

The Privacy-Utility Tradeoff of Robust Local Differential Privacy

Milan Lopuhaä-Zwakenberg and Jasper Goseling

Abstract

We consider data release protocols for data $X = (S, U)$, where S is sensitive; the released data Y contains as much information about X as possible, measured as $I(X; Y)$, without leaking too much about S . We introduce the Robust Local Differential Privacy (RLDP) framework to measure privacy. This framework relies on the underlying distribution of the data, which needs to be estimated from available data. Robust privacy guarantees are ensuring privacy for all distributions in a given set \mathcal{F} , for which we study two cases: when \mathcal{F} is the set of all distributions, and when \mathcal{F} is a confidence set arising from a χ^2 test on a publicly available dataset. In the former case we introduce a new release protocol which we prove to be optimal in the low privacy regime. In the latter case we present four algorithms that construct RLDP protocols from a given dataset. One of these approximates \mathcal{F} by a polytope and uses results from robust optimisation to yield high utility release protocols. However, this algorithm relies on vertex enumeration and becomes computationally inaccessible for large input spaces. The other three algorithms are low-complexity and build on randomised response. Experiments verify that all four algorithms offer significantly improved utility over regular LDP.

I. INTRODUCTION

We consider the setting in which users have data $X = (S, U)$ that a data aggregator is interested in, but users do not wish to disclose information about sensitive data S . Therefore, users release an obfuscated version Y of X , such that Y contains as much information about X as possible, measured as $I(X; Y)$, without leaking too much about S . This scenario and closely related ones have been studied in, for instance, [1]–[8].

This paper introduces a form of local differential privacy (LDP) [9] to measure the amount of information that Y leaks on S . The following version of ε -LDP was introduced in [10]:

$$\mathbb{P}(Y = y|S = s) \leq e^\varepsilon \mathbb{P}(Y = y|S = s'), \quad (1)$$

for all y, s and s' . Note, that this condition is less strict than $\mathbb{P}(Y = y|X = x) \leq e^\varepsilon \mathbb{P}(Y = y|X = x')$ as would be used in ordinary LDP. Also note that (1) relies on the distribution $P_X = P_{S,U}$. From these observations it follows that this privacy definition enables higher utility of the released data Y at the expense of not being completely ‘distribution free’ as would be the case for ordinary LDP.

In [10] condition (1) is studied for the case of known P_X . This is a strong assumption, since users will need to estimate P_X . When an attacker has better knowledge of P_X than the user, it follows from the odds-ratio interpretation of differential privacy [11] that sufficient privacy is not guaranteed in such a scenario.

In this paper we, therefore, provide stronger privacy guarantees. In particular, we introduce robustness constraints, which say that privacy should not just hold for one P_X , but for a set \mathcal{F} of these. As a result we guarantee privacy against attackers with (at least) reasonable estimates of P_X , without sacrificing utility to protect against attackers with no or unreliable information on P_X . We refer to the resulting privacy framework as Robust Local Differential Privacy (RLDP).

We consider two cases for \mathcal{F} . In the first case, we let \mathcal{F} be the set of *all* probability distributions \mathcal{F} . We show that in this case, privacy w.r.t. S is very similar, but not equivalent, to privacy w.r.t. X . We introduce a new privacy protocol that exploits the small difference that remains between these two definitions and show that this protocol is optimal in the low privacy regime.

In the second case, we assume that there is publicly available data from n users, which allows the aggregator and the users to estimate \hat{P}_X . The set \mathcal{F} consists of those P that are close enough to \hat{P} so that the difference is not statistically significant for a chosen significance level α ; this choice of \mathcal{F} is common in statistical optimisation. Here, we introduce three protocols and study their privacy and utility.

A. Contributions

In addition to introducing the RLDP privacy framework, the main contributions of this paper are as follows.

We consider the setting where $\mathcal{F} = \mathcal{P}_X$. In this setting:

- We introduce a protocol SRR based on the classic Randomized Response protocol [12]. We show that SRR maximises mutual information in the low privacy regime.

We consider the setting where \mathcal{F} is a χ^2 confidence set around \hat{P}_X . In this setting:

- We approximate \mathcal{F} by an enveloping polytope. We then use techniques from robust optimisation [13]–[15] to characterize the protocol that is optimal over this polytope. The resulting lower bound on utility demonstrates the advantage of RLDP over ordinary LDP. A drawback of this method is that it relies on vertex enumeration and is, therefore, computationally unfeasible for large alphabets.
- Therefore, we introduce two low-complexity data release mechanisms: i) Independent Reporting (IR), in which S and U are reported through separate LDP protocols, and ii) Conditional Reporting (CR), in which first S is obfuscated, and either a slightly obfuscated U or a randomly drawn U' is returned, depending on whether the obfuscated S is ‘correct’.

- For both mechanisms we characterize the conditions that underlying LDP protocols have to satisfy in order to ensure RLDP. Furthermore, while both mechanisms can incorporate any LDP protocol, we show that it is optimal to use Randomised Response [12]. This drastically reduces the search space and allows us to find the optimal SR and CR mechanisms using one-dimensional optimisation.

We demonstrate the improved utility of RLDP over LDP with numerical experiments. In particular we provide results for both synthetic datasets as well as real-world census data.

B. Related work

Disclosing information in a privacy-preserving way is one of the main challenges in official statistics [16], [17]. The setting considered in the current paper is closely connected to disclosing a table with micro-data where each record in the table is released independently of the other records. This approach to disclosing micro-data was studied in [1] by considering expected error as the utility measure and mutual information as the privacy measure. The resulting optimization problem corresponds to the traditional rate-distortion problem.

The version of the problem in which both utility and privacy are measured using mutual information is known as the privacy funnel and was studied first in [3]. The dual problem of the privacy funnel, in which utility is maximized w.r.t. a privacy constraint was studied in [5]. The privacy funnel and its dual are intimately related to the information bottleneck problem [18], which seeks to optimise compression while retaining relevant information. Multiple approaches to optimising privacy funnel also work for the information bottleneck and vice versa [6], [7].

In [4] a version of this problem is studied in which privacy leakage is measured through the improved statistical inference by an attacker after seeing the disclosed information. This measure is formulated through a general cost function, with mutual information resulting as a special case. Perfect privacy, which demands the output to be independent of the sensitive data, is studied in [19], and methods are given to find optimal protocols in this setting. In [20] the maximal leakage measure with a clear operational interpretation is defined. In [21] this measure is generalized to a parametrized measure, enabling to interpolate between maximal leakage and mutual information. A multitude of other privacy frameworks and leakage measures exist. We refer to [22] for an overview and restrict the remainder of this section to local differential privacy and robustness, which are most closely related to our work.

In this paper we consider measures based on Local Differential Privacy (LDP) [9], [23]. In this setting, several privacy protocols exist, including Randomised Response [12] and Unary Encoding [24]. Optimal LDP protocols under a variety of utility metrics, including mutual information, are found in [2]. A variation of LDP is proposed in [10] for the case of disclosing $X = (S, U)$, where

only S is sensitive. The privacy metrics given there fit into a general framework called pufferfish privacy [11]. In [8] a general class of privacy metrics called *average information leakage* is introduced in this setting, and it is shown that LDP implies privacy under these metrics.

In all the above work the privacy protocol is derived from the (estimated) distribution $P_{S,X}$. In most cases an analysis of robustness/sensitivity with respect to this estimate is not present. An exception is [4] in which one of the contributions is to quantify the impact of mismatched priors, i.e. the impact of not knowing $P_{S,X}$ exactly. A bound on the resulting level of privacy is derived in terms of the total variational distance between the actual and the estimated $P_{S,X}$. The behaviour of privacy and utility metrics under robustness are studied in [25], [26]. For a wide variety of privacy and utility metrics, they give bounds on the utility loss that occurs when robustness is added to the requirements. In both cases, robustness is defined by looking at an ℓ_1 -ball around the observed empirical distribution. One can also define robustness in other ways, such as by KL-divergence [27], χ^2 -divergence [15], or a general f -divergence [28].

Another line of work builds on recent advances in generative adversarial networks [29]. In [30], [31] the generative adversarial framework is used to provide release protocols that do not use explicit expressions for P_X . Even though it is not explicitly addressed in [30], [31], it is expected that the generalization properties of networks will provide a form of robustness. Closely related approaches are used in the area of face recognition, [32], [33] with the aim of preventing biometric profiling [34]. In [32], [33], however, the leakage measures that are used do not seem to have an operational interpretation.

C. Overview of paper

The structure of this paper is as follows. In Section II we describe the model in detail. In Section III we consider the case that $\mathcal{F} = \mathcal{P}_X$. In Section IV we study the case that \mathcal{F} is a confidence set, and we prove several properties of \mathcal{F} that will be useful in the following sections. In Section V we introduce PolyOpt, an algorithm that finds high utility protocols through approximating \mathcal{F} by an enveloping polytope. In Section VI we discuss Independent Reporting, its privacy and utility, and we show how the optimal IR-protocol can be found using low-dimensional optimisation. In Section VII we do the same for Conditional Reporting. In Section VIII we evaluate the discussed methods experimentally. Finally, in Section IX we provide a discussion of our results and provide an outlook on future work.

II. MODEL AND PRELIMINARIES

The dataspace is $\mathcal{X} = \mathcal{S} \times \mathcal{U}$, where \mathcal{S} and \mathcal{U} are finite sets. We write $|\mathcal{S}| =: a_1$, $|\mathcal{U}| =: a_2$, and $|\mathcal{X}| = a_1 a_2 =: a$. New data items $X = (S, U)$ are drawn from a probability distribution P^*

in $\mathcal{P}_{\mathcal{X}}$, the space of probability distributions on \mathcal{X} . The user's aim is to create a release protocol \mathcal{Q} such that $Y = \mathcal{Q}(X)$ contains as much information about X as possible, while not leaking too much information about S . Protocol \mathcal{Q} is a probabilistic map, that we represent by a left stochastic matrix $(Q_{y|x})_{y \in \mathcal{Y}, x \in \mathcal{X}}$, and we write $|\mathcal{Y}| = b$. Often, we identify $\mathcal{Y} = \{1, \dots, b\}$, and likewise for other sets.

The distribution P^* is not known exactly. Instead it is known only that $P^* \in \mathcal{F}$ for some set of possible distributions $\mathcal{F} \subset \mathcal{P}_{\mathcal{X}}$, where $\mathcal{P}_{\mathcal{X}}$ denotes the probability simplex over \mathcal{X} . We give various examples of such \mathcal{F} below. The uncertainty set \mathcal{F} captures our uncertainty about P^* . The idea is that we guarantee privacy for all $P \in \mathcal{F}$. We will denote this as robust local differential privacy (RLDP).

Definition 1. Let $\varepsilon \geq 0$ and $\mathcal{F} \subset \mathcal{P}_{\mathcal{X}}$. We say that \mathcal{Q} satisfies $(\varepsilon, \mathcal{F})$ -RLDP if for all $s, s' \in \mathcal{S}$, all $y \in \mathcal{Y}$, and all $P \in \mathcal{F}$ we have

$$\mathbb{P}_{X \sim P}(Y = y | S = s) \leq e^\varepsilon \mathbb{P}_{X \sim P}(Y = y | S = s'). \quad (2)$$

Note that we use the notation $\mathbb{P}_{X \sim P}(\bullet)$ to emphasize that X is distributed according to P . If no confusion can arise, we will often leave out the subscript $X \sim P$ to improve readability.

We consider various forms of uncertainty on P^* , as captured by \mathcal{F} :

- 1) Nothing is known about P^* . In this case $\mathcal{F} = \mathcal{P}_{\mathcal{X}}$. Regarding privacy, this is the 'safest' choice.
- 2) We suppose there is a database $\vec{x} = (x_1, \dots, x_n)$ accessible to the user, where each $x_i = (s_i, u_i)$ is drawn independently from P^* . Based on this, the user produces an estimate \hat{P} of P . Fix a significance level α : we let \mathcal{F} be the $(1 - \alpha)$ -confidence interval for P in a χ^2 -test, i.e.

$$\mathcal{F} = \left\{ P : \sum_x \frac{(\hat{P}_x - P_x)^2}{P_x} \leq B := \frac{F_{\#\mathcal{X}-1}^{-1}(1 - \alpha)}{n} \right\}, \quad (3)$$

where F_d is the cdf of the χ^2 -distribution with d degrees of freedom. At times, it will be convenient to express this as

$$\mathcal{F} = \left\{ P : \sum_x \frac{\hat{P}_x^2}{P_x} \leq B + 1 \right\}. \quad (4)$$

This situation is expressed in Figure 1.

Another option would be to have \mathcal{F} be a singleton, i.e. to assume that P is known. This setting is studied in [10].

For completeness we give the definition of LDP.

Definition 2. Let $\varepsilon \geq 0$. We say that $\mathcal{Q}: \mathcal{X} \rightarrow \mathcal{Y}$ satisfies ε -LDP if for all $x, x' \in \mathcal{X}$ and all $y \in \mathcal{Y}$ we have

$$\mathbb{P}(Y = y | X = x) \leq e^\varepsilon \mathbb{P}(Y = y | X = x'). \quad (5)$$

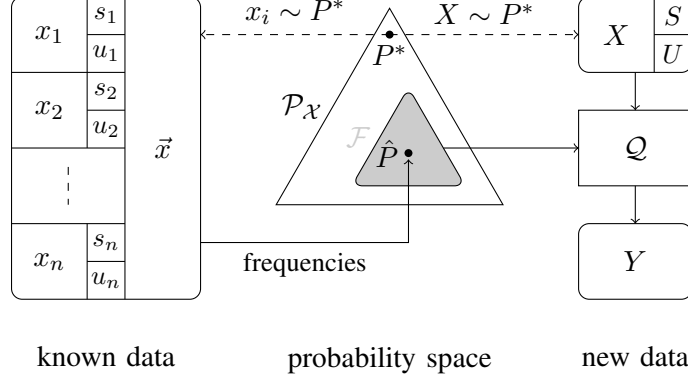


Fig. 1: An overview of the setting of this paper when \mathcal{F} is a confidence set based on a dataset \vec{x} .

In Sections VI and VII, we build RLDP protocols from regular LDP protocols. To establish the privacy guarantees of these protocols, we will need the following lemma that relates LDP to the ℓ_1 -distance of probability distributions.

Lemma 3. *Let $\mathcal{Q}: \mathcal{X} \rightarrow \mathcal{Y}$ be an ε -LDP protocol. Then for all $y \in \mathcal{Y}$ and all $P, P' \in \mathcal{P}_{\mathcal{X}}$ we have*

$$\frac{\mathbb{P}_{X \sim P}(\mathcal{Q}(X) = y)}{\mathbb{P}_{X \sim P'}(\mathcal{Q}(X) = y)} \leq 1 + \frac{e^\varepsilon - 1}{2} \|P - P'\|_1. \quad (6)$$

Proof. Let $Q_y^{\max} = \max_x Q_{y|x}$ and $Q_y^{\min} = \min_x Q_{y|x}$; note that $Q_y^{\max} \leq e^\varepsilon Q_y^{\min}$. Furthermore, $\mathbb{P}_{X \sim P}(\mathcal{Q}(X) = y) = \sum_{x \in \mathcal{X}} Q_{y|x} P_x$ and $\mathbb{P}_{X \sim P'}(\mathcal{Q}(X) = y) = \sum_{x \in \mathcal{X}} Q_{y|x} P'_x$, hence

$$\mathbb{P}_{X \sim P}(\mathcal{Q}(X) = y) - \mathbb{P}_{X \sim P'}(\mathcal{Q}(X) = y) \quad (7)$$

$$= \sum_{x: P_x \geq P'_x} Q_{y|x} (P_x - P'_x) - \sum_{x: P'_x > P_x} Q_{y|x} (P'_x - P_x) \quad (8)$$

$$\leq \frac{Q_y^{\max}}{2} \|P - P'\|_1 - \frac{Q_y^{\min}}{2} \|P - P'\|_1 \quad (9)$$

$$\leq \frac{(e^\varepsilon - 1) Q_y^{\min}}{2} \|P - P'\|_1 \quad (10)$$

$$\leq \frac{(e^\varepsilon - 1) \mathbb{P}_{X \sim P'}(\mathcal{Q}(X) = y)}{2} \|P - P'\|_1, \quad (11)$$

from which the lemma directly follows. \square

Next to a privacy leakage measure we need to define a utility measure. Throughout this paper, we follow the original Privacy Funnel [3] and its LDP counterpart [10] in taking mutual information $I(X; Y)$ as a utility measure. As is argued in [3], mutual information arises naturally when minimising log loss distortion in the Privacy Funnel scenario.

The value of $I(X; Y)$ depends on \mathcal{Q} and on the probability distribution on \mathcal{X} . As this is unknown, we consider two possibilities:

- 1) One can take $I_{X \sim \hat{P}}(X; Y)$, abbreviated to $I_{\hat{P}}(X; Y)$;

2) One can consider $\min_{P \in \mathcal{F}} I_P(X; Y)$.

Throughout this paper, all results will be proven for general P . Furthermore, it will turn out that many protocols we find will not depend on P . In the experiments of Section VIII, we focus on $I_{\hat{P}}(X; Y)$, although we also investigate the effect of P on utility by comparing $I_{P^*}(X; Y)$ to $I_{\hat{P}}(X; Y)$ in Section VIII-F.

III. MAXIMAL \mathcal{F}

In this section, we consider the case where \mathcal{F} is maximal, i.e. $\mathcal{F} = \mathcal{P}_{\mathcal{X}}$. We show that in this situation, RLDP is almost equivalent to LDP. However, it is not completely equivalent, and we use this to describe a version of Generalised Randomised Response (GRR) that exploits the difference between RLDP and LDP. We show that this new protocol is optimal in the low privacy regime (i.e. $\varepsilon \gg 0$), similar to how GRR is the optimal LDP-protocol in the low privacy regime [2]. The following Proposition gives a characterisation of $(\varepsilon, \mathcal{P}_{\mathcal{X}})$ -RLDP.

Proposition 4. *\mathcal{Q} satisfies $(\varepsilon, \mathcal{P}_{\mathcal{X}})$ -RLDP if and only if for all $y \in \mathcal{Y}$ and $(s, u), (s', u') \in \mathcal{X}$ with $s \neq s'$ one has*

$$\frac{Q_{y|s,u}}{Q_{y|s',u'}} \leq e^\varepsilon. \quad (12)$$

Proof. Suppose that \mathcal{Q} satisfies $(\varepsilon, \mathcal{F})$ -RLDP w.r.t. $\mathcal{P}_{\mathcal{X}}$. Let $(s, u), (s', u') \in \mathcal{X}$ with $s \neq s'$. Let P be given by

$$P_x = \begin{cases} \frac{1}{2}, & \text{if } x \in \{(s, u), (s', u')\}, \\ 0, & \text{otherwise.} \end{cases} \quad (13)$$

Then

$$\frac{Q_{y|s,u}}{Q_{y|s',u'}} = \frac{\mathbb{P}(\mathcal{Q}(X) = y | S = s)}{\mathbb{P}(\mathcal{Q}(X) = y | S = s')} \leq e^\varepsilon. \quad (14)$$

On the other hand, suppose that $\frac{Q_{y|s,u}}{Q_{y|s',u'}} \leq e^\varepsilon$ for all $s \neq s'$ and u, u' . Then for all $s \neq s'$ and P we have

$$\frac{\mathbb{P}(\mathcal{Q}(X) = y | S = s)}{\mathbb{P}(\mathcal{Q}(X) = y | S = s')} = \frac{\sum_u Q_{y|s,u} P_{u|s}}{\sum_{u'} Q_{y|s',u'} P_{u'|s'}} \leq e^\varepsilon. \quad (15)$$

Hence, \mathcal{Q} satisfies $(\varepsilon, \mathcal{P}_{\mathcal{X}})$ -RLDP w.r.t. \mathcal{F} . \square

The proposition demonstrates that RLDP is very similar to LDP. The difference is that the condition “for all $x, x' \in \mathcal{X}$ ” from Definition 2 is relaxed to only those x and x' for which $s \neq s'$. We will exploit this difference. Recall that Generalised Randomised Response [12] is the privacy protocol $\text{GRR}^\varepsilon: \mathcal{X} \rightarrow \mathcal{X}$ given by

$$\text{GRR}_{y|x}^\varepsilon = \begin{cases} \frac{e^\varepsilon}{e^\varepsilon + a - 1} & \text{if } x = y, \\ \frac{1}{e^\varepsilon + a - 1} & \text{otherwise.} \end{cases} \quad (16)$$

This protocol has been designed such that $\frac{\text{GRP}_{y|x}^\varepsilon}{\text{GRR}_{y|x'}^\varepsilon} = e^{\pm\varepsilon}$ for $x \neq x'$, the maximal fractional difference that ε -LDP allows. We will see that for RLDP we can go up to a difference of $e^{\pm 2\varepsilon}$ if $x = (s, u)$ and $x' = (s, u')$, as we typically only need to satisfy

$$Q_{y|s,u} \leq e^\varepsilon Q_{y|s',u'} \leq e^{2\varepsilon} Q_{y|s,u'}. \quad (17)$$

We capture the intuition from necessary condition (17) in a new protocol called *Secret Randomised Response (SRR)*.

Definition 5 (Secret Randomised Response (SRR)). *Let $\varepsilon > 0$. Then the release protocol $\text{SRR}^\varepsilon : \mathcal{X} \rightarrow \mathcal{X}$ is given by*

$$\text{SRR}_{s',u'|s,u}^\varepsilon = \begin{cases} \frac{e^\varepsilon}{e^\varepsilon + e^{-\varepsilon}(a_2-1) + a - a_2}, & \text{if } (s', u') = (s, u), \\ \frac{e^{-\varepsilon}}{e^\varepsilon + e^{-\varepsilon}(a_2-1) + a - a_2}, & \text{if } s' = s \text{ and } u' \neq u, \\ \frac{1}{e^\varepsilon + e^{-\varepsilon}(a_2-1) + a - a_2}, & \text{if } s' \neq s, \end{cases} \quad (18)$$

The next result demonstrates that the necessary condition (17) is, in the case of SRR, also sufficient.

Lemma 6. *SRR satisfies $(\varepsilon, \mathcal{P}_{\mathcal{X}})$ -RLDP.*

Proof. It can be directly verified that for all $s \neq s', u, u'$ and y we have $\frac{\text{SRR}_{y|s,u}^\varepsilon}{\text{SRR}_{y|s',u'}^\varepsilon} \in \{e^{-\varepsilon}, 1, e^\varepsilon\}$, from which $(\varepsilon, \mathcal{P}_{\mathcal{X}})$ -RPP follows. \square

As for utility, note that the robust utility metric $\min_{P \in \mathcal{F}} I_P(X; Y)$ is not useful if $\mathcal{F} = \mathcal{P}_{\mathcal{X}}$, since by considering a degenerate P it follows immediately that $I_P(X; Y) = 0$ for every Q . However, SRR is optimal in the following sense:

Theorem 7. *For every P , there is a $\varepsilon_0 \geq 0$ such that for all $\varepsilon \gg \varepsilon_0$ such that SRR is the $(\varepsilon, \mathcal{P}_{\mathcal{X}})$ -RLDP protocol maximising $I_P(X; Y)$.*

The proof of this theorem follows along the same lines as the proof of Theorem 14 of [2], in which it is proven that GRR is the optimal LDP protocol for ε large enough. The proof is presented in Appendix A-A.

IV. PROPERTIES OF THE DOMAIN \mathcal{F}

From this point onwards we consider \mathcal{F} to be of the form in (3). Before we introduce new algorithms in Sections V–VII, we need some technical results on properties of \mathcal{F} . First some notation: for $u \in \mathcal{U}$

and $s \in \mathcal{S}$, we write

$$P_u = \sum_s P_{u,s}, \quad (19)$$

$$P_s = \sum_u P_{u,s}, \quad (20)$$

$$P_{u|s} = \frac{P_{u,s}}{P_s}, \quad (21)$$

$$P_{\mathcal{U}|s} = (P_{u|s})_{u \in \mathcal{U}} \in \mathcal{P}_{\mathcal{U}}. \quad (22)$$

The following lemma states that for every s , the image of \mathcal{F} under the projection $P \mapsto P_{\mathcal{U}|s}$ is again of the form in (3).

Lemma 8. *Let $s \in \mathcal{S}$ such that $\hat{P}_s > 0$. Let $\mathcal{F}_{\mathcal{U}|s}$ be the projection of \mathcal{F} onto $\mathcal{P}_{\mathcal{U}}$ via the map $P \mapsto P_{\mathcal{U}|s} \in \mathcal{P}_{\mathcal{U}}$. Define $B_s := \frac{(\sqrt{B+1} + \hat{P}_s - 1)^2}{\hat{P}_s^2} - 1$. Then*

$$\mathcal{F}_{\mathcal{U}|s} = \left\{ R \in \mathcal{P}_{\mathcal{U}} : \sum_u \frac{(\hat{P}_{u|s} - R_u)^2}{R_u} \leq B_s \right\}. \quad (23)$$

Proof. For $P \in \mathcal{F}$ and $s \in \mathcal{S}$ one has, using the definition of \mathcal{F} in (4),

$$\frac{\hat{P}_s^2}{P_s} \sum_u \frac{\hat{P}_{u|s}^2}{P_{u|s}} = \sum_u \frac{\hat{P}_{s,u}^2}{P_{s,u}} \quad (24)$$

$$\leq B + 1 - \sum_{s' \neq s} \sum_u \frac{\hat{P}_{s',u}^2}{P_{s',u}} \quad (25)$$

$$= B + 1 - \frac{(1 - \hat{P}_s)^2}{1 - P_s} \sum_{s' \neq s} \sum_u \frac{\hat{P}_{s',u| \neg s}^2}{P_{s',u| \neg s}}, \quad (26)$$

where for $s' \in \mathcal{S} \setminus \{s\}$ and $u \in \mathcal{U}$ we define $P_{s',u| \neg s} = \frac{P_{u,s'}}{1 - P_s}$. These form a probability distribution on $(\mathcal{S} \setminus \{s\}) \times \mathcal{U}$. As such we have

$$\sum_{s' \neq s} \sum_u \frac{\hat{P}_{u,s'| \neg s}^2}{P_{u,s'| \neg s}} = 1 + \sum_{s' \neq s} \sum_u \frac{(P_{u,s'| \neg s} - \hat{P}_{u,s'| \neg s})^2}{P_{u,s'| \neg s}} \geq 1. \quad (27)$$

It follows that

$$\sum_u \frac{\hat{P}_{u|s}^2}{P_{u|s}} \leq \frac{P_s}{\hat{P}_s^2} \left(B + 1 - \frac{(1 - \hat{P}_s)^2}{1 - P_s} \right). \quad (28)$$

We find the maximum of the right hand side by differentiating with respect to P_s , for which we get

$$\frac{B + 1}{\hat{P}_s^2} - \frac{(1 - \hat{P}_s)^2}{\hat{P}_s^2 (1 - P_s)^2}. \quad (29)$$

Setting this equal to 0 and solving w.r.t. P_s , we find that the maximum is attained at $P_s = 1 - \frac{1 - \hat{P}_s}{\sqrt{B+1}}$.

Substituting this, we find

$$\frac{P_s}{\hat{P}_s^2} \left(B + 1 - \frac{(1 - \hat{P}_s)^2}{1 - P_s} \right) \leq \frac{(\sqrt{B+1} - 1 + \hat{P}_s)^2}{\hat{P}_s^2} = B_s + 1, \quad (30)$$

hence $\sum_u \frac{(P_{u|s} - \hat{P}_{u|s})^2}{P_{u|s}} \leq B_s$; this shows the inclusion “ \subset ” in (23). On the other hand, suppose that $R \in \mathcal{P}_{\mathcal{U}}$ satisfies $\sum_u \frac{\hat{P}_{u|s}^2}{R_u} \leq B_s + 1$. Let $c = 1 - \frac{1 - \hat{P}_s}{\sqrt{B+1}}$, and define $P \in \mathcal{P}_{\mathcal{X}}$ by

$$P_{u,s'} = \begin{cases} cR_u, & \text{if } s' = s, \\ \frac{\hat{P}_{u,s'}}{\sqrt{B+1}} & \text{otherwise.} \end{cases} \quad (31)$$

Then $P_{\mathcal{U}|s} = R$, and

$$\sum_{u,s'} \frac{\hat{P}_{u,s'}^2}{P_{u,s'}} = \sum_u \frac{\hat{P}_{u,s}^2}{cR_u} + \sum_u \sum_{s' \neq s} \sqrt{B+1} \hat{P}_{u,s'} \quad (32)$$

$$= \frac{\hat{P}_s^2 \sqrt{B+1}}{\sqrt{B+1} - 1 + \hat{P}_s} \sum_u \frac{\hat{P}_{u|s}^2}{R_u} + \sqrt{B+1}(1 - \hat{P}_s) \quad (33)$$

$$\leq \frac{\hat{P}_s^2 \sqrt{B+1}}{\sqrt{B+1} - 1 + \hat{P}_s} \cdot \frac{(\sqrt{B+1} - 1 + \hat{P}_s)^2}{\hat{P}_s^2} + \sqrt{B+1}(1 - \hat{P}_s) \quad (34)$$

$$= B + 1, \quad (35)$$

hence $P \in \mathcal{F}$. This shows the inclusion “ \supset ” in (23). \square

This lemma implies that many results which hold for \mathcal{F} also hold for $\mathcal{F}_{\mathcal{U}|s}$. For what follows, we need Lemma 9 and Proposition 10 that are given next. Lemma 9 gives tight bounds on P_x given \hat{P} and B . Will use this in Section V to describe polyhedral approximations of \mathcal{F} and the \mathcal{F}_s , which we will use in turn to obtain useful lower bounds on the utility that can be obtained under RLDP.

Lemma 9. *Let $x \in \mathcal{X}$. Then*

$$\min_{P \in \mathcal{F}} P_x = \frac{B + 2\hat{P}_x - \sqrt{B^2 + 4B\hat{P}_x - 4B\hat{P}_x^2}}{2B + 2}, \quad (36)$$

$$\max_{P \in \mathcal{F}} P_x = \frac{B + 2\hat{P}_x + \sqrt{B^2 + 4B\hat{P}_x - 4B\hat{P}_x^2}}{2B + 2}. \quad (37)$$

Proof. Evidently the minimum and maximum exist and are attained on the boundary, i.e. for P satisfying $\sum_{x'} \frac{\hat{P}_{x'}^2}{P_{x'}} = B + 1$. Thus for finding both the minimum and the maximum we have to find the stationary points of

$$P_x + \lambda \left(\sum_{x'} \frac{\hat{P}_{x'}^2}{P_{x'}} - B - 1 \right) + \mu \left(\sum_{x'} P_{x'} - 1 \right). \quad (38)$$

Taking derivatives with respect to all $P_{x'}$, we find

$$1 + \mu - \lambda \frac{\hat{P}_x^2}{P_x^2} = 0, \quad (39)$$

$$\forall x' \neq x : \mu - \lambda \frac{\hat{P}_{x'}^2}{P_{x'}^2} = 0. \quad (40)$$

It follows that for $x' \neq x$, we have $P_{x'} = c\hat{P}_{x'}$, with $c = \sqrt{\frac{\lambda}{\mu}}$. Since $\sum_{x'} P_{x'} = \sum_{x'} \hat{P}_{x'} = 1$, hence $c = \frac{1-P_x}{1-\hat{P}_x}$. Substituting this in the boundary constraint yields

$$\frac{\hat{P}_x^2}{P_x} + \frac{(1-\hat{P}_x)^2}{1-P_x} = B + 1. \quad (41)$$

Solving this for P_x gives us

$$P_x = \frac{B + 2\hat{P}_x \pm \sqrt{B^2 + 4B\hat{P}_x - 4B\hat{P}_x^2}}{2B + 2}, \quad (42)$$

giving both the minimum and maximum. \square

The following Proposition gives a bound on $\|P - \hat{P}\|_1$ in terms of \hat{P} and B , which is tight for $B \geq 1$. This is an essential ingredient to the explicit privacy protocols introduced in Sections VI and VII. The proof is rather long and technical, so we present it in Appendix A-B.

Proposition 10. *Let B and $\hat{P} \in \mathcal{P}_{\mathcal{X}}$ be given.*

1) *Suppose $B \geq 1$. Let $x_{\min} \in \arg \min_{x \in \mathcal{X}} \hat{P}_x$. Then*

$$\max_{P \in \mathcal{F}} \|P - \hat{P}\|_1 = \frac{B - 2B\hat{P}_{x_{\min}} + \sqrt{B^2 + 4B\hat{P}_{x_{\min}} - 4B\hat{P}_{x_{\min}}^2}}{B + 1}. \quad (43)$$

2) *Suppose $B < 1$. Then $\max_{P \in \mathcal{F}} \|P - \hat{P}\|_1 \leq \sqrt{B}$.*

V. POLYHEDRAL APPROXIMATION: POLYOPT

Our first method to find RLDP protocols for when \mathcal{F} is a confidence interval from a χ^2 test relies on optimising $I_P(X; Y)$ over protocols that satisfy a more stringent privacy constraint; this yields a lower bound on the maximal $I_P(X; Y)$. More concretely, we consider protocols that satisfy (2) for all P for which $P_{\mathcal{U}|s} \in \mathcal{D}_{\mathcal{U}|s}$, where each $\mathcal{D}_{\mathcal{U}|s}$ is a polyhedron containing the set $\mathcal{F}_{\mathcal{U}|s}$ from Lemma 8. All $P \in \mathcal{F}$ certainly satisfy this condition. For each s, u , let $P_{u|s}^{\min} = \inf_{P \in \mathcal{F}} P_{u|s}$: an explicit formula is given in Lemma 9. When each $\mathcal{D}_{\mathcal{U}|s}$ is the simplex $\{R : \forall u R_u \geq P_{u|s}^{\min}\}$, robust optimisation for polytopes [13] yields the following result. Let Γ be the convex cone consisting of all $T \in \mathbb{R}_{\geq 0}^{\mathcal{X}}$ satisfying

$$\forall s_1, s_2, u_1, u_2 : T_{s_1, u_1} - e^\varepsilon T_{s_2, u_2} + \sum_u P_{u|s_1}^{\min} (T_{s_1, u} - T_{s_1, u_1}) - e^\varepsilon \sum_u P_{u|s_2}^{\min} (T_{s_2, u} - T_{s_2, u_2}) \leq 0. \quad (44)$$

Theorem 11. *Let \mathcal{Q} be a privacy protocol such that for all y we have $Q_y \in \Gamma$. Then \mathcal{Q} satisfies $(\varepsilon, \mathcal{F})$ -RLDP.*

Theorem 12. Let $\hat{\Gamma}$ be polytope given by $\{T \in \Gamma : \sum_x T_x = 1\}$. Let \mathcal{V} be the set of vertices of $\hat{\Gamma}$. For $v \in \mathcal{V}$, define

$$\mu^1(v) = \sum_x v_x \hat{P}_x \log \frac{v_x}{\sum_{x'} v_{x'} \hat{P}_{x'}}, \quad (45)$$

$$\mu^2(v) = \min_{P \in \mathcal{F}} \sum_x v_x P_x \log \frac{v_x}{\sum_{x'} v_{x'} P_{x'}}. \quad (46)$$

For $i = 1, 2$, let $\hat{\theta}^i$ be the solution to the optimisation problem

$$\text{maximise}_{\theta} \sum_{v \in \mathcal{V}} \theta_v \mu^i(v) \quad (47)$$

$$\text{satisfying } \theta \in \mathbb{R}_{\geq 0}^{\mathcal{V}},$$

$$\sum_v \theta_v v = 1_{\mathcal{X}}.$$

Let the privacy protocol \mathcal{Q}^i be given by $\mathcal{Y}^i = \{v \in \mathcal{V} : \hat{\theta}_v^i > 0\}$ and $Q_{v|x}^i = \hat{\theta}_v^i v_x$. Then:

- 1) The protocol \mathcal{Q}^1 maximises $I_{\hat{P}}(X; Y)$ among all protocols satisfying the condition of Theorem 11. One has $|\mathcal{Y}^i| \leq a$.
- 2) Let $L = \sum_{v \in \mathcal{V}} \hat{\theta}_v^2 \mu^2(v)$. Then \mathcal{Q}^2 satisfies $\inf_{P \in \mathcal{F}} I_P(X; Y) \geq L$.

Together, these two theorems show, if we can solve a vertex enumeration problem, that we can find a protocol \mathcal{Q}^1 that maximises $I_{\hat{P}}(X; Y)$ among a subset of all $(\varepsilon, \mathcal{F})$ -RLDP \mathcal{Q} , a lower bound for the achievable $\min_P I_P(X; Y)$, and a protocol \mathcal{Q}^2 that exceeds this bound.

In Theorem 12, to calculate $\mu^2(v)$ one needs to take the minimum over all $P \in \mathcal{F}$. To approximate this, one may replace \mathcal{F} by a polyhedron containing it; the minimum is then attained at one of its vertices.

Before we prove Theorem 11, we need an intermediate result. For a privacy protocol \mathcal{Q} and a $y \in \mathcal{Y}$, we let Q_y be the vector $(Q_{y|x})_x \in \mathbb{R}^{\mathcal{X}}$. Furthermore, for $s_1, s_2 \in \mathcal{S}$, let $B^{s_1, s_2} \in \mathbb{R}^{\mathcal{X} \times \mathcal{X}}$ be given by

$$B_{s, u; s', u'}^{s_1, s_2} = \begin{cases} 1, & \text{if } u = u' \text{ and } s = s' = s_1, \\ -\varepsilon, & \text{if } u = u' \text{ and } s = s' = s_2, \\ 0, & \text{otherwise.} \end{cases} \quad (48)$$

Lemma 13. Let $\mathcal{D} \subset (\mathcal{P}_{\mathcal{U}})^{\mathcal{S}} \subset \mathbb{R}^{\mathcal{X}}$ be a polyhedron such that for every $P \in \mathcal{F}$ one has $(P_{\mathcal{U}|s})_{s \in \mathcal{S}} \in \mathcal{D}$. Let \mathcal{D} be given by the equations $DR + d \geq 0$ and $ER + e = 0$, for matrices D and E , vectors d and e , and $R \in \mathbb{R}^{\mathcal{S} \times \mathcal{U}}$. Let \mathcal{Q} be a privacy protocol such that for all $y \in \mathcal{Y}$ and $s_1, s_2 \in \mathcal{S}$ there

exist z, w such that

$$D^T z + E^T w = -(B^{s_1, s_2})^T Q_y, \quad (49)$$

$$z \geq 0, \quad (50)$$

$$d^T z + e^T w \leq 0. \quad (51)$$

Then Q satisfies ε -RLDP w.r.t. \mathcal{F} .

Proof. For $y \in \mathcal{Y}$ and $s \in \mathcal{S}$, write $Q_{y,s} := (Q_{y|s,u})_u \in \mathbb{R}^{\mathcal{U}}$, and $Q_y := (Q_{y|s,u})_{s,u} \in \mathbb{R}^{\mathcal{X}}$. We can then formulate ε -RLDP as

$$\forall y, s_1, s_2 : \max_{P \in \mathcal{F}} P_{\mathcal{U}|s_1}^T Q_{y,s_1} - \varepsilon P_{\mathcal{U}|s_2}^T Q_{y,s_2} \leq 0. \quad (52)$$

Set $\mathcal{G} = \prod_s \mathcal{F}_{\mathcal{U}|s}$. Then \mathcal{D} satisfies the conditions of the Lemma if and only if $\mathcal{G} \subset \mathcal{D}$. In particular, the following condition implies (52):

$$\forall y, s_1, s_2 : \max_{R \in \mathcal{D}} R_{s_1}^T Q_{y,s_1} - \varepsilon R_{s_2}^T Q_{y,s_2} \leq 0. \quad (53)$$

Using the matrices B^{s_1, s_2} , we can rewrite (53) as

$$\forall y, s_1, s_2 : \max_{R \in \mathcal{D}} ((B^{s_1, s_2})^T Q_y)^T R \leq 0. \quad (54)$$

Now fix s_1, s_2, y . By dualising we have

$$\max_{R \in \mathcal{D}} ((B^{s_1, s_2})^T Q_y)^T R = \min_{\substack{z, w: \\ D^T z + E^T w = -(B^{s_1, s_2})^T Q_y \\ z \geq 0}} d^T z + e^T w. \quad (55)$$

It follows that Q satisfies ε -RLDP if for each y, s_1, s_2 there exist $z \geq 0$ and w satisfying $D^T z + E^T w = -(B^{s_1, s_2})^T Q_y$ such that $d^T z + e^T w \leq 0$. \square

Proof of Theorem 11. Define $\mathcal{D}_{\mathcal{U}|s} = \{R \in \mathcal{P}_{\mathcal{S}} : \forall u R_u \geq P_{u|s}^{\min}\}$, and let $\mathcal{D} = \prod_s \mathcal{D}_{\mathcal{U}|s}$. This satisfies the conditions of Lemma 13. One checks that in this case we have $D = \text{id}_{\mathcal{X}}$, $d \in \mathbb{R}^{\mathcal{X}}$ is given by $d_{s,u} = -P_{u|s}^{\min}$, $E \in \mathbb{R}^{\mathcal{S} \times \mathcal{X}}$ is given by $E_{s;u',s'} = \delta_{s=s'}$, and $e = -1_{\mathcal{S}}$. This also means that $z \in \mathbb{R}^{\mathcal{X}}$ and $w \in \mathbb{R}^{\mathcal{S}}$. It follows from these descriptions that

$$D^T z = z, \quad (56)$$

$$(E^T w)_{s,u} = w_s, \quad (57)$$

$$((B^{s_1, s_2})^T Q_y)_{s,u} = \begin{cases} Q_{y|s_1,u}, & \text{if } s = s_1, \\ -\varepsilon Q_{y|s_2,u}, & \text{if } s = s_2, \\ 0, & \text{otherwise.} \end{cases} \quad (58)$$

It follows that $D^T z + E^T w = -(B^{s_1, s_2})^T Q_y$ can be rewritten as

$$z_{s,u} = \begin{cases} -Q_{y|s_1,u} - w_{s_1}, & \text{if } s = s_1, \\ e^\varepsilon Q_{y|s_2,u} - w_{s_2}, & \text{if } s = s_2, \\ -w_s & \text{otherwise.} \end{cases} \quad (59)$$

Eliminating z from (50) and(51), we get

$$-\sum_s \left(1 - \sum_u P_{u|s}^{\min}\right) w_s + \sum_u Q_{y|s_1,u} P_{u|s_1}^{\min} - e^\varepsilon \sum_u Q_{y|s_2,u} P_{u|s_2}^{\min} \leq 0, \quad (60)$$

$$\forall u : w_{s_1} \leq -Q_{y|s_1,u}, \quad (61)$$

$$\forall u : w_{s_2} \leq e^\varepsilon Q_{y|s_2,u}, \quad (62)$$

$$\forall s \neq s_1, s_2 : w_s \leq 0. \quad (63)$$

Since $\sum_u P_{u|s}^{\min} \leq 1$ for all s , it follows that the left hand side of (60) is minimal if each w_s attains its maximal value, subject to the constraints (61–63). It follows that the minimum of the left hand side is equal to

$$\begin{aligned} & \left(1 - \sum_u P_{u|s_1}^{\min}\right) \left(\max_{u_1} Q_{y|u_1, s_1}\right) - e^\varepsilon \left(1 - \sum_u P_{u|s_2}^{\min}\right) \left(\min_{u_2} Q_{y|u_2, s_2}\right) \\ & + \sum_u Q_{y|s_1,u} P_{u|s_1}^{\min} - e^\varepsilon \sum_u Q_{y|s_2,u} P_{u|s_2}^{\min}. \end{aligned} \quad (64)$$

This is nonpositive if and only if it is nonpositive for all choices of u_1 and u_2 ; but this is true precisely if $Q_y \in \Gamma$. \square

The proof of Theorem 12 is analogous to the proof of Theorem 4 of [2]. It is presented in Appendix A-C. The algorithm that produces \mathcal{Q}^1 from \hat{P} and ε will be referred to as *PolyOpt* in the remainder of this paper.

Remark 14. *A simplex is not the only possible choice for $\mathcal{D}_{\mathcal{U}|s}$. In general, we can make $\mathcal{D}_{\mathcal{U}|s}$ closer to $\mathcal{F}_{\mathcal{U}|s}$ by adding more defining hyperplanes. Doing this allows more \mathcal{Q} to satisfy Theorem 11, and in turn increase the utility of the \mathcal{Q} we find via Theorem 12. However, since Γ is related to the $\mathcal{D}_{\mathcal{U}|s}$ via duality, adding extra constraints to the $\mathcal{D}_{\mathcal{U}|s}$ will increase the dimension of Γ through the addition of auxiliary variables. This makes the vertex enumeration problem of Theorem 12 more computationally involved. Thus we have a tradeoff between utility and computational complexity.*

It should be noted that in general the increasing utility found in this way does not approach the optimal utility over all $(\varepsilon, \mathcal{F})$ -RLDP protocols. This is because, as we take increasingly finer $\mathcal{D}_{\mathcal{U}|s}$, we approach the set of \mathcal{Q} that satisfy (2) for all P in $\mathcal{F}' := \{P : \forall s P_{\mathcal{U}|s} \in \mathcal{F}_{\mathcal{U}|s}\}$. Since in general $\mathcal{F} \subsetneq \mathcal{F}'$, the set of $(\varepsilon, \mathcal{F}')$ -RLDP protocols is strictly smaller than the set of $(\varepsilon, \mathcal{F})$ -RLDP protocols.

VI. INDEPENDENT REPORTING

As PolyOpt relies on vertex enumeration, it can be computationally infeasible for larger a . In this section, we consider a class of release protocols which we call *Independent Reporting*. We show that within this class the optimal protocols can be found by finding the maximum of a one-dimensional function. Since the dimension of this optimisation problem does not depend on a , this approach can be used when vertex enumeration is out of reach. As mentioned before we continue to let \mathcal{F} be a confidence set for a χ^2 test.

The basis of IR is to apply two separate LDP protocols \mathcal{R}^1 and \mathcal{R}^2 to S and U , respectively, and output $(\mathcal{R}^1(S), \mathcal{R}^2(U))$. This is described in Protocol 1.

Protocol 1: $\text{IR}_{\mathcal{Q}^1, \mathcal{Q}^2}$ (Independent reporting)

Input : Privacy protocols $\mathcal{R}^1: \mathcal{S} \rightarrow \mathcal{Y}^1$ and $\mathcal{R}^2: \mathcal{U} \rightarrow \mathcal{Y}^2$; $x = (s, u) \in \mathcal{X}$.

Output: Output datum $y \in \mathcal{Y} := \mathcal{Y}^1 \times \mathcal{Y}^2$

Compute $y_1 \leftarrow \mathcal{Q}^1(s)$;

Compute $y_2 \leftarrow \mathcal{Q}^2(u)$;

$y \leftarrow (y_1, y_2)$;

While only S needs to be protected, we also need to apply a privacy protocol to U because of the possible correlation between the two. However, since U only indirectly leaks information about S , we can get away with less strict privacy requirements. This is reflected in the following theorem.

Theorem 15. Let $\varepsilon_1, \varepsilon_2 \in \mathbb{R}_{\geq 0}$. For each s , define B_s as in Lemma 8, and let $u_s \in \mathcal{U}$ be such that $\hat{P}_{u_s|s}$ is minimal. Define

$$d_s := \begin{cases} \frac{B_s(1-2\hat{P}_{u_s|s}) + \sqrt{B_s^2 + 4B_s\hat{P}_{u_s|s} - 4B_s\hat{P}_{u_s|s}^2}}{B_s+1}, & \text{if } B_s \geq 1; \\ \sqrt{B_s}, & \text{if } B_s \leq 1. \end{cases} \quad (65)$$

Furthermore, define

$$d := \min \left\{ 2, \max_s(2d_s) + \max_{s,s'} \|\hat{P}_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1 \right\}. \quad (66)$$

Let $\delta_2 = \log \left(1 + \frac{2(e^{\varepsilon_2}-1)}{d} \right)$. Suppose that \mathcal{R}^1 is ε_1 -LDP and that \mathcal{R}^2 is δ_2 -LDP. Then IR is $(\varepsilon_1 + \varepsilon_2, \mathcal{F})$ -RLDP.

Proof. We start by showing that d is an upper bound for $\|P_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1$. If $d = 2$, this is certainly the case. Suppose $d = \max_s(2d_s) + \max_{s,s'} \|\hat{P}_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1$. It follows from Lemma 10 that for each

$P \in \mathcal{F}$ and each $s \in \mathcal{S}$ we have $\|P_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s}\|_1 \leq d_s$. Hence, for all $s, s' \in \mathcal{S}$ and $P \in \mathcal{F}$ we have

$$\|P_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1 \leq \|P_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s}\|_1 + \|\hat{P}_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1 + \|\hat{P}_{\mathcal{U}|s'} - P_{\mathcal{U}|s'}\|_1 \quad (67)$$

$$\leq d_s + d_{s'} + \|\hat{P}_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1 \quad (68)$$

$$\leq d. \quad (69)$$

Combining Lemma 3 with the fact that $\varepsilon_2 = \log\left(1 + \frac{d(e^{\delta_2} - 1)}{2}\right)$, it follows that for every $y_2 \in \mathcal{Y}_2$ we have

$$\frac{\mathbb{P}_P(\mathcal{R}^2(U) = y_2 | S = s)}{\mathbb{P}_P(\mathcal{R}^2(U) = y_2 | S = s')} \leq 1 + \frac{e^{\delta_2} - 1}{2} \|P_{\mathcal{U}|s} - P_{\mathcal{U}|s'}\|_1 \quad (70)$$

$$\leq 1 + \frac{d(e^{\delta_2} - 1)}{2} \quad (71)$$

$$= e^{\varepsilon_2}. \quad (72)$$

Since \mathcal{R}^1 is ε_1 -LDP, it follows that for every $y_1 \in \mathcal{Y}_1$ and every $y_2 \in \mathcal{Y}_2$ we have

$$\frac{\mathbb{P}(\mathcal{R}^1(S) = y_1, \mathcal{R}^2(U) = y_2 | S = s)}{\mathbb{P}(\mathcal{R}^1(S) = y_1, \mathcal{R}^2(U) = y_2 | S = s')} \leq e^{\varepsilon_1 + \varepsilon_2}, \quad (73)$$

which shows that $\text{IR}_{\mathcal{R}^1, \mathcal{R}^2}$ is $(\varepsilon_1 + \varepsilon_2, \mathcal{F})$ -RLDP. \square

The more independent S and U are, the smaller $\max_{s, s'} \|\hat{P}_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1$ will be. Theorem 15 then tells us that for more independent S and U , the privacy requirements on \mathcal{R}^2 will be less strict, resulting in better utility. The utility of IR is described by the following theorem:

Theorem 16. *For any $P \in \mathcal{P}_{\mathcal{X}}$ one has*

$$I_P(\text{IR}_{\mathcal{R}^1, \mathcal{R}^2}(X); X) = I_P(\mathcal{R}^1(S); S) + I_P(\mathcal{R}^2(U); U | \mathcal{R}^1(S)). \quad (74)$$

Proof. Since $\mathcal{R}^1(S)$ and U are independent given S , and $\mathcal{R}^2(U)$ and S are independent given U , we have

$$I_P(\text{IR}_{\mathcal{R}^1, \mathcal{R}^2}(X); X) = I_P(\mathcal{R}^1(S), \mathcal{R}^2(U); U, S) \quad (75)$$

$$= I_P(\mathcal{R}^1(S); U, S) + I_P(\mathcal{R}^2(U); U, S | \mathcal{R}^1(S)) \quad (76)$$

$$= I_P(\mathcal{R}^1(S); S) + I_P(\mathcal{R}^2(U); U | \mathcal{R}^1(S)). \square \quad (77)$$

Given an $\varepsilon \geq 0$, we can use these theorems to find $(\varepsilon, \mathcal{F})$ -RLDP protocols. Per Theorem 15, it suffices to take a ε_2 , and use a ε_1 -LDP protocol \mathcal{R}^1 and a δ_2 -LDP protocol \mathcal{R}^2 , where $\varepsilon_1 = \varepsilon - \varepsilon_2$ and δ_2 is as in Theorem 15. We want to choose ε_2 , \mathcal{R}^1 and \mathcal{Q}^2 in such a way that we maximise the expression in Theorem 16. For ε large enough, the \mathcal{R}^1 that maximises $I_P(\mathcal{R}^1(S); S)$ is GRR. Furthermore, since

$$I_P(\mathcal{R}^2(U); U | \mathcal{R}^1(S)) = \mathbb{E}_r [I_{U \sim P_U | \mathcal{R}^1(S)=r}(\mathcal{R}^2(U); U)], \quad (78)$$

and GRR maximises $I(\mathcal{R}^2(U); U)$ for large enough ε for any distribution of U , we should take \mathcal{R}^2 to be GRR as well; this is true regardless of the value of P . We are left with only the unknown ε_2 , hence to maximise the mutual information of IR for a given P we have to solve a one-dimensional optimisation problem.

VII. CONDITIONAL REPORTING

From Theorem 15 it is clear that in IR we can afford a larger privacy budget to \mathcal{R}^2 if S and U are only weakly correlated. When S and R are closer related, however, the difference between δ_2 and ε_2 will be small, and IR cannot offer any advantage over general LDP protocols. To this end, we introduce two other protocols that fall under the umbrella term *Conditional Reporting*. In both these protocols, we apply an established privacy protocol \mathcal{R}^1 to S . Furthermore, we return U (with a small perturbation) if \mathcal{R}^1 returns a ‘correct’ response and a random U otherwise. We will see that the noise on U depends on the size of the feasible set \mathcal{F} rather than on the correlation between S and U .

A. GRR-CR

For the first CR protocol, GRR-CR, we first need to specify a parameter ε_1 and, for each $s \in \mathcal{S}$, a privacy protocol $\mathcal{R}^s: \mathcal{U} \rightarrow \mathcal{Y}_s$, where each \mathcal{Y}_s is a finite set. To apply it to an input datum $(s, u) \in \mathcal{X}$, we first apply GRR with parameter ε_1 to s ; call the outcome \tilde{S} . If $\tilde{S} = s$, we apply \mathcal{R}^s to u , and we output $(s, \mathcal{R}^s(u))$. If $\tilde{S} \neq s$, we draw a random $\tilde{U} \in \mathcal{U}$ from the probability distribution $\hat{P}_{\mathcal{U}|\tilde{S}}$, and we output $(\tilde{S}, \mathcal{R}^{\tilde{S}}(\tilde{U}))$. This protocol is described in Protocol 2.

Protocol 2: GRR-CR

Input : Privacy parameter ε_1 ; For every $s \in \mathcal{S}$, a privacy protocol $\mathcal{Q}^s: \mathcal{U} \rightarrow \mathcal{Y}_s$; input datum

$$x = (s, u) \in \mathcal{X}$$

Output: Output datum $Y \in \mathcal{S} \times \bigcup_{s \in \mathcal{S}} \mathcal{Y}_s$

Take $\tilde{S} \leftarrow \text{GRR}^{\varepsilon_1}(s) \in \mathcal{S}$;

if $\tilde{S} = s$ **then**

 | Compute $Y \leftarrow (s, \mathcal{R}^s(u))$;

else

 | Sample $\tilde{U} \in \mathcal{U}$ with $\mathbb{P}(\tilde{U} = u') = \hat{P}_{u'|\tilde{S}}$;

 | Compute $Y \leftarrow (\tilde{S}, \mathcal{R}^{\tilde{S}}(\tilde{U}))$;

end

Output Y ;

Although we have already obfuscated S via GRR, we still need to obfuscate U and \tilde{U} via $\mathcal{R}^{\tilde{S}}$ for the following reason. Suppose we omit this last step, and instead return (\tilde{S}, \tilde{U}) , with $\tilde{U} = u$ if $\tilde{S} = s$. From the viewpoint of an attacker, given \tilde{S} , the random variable \tilde{U} is drawn from the distribution $P_{\mathcal{U}|\tilde{S}}$ if $\tilde{S} = s$, and from the distribution $\hat{P}_{\mathcal{U}|\tilde{S}}$ otherwise. In the LDP model the attacker may collude with an arbitrary amount of users, and as such we may assume that they have access to the real distribution $P \in \mathcal{P}_{\mathcal{X}}$. Under this assumption, the output \tilde{U} contains information about whether it was drawn from $P_{\mathcal{U}|\tilde{S}}$ or $\hat{P}_{\mathcal{U}|\tilde{S}}$, and hence whether $S = \tilde{S}$ or not. To prevent this leakage, we have to mask \tilde{U} with the privacy protocol $\mathcal{R}^{\tilde{S}}$. As the following theorem shows, the privacy level that is needed for \mathcal{R}^s depends on $\|\hat{P}_{\mathcal{U}|s} - P_{\mathcal{U}|s}\|_1$, which explains why we need a different protocol \mathcal{R}^s for every s .

Theorem 17. *Let $\varepsilon_1, \varepsilon_2 \in \mathbb{R}_{\geq 0}$. For every $s \in \mathcal{S}$ define d_s as in (65), define $\delta_s := \log\left(1 + \frac{2(e^{\varepsilon_2} - 1)}{d_s}\right)$, and let \mathcal{Q}^s satisfy δ_s -LDP. Then Algorithm 2 satisfies $(\varepsilon_1 + \varepsilon_2, \mathcal{F})$ -RLDP.*

Proof. For $s \in \mathcal{S}$, $u \in \mathcal{U}$, and $y \in \mathcal{Y}_s$, let $R_{y|u}^s = \mathbb{P}(\mathcal{Q}^s(u) = y)$. Then for every $\tilde{s}, s \in \mathcal{S}$ and every $y \in \mathcal{Y}_{\tilde{s}}$ we have

$$\mathbb{P}(\tilde{S} = \tilde{s}, \mathcal{R}^{\tilde{s}}(U) = y | S = s) = \begin{cases} \frac{e^{\varepsilon_1} \sum_u R_{y|u}^{\tilde{s}} P_{u|\tilde{s}}}{e^{\varepsilon_1} + a_1 - 1}, & \text{if } s = \tilde{s}, \\ \frac{\sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}}{e^{\varepsilon_1} + a_1 - 1}, & \text{if } s \neq \tilde{s}. \end{cases} \quad (79)$$

It follows that for every $s' \in \mathcal{S}$ we have

$$\frac{\mathbb{P}(\tilde{S} = \tilde{s}, \mathcal{R}^{\tilde{s}}(U) = y | S = s)}{\mathbb{P}(\tilde{S} = \tilde{s}, \mathcal{R}^{\tilde{s}}(U) = y | S = s')} = \begin{cases} 1, & \text{if } s = s' \\ e^{\varepsilon_1} \frac{\sum_u R_{y|u}^{\tilde{s}} P_{u|\tilde{s}}}{\sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}}, & \text{if } s = \tilde{s} \neq s' \\ e^{-\varepsilon_1} \frac{\sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}}{\sum_u R_{y|u}^{\tilde{s}} P_{u|\tilde{s}}}, & \text{if } s \neq \tilde{s} = s' \\ 1, & \text{if } s \neq \tilde{s} \neq s' \end{cases} \quad (80)$$

$$\leq e^{\varepsilon_1} \max \left\{ \frac{\sum_u R_{y|u}^{\tilde{s}} P_{u|\tilde{s}}}{\sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}}, \frac{\sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}}{\sum_u R_{y|u}^{\tilde{s}} P_{u|\tilde{s}}} \right\}. \quad (81)$$

Since $\|P_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s}\|_1 \leq d_s$, we find by Lemma 3 that

$$\frac{\sum_u R_{y|u}^{\tilde{s}} P_{u|\tilde{s}}}{\sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}} = \frac{\mathbb{P}_{\tilde{U} \sim P_{\mathcal{U}|\tilde{s}}}(\tilde{Q}^{\tilde{s}}(\tilde{U}) = y)}{\mathbb{P}_{\tilde{U} \sim \hat{P}_{\mathcal{U}|\tilde{s}}}(\tilde{Q}^{\tilde{s}}(\tilde{U}) = y)} \quad (82)$$

$$\leq 1 + \frac{\|P_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s}\|_1 (e^{\delta_s} - 1)}{2} \quad (83)$$

$$\leq 1 + \frac{d_s (e^{\delta_s} - 1)}{2} \quad (84)$$

$$= e^{\varepsilon_2}. \quad (85)$$

The same holds analogously for $\frac{\sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}}{\sum_u R_{y|u}^{\tilde{s}} P_{u|\tilde{s}}}$, and it follows that GRR-CR satisfies $(\varepsilon_1 + \varepsilon_2, \mathcal{F})$ -RLDP w.r.t. \mathcal{F} . \square

As we can see, the privacy level of \mathcal{R}^s only depends on $\max_P \|P_{U|s} - \hat{P}_{U|s}\|_1$. This makes GRR-CR an attractive protocol if this is small, which happens if either the number of known data points n is large or if α is small.

On the side of utility, we have the following:

Theorem 18. *For any P one has*

$$I_P(\text{GRR-CR}(X); X) = I_P(\text{GRR}^{\varepsilon_1}(S); S) + \frac{e^{\varepsilon_1}}{e^{\varepsilon_1} + a_1 - 1} I_P(\mathcal{R}^S(U); U|S). \quad (86)$$

Proof. One has

$$\begin{aligned} I_P(\text{GRR-CR}(X); X) &= I_P(\tilde{S}; S) + I_P(\mathcal{R}^{\tilde{S}}(U); S|\tilde{S}) \\ &\quad + I_P(\tilde{S}; U|S) + I_P(\mathcal{R}^{\tilde{S}}(U); U|S, \tilde{S}). \end{aligned} \quad (87)$$

Note that $\mathbb{P}(\mathcal{Q}^{\tilde{S}}(U) = y|S = s, \tilde{S} = \tilde{s}) = \sum_u R_{y|u}^{\tilde{s}} \hat{P}_{u|\tilde{s}}$. This does not depend on s , hence S and $\mathcal{R}^{\tilde{S}}(U)$ are independent given \tilde{S} , and $I(\mathcal{R}^{\tilde{S}}(U); S|\tilde{S}) = 0$. Furthermore, \tilde{S} and U are independent given S , hence $I_P(\tilde{S}; U|S) = 0$. For the last term we have

$$I_P(\mathcal{R}^{\tilde{S}}(U); U|S, \tilde{S}) = \sum_{s, \tilde{s}} \mathbb{P}(\tilde{S} = \tilde{s}|S = s) \mathbb{P}(S = s) I_P(\mathcal{R}^{\tilde{s}}(U); U|S = s, \tilde{S} = \tilde{s}). \quad (88)$$

We know that $\mathcal{R}^s(U)$ and U are independent given S and \tilde{S} if $S \neq \tilde{S}$, hence in the summation above only terms with $s = \tilde{s}$ matter; hence this is equal to

$$\sum_s \mathbb{P}(\tilde{S} = s|S = s) \mathbb{P}(S = s) I_P(\mathcal{R}^s(U); U|S = \tilde{S} = s) = \frac{e^{\varepsilon_1}}{e^{\varepsilon_1} + a_1 - 1} I_P(\mathcal{R}^S(U); U|S). \quad (89)$$

The theorem now follows from putting this all together. \square

Compared to Theorem 16, we see that if we take $\mathcal{R}^s = \mathcal{R}^2$ for every s , then GRR-CR typically has a lower utility than IR. However, the advantage of GRR-CR is that \mathcal{R}^s can typically be chosen with more relaxed privacy conditions than \mathcal{R}^2 , which will increase the utility again.

Since $I_P(\mathcal{R}^S(U); U|S) = \sum_s \mathbb{P}(S = s) I_{U \sim P_{U|s}}(\mathcal{R}^s(U); U)$, Theorem 18 tells us that we want to choose each \mathcal{R}^s to be the δ_s -LDP protocol for which $I_{U \sim P_{U|s}}(\mathcal{R}^s(U); U)$ is maximised. For big enough δ_s , this is GRR. As was the case with IR, it is a one-dimensional optimisation problem to optimise for $I_{\hat{P}}(X; Y)$.

B. UE-CR

The second CR protocol we introduce originates from Unary Encoding (UE) [24]. UE is a protocol $\text{UE}^{\kappa, \lambda}: \mathcal{S} \rightarrow 2^{\mathcal{S}}$ given by parameters $0 < \lambda \leq \kappa < 1$ that for an input s outputs a binary vector $(E_{s'})_{s' \in \mathcal{S}}$, where each coefficient is an independent Bernoulli variable with

$$\mathbb{P}(E_{s'} = 1) = \begin{cases} \kappa, & \text{if } s = s', \\ \lambda, & \text{if } s \neq s'. \end{cases} \quad (90)$$

This protocol satisfies ε -LDP for $\varepsilon \geq \log \frac{\kappa(1-\lambda)}{\lambda(1-\kappa)}$. Popular choices for (κ, λ) are $(\frac{e^{\varepsilon/2}}{e^{\varepsilon/2}+1}, \frac{1}{e^{\varepsilon/2}+1})$, $(\frac{1}{2}, \frac{1}{e^{\varepsilon}+1})$, and $(\frac{e^{\varepsilon}}{e^{\varepsilon}+1}, \frac{1}{2})$ [24]. It will be convenient for us to consider the output of UE as a subset of \mathcal{S} , rather than a binary vector.

To apply UE-CR to a $(s, u) \in \mathcal{X}$, we first fix parameters κ, λ , and a privacy protocol $\mathcal{R}^s: \mathcal{U} \rightarrow \mathcal{Y}^s$ for every s . We perform UE on s , yielding a subset $\tilde{S} \subset \mathcal{S}$. For every $s' \in \mathcal{S}$, we output a $Y_{s'} \in \mathcal{Y}^{s'}$ as follows: if $s' = s$, we take $Y_s = \mathcal{R}^s(u)$. If $s' \neq s$, we draw a $\tilde{U} \in \mathcal{U}$ with probability distribution $\hat{P}_{\mathcal{U}|s'}$, and we take $Y_{s'} = \mathcal{R}^{s'}(\tilde{U})$. Finally, we output $(\tilde{S}, (Y_{s'})_{s' \in \tilde{S}})$. This is described in Protocol 3.

Protocol 3: UE-CR

Input : Parameters $0 \leq \lambda \leq \kappa \leq 1$; For every $s \in \mathcal{S}$, a privacy protocol $\mathcal{R}^s: \mathcal{U} \rightarrow \mathcal{Y}_s$; a

probability distribution \hat{P} on \mathcal{X} ; input datum $x = (s, u) \in \mathcal{X}$

Output: Output datum $(\tilde{S}, (Y_{s'})_{s' \in \tilde{S}})$ with $\tilde{S} \subset \mathcal{S}$ and $Y_{s'} \in \mathcal{Y}^{s'}$ for each s'

Take $\tilde{S} \leftarrow \text{UE}^{\kappa, \lambda}(s) \subset \mathcal{S}$;

for $s' \in \tilde{S}$ **do**

if $s' = s$ **then**

$Y_s \leftarrow \mathcal{R}^s(u)$;

else

 Sample $\tilde{U} \sim \hat{P}_{\mathcal{U}|s}$;

$Y_{s'} \leftarrow \mathcal{R}^{s'}(\tilde{U})$;

end

end

Output $(\tilde{S}, (Y_{s'})_{s' \in \tilde{S}})$;

As for GRR-CR, the privacy protocols $\mathcal{R}^{s'}$ are needed to obfuscate the difference between $\hat{P}_{\mathcal{U}|s}$ and $P_{\mathcal{U}|s}$. The precise privacy requirements for the \mathcal{R}^s are given in the theorem below.

Theorem 19. *Let $\varepsilon_1, \varepsilon_2 \in \mathbb{R}_{\geq 0}$. For every s , let δ_s be as in Theorem 15, and assume \mathcal{R}^s has δ_s -LDP and that $\frac{\kappa(1-\lambda)}{\lambda(1-\kappa)} \leq e^{\varepsilon_1}$. Then UE-CR satisfies $(\max\{\varepsilon_1 + \varepsilon_2, 2\varepsilon_2\}, \mathcal{F})$ -SLDP w.r.t. \mathcal{F} .*

Proof. For $s \in \mathcal{S}$, define $T_s := \sum_u R_{y_s|u}^s P_{u|s}$ and $\hat{T}_s := \sum_u R_{y_s|u}^s \hat{P}_{u|s}$. One has

$$\mathbb{P}(Y = (\tilde{s}, (y_{s'})_{s' \in \tilde{s}}) | S = s) \tag{91}$$

$$= \begin{cases} \kappa \lambda^{|\tilde{s}|-1} (1-\lambda)^{a_1-|\tilde{s}|} T_s \prod_{s' \in \tilde{s} \setminus \{s\}} \hat{T}_{s'}, & \text{if } s \in \tilde{s}, \\ \lambda^{|\tilde{s}|} (1-\kappa) (1-\lambda)^{a_1-|\tilde{s}|-1} \prod_{s' \in \tilde{s}} \hat{T}_{s'}, & \text{if } s \notin \tilde{s}. \end{cases} \tag{92}$$

It follows that

$$\frac{\mathbb{P}(Y = (\tilde{s}, (y_{s'})_{s' \in \tilde{s}}) | S = s)}{\mathbb{P}(Y = (\tilde{s}, (y_{s'})_{s' \in \tilde{s}}) | S = s')} \quad (93)$$

$$= \begin{cases} 1, & \text{if } s = s' \text{ or } s, s' \notin \tilde{s}, \\ \frac{\kappa(1-\lambda)T_s}{\lambda(1-\kappa)\hat{T}_s}, & \text{if } s \in \tilde{s} \not\ni s', \\ \frac{\lambda(1-\kappa)\hat{T}_{s'}}{\kappa(1-\lambda)T_{s'}}, & \text{if } s \notin \tilde{s} \ni s', \\ \frac{T_s \hat{T}_{s'}}{\hat{T}_s T_{s'}}, & \text{if } s, s' \in \tilde{s} \text{ and } s \neq s'. \end{cases} \quad (94)$$

By Lemma 3, one has $\frac{T_s}{\hat{T}_s}, \frac{\hat{T}_s}{T_s} \leq 1 + \frac{e^{\delta s} - 1}{2} \|P_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s}\|_1 \leq e^{\varepsilon_2}$ for each $s \in \mathcal{S}$. Since $\frac{\lambda(1-\kappa)}{\kappa(1-\lambda)} \leq \frac{\kappa(1-\lambda)}{\lambda(1-\kappa)} \leq e^{\varepsilon_1}$, it follows that

$$\frac{\mathbb{P}(Y = (\tilde{s}, (y_{s'})_{s' \in \tilde{s}}) | S = s)}{\mathbb{P}(Y = (\tilde{s}, (y_{s'})_{s' \in \tilde{s}}) | S = s')} \leq \max \{e^{\varepsilon_1 + \varepsilon_2}, e^{2\varepsilon_2}\}, \quad (95)$$

which proves the SLDP. \square

This theorem shows that, similar to GRR-CR, the privacy requirements on the \mathcal{R}^s become less strict as $P_{\mathcal{U}|s}$ and $\hat{P}_{\mathcal{U}|s}$ are closer. As for utility, we find the following theorem:

Theorem 20. *For any P one has*

$$\mathbb{I}_P(\text{UE-CR}(X); X) = \mathbb{I}_P(\text{UE}(S); S) + \kappa \mathbb{I}_P(\mathcal{R}^S(U); U|S). \quad (96)$$

Proof. One has

$$\mathbb{I}_P(\text{UE-CR}(X); X) = \mathbb{I}_P(\tilde{S}, (Y_{s'})_{s' \in \tilde{S}}; S, U) \quad (97)$$

$$= \mathbb{I}_P(\tilde{S}; S, U) + \mathbb{I}_P((Y_{s'})_{s' \in \tilde{S}}; S|\tilde{S}) + \mathbb{I}_P((Y_{s'})_{s' \in \tilde{S}}; U|S, \tilde{S}). \quad (98)$$

Similar to the proof of Theorem 18 we have that given \tilde{S} , the random variables $(Y_{s'})_{s' \in \tilde{S}}$ and S are independent, hence $\mathbb{I}_P((Y_{s'})_{s' \in \tilde{S}}; S|\tilde{S}) = 0$. Furthermore, \tilde{S} and U are independent given S , hence $\mathbb{I}_P(\tilde{S}; S, U) = \mathbb{I}_P(\tilde{S}; S)$. Furthermore, we can write

$$\mathbb{I}_P((Y_{s'})_{s' \in \tilde{S}}; U|S, \tilde{S}) = \mathbb{E}_{s, \tilde{s}} \left[\mathbb{I}_P((Y_{s'})_{s' \in \tilde{S}}; U|S = s, \tilde{S} = \tilde{s}) \right]. \quad (99)$$

If $s \notin \tilde{s}$, then U and $(Y_{s'})_{s' \in \tilde{S}}$ are independent given $S = s$ and $\tilde{S} = \tilde{s}$. If $s \in \tilde{s}$, then U and $Y_{s'}$ are independent given $S = s$ and $\tilde{S} = \tilde{s}$, for $s' \neq s$. It follows that

$$\mathbb{I}_P((Y_{s'})_{s' \in \tilde{S}}; U|S = s, \tilde{S} = \tilde{s}) = \begin{cases} \mathbb{I}_P(Y_s; U|S = s), & \text{if } s \in \tilde{s}, \\ 0, & \text{otherwise.} \end{cases} \quad (100)$$

From this we conclude that

$$\mathbb{I}_P((Y_{s'})_{s' \in \tilde{S}}; U|S, \tilde{S}) = \sum_s \mathbb{P}(s \in \tilde{S}|S = s) \mathbb{P}(S = s) \mathbb{I}_P(Y_s; U|S = s) \quad (101)$$

$$= \kappa \mathbb{I}(Y_S; U|S). \quad (102)$$

As before, we can conclude from this that we should let all \mathcal{R}^s be GRR. This leaves us to finding ε_2 , κ and λ , which is a three-dimensional optimisation problem.

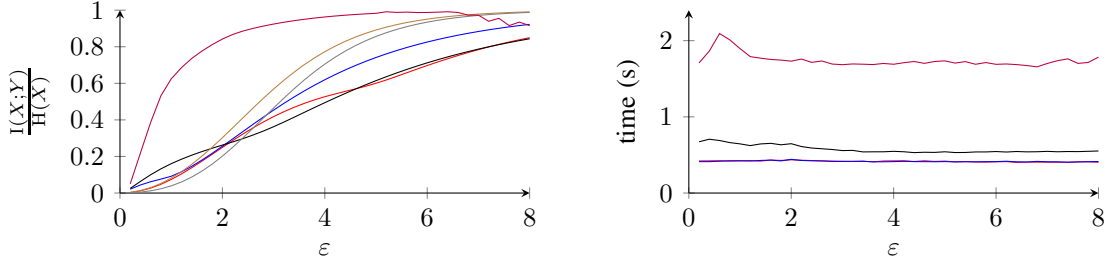


Fig. 2: Experiments on synthetic data for $a_1 = a_2 = 3$ (— *PolyOpt*, — *IR*, — *GRR-CR*, — *UE-CR*, — *SRR*, — *GRR*).

VIII. EXPERIMENTS

In order to test the feasibility of the different methods we perform several experiments, both on synthetic and real data. Throughout, we take $\alpha = 0.05$ unless stated otherwise.

Throughout the experiments, we use $I_{\hat{P}}(X;Y)$ as a utility metric, occasionally normalised by dividing by $H(X)$. We use this rather than $I_{P^*}(X;Y)$, as the aggregator only has access to the former. In fact, while P^* is known for the synthetic data, this is not the case for real data, so we cannot even use $I_{P^*}(X;Y)$ as a utility metric.

The outline of the remainder of this section is as follows. In Section VIII-A we focus on the *PolyOpt* method. In Section VIII-B we consider the impact of a_1 and a_2 . In Section VIII-C we analyse the optimal value of the parameters of our protocols. In Section VIII-D we investigate the difference between *IR* and *GRR-CR*. In Section VIII-E we analyze the role of n and α . In Section VIII-F, we investigate the difference $I_{\hat{P}}(X;Y) - I_{P^*}(X;Y)$ for synthetic data, to evaluate the robustness of the utility metric. Finally, in Section VIII-G we consider real data.

A. *PolyOpt*

We first perform experiments to test the utility of the *PolyOpt* method introduced in Section V. We perform numerical experiments on synthetic data. For $a_1 = a_2 = 3$, we draw 200 distributions from the Jeffreys prior on the space of probability distributions on \mathcal{X} . For each distribution, we draw $n = 1000$ items from this distribution, and we demand robustness w.r.t. this observed distribution. For each observed distribution, for $\varepsilon \in [0.2, 8]$, and for each protocol of *PolyOpt*, *IR*, *GRR-CR*, *UE-CR* and *SRR*¹, we calculate the normalised utility $\frac{I_{\hat{P}}(X;Y)}{H(X)}$, which we average over all distributions. As a reference we perform the same analysis on *GRR*, the LDP protocol that maximises mutual information for large ε . Since *GRR* satisfies ε -LDP, it certainly satisfies $(\varepsilon, \mathcal{F})$ -RLDP for any \mathcal{F} .

¹for *SRR* we ignore the value of α .

The results are in Figure 2. As we can see, PolyOpt significantly outperforms the other methods, although the optimisation we used (we used Matlab, specifically the MPT3 toolbox) becomes more inaccurate at larger ε . However, the downside of Polyopt lies in its computation time, which is significantly higher than that of other methods. All experiments were conducted on a PC with Intel Core i7-7700HQ 2.8GHz and 32GB memory. As can be observed in Figure 2 for larger a the computation time increases dramatically: for $a_1 = 3$, $a_2 = 4$ the computation time is 72s on average, and for $a_1 = a_2 = 4$ we terminated the computation when it was still running after 12 hours.

In general, if the user has enough computation power to use the PolyOpt method, then this is recommended, because it clearly outperforms all other protocols. However, it is possible that this is computationally unfeasible. For most of our other experiments, we assume that this is the case, and we study the utility of the other methods.

B. Synthetic data

We perform the same procedure as before, but for different a_1, a_2 . The results are in Figures 3 and 4. As can be seen, for ε large enough, SRR is the best protocol, which is remarkable as it has the strictest privacy requirement. The larger a_1 and a_2 are, the larger ε has to be for SRR to become the preferred method. We see that IR and GRR-CR perform more or less similar. For small a_1 and a_2 , we see that UE-CR outperforms these; for high ε , on the other hand, UE-CR is the worse choice. This is understandable considering the fact that UE yields less mutual information between input and output than GRR [10].

Looking at GRR, we see that it performs slightly worse than SRR across all ε . This is expected behaviour since the protocols are very similar, but SRR is better tailored to the Privacy Funnel scenario.

C. Optimal parameter settings

We plot the values of $\varepsilon_2/\varepsilon$ for IR and GRR-CR, and κ and λ for UE-CR, to get insight into the ideal parameter settings. We take $a_1 = a_2 = 5$, and we draw two distributions from the Jeffreys prior ($n = 1000$). We also draw 200 distributions, and take the average parameter settings over these distributions. These graphs are depicted in Figure 5.

As we can see from the samples, for low ε it is optimal to take either $\varepsilon_2 = 0$ or $\varepsilon_2 = \varepsilon$; in IR and GRR-CR this means to use all the privacy budget for either transmitting S or U . Note that when the whole privacy budget is spent on S , then IR and GRR-CR are the same protocol; this explains why they behave so similar for low ε . Furthermore, we see that for GRR-CR it is beneficial for any P to spend the entire privacy budget on U for $\varepsilon < 1$, and on S for slightly higher $1 < \varepsilon < 3$. By

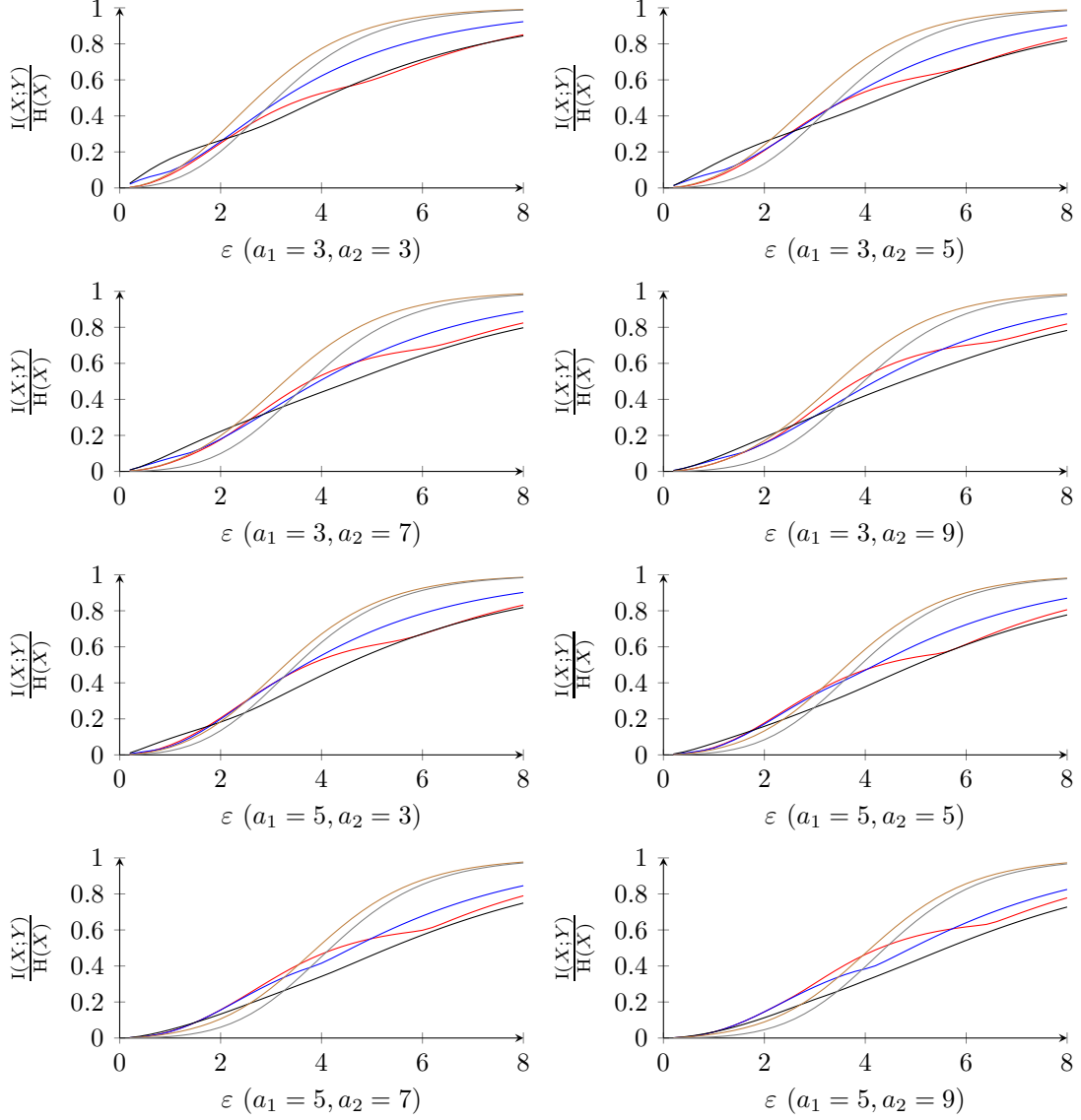


Fig. 3: Experiments on synthetic data, $a_1 \in \{3, 5\}$, $a_2 \in \{3, 5, 7, 9\}$ (— IR , — $GRR-CR$, — $UE-CR$, — SRR , — GRR).

contrast, it depends on P whether the privacy budget of IR is to be spent on S or on U for low ε . This explains why the average value of $\varepsilon_2/\varepsilon$ is close to 0.5 regardless of the value of ε for IR .

For $UE-CR$ we also see that the privacy budget is spent only on one of the two components: for low ε , it is optimal to take $\kappa = \lambda = 1$, which means there is no information leakage about S . It is only when ε grows larger that it becomes optimal to divide the privacy budget among S and U . The point where such a division is optimal depends on the distribution.

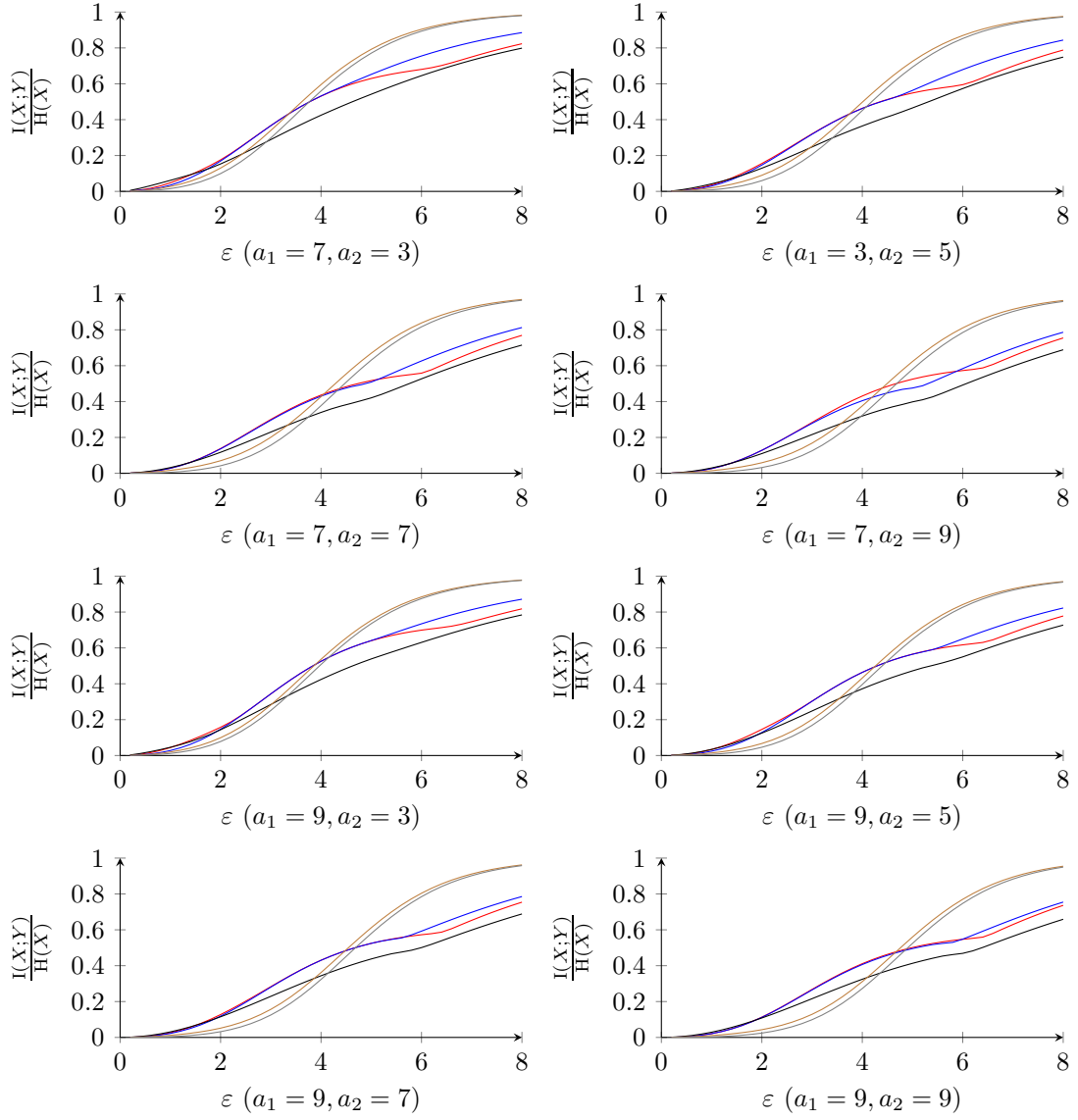


Fig. 4: Experiments on synthetic data, $a_1 \in \{7, 9\}$, $a_2 \in \{3, 5, 7, 9\}$ (— IR, — GRR-CR, — UE-CR, — SRR, — GRR).

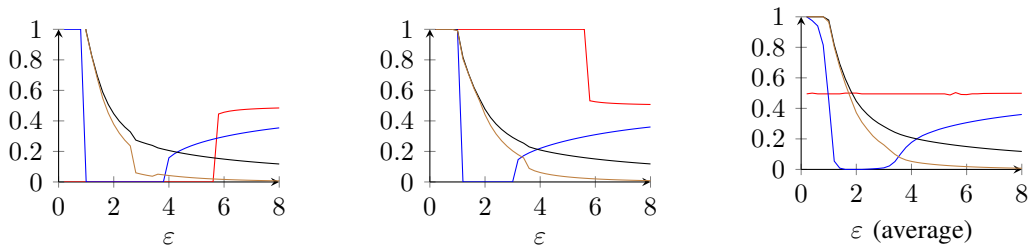


Fig. 5: Parameter values for IR and CR ($a_1 = a_2 = 5$) for two distributions and the average over 200 distributions (— $\varepsilon_2/\varepsilon$ (IR), — $\varepsilon_2/\varepsilon$ (GRR-CR), — κ (UE-CR), — λ (UE-CR)).

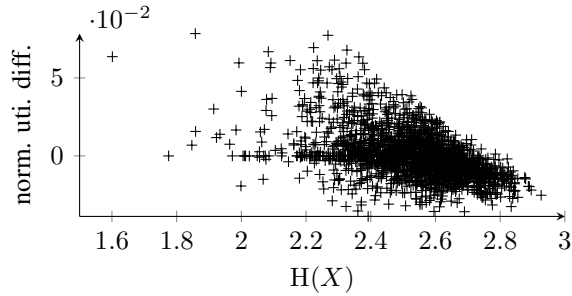


Fig. 6: Difference in $\frac{I(X;Y)}{H(X)}$ between SR and GRR-CR, plotted against $H(X)$, for 2000 distributions drawn from the Jeffreys prior for $a_1 = a_2 = 5$; $\varepsilon = 4$.

D. IR vs GRR-CR

We also try to find out what causes the difference between IR and GRR-CR. In Figure 6 we plot the normalised difference in utility between these two protocols against $H(X)$, for 2000 randomly generated probability distributions (with $a_1 = a_2 = 5$ and $\varepsilon = 4$). As one can see, there are many distributions where the two have equal utility, which is caused by the fact that the two protocols coincide when the whole privacy budget is allotted to S . Among the other distributions, however, we see a downward trend signifying that GRR-CR outperforms IR for large $H(X)$.

E. Role of n and α

We also vary n and α , which were taken to be 1000 and 0.05 before, respectively. Taking larger n and smaller α have the same effect, namely that \mathcal{F} is smaller. As can be seen from Figure 7, taking larger n has no effect on SRR and GRR, as they do not depend on \mathcal{F} . For IR, GRR-CR and UE-CR, we see that the larger n is, the better utility they provide. This is more pronounced for GRR-CR and UE-CR than it is for SR, which can be explained from the fact that the privacy parameter δ_2 from Theorem 15 does not only depend on the size of \mathcal{F} , but also on $\max_{s,s'} \|\hat{P}_{\mathcal{U}|s} - \hat{P}_{\mathcal{U}|s'}\|_1$. As such, the increase in utility that comes from reducing \mathcal{F} is more limited than with CR.

We also look at the effect of α on the utility of IR, GRR-CR, and UE-CR. As mentioned before, the smaller α , the larger \mathcal{F} , and the less utility the protocols will provide. This is reflected in Figure 8, where we see that having a smaller α reduces the utility of GRR-CR and UE-CR (we take $a_1 = a_2 = 5$ and $n = 1000$, and take the average over 200 distributions). For IR there is no difference at all: this is because the maximum $d = 2$ is obtained in Theorem 15, at which point the protocol is not affected by changing α . For CR the loss of utility caused by changing α is rather small in absolute terms, but becomes important for low ε , as a change from α from 0.1 to 0.0001 causes an average utility loss of 36% for GRR-CR, and 47% for UE-CR for $\varepsilon = 0.1$.

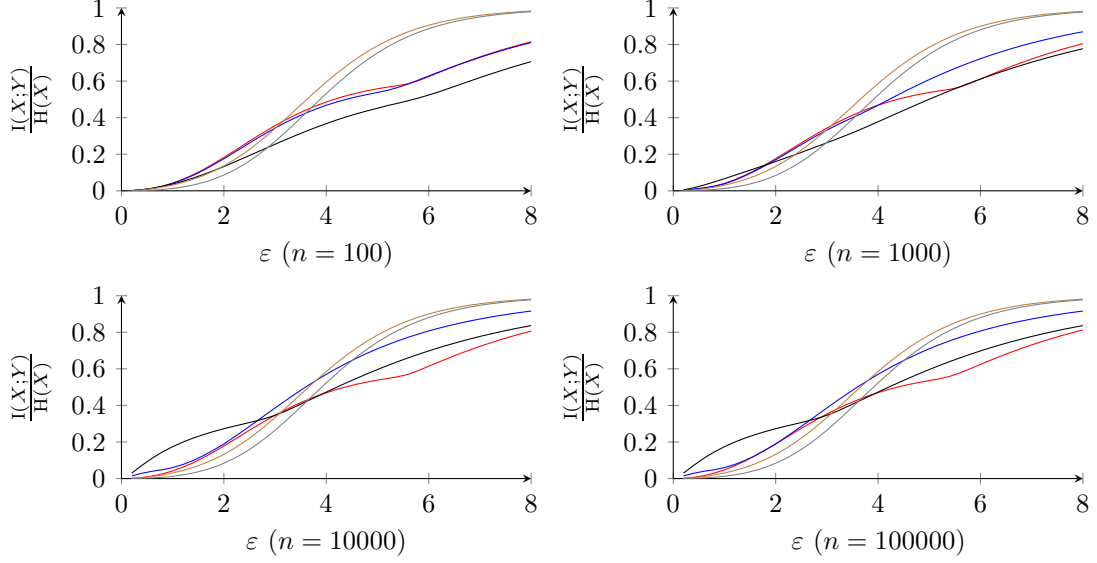


Fig. 7: Experiments on synthetic data with n changing, $a_1 = a_2 = 5$ (— IR, — GRR-CR, — UE-CR, — SRR, — GRR).

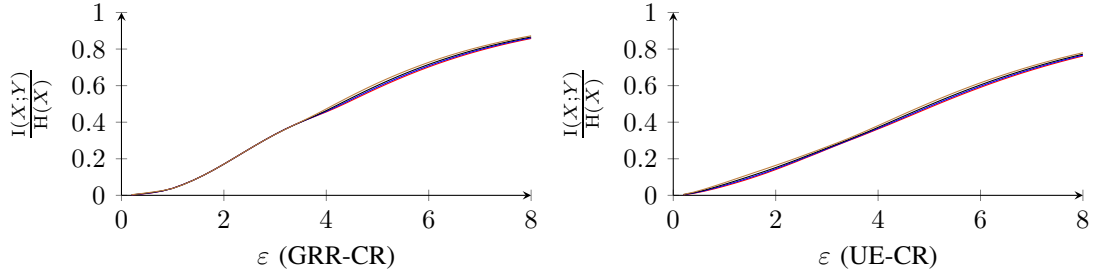


Fig. 8: Experiments on synthetic data with α changing, $a_1 = a_2 = 5$ (— $\alpha = 10^{-4}$, — $\alpha = 10^{-3}$, — $\alpha = 10^{-2}$, — $\alpha = 10^{-1}$).

F. Robustness of utility

We also consider the robustness of the utility by comparing the ‘true utility’ $I_{P^*}(X; Y)$ to $I_{\hat{P}}(X; Y)$, the latter of which is maximised in IR and CR. The results (for $a_1 = a_2 = 5$) are in Figure 9. As one can see, the true utility is on average often actually higher than the optimised utility, especially for small n . Furthermore, the difference between the two utilities rapidly becomes negligible for larger n . We conclude that IR and CR produce robust utility results.

G. Adult dataset

We also perform numerical experiments on the adult-dataset ($n = 32561$) [35], which contains demographic data from the 1994 US census. Some examples, where we use different categorical attributes from the dataset as S and U , are depicted in Figure 10. To compare them to the synthetic

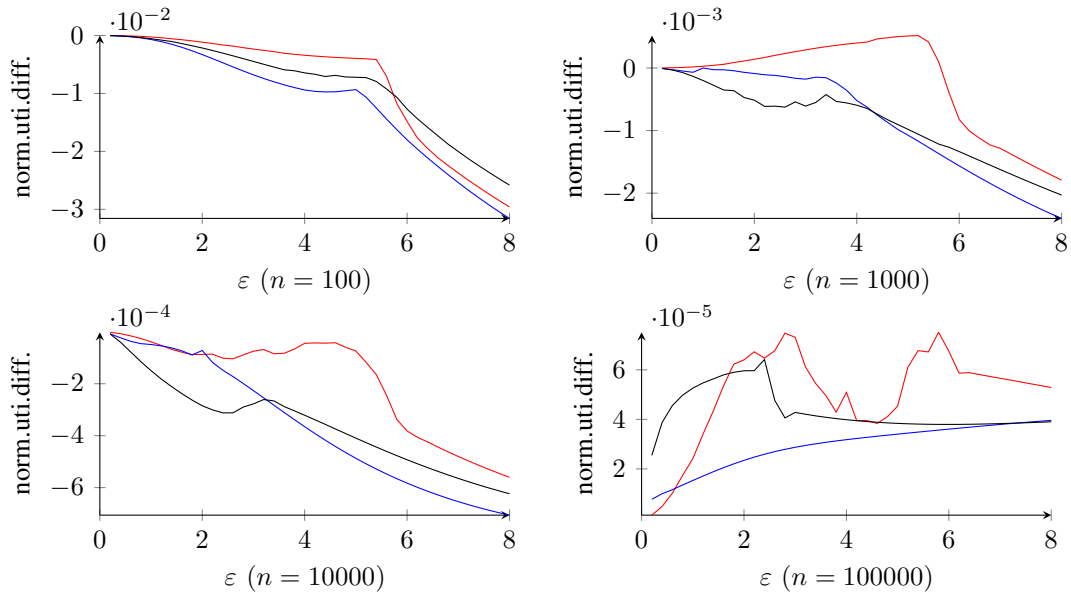


Fig. 9: The average value of $\frac{I_{\hat{P}}(X;Y) - I_{P^*}(X;Y)}{H(X)}$, for 200 randomly generated distributions, with $a_1 = a_2 = 5$ (— IR, — GRR-CR, — UE-CR).

data, we also perform experiments on synthetic data with the same a_1, a_2 as in the experiments in Figure 10; these experiments are in Figure 11. As we can see, the relative behaviour of the methods on the real data and the synthetic data of the same dimension align closely. The largest difference is the fact that IR outperforms GRR-CR for the synthetic data for $a_1 = 6, a_2 = 42$, but this is because the relative performance of IR and GRR-CR is distribution-specific, as we have seen in Section VIII-D. The close correspondence between the synthetic and real-data experiments lends additional validity to the experiments on synthetic data in the rest of this section.

IX. CONCLUSION AND FUTURE WORK

In this paper, we presented a number of algorithms that, given a desired privacy level ϵ , an estimated distribution \hat{P} , and a set of probability distributions \mathcal{F} of a specified form, return a release protocol that aim to maximise the mutual information between input and output, while satisfying privacy w.r.t. a given sensitive part of the data, for all distributions in \mathcal{F} . In the case that $\mathcal{F} = \mathcal{P}_{\chi}$, we have introduced SRR, which we have shown to be optimal for this \mathcal{F} in the low privacy regime, irrespective of the actual probability distribution. Furthermore, experiments show that in the low privacy regime SRR outperforms most of our other algorithms, even though these have smaller \mathcal{F} . The privacy level at which SRR overtakes the other algorithms in utility is lower for larger input spaces and smaller \mathcal{F} . However, in the high privacy regime the other algorithms offer significantly better utility. This shows the validity of using confidence sets in the RLDP framework.

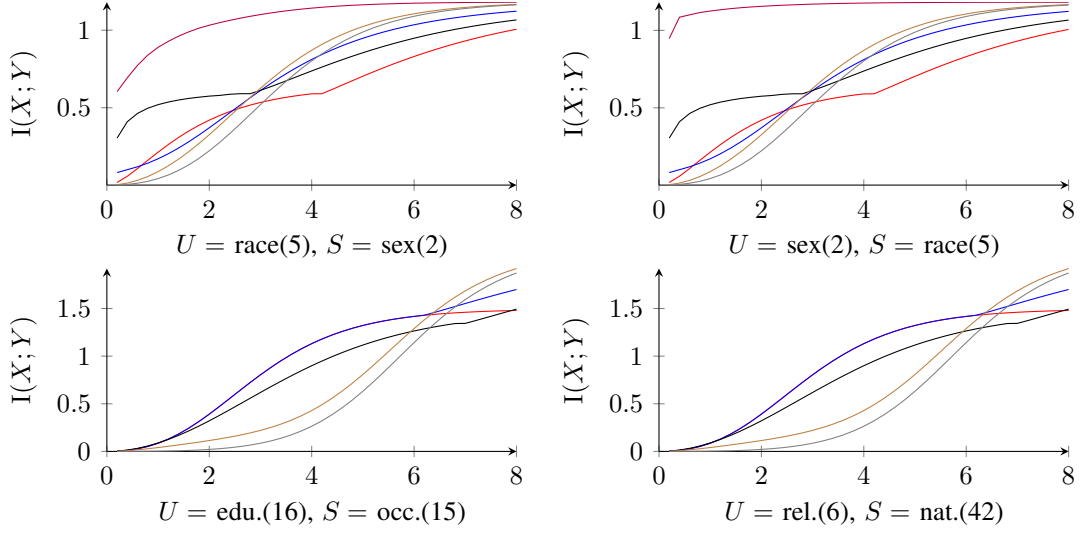


Fig. 10: Experiments on the categories *sex*, *race*, *education*, *occupation*, *relation* and *native-country* of the adult-dataset. Numbers between brackets indicate a_1 and a_2 (— PolyOpt, — IR, — GRR-CR, — UE-CR, — SRR, — GRR).

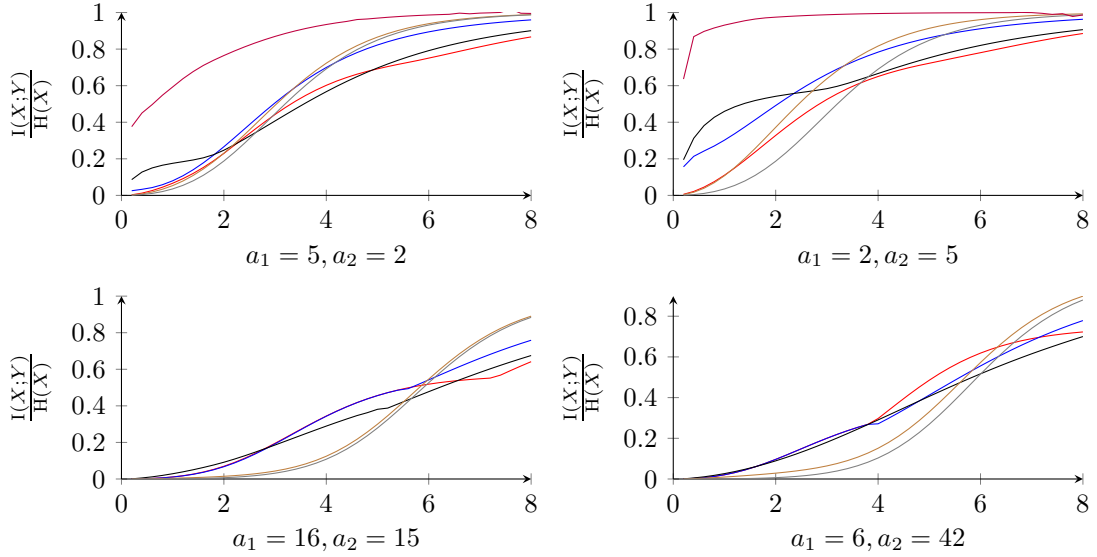


Fig. 11: Experiments on synthetic data, with a_1, a_2, n as in the adult-dataset. Horizontal axis is ϵ -LDP; vertical axis is $I(X;Y)/H(X)$ (— PolyOpt, — IR, — GRR-CR, — UE-CR, — SRR, — GRR).

In the case that \mathcal{F} is a confidence set around \hat{P} , arising from a χ^2 -test with given confidence level, we offer multiple algorithms. One of these, PolyOpt, offers significantly higher utility, especially in the high privacy regime. However, it relies on vertex enumeration, making it computationally infeasible for larger input spaces. The other 3 algorithms, SR, GRR-CR and UE-CR, rely on processing the sensitive and non-sensitive data separately. These algorithms rely on low-dimensional optimisation, independent of the size of the input space, allowing these to be used when PolyOpt is outside the computational capabilities. Of these protocols, UE-CR is the best option when either \mathcal{F} or the input space is small. SR and GRR-CR perform similar in the high privacy regime, with GRR-CR performing better for input distributions with large probability.

Our results suggest several avenues for future research. First, one may want to incorporate not only robustness in privacy, but also in utility, i.e. to find the protocol \mathcal{Q} that maximises $\min_{P \in \mathcal{F}} I_P(X; Y)$. An obstacle for this is that $I_P(X; Y)$ is concave in P , which makes finding its minimum over \mathcal{F} difficult. Second, instead of looking at the situation where X splits into a sensitive part S and a non-sensitive part U , one can consider the more general case that X is correlated with the sensitive data S . This is already done in work on the privacy funnel, but this generally does not incorporate robustness. Furthermore, the utility of IR and CR might be improved in the high privacy regime by incorporating other LDP protocols than GRR. It is shown in [2] that GRR is the optimal LDP protocol for high ε , but for low ε the optimum typically takes a different form. One obstacle in incorporating this is that these optima depend on P^* , which is inaccessible in the RLDP framework.

ACKNOWLEDGEMENTS

This work was supported by NWO grant 628.001.026.

REFERENCES

- [1] D. Rebollo-Monedero, J. Forne, and J. Domingo-Ferrer, "From t-closeness-like privacy to postrandomization via information theory," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 11, pp. 1623–1636, 2010.
- [2] P. Kairouz, S. Oh, and P. Viswanath, "Extremal mechanisms for local differential privacy," *arXiv:1407.1338*, 2014.
- [3] A. Makhdoumi, S. Salamatian, N. Fawaz, and M. Médard, "From the information bottleneck to the privacy funnel," in *2014 IEEE Information Theory Workshop (ITW 2014)*, IEEE, 2014, pp. 501–505.
- [4] S. Salamatian, A. Zhang, F. du Pin Calmon, S. Bhamidipati, N. Fawaz, B. Kveton, P. Oliveira, and N. Taft, "Managing your private and public data: Bringing down inference attacks against your privacy.," *J. Sel. Topics Signal Processing*, vol. 9, no. 7, pp. 1240–1255, 2015.

- [5] S. Asoodeh, M. Diaz, F. Alajaji, and T. Linder, “Information extraction under privacy constraints,” *Information*, vol. 7, no. 1, p. 15, 2016.
- [6] S. Kung, “A compressive privacy approach to generalized information bottleneck and privacy funnel problems,” *Journal of the Franklin Institute*, vol. 355, no. 4, pp. 1846–1872, 2018.
- [7] N. Ding and P. Sadeghi, “A submodularity-based agglomerative clustering algorithm for the privacy funnel,” *arXiv:1901.06629*, 2019.
- [8] S. Salamatian, F. P. Calmon, N. Fawaz, A. Makhdoumi, and M. Médard, “Privacy-utility tradeoff and privacy funnel,” 2020, Preprint.
- [9] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, “Local privacy and statistical minimax rates,” in *Foundations of Computer Science (FOCS), 2013 IEEE 54th Annual Symposium on*, IEEE, 2013, pp. 429–438.
- [10] M. Lopuhaä-Zwakenberg, H. Tong, and B. Škorić, “Data sanitisation for the privacy funnel with differential privacy guarantees,” *arXiv:2008.13151*, 2020.
- [11] D. Kifer and A. Machanavajjhala, “Pufferfish: A framework for mathematical privacy definitions,” *ACM Transactions on Database Systems (TODS)*, vol. 39, no. 1, pp. 1–36, 2014.
- [12] S. L. Warner, “Randomized response: A survey technique for eliminating evasive answer bias,” *Journal of the American Statistical Association*, vol. 60, no. 309, pp. 63–69, 1965.
- [13] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski, *Robust optimization*. Princeton University Press, 2009, vol. 28.
- [14] A. Ben-Tal, D. Den Hertog, and J.-P. Vial, “Deriving robust counterparts of nonlinear uncertain inequalities,” *Mathematical programming*, vol. 149, no. 1-2, pp. 265–299, 2015.
- [15] D. Bertsimas, V. Gupta, and N. Kallus, “Data-driven robust optimization,” *Mathematical Programming*, vol. 167, no. 2, pp. 235–292, 2018.
- [16] L. Willenborg and T. De Waal, *Elements of statistical disclosure control*. Springer Science & Business Media, 2012, vol. 155.
- [17] A. Hundepool, J. Domingo-Ferrer, L. Franconi, S. Giessing, E. S. Nordholt, K. Spicer, and P.-P. De Wolf, *Statistical disclosure control*. John Wiley & Sons, 2012.
- [18] N. Tishby, F. C. Pereira, and W. Bialek, “The information bottleneck method,” *arXiv:physics/0004057*, 2000.
- [19] B. Rassouli and D. Gunduz, “On perfect privacy,” in *2018 IEEE International Symposium on Information Theory (ISIT)*, IEEE, 2018, pp. 2551–2555.
- [20] I. Issa, A. B. Wagner, and S. Kamath, “An operational approach to information leakage,” *IEEE Transactions on Information Theory*, vol. 66, no. 3, pp. 1625–1657, 2019.

- [21] J. Liao, O. Kosut, L. Sankar, and F. du Pin Calmon, “Tunable measures for information leakage and applications to privacy-utility tradeoffs,” *IEEE Transactions on Information Theory*, vol. 65, no. 12, pp. 8043–8066, 2019.
- [22] I. Wagner and D. Eckhoff, “Technical privacy metrics: A systematic survey,” *ACM Computing Surveys (CSUR)*, vol. 51, no. 3, pp. 1–38, 2018.
- [23] S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith, “What can we learn privately?” *SIAM Journal on Computing*, vol. 40, no. 3, pp. 793–826, 2011.
- [24] T. Wang, J. Blocki, N. Li, and S. Jha, “Locally differentially private protocols for frequency estimation,” in *26th {USENIX} Security Symposium ({USENIX} Security 17)*, 2017, pp. 729–745.
- [25] H. Wang, M. Diaz, F. P. Calmon, and L. Sankar, “The utility cost of robust privacy guarantees,” in *2018 IEEE International Symposium on Information Theory (ISIT)*, IEEE, 2018, pp. 706–710.
- [26] M. Diaz, H. Wang, F. P. Calmon, and L. Sankar, “On the robustness of information-theoretic privacy measures and mechanisms,” *IEEE Transactions on Information Theory*, vol. 66, no. 4, pp. 1949–1978, 2019.
- [27] Z. Wang, P. W. Glynn, and Y. Ye, “Likelihood robust optimization for data-driven problems,” *Computational Management Science*, vol. 13, no. 2, pp. 241–261, 2016.
- [28] J. Duchi, P. Glynn, and H. Namkoong, “Statistics of robust optimization: A generalized empirical likelihood approach,” *arXiv:1610.03425*, 2016.
- [29] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [30] C. Huang, P. Kairouz, X. Chen, L. Sankar, and R. Rajagopal, “Context-aware generative adversarial privacy,” *Entropy*, vol. 19, no. 12, p. 656, 2017.
- [31] A. Tripathy, Y. Wang, and P. Ishwar, “Privacy-preserving adversarial networks,” in *2019 57th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, IEEE, 2019, pp. 495–505.
- [32] V. Mirjalili, S. Raschka, A. Namboodiri, and A. Ross, “Semi-adversarial networks: Convolutional autoencoders for imparting privacy to face images,” in *2018 International Conference on Biometrics (ICB)*, IEEE, 2018, pp. 82–89.
- [33] B. Bortolato, M. Ivanovska, P. Rot, J. Križaj, P. Terhörst, N. Damer, P. Peer, and V. Štruc, “Learning privacy-enhancing face representations through feature disentanglement,” in *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)(FG)*, 2020, pp. 45–52.

- [34] J. I. Stoker, H. Garretsen, and L. J. Spreeuwiers, “The facial appearance of ceos: Faces signal selection but not performance,” *PloS one*, vol. 11, no. 7, e0159950, 2016.
- [35] D. Dua and C. Graff, *UCI machine learning repository*, 2017. [Online]. Available: <http://archive.ics.uci.edu/ml/d>

APPENDIX A

PROOFS

A. Proof of Theorem 7

We follow the proof of Theorem 14 in [2]. For $C \in \mathbb{R}_{\geq 0}^{\mathcal{X}}$, define

$$\mu(C) = \sum_x P_x C_x \log \frac{C_x}{\sum_{x'} P_{x'} C_{x'}}. \quad (103)$$

Then the utility of a protocol $\mathcal{Q}: \mathcal{X} \rightarrow \mathcal{Y}$ is given by $I_P(X; Y) = \sum_y \mu(Q_{y|\bullet})$. Furthermore, μ is a sublinear function in the sense of [2, Definition 1].

We fix an $\varepsilon > 0$. Furthermore, let $\mathcal{C} \subset \mathbb{R}_{\geq 0}^{\mathcal{X}}$ be the positive cone defined by the inequalities of the following form, for $s, s' \in \mathcal{S}$ with $s \neq s'$ and $u, u' \in \mathcal{U}$:

$$C_{s,u} \leq e^\varepsilon C_{s',u'}. \quad (104)$$

Then a protocol \mathcal{Q} satisfies ε -SLDP if and only if each $Q_{y|\bullet}$ is an element of \mathcal{C} . Furthermore, \mathcal{C} is spanned (as a cone) by the set

$$\mathcal{V} = \left\{ v \in \mathbb{R}^{\mathcal{X}} : \begin{array}{l} \forall x: v_x \in \{1, e^\varepsilon\}, \\ |\{s: \exists u \text{ s.t. } v_{s,u} = e^\varepsilon\}| \geq 2 \end{array} \right\} \cup \left\{ v \in \mathbb{R}^{\mathcal{X}} : \exists s \text{ s.t. } \begin{array}{l} \forall u: v_{s,u} \in \{e^{-\varepsilon}, e^\varepsilon\}; \\ \forall s' \neq s, \forall u: v_{s',u} = 1 \end{array} \right\}. \quad (105)$$

Let \mathcal{D} be the polytope spanned by \mathcal{V} . If \mathcal{Q} satisfies ε -SLDP, then every column $Q_{y|\bullet}$ is of the form $\theta_y \cdot d_y$, where $d_y \in \mathcal{D}$ and $\theta_y \in \mathbb{R}_{\geq 0}$ are such that $\sum_y \theta_y d_y = 1_{\mathcal{X}}$. Analogous to the proof of Theorems 2 and 4 in [2, Section 7], one proves that the optimal \mathcal{Q} is found by taking $b = a$, and taking $d_y \in \mathcal{V}$ for all d . Since

$$I(X; Y) = \sum_y \mu(Q_{y|\bullet}) = \sum_y \theta_y \mu(d_y) \quad (106)$$

we can find the optimal \mathcal{Q} by solving the following optimisation problem, where m is the vector $(\mu(v))_{v \in \mathcal{V}}$, and where $A \in \mathbb{R}^{\mathcal{X} \times \mathcal{V}}$ is the matrix whose v -th column is v :

$$\begin{aligned} & \text{maximise}_{\theta \in \mathbb{R}^{\mathcal{V}}} m \cdot \theta \\ & \text{such that } A \cdot \theta = 1_{\mathcal{X}}, \\ & \theta \geq 0. \end{aligned}$$

From here, we follow [2, Section 9.5]. The dual to the above problem is

$$\begin{aligned} & \text{minimise}_{\alpha \in \mathbb{R}^{\mathcal{X}}} (1_{\mathcal{X}}) \cdot \alpha \\ & \text{such that } A^T \cdot \alpha \geq m, \\ & \alpha \geq 0. \end{aligned}$$

By duality we have $\max_{\theta} m \cdot \theta = \min_{\alpha} (1_{\mathcal{X}}) \cdot \alpha$. We describe α^* and θ^* , depending on ε , such that for ε large enough one has $A^{\top} \cdot \alpha^* \geq m$, such that $m \cdot \theta^* = (1_{\mathcal{X}}) \cdot \alpha^*$ and $A\theta^* = 1_{\mathcal{X}}$, and such that θ^* corresponds to SGRR, i.e. for each $y \in \mathcal{Y} = \mathcal{X}$ there is a $\hat{v}_y \in \mathcal{V}$ such that $Q_{y|\bullet} = \theta_y^* v_y$. Together, this proves that SGRR is optimal for $\varepsilon \gg 0$.

More concretely, for $y = (s, u) \in \mathcal{X}$, define \hat{v}_y by

$$(\hat{v}_y)_{s',u'} = \begin{cases} e^{\varepsilon}, & \text{if } (s', u') = (s, u), \\ e^{-\varepsilon}, & \text{if } s' = s \text{ and } u' \neq u, \\ 1, & \text{if } s' \neq s, \end{cases} \quad (107)$$

and let $\theta^* \in \mathbb{R}^{\mathcal{V}}$ be given by

$$\theta_v^* = \begin{cases} \frac{1}{e^{\varepsilon} + e^{-\varepsilon}(a_2 - 1) + a - a_2}, & \text{if there is a } y \in \mathcal{X} \text{ such that } v = \hat{v}_y, \\ 0, & \text{otherwise;} \end{cases} \quad (108)$$

Then SRR satisfies $Q_{y|\bullet} = \theta_{\hat{v}_y}^* \hat{v}_y$ for all $y \in \mathcal{X}$, and also

$$(A\theta^*)_x = \sum_v A_{x,v} \theta_v^* \quad (109)$$

$$= \sum_v v_x \theta_v^* \quad (110)$$

$$= \frac{\sum_y (\hat{v}_y)_x}{e^{\varepsilon} + e^{-\varepsilon}(a_2 - 1) + a - a_2} \quad (111)$$

$$= 1, \quad (112)$$

which shows that $A\theta^* = 1_{\mathcal{X}}$. Furthermore, define $\alpha^* \in \mathbb{R}^{\mathcal{X}}$ by

$$\alpha_{s,u}^* = c_1 \mu(\hat{v}_{s,u}) + c_2 \sum_{u' \neq u} \mu(\hat{v}_{s,u'}) + c_3 \sum_{\substack{s' \neq s, \\ u'}} \mu(\hat{v}_{s',u'}), \quad (113)$$

where

$$c_1 = \frac{-(a_2 - 2)(a_2 - 1) + (a - a_2 + 1)(a_2 - 2)e^{\varepsilon} + (a - 2a_2 + 1)e^{2\varepsilon} + e^{3\varepsilon}}{(e^{\varepsilon} - 1)(e^{\varepsilon} + 1)(e^{\varepsilon} - a_2 + 1)(e^{\varepsilon} + (a_2 - 1)e^{-\varepsilon} + a - a_2)}, \quad (114)$$

$$c_2 = \frac{a_2 - 1 + (a - a_2 + 1)e^{\varepsilon}}{(e^{\varepsilon} - 1)(e^{\varepsilon} + 1)(e^{\varepsilon} - a_2 + 1)(e^{\varepsilon} + (a_2 - 1)e^{-\varepsilon} + a - a_2)}, \quad (115)$$

$$c_3 = \frac{-e^{2\varepsilon}}{(e^{\varepsilon} - 1)(e^{\varepsilon} - a_2 + 1)(e^{\varepsilon} + (a_2 - 1)e^{-\varepsilon} + a - a_2)}. \quad (116)$$

One readily calculates that for all x we have

$$m \cdot \theta^* = (1_{\mathcal{X}}) \cdot \alpha^* = \frac{1}{e^{\varepsilon} + e^{-\varepsilon}(a_2 - 1) + a - a_2} \sum_x \mu(\hat{v}_x), \quad (117)$$

$$\hat{v}_x \cdot \alpha^* = m_{\hat{v}_x} = \mu(\hat{v}_x). \quad (118)$$

It remains to be shown that α^* satisfies the dual problem for $\varepsilon \gg 0$, i.e. $A^T \alpha \geq m$ for ε large enough. To this end, for $v \in \mathcal{V}$, set

$$F_v = \{x \in \mathcal{X} : v_x = e^\varepsilon\}, \quad (119)$$

$$G_v = \{x \in \mathcal{X} : v_x = 1\}, \quad (120)$$

$$H_v = \{x \in \mathcal{X} : v_x = e^{-\varepsilon}\}, \quad (121)$$

Then $\#F_v \geq 1$ for all v , and $\#F_v = 1$ if and only if there exist s, u such that $v = \hat{v}_{s,u}$. We write $P_{F_v} = \sum_{x \in F_v} P_x$ and likewise for G_v, H_v . For large ε we have

$$m_v = \mu(v) = e^\varepsilon \sum_{x \in F_v} P_x \log \frac{1}{P_{F_v} + e^{-\varepsilon} P_{G_v} + e^{-2\varepsilon} P_{H_v}} \quad (122)$$

$$+ \sum_{x \in G_v} P_x \log \frac{1}{e^\varepsilon P_{F_v} + P_{G_v} + e^{-\varepsilon} P_{H_v}} \quad (123)$$

$$+ e^{-\varepsilon} \sum_{x \in H_v} P_x \log \frac{1}{e^{2\varepsilon} P_{F_v} + e^\varepsilon P_{G_v} + P_{H_v}} \quad (124)$$

$$= (-P_{F_v} \log P_{F_v}) e^\varepsilon + \mathcal{O}(\varepsilon) \quad (125)$$

and furthermore

$$c_1 = e^{-\varepsilon} + \mathcal{O}(e^{-2\varepsilon}), \quad (126)$$

$$c_2, c_3 = \mathcal{O}(e^{-2\varepsilon}), \quad (127)$$

$$\alpha_x^* = c_1 \mu(\hat{v}_x) + (c_2 + c_3) \mathcal{O}(e^\varepsilon) \quad (128)$$

$$= -P_x \log P_x + \mathcal{O}(\varepsilon e^{-\varepsilon}), \quad (129)$$

$$v^T \alpha^* = \left(- \sum_{x \in F_v} P_x \log P_x \right) e^\varepsilon + \mathcal{O}(\varepsilon). \quad (130)$$

For $|F_v| \geq 2$ one has $P_{F_v} \log P_{F_v} > \sum_{x \in F_v} P_x \log P_x$. This means that if v is not of the form \hat{v}_x , one has $v^T \alpha^* \geq m_v$ for ε large enough. Together with (118) this shows that $A^T \alpha^* \geq m$ for ε large enough; this concludes the proof.

B. Proof of Proposition 10

Before we can prove this proposition, we need the following auxiliary lemma.

Lemma 21. *Let $B \in \mathbb{R}_{\geq 1}$. Then the function $g: [0, 1] \rightarrow \mathbb{R}$ given by*

$$g(x) = \frac{B(1-2x) + \sqrt{B(B+4x(1-x))}}{B+1} \quad (131)$$

is nonincreasing.

Proof. Since $B \geq 1$ we have for all $x \in [0, 1]$ that $B^2 - 1 + 4(B + 1)x(1 - x) \geq 0$. Rearranging terms, it follows that

$$B(B + 4x(1 - x)) \geq 1 - 4x + 4x^2, \quad (132)$$

hence $\sqrt{B(B + 4x(1 - x))} \geq 1 - 2x$. Using this, one calculates

$$g'(x) = \frac{2B \left(1 - 2x - \sqrt{B(B + 4x(1 - x))}\right)}{(B + 1)\sqrt{B(B + 4x(1 - x))}} \leq 0, \quad (133)$$

hence g is nonincreasing. \square

Proof of Proposition 10. The distribution P maximising $\|P - \hat{P}\|_1$ is located on the boundary of \mathcal{F} , hence $\sum_x \frac{\hat{P}_x^2}{P_x} = B + 1$. We define sets

$$\mathcal{X}_1 := \left\{x \in \mathcal{X} : P_x \geq \hat{P}_x > 0\right\}, \quad (134)$$

$$\mathcal{X}_2 := \left\{x \in \mathcal{X} : P_x < \hat{P}_x\right\}, \quad (135)$$

$$\mathcal{X}_3 := \left\{x \in \mathcal{X} : \hat{P}_x = 0\right\}. \quad (136)$$

Note that \mathcal{X}_2 and $\mathcal{X}_1 \cup \mathcal{X}_3$ both have to be nonempty. Then

$$\|P - \hat{P}\|_1 = \sum_{x \in \mathcal{X}_1} (P_x - \hat{P}_x) + \sum_{x \in \mathcal{X}_2} (\hat{P}_x - P_x) + \sum_{x \in \mathcal{X}_3} P_x. \quad (137)$$

We can find the P maximising this, subject to the constraints $\sum_x \frac{\hat{P}_x^2}{P_x} = B + 1$ and $\sum_x P_x = 1$, by finding critical points of the Lagrange multiplier expression

$$\sum_{x \in \mathcal{X}_1} (P_x - \hat{P}_x) + \sum_{x \in \mathcal{X}_2} (\hat{P}_x - P_x) + \sum_{x \in \mathcal{X}_3} P_x + \lambda \left(\sum_x \frac{\hat{P}_x^2}{P_x} - B - 1 \right) + \mu \left(\sum_x P_x - 1 \right). \quad (138)$$

Differentiating this with respect to P_x for $x \in \mathcal{X}_1, \mathcal{X}_2, \mathcal{X}_3$, respectively, we get

$$\forall x \in \mathcal{X}_1 : 1 - \frac{\lambda \hat{P}_x^2}{P_x^2} + \mu = 0, \quad (139)$$

$$\forall x \in \mathcal{X}_2 : -1 - \frac{\lambda \hat{P}_x^2}{P_x^2} + \mu = 0, \quad (140)$$

$$\forall x \in \mathcal{X}_3 : 1 + \mu = 0. \quad (141)$$

If $\mathcal{X}_1, \mathcal{X}_2, \mathcal{X}_3$ are all nonempty, then (141) implies $\mu = -1$, so from (139) we get $\lambda = 0$; however, then (140) leads to a contradiction. Hence either \mathcal{X}_1 or \mathcal{X}_3 is empty; we will discuss these cases separately.

Suppose $\mathcal{X}_1 = \emptyset$; then $\sum_{x \in \mathcal{X}_2} \hat{P}_x = 1$. From (140) and (141) we find $P_x = \sqrt{-\frac{\lambda}{2}} \hat{P}_x$ for all $x \in \mathcal{X}_2$. Writing $c_- := \sqrt{-\frac{\lambda}{2}}$, we know that c_- should satisfy

$$\frac{1}{c_-} = \sum_{x \in \mathcal{X}_2} \frac{\hat{P}_x}{c_-} = \sum_{x \in \mathcal{X}_2} \frac{\hat{P}_x^2}{P_x} = B + 1. \quad (142)$$

Hence $c_- = \frac{1}{B+1}$, and

$$\|P - \hat{P}\|_1 = 2 \sum_{x \in \mathcal{X}_2} (\hat{P}_x - P_x) \quad (143)$$

$$= 2(1 - c_-) \sum_{x \in \mathcal{X}_2} \hat{P}_x \quad (144)$$

$$= \frac{2B}{B+1}, \quad (145)$$

which is indeed the formula in (43) when $P_{x_{\min}} = 0$. Furthermore, by the AM-GM inequality we have $\frac{2B}{B+1} = \sqrt{B} \frac{2}{\sqrt{B} + \frac{1}{\sqrt{B}}} \leq \sqrt{B}$, which shows that $\|P - \hat{P}\|_1 \leq \sqrt{B}$ in general, and in particular for $B \leq 1$.

Now suppose $\mathcal{X}_3 = \emptyset$. In that case we find from (139) that for $x \in \mathcal{X}_1$ we have $P_x = \sqrt{\frac{\lambda}{\mu+1}} \hat{P}_x$, while for $x \in \mathcal{X}_2$ we have $P_x = \sqrt{\frac{\lambda}{\mu-1}} \hat{P}_x$. Setting $c_+ := \sqrt{\frac{\lambda}{\mu+1}}$ and $c_- := \sqrt{\frac{\lambda}{\mu-1}}$, then

$$\hat{P}_{\mathcal{X}_1} c_+ + (1 - \hat{P}_{\mathcal{X}_1}) c_- = \sum_{x \in \mathcal{X}_1} c_+ \hat{P}_x + \sum_{x \in \mathcal{X}_2} c_- \hat{P}_x \quad (146)$$

$$= 1, \quad (147)$$

$$\frac{\hat{P}_{\mathcal{X}_1}}{c_+} + \frac{1 - \hat{P}_{\mathcal{X}_1}}{c_-} = \sum_{x \in \mathcal{X}_1} \frac{\hat{P}_x}{c_+} + \sum_{x \in \mathcal{X}_2} \frac{\hat{P}_x}{c_-} \quad (148)$$

$$= \sum_{x \in \mathcal{X}_1} \frac{\hat{P}_x^2}{P_x} + \sum_{x \in \mathcal{X}_2} \frac{\hat{P}_x^2}{P_x} \quad (149)$$

$$= B + 1. \quad (150)$$

Jointly solving (147) and (150) we find

$$c_+ = \frac{B + 2\hat{P}_{\mathcal{X}_1} \pm \sqrt{B^2 + 4B\hat{P}_{\mathcal{X}_1} - 4B\hat{P}_{\mathcal{X}_1}^2}}{2(B+1)\hat{P}_{\mathcal{X}_1}}, \quad (151)$$

$$c_- = \frac{B + 2(1 - \hat{P}_{\mathcal{X}_1}) \mp \sqrt{B^2 + 4B\hat{P}_{\mathcal{X}_1} - 4B\hat{P}_{\mathcal{X}_1}^2}}{2(B+1)(1 - \hat{P}_{\mathcal{X}_1})}. \quad (152)$$

By definition of \mathcal{X}_1 and \mathcal{X}_2 we know that $c_+ \geq 1$ and $c_- < 1$; this only occurs if c_+ is the "+" solution while c_- is the "-" solution. It follows that

$$\|P - \hat{P}\|_1 = \hat{P}_{\mathcal{X}_1}(c_+ - 1) + (1 - \hat{P}_{\mathcal{X}_1})(1 - c_-) \quad (153)$$

$$= \frac{B - 2B\hat{P}_{\mathcal{X}_1} + \sqrt{B^2 + 4B\hat{P}_{\mathcal{X}_1} - 4B\hat{P}_{\mathcal{X}_1}^2}}{B+1}. \quad (154)$$

It follows that the P maximising $\|P - \hat{P}\|_1$ is obtained by finding the (nonempty) subset $\mathcal{X}_1 \subset \mathcal{X}$ that maximises (154). By Lemma 21, this is when $\mathcal{X}_1 = \{x_{\min}\}$ if $B \geq 1$, proving (43). If $B < 1$, the optimal \mathcal{X}_1 is unfortunately harder to determine. However, we can still find an upper bound, by

finding the value of $\hat{P}_{\mathcal{X}_1}$ that maximises (154). We find the maximum by taking the derivative with respect to $\hat{P}_{\mathcal{X}_1}$, and we have to solve

$$\frac{-2B}{B+1} + \frac{2B - 4B\hat{P}_{\mathcal{X}_1}}{(B+1)\sqrt{B^2 + 4B\hat{P}_{\mathcal{X}_1} - 4B\hat{P}_{\mathcal{X}_1}^2}} = 0, \quad (155)$$

which leads to $\hat{P}_{\mathcal{X}_1} = \frac{1-\sqrt{B}}{2}$. Substituting this in (154), we find

$$\|P - \hat{P}_1\| \leq \frac{B - B(1 - \sqrt{B}) + \sqrt{B^2 + 2B(1 - \sqrt{B}) - B(1 - \sqrt{B})^2}}{B + 1} \quad (156)$$

$$= \sqrt{B}. \square \quad (157)$$

C. Proof of Theorem 12

This is essentially analogous to the proof of Theorem 4 in [2]; the main difference is that the equivalent of $\hat{\Gamma}$ is a hypercube, so there a vertex enumeration step is not needed. Let \mathcal{Q} be a protocol such that $Q_y \in \Gamma$ for all y ; then there exist $\alpha_y \in \mathbb{R}_{\geq 0}$, $\gamma_y \in \hat{\Gamma}$ such that $Q_y = \alpha_y \gamma_y$. One has

$$I_{\hat{P}}(X; Y) = \sum_y \mu^1(Q_y) = \sum_y \alpha_y \mu^1(\gamma_y). \quad (158)$$

Since $\hat{\Gamma}$ is the convex hull of \mathcal{V} , we can write $\gamma_y = \sum_v \lambda_{y,v} v$ for suitable constants $\lambda_{y,v}$. Define $\theta \in \mathbb{R}_{\geq 0}^{\mathcal{V}}$ by $\theta_v = \sum_y \lambda_{y,v} \alpha_y$. Then

$$\sum_v \theta_v v = \sum_y Q_y = 1_{\mathcal{X}}. \quad (159)$$

As such, the matrix $Q' \in \mathbb{R}^{\mathcal{V} \times \mathcal{X}}$ defined by $Q'_v = \theta_v v$ defines a privacy protocol \mathcal{Q}' . One has

$$I_{\hat{P}}(X; \mathcal{Q}'(X)) = \sum_v \mu^1(Q'_v) \quad (160)$$

$$= \sum_v \theta_v \mu^1(v) \quad (161)$$

$$= \sum_y \alpha_y \sum_v \lambda_{y,v} \mu^1(v) \quad (162)$$

$$\geq \sum_y \alpha_y \mu^1\left(\sum_v \lambda_{y,v} v\right) \quad (163)$$

$$= I_{\hat{P}}(X; \mathcal{Q}(X)), \quad (164)$$

where we use the fact that μ^1 is convex. This shows that the Q_y of the optimal protocol satisfying Theorem 11 are all of the form $\theta_v \cdot v$; hence (47) yields the optimal protocol. For \mathcal{Q}^2 , note that

$$\inf_P (X; \mathcal{Q}^2(X)) = \inf_P \sum_v \hat{\theta}_v^2 \sum_x v_x P_x \log \frac{v_x}{\sum_{x'} v_{x'} P_{x'}} \quad (165)$$

$$\geq \sum_v \hat{\theta}_v^2 \inf_P \sum_x v_x P_x \log \frac{v_x}{\sum_{x'} v_{x'} P_{x'}} \quad (166)$$

$$= \sum_v \hat{\theta}_v^2 \mu^2(v). \quad (167)$$