \text{diag}(Q \cdot p)^{-1} Q \in \mathbb{R}^{a \times a}, \quad (36)$$

$$E_x := \text{diag}(w_x) - w_x w_x^T \in \mathbb{R}^{(b-1) \times (b-1)}, \quad (37)$$

$$G_p := \sum_{x=1}^a p_x E_x \in \mathbb{R}^{(b-1) \times (b-1)}. \quad (38)$$

Furthermore, we define constants  $\gamma_\mu(\mathcal{Q})$ ,  $\delta_\mu(\mathcal{Q})$  by

$$\gamma_\mu(\mathcal{Q}) = \frac{a-1}{2} \log(2\pi e) - \frac{1}{2} \mathbb{E}_P \log \det D_P, \quad (39)$$

$$\delta_\mu(\mathcal{Q}) = \gamma_\mu(\mathcal{Q}) + \frac{1}{2} \mathbb{E}_P \log \frac{\det G_P}{\prod_{y=1}^b (Q \cdot P)_y}. \quad (40)$$

While the matrix  $G_P$  depends on the choice of enumeration of  $\mathcal{B}$ , the resulting constant  $\delta_\mu(\mathcal{Q})$  does not. The introduction of these constants allows us to formulate the following Theorem:

**Theorem 4.2.** *It holds that*

$$\lim_{n \rightarrow \infty} \mathbb{h}(P|\bar{Y}) + \frac{a-1}{2} \log n = \gamma_\mu(\mathcal{Q}), \quad (41)$$

$$\lim_{n \rightarrow \infty} \mathbb{H}(F|\bar{Y}) - \frac{a-1}{2} \log n = \delta_\mu(\mathcal{Q}). \quad (42)$$

*Proof.* Framed in the language of this section, and applying our assumption that  $Q$  has rank  $a$ , [19, Thm. 6.7.1] states that

$$\lim_{n \rightarrow \infty} \mathbb{I}(\bar{Y}; P) - \frac{a-1}{2} \log n = -\gamma_\mu(\mathcal{Q}) + \mathbb{h}(P). \quad (43)$$

Since  $\mathbb{h}(P|\bar{Y}) = \mathbb{h}(P) - \mathbb{I}(\bar{Y}; P)$ , this proves (41). Similarly, we have  $\mathbb{H}(F|\bar{Y}) = \mathbb{H}(T|\bar{Y}) = \mathbb{H}(T) - \mathbb{I}(\bar{Y}; T) = \mathbb{H}(T) - \mathbb{I}(S; T)$ , where  $S$  and  $T$  are as in (3) and (4). The last equation holds because given  $T$ ,  $S$  is a sufficient statistic for  $\bar{Y}$ . We start by describing the limit behaviour of  $\mathbb{I}(S; T) = \mathbb{H}(S|P) + \mathbb{I}(S; P) - \mathbb{H}(S|T)$ . By [19, Lem. C.3], we have

$$\lim_{n \rightarrow \infty} \mathbb{H}(S|P) - \frac{b-1}{2} \log n = \frac{b-1}{2} \log(2\pi e) + \frac{1}{2} \sum_{y=1}^b \mathbb{E}_P[(Q \cdot P)_y] \quad (44)$$

Furthermore,  $\mathbb{I}(\bar{Y}; P) = \mathbb{I}(S; P)$ , so by (43) we get

$$\lim_{n \rightarrow \infty} \mathbb{I}(S; P) - \frac{a-1}{2} \log n = -\gamma_\mu(\mathcal{Q}) + \mathbb{h}(P). \quad (45)$$

It remains to study  $\mathbb{H}(S|T)$ . Let  $S_{y|x}$  be as in (5), and let  $S_{\bullet|x} = (S_{1|x}, \dots, S_{b|x}) \in \mathbb{Z}_{\geq 0}^b$ . Then  $S_{\bullet|x}$  follows a multinomial distribution with  $T_x$  samples and probability vector  $(Q_{1|x}, \dots, Q_{b|x})$ . Let  $S'_x := (S_{1|x}, \dots, S_{b-1|x}) \in \mathbb{Z}^{b-1}$ . By the multivariate de Moivre–Laplace theorem [23], we know that as  $n$  goes to infinity,  $S'_x$  can be approximated by the discretisation of a multivariate normal distribution with mean  $T_x w_x$  and covariance matrix  $T_x E_x$ . Applying [11, Lem. 1.1] to the matrix  $\text{diag}(w_x)$  and the vectors  $w_x$  and  $-w_x$ , we find

$$\det E_x = \prod_{y=1}^b Q_{y|x}. \quad (46)$$

Since we assumed in Section 2.4 that each  $Q_{y|x}$  is strictly positive, this means that  $E_x$  is nonsingular, and hence the associated multivariate normal distribution is nonsingular. Let  $S' = (S_1, \dots, S_{b-1})$ ; then  $S' = \sum_{x=1}^a S'_x$ , so for large  $n$ , given  $T$ , the random variable  $S'$  can be approximated by a multivariate normal variable  $N'$  with mean  $\sum_{x=1}^a T_x w_x$  and covariance matrix  $\sum_{x=1}^a T_x E_x$ . Using the known formula for the differential entropy of a multivariate normal distribution [8], it follows that for large  $n$

$$\mathbb{H}(S|T) = \mathbb{H}(S'|T) \quad (47)$$

$$\approx \mathbb{h}(N'|T) \quad (48)$$

$$\approx \mathbb{E}_T \left[ \frac{b-1}{2} \log(2\pi e) + \frac{1}{2} \log \det \left( \sum_{x=1}^a T_x E_x \right) \right]. \quad (49)$$

Here  $\approx$  means ‘the difference goes to 0 as  $n \rightarrow \infty$ ’. Since  $n^{-1}T \approx P$  for large  $n$ , we find

$$\mathbb{H}(S|T) \approx \frac{b-1}{2} \log(2\pi e n) + \mathbb{E}_P \left[ \frac{1}{2} \log \det G_P \right]. \quad (50)$$

Combining (44), (45) and (50) now finishes the proof of (42).  $\square$

Combining Theorems 4.1 and 4.2 allows us bound the MSE in terms of  $\gamma_\mu(\mathcal{Q})$  and  $\delta_\mu(\mathcal{Q})$ .

**Corollary 4.3.** *One has*

$$\liminf_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}}) \geq \frac{a}{2\pi e} e^{\frac{2}{a-1} \gamma_\mu(\mathcal{Q})}, \quad (51)$$

$$\liminf_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{freq}}(\mathcal{Q}, \Phi_{\text{opt}}) \geq \frac{a}{2\pi e} e^{\frac{2}{a-1} \delta_\mu(\mathcal{Q})}. \quad (52)$$

*Proof.* One has  $\lim_{n \rightarrow \infty} n \frac{a}{2\pi e} e^{\frac{2}{a-1} \mathbb{h}(P|\bar{Y})} = \frac{a}{2\pi e} e^{\frac{2}{a-1} \gamma_\mu(\mathcal{Q})}$  by Theorem 4.2. On the other hand, by Theorem 4.1, one has

$$\liminf_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}}) \geq \liminf_{n \rightarrow \infty} n \frac{a}{2\pi e} e^{\frac{2}{a-1} \mathbb{h}(P|\bar{Y})}. \quad (53)$$

Combining these two statements proves (51). Equation (52) can be proven analogously.  $\square$

We should interpret this Corollary as stating that in the best case we have  $\text{MSE}_\mu^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}}), \text{MSE}_\mu^{\text{freq}}(\mathcal{Q}, \Phi_{\text{opt}}) = \Omega(n^{-1})$  for fixed  $\mathcal{Q}$  and  $a$ , and we can bound the constants involved.

### 4.3. Accuracy bounds from $\varepsilon$ -LDP

While the constants  $\gamma_\mu(\mathcal{Q})$  and  $\delta_\mu(\mathcal{Q})$  do not depend on  $n$ , they still involve integration over  $\mathcal{P}_A$ , and as such can be computationally difficult for large  $a$ . However, it is possible to give lower bounds for these constants that are independent of  $\mu$ , and whose  $\mathcal{Q}$ -dependence only appears in the privacy parameter  $\varepsilon$ .

**Theorem 4.4.** *Suppose  $\mathcal{Q}$  satisfies  $\varepsilon$ -LDP. Then*

$$\gamma_\mu(\mathcal{Q}) \geq (a-1) \log \frac{\sqrt{2\pi e}}{e^\varepsilon - 1}, \quad (54)$$

$$\delta_\mu(\mathcal{Q}) \geq (a-1) \log \frac{\sqrt{2\pi e}}{e^\varepsilon - 1} - \frac{b\varepsilon}{2}. \quad (55)$$

*Proof.* In the terminology of the present paper,  $\mathbb{U}_\mu^{\text{as}}(\mathcal{Q})$  of [19, Def. 6.5] equals  $-\frac{1}{a-1} \gamma_\mu(\mathcal{Q})$ , and  $\mathbb{S}^{\text{wc}}(\mathcal{Q})$  of [19, Def. 5.1] equals  $e^{-\varepsilon'}$ , where  $\varepsilon' \leq \varepsilon$  is minimal such that  $\mathcal{Q}$  satisfies  $\varepsilon'$ -LDP. In this terminology, [19, Prop. 8.1] tells us that

$$\gamma_\mu(\mathcal{Q}) \geq (a-1) \log \frac{\sqrt{2\pi e}}{e^{\varepsilon'} - 1} \geq (a-1) \log \frac{\sqrt{2\pi e}}{e^\varepsilon - 1}, \quad (56)$$

proving (54). As for (55), by the definition of  $\varepsilon$ -LDP, for every  $x \in \mathcal{A}$  and every  $p \in \mathcal{P}_A$  we have  $\prod_{y=1}^b Q_{y|x} \geq e^{-b\varepsilon} \prod_{y=1}^b (Q \cdot p)_y$ . Let  $E_x$  be as in (37). By (46) we get

$$\det E_x = \prod_{y=1}^b Q_{y|x} \geq e^{-b\varepsilon} \prod_{y=1}^b (Q \cdot p)_y \quad (57)$$

for every  $x$  and  $p$ . Since  $\log \det$  is concave on the space of positive symmetric matrices, we find for every  $p$  that

$$\log \det G_p \geq \sum_x p_x \log \det E_x \geq -b\varepsilon + \sum_{y=1}^b \log(Q \cdot p)_y. \quad (58)$$

Combined with the definition of  $\delta_\mu(Q)$  this now directly proves (55).  $\square$

As a corollary of this Theorem, we find a lower bound for the MSE in terms of  $\varepsilon$ . We write  $f(\varepsilon) \succeq g(\varepsilon)$  if  $\liminf_{\varepsilon \rightarrow 0} \frac{f(\varepsilon)}{g(\varepsilon)} \geq 1$ . This Theorem follows directly from substituting the bounds for  $\gamma_\mu(Q)$  and  $\delta_\mu(Q)$  from Theorem 4.4 into Corollary 4.3.

**Theorem 4.5.** *Let  $\Pi_{\text{opt}}$  and  $\Phi_{\text{opt}}$  be as in Theorem 4.1. Then for every  $\mu$  one has*

$$\liminf_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{distr}}(Q, \Pi_{\text{opt}}) \geq \frac{a}{(e^\varepsilon - 1)^2}, \quad (59)$$

$$\liminf_{n \rightarrow \infty} \text{MSE}_\mu^{\text{freq}}(Q, \Phi_{\text{opt}}) \geq \frac{ae^{-\frac{b}{2(a-1)\varepsilon}}}{(e^\varepsilon - 1)^2}. \quad (60)$$

For  $n \in \mathbb{Z}_{>0}$  and  $\varepsilon \in \mathbb{R}_{>0}$ , let  $(Q_{n,\varepsilon}, \Pi_{\text{opt}})$  be the pair of an  $\varepsilon$ -LDP privacy protocol and an estimator for  $P$  minimising (7), and let  $(Q'_{n,\varepsilon}, \Phi_{\text{opt}})$  be the pair of an  $\varepsilon$ -LDP privacy protocol and an estimator for  $T$  minimising (8). Then as  $\varepsilon \rightarrow 0$ , we have

$$\liminf_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{distr}}(Q_{n,\varepsilon}, \Pi_{\text{opt}}) \succeq \frac{a}{\varepsilon^2}, \quad (61)$$

$$\liminf_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{freq}}(Q'_{n,\varepsilon}, \Phi_{\text{opt}}) \succeq \frac{a}{\varepsilon^2}. \quad (62)$$

Note that these results are strictly better than what is known in the literature for the case  $n \rightarrow \infty$ , to the best of our knowledge: (59) improves the result in the proof of Proposition 6 of [14] in three ways: First, our result does not just give a bound for the WMSE, but also for the MSE. Also, we improve the lower bound by a factor 64. Furthermore, the result of [14] only holds for  $\varepsilon \leq 1$ . The downside of our result, however, is that it only holds for the limit case  $n \rightarrow \infty$ .

As for the results for  $\Phi_n$ , it follows from results in [4] that the optimal  $(Q, \Phi_{n,\varepsilon})$  satisfies  $\text{WMSE}^{\text{freq}}(Q', \Phi_{n,\varepsilon}) = \Omega(\frac{a}{n\varepsilon^2})$  for small  $\varepsilon$ , but the authors do not make a statement about the constants involved. Equation (60) improves upon this by giving a quantitative lower bound for the MSE, which itself is a lower bound for the WMSE. Also, the OUE and OLH protocols from [26] performs as  $\approx \frac{4a}{n\varepsilon^2}$  for large  $n$  and small  $\varepsilon$ , while our lower bound is of the form  $\frac{a}{n\varepsilon^2}$ . Therefore, our bound is quite near to what is possible in practice.

The bound in (60) seems to imply that when looking for optimal  $Q$ , we should take  $b$  as large as possible. This seemingly contradicts results in [17], where it is found that taking  $b = a$  is always sufficient. There are two possible explanations for this. First, it is possible that our bounds are not sharp enough to accurately detect the dependence on  $b$ ; this is especially probable since (60) is only the latest in a chain of inequalities. However, the discrepancy between our results and those of [17] can also be caused by the fact that different utility metrics were being used. Whereas we focus on asymptotically

many users, the utility metric in [17] looks at the KL-divergence between the probability distributions induced by the private datum of one user. It is a possibility that the optimal protocols for these different metrics do not coincide. Overall, the  $b$ -dependence of the estimation error is difficult to assess, as there exist protocols achieving the best known MSE  $\frac{4a}{n\varepsilon^2}$  for both  $b = 2a$  [1] and  $b = 2^a$  [14].

#### 4.4. Tightness for $a = 2$

In this section, we show that if  $a = 2$  and as  $\varepsilon \rightarrow 0$ , the bounds in Theorem 4.5 are tight. For this, we recall the Randomised Response protocol [27] for  $a = 2$ , which, for an  $\varepsilon > 0$ , is the LDP protocol  $\text{RR}_\varepsilon: \{1, 2\} \rightarrow \{1, 2\}$  given by the matrix

$$\begin{pmatrix} \frac{e^\varepsilon}{e^\varepsilon+1} & \frac{1}{e^\varepsilon+1} \\ \frac{1}{e^\varepsilon+1} & \frac{e^\varepsilon}{e^\varepsilon+1} \end{pmatrix}. \quad (63)$$

Note that  $\text{RR}_\varepsilon$  satisfies  $\varepsilon$ -LDP. Let  $s_1, s_2$  be as in (4). As estimators for  $P$  and  $T$ , respectively, we use the maps  $\Pi_{\text{RR}}, \Phi_{\text{RR}}: \{1, 2\}^n \rightarrow \mathbb{R}^2$  given by

$$\Pi_{\text{RR}}(\vec{y}), \Phi_{\text{RR}}(\vec{y}) = \left( \frac{(e^\varepsilon + 1)s_1 - 1}{n(e^\varepsilon - 1)}, \frac{(e^\varepsilon + 1)s_2 - 1}{n(e^\varepsilon - 1)} \right). \quad (64)$$

These are unbiased estimators for  $P$  and  $F$  [27].

**Proposition 4.6.** *For any prior  $\mu$  of  $P$  one has*

$$\text{MSE}_\mu^{\text{distr}}(\text{RR}_\varepsilon, \Pi_{\text{RR}}) = \frac{2}{n} \left( \frac{e^\varepsilon}{(e^\varepsilon - 1)^2} + \mathbb{E}_P[P_1 P_2] \right), \quad (65)$$

$$\text{MSE}_\mu^{\text{freq}}(\text{RR}_\varepsilon, \Phi_{\text{RR}}) = \frac{2e^\varepsilon}{n(e^\varepsilon - 1)^2}. \quad (66)$$

*Proof of Proposition 4.6.* Since  $\Pi_{\text{RR}}$  is an unbiased estimator of  $P$  and  $\Pi_{\text{RR}}(\vec{Y})_1 + \Pi_{\text{RR}}(\vec{Y})_2 = 1$ , we have

$$\begin{aligned} \text{MSE}_\mu^{\text{distr}}(\text{RR}_\varepsilon, \Pi_{\text{RR}}) &= \mathbb{E}_p \text{Var}(\Pi_{\text{RR}}(\vec{Y})_1 | P = p) \\ &\quad + \mathbb{E}_p \text{Var}(\Pi_{\text{RR}}(\vec{Y})_2 | P = p) \\ &= 2\mathbb{E}_p \text{Var}(\Pi_{\text{RR}}(\vec{Y})_1 | P = p). \end{aligned} \quad (68)$$

For a given  $P = p$ , we know that  $S_1$  is binomially distributed with  $n$  samples and probability  $\frac{1 + (e^\varepsilon - 1)p_1}{e^\varepsilon + 1}$ . Substituting this into (64) yields

$$\begin{aligned} \text{Var}(\Pi_{\text{RR}}(\vec{Y})_1 | P = p) &= \frac{(e^\varepsilon + 1)^2}{n^2(e^\varepsilon - 1)^2} \text{Var}(S_1 | P = p) \\ &= \frac{(e^\varepsilon + 1)^2}{n^2(e^\varepsilon - 1)^2} \cdot n \cdot \frac{1 + (e^\varepsilon - 1)p_1}{e^\varepsilon + 1} \cdot \frac{e^\varepsilon - (e^\varepsilon - 1)p_1}{e^\varepsilon + 1} \end{aligned} \quad (69)$$

$$= \frac{(e^\varepsilon + 1)^2}{n^2(e^\varepsilon - 1)^2} \cdot n \cdot \frac{1 + (e^\varepsilon - 1)p_1}{e^\varepsilon + 1} \cdot \frac{e^\varepsilon - (e^\varepsilon - 1)p_1}{e^\varepsilon + 1} \quad (70)$$

$$= \frac{e^\varepsilon + (e^\varepsilon - 1)^2 p_1 (1 - p_1)}{n(e^\varepsilon - 1)^2}. \quad (71)$$

Equation (65) follows directly from this. Since  $\Phi_{\text{RR}}$  is an unbiased estimator of  $F = n^{-1}T$ , we analogously find that we only need to determine  $\text{Var}(\Phi_{\text{RR}}(\vec{Y})_1 | T = t)$ . Let  $S_{1|1}, S_{1|2}$  be as in (5). Then  $S_1 = S_{1|1} + S_{1|2}$ , and  $S_{1|1}$  is binomially distributed with  $t$  samples and probability  $\frac{e^\varepsilon}{e^\varepsilon + 1}$ , while  $S_{1|2}$  is binomially distributed with  $n - t$



samples and probability  $\frac{1}{e^\varepsilon + 1}$ . Since  $S_{1|1}$  and  $S_{1|2}$  are independent given  $T$ , it follows from (64) that

$$\begin{aligned} & \text{Var}(\Phi_{\text{RR}}(\vec{Y})_1 | T = t) \\ &= \frac{(e^\varepsilon + 1)^2}{(e^\varepsilon - 1)^2} (\text{Var}(S_{1|1} | T = t) + \text{Var}(S_{1|2} | T = t)) \end{aligned} \quad (72)$$

$$= \frac{(e^\varepsilon + 1)^2}{(e^\varepsilon - 1)^2} \left( \frac{te^\varepsilon}{(e^\varepsilon + 1)^2} + \frac{(n-t)e^\varepsilon}{(e^\varepsilon + 1)^2} \right) \quad (73)$$

$$= \frac{ne^\varepsilon}{(e^\varepsilon - 1)^2}. \quad (74)$$

It follows that

$$\text{MSE}_\mu^{\text{tally}}(\text{RR}_\varepsilon, \Phi_{\text{RR}}) = 2\mathbb{E}_t \text{Var}(\Phi_{\text{RR}}(\vec{Y})_1 | T = t) = \frac{2ne^\varepsilon}{(e^\varepsilon - 1)^2} \quad (75)$$

□

This Proposition yields the following Corollary. Below, we use  $f(\varepsilon) \sim g(\varepsilon)$  to denote  $\lim_{\varepsilon \rightarrow 0} \frac{f(\varepsilon)}{g(\varepsilon)} = 1$ .

**Corollary 4.7.** *Let  $a = 2$ , and let  $\mu$  be given. Let  $(\mathcal{Q}_{n,\varepsilon}, \Pi_{\text{opt}})$  and  $(\mathcal{Q}'_{n,\varepsilon}, \Phi_{\text{opt}})$  be as in Theorem 4.5. Then as  $\varepsilon \rightarrow 0$  we have*

$$\lim_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{distr}}(\mathcal{Q}_{n,\varepsilon}, \Pi_{n,\varepsilon}) \sim \frac{2}{\varepsilon^2}, \quad (76)$$

$$\lim_{n \rightarrow \infty} n \text{MSE}_\mu^{\text{freq}}(\mathcal{Q}'_{n,\varepsilon}, \Phi_{n,\varepsilon}) \sim \frac{2}{\varepsilon^2}. \quad (77)$$

*Proof.* A lower bound (of the behaviour in  $\varepsilon$  as  $\varepsilon \rightarrow 0$ ) is provided by Theorem 4.5, while an upper bound is provided by Proposition 4.6. □

In other words, the bounds in Theorem 4.5 become tight when  $\varepsilon \rightarrow 0$ , provided that  $a = 2$ . This shows that our bounds in Table 1 give the best possible coefficient for  $\frac{a}{n\varepsilon^2}$  that holds for all  $a$ .

## 5. MLE experiments

**Synthetic dataset.** We perform synthetic experiments to see how the estimator  $\Pi_{\text{MLE}}$  from Section 3.1 performs in practice. We apply the MLE to two well-established privacy protocols, Randomised Response (RR) [27] and Unary Encoding (UE) [26]; the latter is one of the protocols that obtains the known optimal error  $\frac{4a}{n\varepsilon^2}$  for  $\hat{P}$  and  $\hat{F}$ . Both of these protocols are parametrised by their LDP parameter  $\varepsilon$ . We take  $\varepsilon \in [0.2, 2]$ . Since  $b = 2^a$  for UE, the number of summands in (22) grows too large to handle for large  $a$ . Therefore we take  $a = 10$  for UE. For RR  $a = b$ , so we do not have this problem, and we take the more general setting  $a = 1024$ . We consider both large number of samples ( $n \gg a$ ) and moderate number of samples ( $n \approx 10a$ ) scenarios. For  $\mu$ , we take the Jeffreys prior on  $\mathcal{P}_A$ , i.e. the symmetric Dirichlet distribution with parameter  $-\frac{1}{2}$ . We take this prior because it is noninformative, as its definition does not depend on the parametrisation of  $\mathcal{P}_A$ . We draw  $P$  from  $\mu$  100 times and generate a dataset of  $n$  users from it. We then randomise the data via the LDP protocol, and perform MLE on the outcome to produce an estimate  $\hat{P}$  of  $P$ ; we furthermore use  $\hat{F} = \hat{P}$  as an estimator of  $F$ . Finally, we measure  $\|P - \hat{P}\|_2^2$  and  $\|F - \hat{F}\|_2^2$ , and average this over all samples to determine the MSE. To calculate

the MLE, we use projected gradient descent to solve (22) for UE. For RR, there is a direct method to find  $\hat{P}$ , see Appendix A.

We compare the MLE estimator to two other estimators: First, the baseline Frequency Oracles (FO), which are affine transformations of  $S$  used to produce unbiased estimators of  $P$  and  $F$ . Second, we look at Norm-Sub, which was found in [25] to be postprocessing method of the FO outcome that gives the best MSE, among a wide selection of considered postprocessing methods.

Figures 2 (RR) and 3 (UE) show the experimental results. The simulation results show that both Norm-sub and MLE post-processing can elevate the accuracy of the results. For RR, we see that Norm-Sub and MLE give the same accuracy when  $n \gg a$ . This is not unexpected, as the results in [25] show that for frequency estimation, Norm-Sub gives similar accuracy to MLE, and for RR the MLEs defined here and in that paper are equivalent. However, when  $n \approx 10a$ , MLE results are better than Norm-Sub, suggesting that when we use RR and  $n$  is comparable to  $a$ , MLE can give better accuracies. For UE, we see that MLE gives more accurate frequency estimations than Norm-sub when  $\varepsilon > 1.5$  when  $n \gg a$ . For a moderate number of samples cases, the results of MLE and Norm-Sub are similar.

**Real world datasets.** To show that the MLE estimator can help to improve the estimation on real world datasets, we also apply our method on two real world datasets: the taxi dataset [21] and the adult dataset [12]. We use discretized taxi distance ( $a = 1000$ , valid samples  $n = 434195$ ) in the taxi dataset and age ( $n = 32561$ ,  $a = 75$ ) in the adult dataset for RR, taxi payment types ( $a = 5$ , valid samples  $n = 359902$ ) and work classes ( $n = 32561$ ,  $a = 8$ ) in adult dataset for UE. We choose different attributes for RR and UE because of the same reason as in synthetic experiments (the output domain size  $b = 2^a$  for UE). Since the true probability  $P$  of real world datasets is unknown, we can only compare the frequency error,  $\|F - \hat{F}\|_2^2$ . All experiments are repeated 100 times and the error means are shown in Figure 4 and 5.

The experiment results in Figure 4 (RR) and 5 (UE) show that, as expected, both Norm-Sub and MLE post-processing can improve the accuracy of the results. MLE and Norm-Sub with RR give similar accuracy on both dataset, while MLE with UE can give more accurate results than Norm-Sub with UE.

## 6. Further research

Although we have shown that our bounds are tight for  $a = 2$ ,  $n \gg 0$  and  $\varepsilon \ll 1$ , it would be interesting to know what frequency estimation accuracy is achievable in more general settings. In particular, it is interesting to know what happens for large  $a$ , and for large  $\varepsilon$  (i.e. in the low privacy domain). For large  $\varepsilon$ , the dependence on  $b$  will probably also come into play, which might be able to help point us towards optimal protocols.

Another useful approach would be to give computationally feasible approximations to  $\hat{\Phi}_n$  and  $\hat{\Pi}_n$ . The formulas of Proposition 3.2, which are presented explicitly in Appendix A, can become too complex for large  $a, b, n$ , and while the MLE approach reduces computational complexity, it is still too involved for large  $a$ . This is impor-

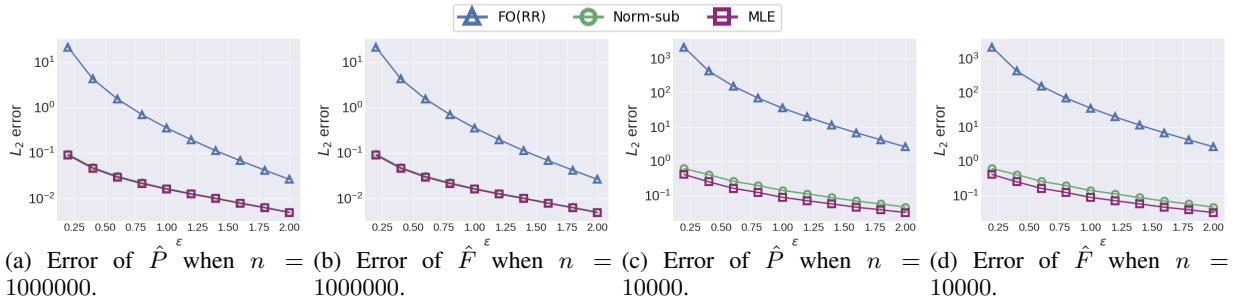


Figure 2: Experiments with RR,  $a = 1024$  and Jeffreys prior  $\mu$ .

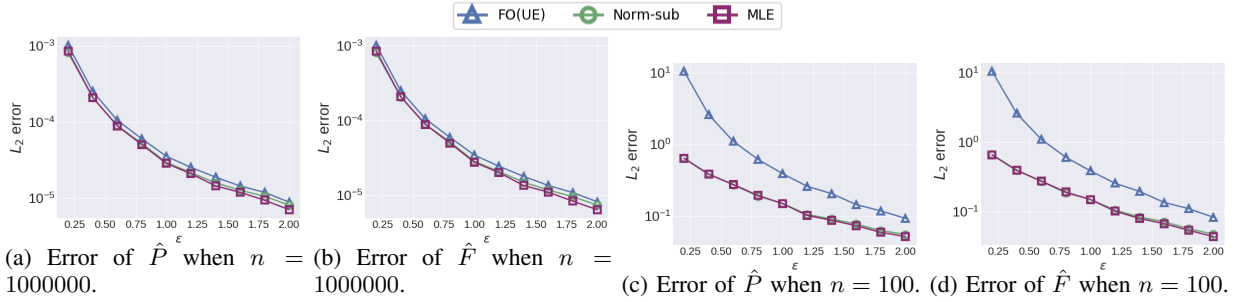


Figure 3: Experiments with UE,  $a = 10$  and Jeffreys prior  $\mu$ .

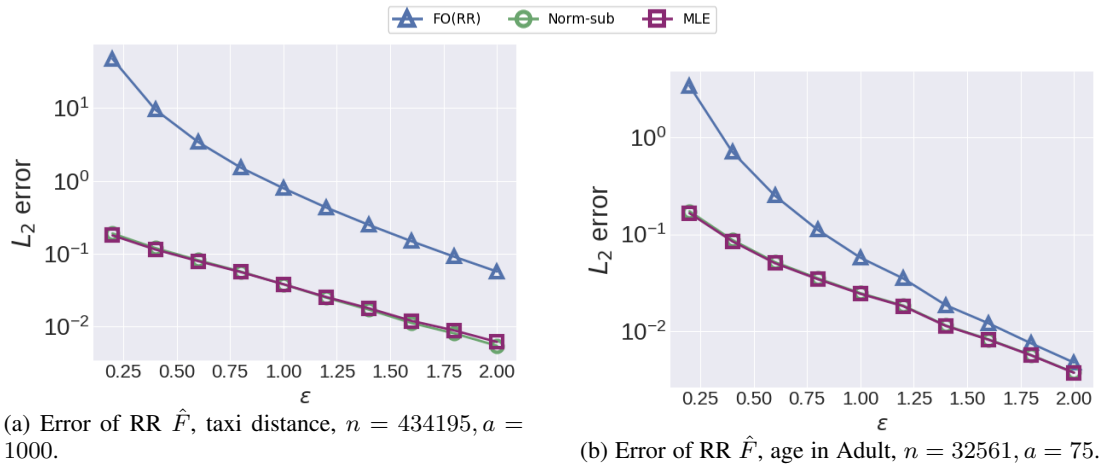


Figure 4: Experiments of RR with real world datasets.

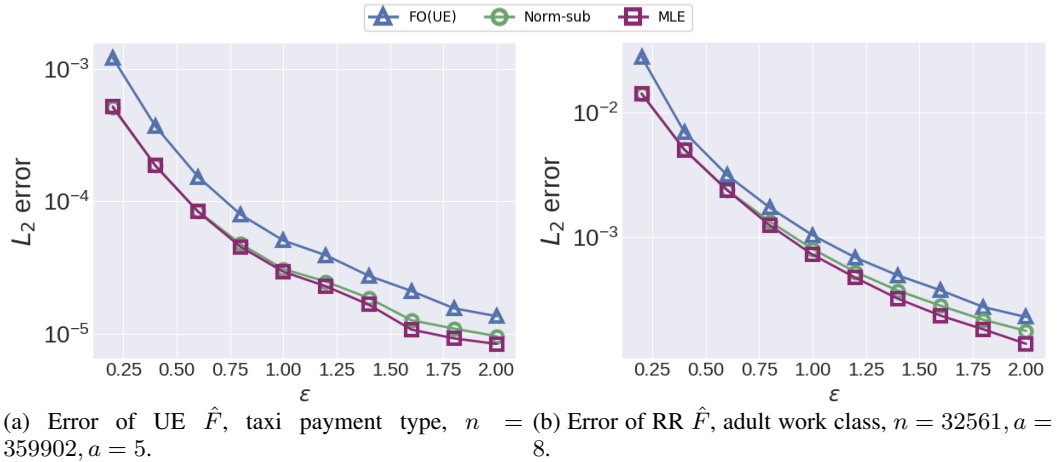


Figure 5: Experiments of UE with real world datasets.

tant, as for small  $a$  the MLE is shown to give accurate frequency estimations; it would be interesting to find a way to extend this approach to larger  $a$ . Furthermore, both our Theorem 3.1, as well as empirical results [16, 25], show that prior knowledge about the distribution can have a significant impact on frequency estimation. The MLE ignores this prior knowledge, hence would be good to find a way to enhance the MLE estimation by taking prior knowledge into account.

By their asymptotic nature, our results mainly concern the case where  $n \gg a$ . However, LDP is also used in situations where  $a \approx n$  [15]. In such situations, one is more interested in identifying the most frequent items rather than estimating the frequency of all items. Different information-theoretical techniques are needed to estimate the utility in this scenario.

## Acknowledgments

This project is supported by NSF grant 1640374, NWO grant 628.001.026, and NSF grant 1931443. We thank the anonymous reviewers for their helpful suggestions.

## References

- [1] Jayadev Acharya, Ziteng Sun, and Huanyu Zhang. “Hadamard response: Estimating distributions privately, efficiently, and with little communication”. In: *The 22nd International Conference on Artificial Intelligence and Statistics*. 2019, pp. 1120–1129.
- [2] Raef Bassily and Adam Smith. “Local, private, efficient protocols for succinct histograms”. In: *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*. 2015, pp. 127–135.
- [3] Michael Betancourt. “Cruising the simplex: Hamiltonian Monte Carlo and the Dirichlet distribution”. In: *AIP Conference Proceedings 31st*. Vol. 1443. 1. American Institute of Physics. 2012, pp. 157–164.
- [4] Jaroslaw Błasiok et al. “Towards instance-optimal private query release”. In: *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM. 2019, pp. 2480–2497.
- [5] Russell CH Cheng. “Random variate generation”. In: *Handbook of Simulation* (1998), pp. 139–172.
- [6] Albert Cheu, Adam Smith, and Jonathan Ullman. “Manipulation attacks in local differential privacy”. In: *arXiv:1909.09630* (2019). Preprint.
- [7] Graham Cormode, Tejas Kulkarni, and Divesh Srivastava. “Answering range queries under local differential privacy”. In: *Proceedings of the VLDB Endowment* 12.10 (2019), pp. 1126–1138.
- [8] Thomas M. Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- [9] Paul Cuff and Lanqing Yu. “Differential privacy as a mutual information constraint”. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. 2016, pp. 43–54.
- [10] Bolin Ding, Janardhan Kulkarni, and Sergey Yekhanin. “Collecting telemetry data privately”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 3571–3580.
- [11] Jiu Ding and Aihui Zhou. “Eigenvalues of rank-one updated matrices with some applications”. In: *Applied Mathematics Letters* 20.12 (2007), pp. 1223–1226.
- [12] Dheeru Dua and Casey Graff. *UCI Machine Learning Repository*. 2017. URL: <http://archive.ics.uci.edu/ml/datasets/Adult>.
- [13] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. “Local privacy and statistical minimax rates”. In: *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*. IEEE. 2013, pp. 429–438.
- [14] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. “Local privacy, data processing inequalities, and statistical minimax rates”. In: *arXiv:1302.3203* (2013). Preprint.
- [15] Úlfar Erlingsson, Vasyli Pihur, and Aleksandra Korolova. “Rappor: Randomized aggregatable privacy-preserving ordinal response”. In: *Proceedings of the 2014 ACM SIGSAC conference on computer and communications security*. 2014, pp. 1054–1067.
- [16] Jinyuan Jia and Neil Zhenqiang Gong. “Calibrate: Frequency estimation and heavy hitter identification with local differential privacy via incorporating prior knowledge”. In: *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*. IEEE. 2019, pp. 2008–2016.
- [17] Peter Kairouz, Sewoong Oh, and Pramod Viswanath. “Extremal mechanisms for local differential privacy”. In: *Advances in neural information processing systems*. 2014, pp. 2879–2887.
- [18] Shiva Prasad Kasiviswanathan et al. “What can we learn privately?” In: *SIAM Journal on Computing* 40.3 (2011), pp. 793–826.
- [19] Milan Lopuhaä-Zwakenberg, Boris Škorić, and Ninghui Li. “Information-theoretic metrics for Local Differential Privacy protocols”. In: *arXiv:1910.07826* (2019). Preprint.
- [20] Helmut Strasser. “The asymptotic equivalence of Bayes and maximum likelihood estimation”. In: *Journal of Multivariate Analysis* 5.2 (1975), pp. 206–226.
- [21] NYC Taxi and Limousine Commission. *TLC Trip Record Data*. 2019. URL: <https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page>.
- [22] Apple Differential Privacy Team. *Learning with Privacy at Scale*. 2017.
- [23] Jerry Alan Veeh. “The multivariate Laplace-De Moivre theorem”. In: *Journal of multivariate analysis* 18.1 (1986), pp. 46–51.
- [24] Ning Wang et al. “PrivTrie: Effective frequent term discovery under local differential privacy”. In: *2018 IEEE 34th International Conference on Data Engineering (ICDE)*. IEEE. 2018, pp. 821–832.
- [25] Tianhao Wang et al. “Locally differentially private frequency estimation with consistency”. In: *arXiv:1905.08320* (2019). Preprint.
- [26] Tianhao Wang et al. “Locally differentially private protocols for frequency estimation”. In: *26th {USENIX} Security Symposium ({USENIX} Security 17)*. 2017, pp. 729–745.

[27] Stanley L. Warner. "Randomized response: A survey technique for eliminating evasive answer bias". In: *Journal of the American Statistical Association* 60.309 (1965), pp. 63–69.

## Appendix

Before we can prove Proposition 3.2, we recall two well-known facts about Dirichlet distributions and distribution mixtures, and prove an auxiliary Proposition.

**Fact A.1.** Let  $P \in \mathcal{P}_{\mathcal{A}}$  be drawn from a Dirichlet distribution with parameter vector  $\gamma \in \mathbb{R}_{\geq 0}^a$ . Then

$$\mathbb{E}[P] = \frac{1}{\sum_{\alpha \in \mathcal{A}} \gamma_{\alpha}} \gamma, \quad (78)$$

$$\text{Var}(P_{\alpha}) = \frac{\mathbb{E}[P_{\alpha}](1 - \mathbb{E}[P_{\alpha}])}{1 + \sum_{\alpha' \in \mathcal{A}} \gamma_{\alpha'}}. \quad (79)$$

**Fact A.2.** Let  $P_1, \dots, P_r$  be continuous random variables in  $\mathbb{R}$ , and let  $w \in \mathbb{R}_{\geq 0}^r$  be such that  $\sum_{j=1}^r w_j = 1$ . Let  $M$  be the mixture of the  $P_j$  with weight vector  $w$ , i.e.  $M = P_j$  with probability  $w_j$ . Then

$$\mathbb{E}[M] = \sum_{j=1}^r w_j \mathbb{E}[P_j], \quad (80)$$

$$\text{Var}(M) = \sum_{j=1}^r w_j (\text{Var}(P_j) + \mathbb{E}[P_j]^2 - \mathbb{E}[M]^2). \quad (81)$$

The following proposition gives explicit, if lengthy, formulas to calculate  $\Pi_{\text{opt}}$  as well as  $\text{MSE}_{\mu}^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}})$ .

**Proposition A.3.** Let  $\Pi_{\text{opt}}$  be as in Theorem 3.1. Let  $\mu$  be a Dirichlet distribution with parameter vector  $\gamma \in \mathbb{R}_{> 0}^a$ . Let  $B$  be the multivariate beta function. For  $\alpha \in \mathcal{A}$ ,  $\vec{x} \in \mathcal{A}^n$  and  $\vec{y} \in \mathcal{B}^n$ , we define

$$C_{\vec{y}} = \sum_{\vec{x} \in \mathcal{A}^n} \frac{B(\gamma + t(\vec{x}))}{B(\gamma)} \prod_i Q_{y_i | x_i}, \quad (82)$$

$$w_{\vec{x} | \vec{y}} = C_{\vec{y}}^{-1} \frac{B(\gamma + t(\vec{x}))}{B(\gamma)} \prod_i Q_{y_i | x_i}, \quad (83)$$

$$m_{\vec{x}, \alpha} = \frac{\gamma_{\alpha} + t_{\alpha}(\vec{x})}{\sum_{\alpha' \in \mathcal{A}} \gamma_{\alpha'} + n}, \quad (84)$$

$$m_{\vec{y}, \alpha} = \sum_{\vec{x} \in \mathcal{A}^n} w_{\vec{x} | \vec{y}} m_{\vec{x}, \alpha}, \quad (85)$$

$$\sigma_{\vec{x}, \alpha}^2 = \frac{m_{\vec{x}, \alpha}(1 - m_{\vec{x}, \alpha})}{\sum_{\alpha' \in \mathcal{A}} \gamma_{\alpha'} + n + 1}. \quad (86)$$

Then the  $\alpha$ -coefficient of  $\Pi_{\text{opt}}(\vec{y})$  is equal to  $m_{\vec{y}, \alpha}$ , and  $\text{MSE}_{\mu}^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}})$  is equal to

$$\sum_{\vec{y} \in \mathcal{B}^n} C_{\vec{y}} \sum_{\alpha \in \mathcal{A}} \left( \sum_{\vec{x} \in \mathcal{A}^n} w_{\vec{x} | \vec{y}} (m_{\vec{x}, \alpha}^2 + \sigma_{\vec{x}, \alpha}^2) - m_{\vec{y}, \alpha}^2 \right). \quad (87)$$

*Proof.* For  $\gamma' \in \mathbb{R}_{\geq 0}^a$ , let  $\Delta_{\gamma'}(p) = \frac{1}{B(\gamma')} \prod_{\alpha} p_{\alpha}^{\gamma'_{\alpha}}$  be the probability density function of the Dirichlet distribution with parameter vector  $\alpha'$ . Let  $\delta_{P | \vec{Y} = \vec{y}}$  be the posterior probability density function of  $P$  given  $\vec{Y} = \vec{y}$ . By [19,

Thm. 9.1] we have (note the definition of  $C_{\vec{y}}$  there differs by a factor  $B(\gamma)$  from the one given here):

$$\delta_{P | \vec{Y} = \vec{y}} = \frac{1}{C_{\vec{y}}} \sum_{\vec{x} \in \mathcal{A}^n} w_{\vec{x} | \vec{y}} \Delta_{\gamma + t(\vec{x})}. \quad (88)$$

In other words,  $P | \vec{Y} = \vec{y}$  is a mixture of Dirichlet distributions  $D_{\vec{x}}$ . These are parametrised by  $\vec{x} \in \mathcal{A}^n$ , have weight  $w_{\vec{x} | \vec{y}}$ , and parameter vector  $\gamma + t(\vec{x})$ ; by Fact A.1 one has  $\mathbb{E}[D_{\vec{x}, \alpha}] = m_{\vec{x}, \alpha}$  and  $\text{Var}(D_{\vec{x}, \alpha}) = \sigma_{\vec{x}, \alpha}^2$ . Fact A.2 now shows us that  $\Pi_{\text{opt}}(\vec{y}) = \mathbb{E}[P_{\alpha} | \vec{Y} = \vec{y}] = m_{\vec{y}, \alpha}$ , proving the first claim of the proposition, and

$$\text{Var}(P_{\alpha} | \vec{Y} = \vec{y}) = \sum_{\vec{x} \in \mathcal{A}^n} w_{\vec{x} | \vec{y}} (m_{\vec{x}, \alpha}^2 + \sigma_{\vec{x}, \alpha}^2) - m_{\vec{y}, \alpha}^2. \quad (89)$$

Furthermore, note that

$$\mathbb{P}(\vec{X} = \vec{x}) = \int \Delta_{\gamma}(p) \mathbb{P}(\vec{X} = \vec{x} | P = p) dp \quad (90)$$

$$= \frac{1}{B(\gamma)} \int \prod_{\alpha} p^{\gamma_{\alpha} + t_{\alpha}(\vec{x})} dp \quad (91)$$

$$= \frac{B(\gamma + t(\vec{x}))}{B(\gamma)} \int \Delta_{\gamma + t(\vec{x})} dp \quad (92)$$

$$= \frac{B(\gamma + t(\vec{x}))}{B(\gamma)}, \quad (93)$$

$$\mathbb{P}(\vec{Y} = \vec{y}) = \sum_{\vec{x}} \mathbb{P}(\vec{Y} = \vec{y} | \vec{X} = \vec{x}) \mathbb{P}(\vec{X} = \vec{x}) \quad (94)$$

$$= \sum_{\vec{x}} \frac{B(\gamma + t(\vec{x}))}{B(\gamma)} \prod_{i=1}^n Q_{y_i | x_i} \quad (95)$$

$$= C_{\vec{y}}. \quad (96)$$

Combining this with (20) and (89), this shows that

$$\begin{aligned} \text{MSE}_{\mu}^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}}) &= \sum_{\vec{y} \in \mathcal{B}^n} \mathbb{P}(\vec{Y} = \vec{y}) \sum_{\alpha \in \mathcal{A}} \text{Var}(P_{\alpha} | \vec{Y} = \vec{y}) \\ &= \sum_{\vec{y} \in \mathcal{B}^n} C_{\vec{y}} \sum_{\alpha \in \mathcal{A}} \left( \sum_{\vec{x} \in \mathcal{A}^n} w_{\vec{x} | \vec{y}} (m_{\vec{x}, \alpha}^2 + \sigma_{\vec{x}, \alpha}^2) - m_{\vec{y}, \alpha}^2 \right). \quad \square \end{aligned} \quad (97)$$

Using these formulas allows us to prove Proposition 3.2.

*Proof of Proposition 3.2.* Let  $T_{\alpha}$ ,  $S_{\beta}$  and  $S_{\beta | \alpha}$  be as in (3,4,5). Note that  $t_{\alpha} = \sum_{\beta} s_{\beta | \bullet}$ , where  $s_{\beta | \bullet} = (s_{\beta | 1}, \dots, s_{\beta | a})$ . Furthermore, for a given  $s$ , define  $\mathcal{S}_s = \{(s_{\beta | \alpha})_{\beta, \alpha} : \forall \beta \sum_{\alpha} s_{\beta | \alpha} = s_{\beta}\}$ . Then

$$\begin{aligned} \Pi_{\text{opt}}(\vec{y})_{\alpha} &= \sum_{\vec{x} \in \mathcal{A}^n} w_{\vec{x} | \vec{y}} m_{\vec{x}, \alpha} \\ &= C_{\vec{y}}^{-1} \sum_{\vec{x} \in \mathcal{A}^n} \frac{\gamma_{\alpha} + t_{\alpha}(\vec{x})}{\sum_{\alpha' \in \mathcal{A}} \gamma_{\alpha'} + n} \cdot \frac{B(\gamma + t(\vec{x}))}{B(\gamma)} \prod_i Q_{y_i | x_i} \end{aligned} \quad (98)$$

$$= C_{\vec{y}}^{-1} \sum_{\vec{x} \in \mathcal{A}^n} \frac{\gamma_{\alpha} + t_{\alpha}(\vec{x})}{\sum_{\alpha' \in \mathcal{A}} \gamma_{\alpha'} + n} \cdot \frac{B(\gamma + t(\vec{x}))}{B(\gamma)} \prod_i Q_{y_i | x_i} \quad (99)$$

$$= C_{\vec{y}}^{-1} \sum_{s_{\bullet} \in \mathcal{S}_s} \left( \prod_{\beta} \binom{s_{\beta}}{s_{\beta | \bullet}} \right) \frac{(\gamma_{\alpha} + t_{\alpha}(\vec{x})) B(\gamma + t(\vec{x}))}{(\sum_{\alpha' \in \mathcal{A}} \gamma_{\alpha'} + n) B(\gamma)} \prod_{\beta, \alpha} Q_{\beta | \alpha}^{s_{\beta | \alpha}}. \quad (100)$$

Here  $\binom{s_{\beta}}{s_{\beta | \bullet}} = \frac{s_{\beta}!}{\prod_{\alpha} s_{\beta | \alpha}!}$  is the multinomial coefficient. Since  $\#\mathcal{S}_s = \mathcal{O}(n^{(a-1)b})$ , we can find  $\Pi_{\text{opt}}(\vec{y})_{\alpha}$ , up

to the scaling factor  $C_{\bar{y}}^{-1}$ , can be found by calculating  $\mathcal{O}(n^{(a-1)b})$  summands. We can then find  $C_{\bar{y}}$  by using the fact that  $\sum_{\alpha} \Pi_{\text{opt}}(\bar{y})_{\alpha} = 1$ .

Analogous to the above, we can similarly show that we need  $\mathcal{O}(n^{(a-1)b})$  to calculate each summand of the form

$$\sum_{\bar{x} \in \mathcal{A}^n} w_{\bar{x}|\bar{y}}(m_{\bar{x},\alpha}^2 + \sigma_{\bar{x},\alpha}^2) - m_{\bar{y},\alpha}^2 \quad (101)$$

in (87). We then need to sum over all possible  $s$ , of which there are  $\mathcal{O}(n^{b-1})$ , leading to a total complexity of  $\mathcal{O}(n^{(a-1)b} \cdot n^{b-1}) = \mathcal{O}(n^{ab-1})$  to calculate  $\text{MSE}_{\mu}^{\text{distr}}(\mathcal{Q}, \Pi_{\text{opt}})$ .  $\square$

We describe an efficient algorithm solving (22) for RR with  $\mathcal{O}(a \log a)$  time complexity in Algorithm 1. This algorithm is justified as follows. Recall from [27] that RR (with privacy parameter  $\varepsilon$ ) is given by the  $a \times a$ -matrix  $Q$  satisfying

$$Q_{y|x} = \frac{1 + (e^{\varepsilon} - 1)\delta_{x=y}}{e^{\varepsilon} + a - 1}. \quad (102)$$

Since  $\mathcal{B} = \mathcal{A}$  for RR, we can rewrite (22) to

$$\begin{aligned} \text{maximise}_p \quad & f(p) = \sum_{y \in \mathcal{A}} s_y \log((Qp)_y) \\ \text{subject to} \quad & p \geq 0, \sum_{x \in \mathcal{A}} p_x = 1. \end{aligned} \quad (103)$$

From this, we obtain the Karush-Kuhn-Tucker (KKT) conditions (with extra variables  $(u_x)_{x \in \mathcal{A}}$  and  $v$ ):

$$\begin{aligned} \forall x \in \mathcal{A}, \quad & \frac{\partial f(p)}{\partial p_x} + u_x - v = 0, \quad (\text{stationarity}) \\ & u_x p_x = 0, \quad (\text{complementary slackness}) \\ & u_x \geq 0, \quad (\text{dual feasibility}) \\ & p_x \geq 0, \sum_{x \in \mathcal{A}} p_x = 1 \quad (\text{primal feasibility}), \end{aligned}$$

in which

$$\frac{\partial f(p)}{\partial p_x} = \sum_{y \in \mathcal{A}} \frac{s_y Q_{y|x}}{\sum_{k \in \mathcal{A}} Q_{y|k} p_k} \quad (104)$$

$$= \frac{(e^{\varepsilon} - 1)s_x}{(e^{\varepsilon} - 1)p_x + 1} + \sum_{y \in \mathcal{A}} \frac{s_y}{(e^{\varepsilon} - 1)p_y + 1}. \quad (105)$$

By stationarity and complementary slackness we find for all  $x \in \mathcal{A}$  that

$$p_x \left( v - \frac{\partial f(p)}{\partial p_x} \right) = 0. \quad (106)$$

By summing all  $x \in \mathcal{A}$ , we find  $v = n$ . Suppose we have found the optimal  $\hat{p}$ , and define  $\mathcal{A}' = \{x : \hat{p}_x > 0\}$ . For  $x, x' \in \mathcal{A}'$ , it follows from (106) that we have

$$\frac{s_x}{(e^{\varepsilon} - 1)\hat{p}_x + 1} = \frac{s_{x'}}{(e^{\varepsilon} - 1)\hat{p}_{x'} + 1}, \quad (107)$$

hence

$$\frac{s_{x'}}{s_x} ((e^{\varepsilon} - 1)\hat{p}_x + 1) = (e^{\varepsilon} - 1)\hat{p}_{x'} + 1. \quad (108)$$

Summing over all  $x' \in \mathcal{A}'$ , we can solve  $\hat{p}_i$  as

$$\hat{p}_x = \frac{\#\mathcal{A}' - \frac{\sum_{x' \in \mathcal{A}'} s_{x'}}{s_x} + (e^{\varepsilon} - 1)}{(e^{\varepsilon} - 1) \left( \frac{\sum_{x' \in \mathcal{A}'} s_{x'}}{s_x} \right)}. \quad (109)$$

Therefore, the problem of finding  $\hat{p}$  reduces to finding  $\mathcal{A}'$ . We determine  $\mathcal{A}'$  by starting with  $\mathcal{A}' = \mathcal{A}$ , and then repeatedly removing the  $x$  with the lowest value of  $s_x$ , until the estimation (109) is no longer nonnegative for all  $x$ .

---

**Algorithm 1:** Exact MLE post-processing for RR

---

**Input :** Total number of users  $n$ ; privacy budget  $\varepsilon$ ; tallies of obfuscated data

$$s = [s_0, \dots, s_{a-1}]$$

**Output:** MLE estimate distribution

$$\hat{p} = [\hat{p}_1, \dots, \hat{p}_{a-1}]$$

$[s'_{(0)}, \dots, s'_{(a-1)}] \leftarrow \text{Sort}(s)$ , such that

$\forall i < j, s'_{(i)} < s'_{(j)}$  ;

$k = 0$  ;

**while**  $a - k + e^{\varepsilon} - 1 - \frac{\sum_{i=k}^{a-1} s'_{(i)}}{s'_{(k)}} < 0$  **do**

    |  $k \leftarrow k + 1$  ;

**end**

**for**  $j = 0, \dots, a - 1$  **do**

    |  $\Phi_j = \frac{s_j}{\sum_{i=k}^{a-1} s'_{(i)}} ;$

    |  $\hat{p}_j = \max\{0, \Phi_j \left( \frac{n-k}{e^{\varepsilon}-1} + 1 \right) - \frac{1}{e^{\varepsilon}-1}\}$

**end**

---