



**HAL**  
open science

## On the Risks of Collecting Multidimensional Data Under Local Differential Privacy

Héber H. Arcolezi, Sébastien Gambs, Jean-François Couchot, Catuscia  
Palamidessi

► **To cite this version:**

Héber H. Arcolezi, Sébastien Gambs, Jean-François Couchot, Catuscia Palamidessi. On the Risks of Collecting Multidimensional Data Under Local Differential Privacy. Proceedings of the VLDB Endowment (PVLDB), 2023, 16 (5), pp.1126 - 1139. 10.14778/3579075.3579086 . hal-04082592

**HAL Id: hal-04082592**

**<https://inria.hal.science/hal-04082592v1>**

Submitted on 26 Apr 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License



# On the Risks of Collecting Multidimensional Data Under Local Differential Privacy

Héber H. Arcolezi

Inria and École Polytechnique (IPP)  
heber.hwang-arcolezi@inria.fr

Jean-François Couchot

Femto-ST Institute, Univ. Bourg. Franche-Comté, CNRS  
jean-francois.couchot@univ-fcomte.fr

Sébastien Gambs

Université du Québec à Montréal, UQAM  
gambs.sebastien@uqam.ca

Catuscia Palamidessi

Inria and École Polytechnique (IPP)  
catuscia@lix.polytechnique.fr

## ABSTRACT

The private collection of multiple statistics from a population is a fundamental statistical problem. One possible approach to realize this is to rely on the local model of differential privacy (LDP). Numerous LDP protocols have been developed for the task of frequency estimation of single and multiple attributes. These studies mainly focused on improving the utility of the algorithms to ensure the server performs the estimations accurately. In this paper, we investigate privacy threats (re-identification and attribute inference attacks) against LDP protocols for multidimensional data following two state-of-the-art solutions for frequency estimation of multiple attributes. To broaden the scope of our study, we have also experimentally assessed five widely used LDP protocols, namely, generalized randomized response, optimal local hashing, subset selection, RAPPOR and optimal unary encoding. Finally, we also proposed a countermeasure that improves both utility and robustness against the identified threats. Our contributions can help practitioners aiming to collect users' statistics privately to decide which LDP mechanism best fits their needs.

### PVLDB Reference Format:

Héber H. Arcolezi, Sébastien Gambs, Jean-François Couchot, and Catuscia Palamidessi. On the Risks of Collecting Multidimensional Data Under Local Differential Privacy. PVLDB, 16(5): 1126 - 1139, 2023. doi:10.14778/3579075.3579086

### PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/hharcolezi/risks-ldp>.

## 1 INTRODUCTION

Private and public organizations regularly collect and analyze digital data about their collaborators, volunteers, clients, etc. However, due to the sensitive nature of this personal data, the collection of users' raw data on a centralized server should be avoided. The distributed version of Differential Privacy (DP) [20–22], known as Local DP (LDP) [19, 29], aims to address such a challenge. Indeed, using an LDP mechanism, a user can sanitize her profile locally

before transmitting it to the server, which leads to strong privacy protection even if the server used for the aggregation is malicious. The LDP model has a close connection with the concept of randomized response [54], which provides “plausible deniability” to users' reports. For this reason, LDP has been already implemented in large-scale systems by Google [24], Microsoft [15] and Apple [46].

A fundamental task under LDP guarantees is frequency estimation [13, 15, 24, 27, 28, 46, 50, 51], in which the data collector estimates the number of users for each possible value of one attribute based on the sanitized data of the users. More recently, a new line of research started investigating *security* [7, 11, 33, 55] and *privacy* [9, 23, 25, 36] threats to LDP protocols (mainly for frequency estimation), which are discussed in detail in Section 7.

In this paper, we further investigate the **privacy threats for the users** when the server aims to perform *frequency estimation of multiple attributes* under LDP guarantees. In this setting [4, 38, 47, 48, 51], the profile of each user is characterized by  $d$  attributes  $\mathcal{A} = \{A_1, A_2, \dots, A_d\}$ , in which each attribute  $A_j$  has a discrete domain of size  $k_j = |A_j|$ , for  $j \in [d]$ . There are  $n$  users  $\mathcal{U} = \{u_1, \dots, u_n\}$ , and each user  $u_i$ , for  $i \in [n]$ , holds a private tuple  $\mathbf{v}^{(i)} = [v_1^{(i)}, v_2^{(i)}, \dots, v_d^{(i)}]$ , in which  $v_j^{(i)}$  represents the value of attribute  $A_j$  in record  $\mathbf{v}^{(i)}$ . Thus, for each attribute  $A_j \in \mathcal{A}$ , for  $j \in [d]$ , the aggregator's goal is to estimate a  $k_j$ -bins histogram.

To the best of our knowledge, for the task considered<sup>1</sup>, there are mainly three solutions for satisfying LDP by randomizing the user's tuple  $\mathbf{v} = [v_1, v_2, \dots, v_d]^2$ , which are described in the following:

- **Splitting (SPL)**. This naive solution directly splits the privacy budget  $\epsilon$  by  $d$  attributes and reports all attributes with  $\frac{\epsilon}{d}$ -LDP, thus incurring a high estimation error [5, 38, 48, 51].
- **Sampling (SMP)**. Instead of splitting the privacy budget, one state-of-the-art solution allows users to randomly sample a single attribute and report it with  $\epsilon$ -LDP [5, 38, 48, 51].
- **Random Sampling Plus Fake Data (RS+FD)** [4]. One of the weakness of the SMP solution is that it discloses the sampled attribute, which might not be fair to all users (e.g., some users will sample age but others will sample sensitive attribute such as disease). The objective of the state-of-the-art RS+FD solution is precisely to enable users to “hide” the sampled attribute (i.e.,  $\epsilon$ -LDP value) by also generating one uniformly random fake data for each non-sampled attribute. Thus, RS+FD creates *uncertainty* on the server-side.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 16, No. 5 ISSN 2150-8097. doi:10.14778/3579075.3579086

<sup>1</sup>This is a different task of joint distribution estimation under LDP guarantees [39, 57].

<sup>2</sup>For simplicity, we omit the index notation  $\mathbf{v}^{(i)}$  and focus on one arbitrary user  $u_i$ .

Focusing on the state-of-the-art solutions SMP and RS+FD, first, we empirically demonstrate through extensive experiments that **the SMP solution is vulnerable to re-identification attacks** when collecting users’ multidimensional data several times with  $\epsilon$  values commonly used by industry nowadays [14, 40]. For instance, assume a user has multiple mobile applications each surveying the user with the SMP solution on different attributes. Another possible scenario is the situation in which the same mobile application is used on a regular basis but surveys users with different attributes. This enables the user to sample a (possibly different) attribute each time, thus resulting in sending their sampled attribute along with their  $\epsilon$ -LDP report. Nevertheless, we show that an adversary who can see every tuple containing  $\langle$ sampled attribute,  $\epsilon$ -LDP report $\rangle$  can construct a partial or complete profile of the user, which can possibly be unique (or in a small anonymity set of  $k$  individuals) in the population considered. Therefore, once the set of  $k$  individuals (referred to as top- $k$  in this paper) is characterized, one can leverage well-known attacks (e.g., homogeneity) [12, 31, 34, 41–44].

More specifically, to attack the SMP solution, our adversarial analysis focuses on the reduced “**plausible deniability**” [17, 54] of using the whole privacy budget  $\epsilon$  to report a single attribute out of  $d$  ones. In this setting, the adversary has a higher chance to infer the users’ true value for each data collection performed. Consequently, in multiple data collections, the adversary can build partial or even sometimes complete profiles of each user, then using it to perform a re-identification attack. However, this depends on the LDP protocol being used as the encoding and randomization vary across them [13, 51]. In our experiments, we have assessed five widely used LDP protocols for frequency estimation (a.k.a. frequency oracle protocols [52, 53]), namely Generalized Randomized Response (GRR) [27, 28], Optimal Local Hashing (OLH) [51], Subset Selection (SS) [49, 56] and two Unary Encoding (UE) protocols (Basic One-time RAPPOR [24] and Optimal UE [51]). To assess the risks of re-identification we have also considered two privacy models, usual LDP and the relaxed version of LDP developed in [36] (the latter in Appendix C of [6]) for measuring re-identification risks.

Secondly, we observe that **since the RS+FD solution generates fake data uniformly at random in [4, 47], it is possible to uncover the sampled attribute of users** in certain conditions. In this context, we evaluated the effectiveness of the RS+FD solution in hiding the sampled attribute to the aggregator by varying the privacy budget  $\epsilon$ , the LDP protocol and the fake data generation procedure. In particular, if the aggregator is able to break RS+FD into the SMP solution, the RS+FD solution might also be subject to the same vulnerability to re-identification attacks on multiple collections. Thus, we have proposed three attack models to uncover the sampled attribute of users using the RS+FD solution and evaluated its risks to re-identification attacks. Lastly, as shown in our results, RS+FD is, to some extent, a natural countermeasure to re-identification attacks due to chaining errors from incorrectly predicting the sampled attribute and user’s value in multiple collections. Building on this, we have designed a stronger countermeasure that adapts RS+FD to generate fake data following non-uniform distributions, **almost fully preventing the inference of the sampled attribute while preserving utility**.

To summarize, this paper makes the following contributions:

- We investigate privacy threats against LDP protocols for multidimensional data following two state-of-the-art solutions for frequency estimation of multiple attributes, SMP [5, 38, 48, 51] and RS+FD [4], providing insightful adversarial analysis to help in LDP protocol selection.
- We demonstrate through extensive experiments that the SMP solution is vulnerable to re-identification attacks due to the disclosure of the sampled attribute and lower “plausible deniability” [17, 54] when using the whole privacy budget to report a single attribute.
- We propose three attack models to predict the sampled attribute of users when collecting multidimensional data with the RS+FD solution with about a 2-20 fold increment over a random guess baseline model.
- We show through empirical results that the RS+FD solution can prevent (to some extent) re-identification attacks
- Finally, we present an adaptation of the RS+FD solution that can serve as a countermeasure to the identified privacy threats while improving both privacy and utility.

**Outline.** In Section 2, we review the LDP privacy model, the LDP protocols and solutions for collecting multidimensional data investigated in this paper. Afterwards, in Section 3, we present the system overview and adversarial setting for both SMP and RS+FD solutions. In Section 4, we present our experimental evaluation and analyze our results before in Section 5 presenting an improvement of the RS+FD as a countermeasure. Next, we provide a general discussion in Section 6. Finally in Section 7, we review related work before concluding with future perspectives of this work in Section 8.

## 2 PRELIMINARIES

This section briefly reviews the LDP model, state-of-the-art LDP frequency estimation protocols and three solutions for multiple attribute frequency estimation under LDP.

### 2.1 Local Differential Privacy

In this paper, we use LDP (Local Differential Privacy) [19, 29] as the privacy model considered, which is formalized as:

**DEFINITION 1 ( $\epsilon$ -LOCAL DIFFERENTIAL PRIVACY).** *A randomized algorithm  $\mathcal{M}$  satisfies  $\epsilon$ -local-differential-privacy ( $\epsilon$ -LDP), where  $\epsilon > 0$ , if for any pair of input values  $v_1, v_2 \in \text{Domain}(\mathcal{M})$  and any possible output  $y$  of  $\mathcal{M}$ :*

$$\Pr[\mathcal{M}(v_1) = y] \leq e^\epsilon \cdot \Pr[\mathcal{M}(v_2) = y]. \quad (1)$$

In essence, LDP guarantees that it is unlikely for the data aggregator to reconstruct the data source regardless of the prior knowledge. The privacy budget  $\epsilon$  controls the privacy-utility trade-off for which lower values of  $\epsilon$  result in tighter privacy protection. Similar to central DP, LDP also has several important properties, such as immunity to post-processing and composability [22].

### 2.2 LDP Frequency Estimation Protocols

In this subsection, we review five state-of-the-art LDP protocols, which enables the aggregator to estimate the frequency of any value  $v_i \in A_j$ , for  $i \in [k_j]$ , under LDP guarantees.

**2.2.1 Generalized Randomized Response.** Randomized response (RR) [54] is the classical technique for achieving LDP, which provides “plausible deniability” for individuals responding to embarrassing (binary) questions in a survey. The Generalized RR (GRR) [27, 28] protocol extends RR to the case of  $k_j \geq 2$  while satisfying  $\epsilon$ -LDP. Given a value  $v_i \in A_j$ , for  $i \in [k_j]$ ,  $GRR(v_i)$  outputs the true value with probability  $p$ , and any other value  $v \in A_j \setminus \{v_i\}$  with probability  $1 - p$ . More formally, the perturbation function is:

$$\forall y \in A_j : \Pr[y = a] = \begin{cases} p = \frac{e^\epsilon}{e^\epsilon + k_j - 1}, & \text{if } a = v \\ q = \frac{1}{e^\epsilon + k_j - 1}, & \text{otherwise,} \end{cases}$$

in which  $y$  is the perturbed value sent to the aggregator. The GRR protocol satisfy  $\epsilon$ -LDP since  $\frac{p}{q} = e^\epsilon$ . To estimate the normalized frequency of  $v_i \in A_j$ , for  $i \in [k_j]$ , one counts how many times  $v_i$  is reported, expressed as  $C(v_i)$ , and then computes [51]:

$$\hat{f}(v_i) = \frac{C(v_i) - nq}{n(p - q)}, \quad (2)$$

in which  $n$  is the total number of users. In [51], it was proven that Eq. (2) is an unbiased estimator (i.e.,  $\mathbb{E}(\hat{f}(v_i)) = f(v_i)$ ).

**2.2.2 Optimal Local Hashing.** Local hashing (LH) protocols can handle a large domain size  $k_j$  by first using hash functions to map an input value to a smaller domain of size  $g_j$  (typically  $g_j \ll k_j$ ), and then applying GRR to the hashed value in the smaller domain.

The authors in [51] have proposed Optimal LH (OLH), which selects  $g_j = e^\epsilon + 1$ . Given a value  $v_i \in A_j$ , for  $i \in [k_j]$ , in OLH, one reports  $\langle H, GRR(H(v_i)) \rangle$  in which  $H$  is randomly chosen from a family of universal hash functions that hash each value in  $A_j$  to  $[g_j] = \{1, \dots, g_j\}$ , which is the domain that  $GRR(\cdot)$  will operate on. The hash values will remain unchanged with probability  $p'$  and switch to a different value in  $[g_j]$  with probability  $q'$ , as:

$$\forall y \in [g_j] : \Pr[y = (H, a)] = \begin{cases} p' = \frac{e^\epsilon}{e^\epsilon + g_j - 1}, & \text{if } a = H(v) \\ q' = \frac{1}{e^\epsilon + g_j - 1}, & \text{otherwise,} \end{cases}$$

in which  $y$  is the hash function and perturbed value sent to the aggregator. From this, the aggregator can obtain the unbiased estimation of  $v_i \in A_j$ , for  $i \in [k_j]$ , with Eq. (2) by setting  $p = p'$  and  $q = \frac{1}{g_j} \cdot p' + \left(1 - \frac{1}{g_j}\right) \cdot q' = \frac{1}{g_j}$  [51].

**2.2.3 Subset Selection.** The main idea of  $\omega$ -Subset Selection ( $\omega$ -SS) [49, 56] is to randomly select  $\omega$  items within the input domain to report a subset of values (i.e.,  $\Omega \subseteq A_j$ ). The user’s true value  $v_i \in A_j$ , for  $i \in [k_j]$ , has higher probability of being included in the subset  $\Omega$ , compared to other values in  $A_j \setminus \{v_i\}$  that are sampled uniformly at random (without replacement). The optimal subset size  $\omega = |\Omega|$  that minimizes the variance is  $\omega = \frac{k_j}{e^\epsilon + 1}$  [49, 56].

Given a value  $v_i \in A_j$ , for  $i \in [k_j]$ , the  $\omega$ -SS protocol starts by initializing an empty subset  $\Omega$ . Afterwards, the true value  $v_i$  is added to  $\Omega$  with probability  $p = \frac{\omega e^\epsilon}{\omega e^\epsilon + k_j - \omega}$ . Finally, it adds values to  $\Omega$  as follows [49, 56]:

- If  $v_i$  has been added to  $\Omega$  in the previous step, then  $\omega - 1$  values are sampled from  $A_j \setminus \{v_i\}$  uniformly at random (without replacement) and are added to  $\Omega$ ;

- If  $v_i$  has not been added to  $\Omega$  in the previous step, then  $\omega$  values are sampled from  $A_j \setminus \{v_i\}$  uniformly at random (without replacement) and are added to  $\Omega$ .

From this, the aggregator can obtain the unbiased estimation of  $v_i \in A_j$ , for  $i \in [k_j]$ , with Eq. (2) by setting  $p = \frac{\omega e^\epsilon}{\omega e^\epsilon + k_j - \omega}$  and  $q = \frac{\omega e^\epsilon (\omega - 1) + (k_j - \omega) \omega}{(k_j - 1)(\omega e^\epsilon + k_j - \omega)}$  [49, 56].

**2.2.4 Unary Encoding Protocols.** Unary encoding (UE) protocols interpret the user’s input  $v_i \in A_j$ , for  $i \in [k_j]$  as a one-hot  $k_j$ -dimensional vector. More specifically,  $B = UE(v_i)$  is a binary vector with only the bit at the position  $v_i$  sets to 1 and the other bits set to 0. One well-known UE-based protocol is the Basic One-time RAPPOR [24], hereafter referred to as symmetric UE (SUE) [51], which randomizes the bits from  $B$  independently with probabilities:

$$\forall i \in [k_j] : \Pr[B'_i = 1] = \begin{cases} p = \frac{e^{\epsilon/2}}{e^{\epsilon/2} + 1}, & \text{if } B_i = 1 \\ q = \frac{1}{e^{\epsilon/2} + 1}, & \text{if } B_i = 0. \end{cases} \quad (3)$$

Afterwards, the client sends  $B'$  to the aggregator. More recently, to minimize the variance of the SUE protocol, the authors in [51] proposed Optimal UE (OUE), which selects probabilities  $p = \frac{1}{2}$  and  $q = \frac{1}{e^\epsilon + 1}$  in Eq. (3) asymmetrically (i.e.,  $p + q \neq 1$ ). The estimation method used in Eq. (2) applies equally to both SUE and OUE protocols, in which both satisfy  $\epsilon$ -LDP for  $\epsilon = \ln\left(\frac{p(1-q)}{(1-p)q}\right)$  [24, 51].

## 2.3 Multidimensional Frequency Estimation

Let  $n$  be the total number of users,  $d \geq 2$  be the total number of attributes,  $\mathbf{k} = [k_1, k_2, \dots, k_d]$  be the domain size of each attribute,  $\mathcal{M}$  be a local randomizer and  $\epsilon$  be the privacy budget. Each user holds a tuple  $\mathbf{v} = [v_1, v_2, \dots, v_d]$ , (i.e., a private discrete value per attribute). The two next subsections describes the SPL, SMP and RS+FD solutions for frequency estimation of multiple attributes.

**2.3.1 Standard Solutions.** Previous works in the local DP setting considered the following approaches [5, 38, 48, 51]:

- **SPL.** On the one hand, due to the sequential composition theorem [22], users can split the privacy budget  $\epsilon$  over the number of attributes  $d$  and send all randomized values  $y_j$ , for  $j \in [d]$ , with  $\frac{\epsilon}{d}$ -LDP to the aggregator (i.e., a tuple  $\mathbf{y} = [y_1, y_2, \dots, y_d]$ ). However, this naïve SPL solution leads to high estimation error [5, 38, 48, 51].
- **SMP.** Instead of splitting the privacy budget  $\epsilon$ , this state-of-the-art solution allows each user to sample a single attribute  $j \in [d]$  at random and uses all the privacy budget to send it with  $\epsilon$ -LDP [5, 38, 48, 51]. In this case, each user tells the aggregator which attribute is sampled, and what is the perturbed value for it ensuring  $\epsilon$ -LDP (i.e.,  $\langle j, y_j \rangle$ ).

**2.3.2 Random Sampling Plus Fake Data (RS+FD).** Because the SMP solution discloses the sampled attribute, one can say that it is not fair to all users (e.g., some users will sample age while others will sample disease). To address this issue, the recently proposed RS+FD [4] solution is composed of two steps, namely local randomization and fake data generation. More precisely, each user samples a unique attribute uniformly at random  $j = \text{Uniform}([d])$  and uses an  $\epsilon$ -LDP protocol to sanitize its value  $v_j$ . Next, for each non-sampled

attribute  $i \in [d] \setminus \{j\}$ , the user generates uniform random fake data following  $A_i$ . Finally, each user sends the (LDP or fake) value of each attribute to the aggregator (*i.e.*, a tuple  $\mathbf{y} = [y_1, y_2, \dots, y_d]$ ). In this manner, the sampling result is not disclosed to the aggregator, thus increasing the *uncertainty*. For this reason, to satisfy  $\epsilon$ -LDP, following the parallel composition theorem [22] and the amplification by sampling result [32], RS+FD utilizes an amplified privacy budget  $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$  for the sampled attribute [4].

With the RS+FD solution, the estimator should remove the bias introduced by the local randomizer  $\mathcal{M}$  and uniform fake data. In [4], the authors used GRR and OUE as LDP protocols within the RS+FD solution, which results in RS+FD[GRR], RS+FD[OUE-z] and RS+FD[OUE-r]. We briefly recall how these three protocols, generalizing OUE to UE as one can select either SUE or OUE (*cf.* Section 2.2.4) as local randomizers [4, 47].

For all three protocols, on the *client-side*, each user randomly samples an attribute  $j$  and uses  $\mathcal{M}$  to sanitize the value  $v_j$  with an amplified privacy parameter  $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$ . Next, the fake data generation procedure and the unbiased estimator for the frequency of each value  $v_i \in A_j$ , for  $i \in [k_j]$ , are as follows:

- **RS+FD[GRR]** [4]. For each non-sampled attribute  $i \in [d] \setminus \{j\}$ , the user generates fake data uniformly at random according to the domain size  $k_i$ . On the *server-side*, the unbiased estimator for this protocol is:  $\hat{f}(v_i) = \frac{C(v_i)dk_j - n(d-1+qk_j)}{nk_j(p-q)}$ , in which  $C(v_i)$  is the number of times  $v_i$  has been reported,  $p = \frac{e^{\epsilon'}}{e^{\epsilon'} + k_j - 1}$  and  $q = \frac{1-p}{k_j - 1}$ .
- **RS+FD[UE-z]** [4]. For each non-sampled attribute  $i \in [d] \setminus \{j\}$ , the user generates fake data by applying an UE protocol to zero-vectors (*i.e.*,  $[0, 0, \dots, 0]$ ) of size  $k_i$ . On the *server-side*, the unbiased estimator for this protocol is:  $\hat{f}(v_i) = \frac{d(C(v_i) - nk_j)}{n(p-q)}$ , in which  $C(v_i)$  is the number of times  $v_i$  has been reported and parameters  $p$  and  $q$  can be selected following the SUE [24] or OUE [51] protocols.
- **RS+FD[UE-r]** [4]. For each non-sampled attribute  $i \in [d] \setminus \{j\}$ , the user generates fake data by applying an UE protocol to one-hot-encoded fake data (uniform at random) of size  $k_i$ . On the *server-side*, the unbiased estimator for this protocol is:  $\hat{f}(v_i) = \frac{C(v_i)dk_j - n[qk_j + (p-q)(d-1) + qk_j(d-1)]}{nk_j(p-q)}$ , in which  $C(v_i)$  is the number of times  $v_i$  has been reported and parameters  $p$  and  $q$  can be selected following the SUE [24] or OUE [51] protocols.

### 3 SYSTEM OVERVIEW & PRIVACY THREATS

Hereafter, we describe the system and adversary models before presenting our adversarial analyses of SMP and RS+FD.

#### 3.1 System Overview

We consider the situation in which a (possibly untrusted) server collects users' multidimensional data  $d \geq 2$  for frequency estimation under  $\epsilon$ -LDP guarantees multiple times. Particularly, in each data collection (*i.e.*, survey), the server can select a different number of attributes. For instance, through a mobile app the server may collect private frequency estimation for different users' demographic data and different application usage (*e.g.*, how much time spent on the

application, preferred widget, etc). Users could be encouraged to share their private data through the exchange of discount coupons, statistics to compare usage with other users, etc. For the sake of simplicity, we assume that the set of users  $\mathcal{U}$  is unique across all surveys, although this can be relaxed in real-life allowing users to opt-in or opt-out of a given survey. We assume that the server uses one of the state-of-the-art LDP solutions (*e.g.*, SMP or RS+FD) to collect one random attribute per user. Thus, we do not consider the SPL solution in our attacks as all attributes would be collected at once, thus resulting in a low level of utility [5, 38, 48, 51]

**Adversary model.** Following the LDP assumptions [19, 29], we assume that the server knows the users' pseudonymized IDs, *but not their private data or their real identity*. This also implies that the server has no knowledge about the real data distributions. However, we assume that the server might have some background knowledge  $\mathcal{D}_{BK}$  coming from public available source, such as Census data [1]. This background knowledge could for instance contain partial or complete profiles of users along with their true identities. Thus, the adversary could be for example the server itself, an attacker who intercepts the communication between the client and the server (*e.g.*, through a man-in-the-middle attack) or a third-party analyst with whom the server may have shared the collected data.

#### 3.2 Attacking SMP: Plausible Deniability and Risks of Re-Identification

**Plausible deniability.** Let  $v_y$  be an embarrassing value of  $A_j = \{v_y, v_n\}$  (*e.g.*, a value "Yes" for an attribute  $A_j$  denoting whether someone cheated on their partner). As long as  $\Pr[\mathcal{M}(v_y) = v_y] < 1$ , the user can deny to have  $A_j = v_y$  [17].

The LDP protocols of Section 2.2 are based on RR [54], which provides "plausible deniability" for users' reports. However, increasing  $\epsilon$  to improve utility of LDP protocols compromises the "plausible deniability" of the users' reports. Indeed, common  $\epsilon$  values used daily by users in high-scale industrial systems nowadays range from small  $\epsilon \leq 1$  to high values  $\epsilon \geq 8$  [14, 40, 45]. Thus, we conduct an adversarial analysis to the SMP solution (*cf.* Section 2.3.1) in which the user randomly samples a single attribute among  $d \geq 2$  ones and uses the whole privacy budget  $\epsilon$  to report it. Consequently, since the whole privacy budget will be allocated to a single attribute, the "plausible deniability" for this attribute will be lower, which can lead an attacker to predict the users' true value as the most likely value after randomization (see details in Section 3.2.1). In this setting, in which many surveys are proposed by the server to the same set of users with possibly different number of attributes (*e.g.*, demographic, preference, application usage, ...), an attacker knowing the tuple (sampled attribute,  $\epsilon$ -LDP report) will be able to profile each user throughout time. Therefore, once a partial or complete profile of the target user is built (see details in Sections 3.2.2, 3.2.3 and 3.2.4), the adversary could use his background knowledge  $\mathcal{D}_{BK}$  to possibly re-identify a user within population [31, 34, 41–44], possibly also inferring all other available attributes. The next four subsections analyze the "plausible deniability" of LDP protocols in single and multiple collections, and describes the proposed re-identification attack models, respectively.

**3.2.1 Plausible Deniability of LDP protocols.** Given a user's true value  $v \in A_j$ , different LDP protocols  $\mathcal{M}$  have different type of

output  $y_i = \mathcal{M}(v, \epsilon)$  [51]. For instance, UE protocols output unary encoded vectors,  $\omega$ -SS outputs a subset  $\Omega$  of  $\omega$  non-encoded values and so on (cf. Section 2.2). Thus, for each user  $u_i \in \mathcal{U}$ , for  $i \in [n]$ , given  $y_i$ , the adversary's goal is to predict  $v_i$ , which is denoted as  $\hat{v}_i$ . The attacker's accuracy (ACC) for LDP protocols is measured by the number of correct predictions  $v = \hat{v}$  over the number of users  $n$ :  $ACC_{FO}(\%) = 100 \cdot \frac{\sum_{i=1}^n f(v_i, \hat{v}_i)}{n}$ , in which  $f(v, \hat{v}) = 1$  if  $v = \hat{v}$  and 0 otherwise. Following the "plausible deniability" intuition and the fact that for all LDP protocols the probability  $p$  of reporting the true value  $v_i$  (or bit  $i$ ) is higher than any other value  $v \in A_j \setminus \{v_i\}$ , we now describe our attack strategy to each LDP protocol. By the time of completing this paper, we learned about a recent work showing that the expectation of our attacks could be analytically formalized with the Bayes adversary of [23]. We believe this work is complementary to our "plausible deniability" attacking interpretation.

**Plausible Deniability of GRR.** Since no specific encoding is used with GRR, the most likely value after randomization is the user's  $u_i$  own true value  $v$ . Thus, an attacker can assume that the reported value  $y$  is the true one (i.e.,  $\hat{v} = y$ ), which gives on expectation an  $ACC_{GRR}(\%) = 100 \cdot \frac{e^\epsilon}{e^\epsilon + k_j - 1}$ .

**Plausible Deniability of OLH.** Since the output of OLH for user  $u_i$  is the hash function  $H_i$  used to hash the user's value  $v$  and the hashed value  $h_i = H_i(v)$ , the most likely value after randomization is one within the subset of all values  $v \in A_j$  that hash to  $h_i$  (i.e.,  $A_{jH} = \{v | v \in A_j, H_i(v) = h_i\}$ ). Thus, the attacker's best guess is a random choice  $\hat{v} = \text{Uniform}(A_{jH})$ . On expectation [23], one achieves:  $ACC_{OLH}(\%) = 100 \cdot \frac{1}{2 \cdot \max(\frac{k_j}{e^\epsilon + 1}, 1)}$ .

**Plausible Deniability of  $\omega$ -SS.** Since the output of  $\omega$ -SS for user  $u_i$  is a set  $\Omega \subseteq A_j$ , the most likely value after randomization is one within the subset  $\Omega$ . Thus, the attacker's best guess is a random choice  $\hat{v} = \text{Uniform}(\Omega)$ . Selecting  $\omega = \frac{k_j}{e^\epsilon + 1}$  in  $\omega$ -SS [49, 56], on expectation [23], one achieves:  $ACC_{\omega\text{-SS}}(\%) = 100 \cdot \frac{e^\epsilon + 1}{2k_j}$ .

**Plausible Deniability of UE protocols.** Since the output of UE protocols for user  $u_i$  is a sanitized unary encoded vector  $B$  of size  $k_j$ , there are three possibilities: 1) a single bit  $b$  in  $B$  is set to 1, in which the attacker's best guess is to predict the bit as the true value as  $\hat{v} = B_b$ ; 2) more than one bit in  $B$  is set to 1, in which the attacker's best guess is a random choice of the bits set to 1 as  $\hat{v} = \text{Uniform}(\{b | b \in [k_j] \text{ if } B_b = 1\})$ ; and 3) no bit in  $B$  is set to 1, in which the attacker's best guess is a random choice of the domain  $\hat{v} = \text{Uniform}(A_j)$ . Therefore, on expectation [23], the attacker's accuracy for SUE is:  $ACC_{SUE}(\%) = 100 \cdot \frac{1}{k_j(e^{\epsilon/2} + 1)} \cdot \frac{e^{\epsilon/2}}{e^{\epsilon/2} + 1} + \sum_{i=1}^{k_j} \frac{e^{\epsilon/2}}{(e^{\epsilon/2} + 1)^i} \cdot \text{Bin}(i - 1; k_j - 1, \frac{1}{e^{\epsilon/2} + 1})$ , in which  $\text{Bin}(\cdot)$  denotes a Binomial distribution with  $k_j - 1$  trials, success probability  $\frac{1}{e^{\epsilon/2} + 1}$  and exactly  $i - 1$  successes. On the other hand, on expectation [23], the attacker's accuracy for OUE is:  $ACC_{OUE}(\%) = 100 \cdot \frac{1}{2k_j} \cdot \frac{e^\epsilon}{e^\epsilon + 1} + \sum_{i=1}^{k_j} \frac{1}{2i} \cdot \text{Bin}(i - 1; k_j - 1, \frac{1}{e^\epsilon + 1})$ .

**3.2.2 Plausible Deniability on Multiple Data Collections: Uniform Privacy Metric.** When collecting multidimensional data  $d \geq 2$  with the SMP solution multiple times, the server could implement that all users sample attributes without replacement. This way, each user will randomly select a new attribute in each data collection (i.e.,

survey), ensuring a *uniform privacy metric across all users*. Since for all LDP protocols the expected  $ACC_{FO}$  depends on  $\epsilon$  and  $k_j$ , our analysis focuses on a generic LDP protocol here. Therefore, depending on the LDP protocol, the expected ACC with uniform privacy metric after #surveys =  $d$ , denoted as  $ACC_{FO}^U$ , now follows:

$$ACC_{FO}^U(\%) = 100 \cdot \prod_{j=1}^d ACC_{FO}(\epsilon, k_j). \quad (4)$$

Since each survey is independent and users sample without replacement, Eq. (4) represents the expected probability of accurately profiling users with exactly  $d$  attributes.

**3.2.3 Plausible Deniability on Multiple Data Collections: Non-Uniform Privacy Metric.** On the other hand, when collecting multidimensional data  $d \geq 2$  with the SMP solution multiple times, the server can allow users to sample attributes with replacements in each data collection. In case of a repeated attribute, the user can report the previous randomized value (a.k.a. memoization [5, 15, 24]). This way, *users will have a non-uniform privacy metric*. Depending on the LDP protocol, the expected ACC with non-uniform privacy metric after #surveys =  $d$ , denoted as  $ACC_{FO}^{NU}$ , now follows:

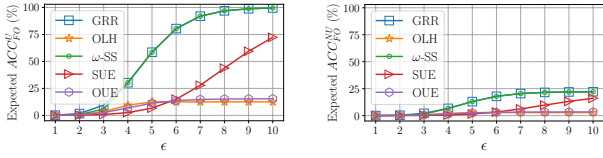
$$ACC_{FO}^{NU}(\%) = 100 \cdot \prod_{j=1}^d \frac{d+1-j}{d} ACC_{FO}(\epsilon, k_j). \quad (5)$$

Since each survey is independent but attributes are sampled with replacement, Eq. (5) denotes the overall adversary's accuracy only considering users that reports a different attribute in each survey (i.e., of accurately profiling users with exactly  $d$  attributes). Thus, in this setting, users can also end-up with partial profiles.

**Analytical analysis of expected ACC.** In Fig. 1, we illustrate the expected  $ACC_{FO}^U$  following Eq. (4) and the  $ACC_{FO}^{NU}$  following Eq. (5) of each LDP protocol with the following parameters (taken from Section 4):  $\epsilon = [1, 2, \dots, 9, 10]$ ,  $d = 3$ ,  $\mathbf{k} = [74, 7, 16]$ , and #surveys =  $d$ . From Fig. 1 (a), one can notice that GRR,  $\omega$ -SS and SUE have the highest attacker's accuracy, which would enable an adversary to accurately infer a **complete profile** after #surveys =  $d$ . Allowing users to have non-uniform privacy metrics in the plot (b), minimizes the attacker's accuracy to infer complete profiles as the probability of selecting different attributes in all  $d$  surveys is  $\frac{d!}{d^d}$ . Note that the expected  $ACC_{FO}$  in both Eqs. (4) and (5) decreases with the #surveys since the probability of accurately inferring the users' true value is independent in each survey.

**3.2.4 Re-Identification Attack Models.** Following the system overview of Section 3.1, we consider two re-identification attack models: **full-knowledge re-identification (FK-RI)** and **partial knowledge re-identification (PK-RI)**, that we detail in the following. The first FK-RI model considers that the attacker has access to the complete background knowledge  $\mathcal{D}_{BK}$  to re-identify the target user. The latter PK-RI model considers that the attacker only has access to a subset  $\mathcal{D}_{PK} \subseteq \mathcal{D}_{BK}$  for her re-identification attack. The re-identification success of both FK-RI and PK-RI models will depend on the results of Sections 3.2.2 and 3.2.3 to accurately profile the target user, which is impacted by the LDP protocol considered.

In particular, after #surveys, the attacker will have a profile  $\mathbf{y}_i$  of at most #surveys sanitized values for the target user  $u_i \in \mathcal{U}$ . The



(a) Uniform privacy metric. (b) Non-uniform privacy metric.

**Figure 1: Analytical attacker’s accuracy when collecting multidimensional data ( $d = 3$ ) with the SMP solution multiple times ( $\#surveys = 3$ ) with attributes’ domain size  $k = [74, 7, 16]$ : (a) uniform privacy metric across users with Eq. (4) and (b) non-uniform privacy metric across users with Eq. (5).**

number of attributes inferred per target user depends on the setting used (*i.e.*, uniform or non-uniform privacy metrics). Therefore, the re-identification attack starts with a *matching algorithm*  $\mathcal{R}$ , which takes as input the sanitized profile  $\mathbf{y}_i$  and the background knowledge  $\mathcal{D}_{BK}$  (or  $\mathcal{D}_{PK}$  for PK-RI), and outputs a score  $c_i \in \mathbb{R}$ . More precisely, the score  $c_i$  measures the distance between the target  $\mathbf{y}_i$  and all samples  $\mathbf{r} \in \mathcal{D}_{BK}$ . Since the LDP protocols from Section 2.2 do not have a notion of “distance” when randomizing a value, when an attribute in  $\mathbf{y}_i \neq \mathbf{r}$  the distance is 1 and 0 otherwise. A smaller distance between  $\mathbf{y}_i$  and a profile in  $\mathcal{D}_{BK}$  indicates that is highly likely that  $\mathbf{y}_i$  has been re-identified through the uniqueness combination of  $\#surveys$  attributes [31, 34, 41–44]. Finally, a *decision algorithm*  $\mathcal{G}$  takes as input the computed distances and outputs a list of top- $k$  possible profiles (or IDs) in  $\mathcal{D}_{BK}$  that corresponds to the target user  $u_i \in \mathcal{U}$ . The attacker’s re-identification accuracy (RID-ACC) is measured by the number of correct re-identification  $u_{id} = \hat{u}_{id}$  over the number of users  $n$ :  $RID-ACC(\%) = 100 \cdot \frac{\sum_{i=1}^n f(u_{id}, \hat{u}_{id})}{n}$ , in which  $f(u_{id}, \hat{u}_{id}) = 1$  if  $u_{id} = \hat{u}_{id}$  and 0 otherwise. The attacker’s RID-ACC depends on the accuracy of partially or completely profiling the target user (*i.e.*, as measured by Eqs. (4) and (5)) and the “uniqueness” of users with respect to the collected attributes (unknown by the server) and in the background knowledge  $\mathcal{D}_{BK}$ .

### 3.3 Attacking RS+FD: Uncovering the Sampled Attribute ( $\rightarrow$ SMP)

Because the objective of the RS+FD solution is to hide the LDP value among fake data [4], discovering the sampled attribute of each user would convert RS+FD into the SMP solution again. Even more, unlike SMP (and SPL), RS+FD utilizes an amplified  $\epsilon' > \epsilon$ , which decreases the “plausible deniability” of the user’s report (*cf.* Section 3.2.1) and could thus be leveraged for re-identification attacks (*cf.* Section 3.2.4) under multiple data collections.

For instance, consider the scenario in which a given user  $u \in \mathcal{U}$  whose sampled attribute is  $t \in [d]$  produces an RS+FD’s output tuple as  $\mathbf{y} = [y_1, y_2, \dots, y_d]$ . In this situation, the **baseline classification model** is just a random guess  $\hat{t} = Uniform([d])$ . In addition, we propose a **classifier learning setting** in which an attacker aims to train a classifier over a learning dataset  $\mathbf{D}_l = \{(\mathbf{y}_i, t_i) \mid i \in [r]\}$  of  $r$  rows and  $c = d + 1$  columns. That is, for each user  $u_i$ ,  $\mathbf{y}_i$  is the output tuple of the RS+FD solution (LDP/fake values, *i.e.*, a full profile of  $d$  attributes) and  $t_i$  is the sampled attribute (target is

a class within  $[d]$ ). **Because the sampled attribute  $t_i$  of users should be unknown to the attacker**, in this work, we propose three settings to build a learning dataset  $\mathbf{D}_l$ , which depends on the attack model. In all these settings, we assume that the attacker has the knowledge of the privacy budget  $\epsilon$  and the LDP protocol used by users with the RS+FD solution. Finally, the attacker’s attribute inference accuracy (AIF-ACC) is measured by the number of correct predictions  $t = \hat{t}$  over the number of users in the testing dataset  $n_t$ :  $AIF-ACC(\%) = 100 \cdot \frac{\sum_{i=1}^{n_t} f(t_i, \hat{t}_i)}{n_t}$ , in which  $f(t, \hat{t}) = 1$  if  $t = \hat{t}$  and 0 otherwise.

**3.3.1 No Knowledge: Training a Classifier Over Synthetic Profiles.** With no knowledge of the real sampled attribute of the  $n$  users  $u \in \mathcal{U}$  and after aggregating users’ LDP data, an attacker could use the estimated frequencies  $\hat{\mathbf{f}} = [\hat{f}_1, \hat{f}_2, \dots, \hat{f}_d]$  to generate  $s$  **synthetic profiles**  $\mathbf{s}_i = [s_{i1}, s_{i2}, \dots, s_{id}]$ , for  $i \in [s]$ , *i.e.*, mimic the real profiles with one value per attribute. Afterwards, for all  $s$  synthetic profiles, the attacker could follow the same protocol used by the real users (*i.e.*, RS+FD with an LDP protocol) to generate the learning set  $\mathbf{D}_l$ . Notice that the attacker has full control over the training set size  $s$ , which can be seen as a trade-off between computational costs (*i.e.*, generating  $s$  synthetic profiles and use as training set) and the attacker’s AIF-ACC. In this **no knowledge (NK)** model, the testing set  $\mathbf{D}_t$  is composed of all the real RS+FD’s sanitized tuples  $\mathbf{y}$  of users  $u \in \mathcal{U}$ , and the objective is to accurately classify their sampled attribute  $t \in [d]$ .

**3.3.2 Partial-Knowledge: Training a Classifier Over Real (Known) Profiles.** This second setting considers the scenario in which the attacker has knowledge about the sampled attribute of  $n_{pk} < n$  real users, *i.e.*, the subset  $\mathcal{U}_{pk} \subset \mathcal{U}^3$ . This setting corresponds in situations in which some users disclose the sampled attribute by preference (*e.g.*, less “sensitive” attributes) or due to security breaches. In this **partial-knowledge (PK)** model, the learning set  $\mathbf{D}_l$  depends on the number of (compromised) profiles  $n_{pk}$  the attacker has access to and the testing set  $\mathbf{D}_t$  has  $n - n_{pk}$  sanitized tuples  $\mathbf{y}$  of users  $u \in \mathcal{U} \setminus \mathcal{U}_{pk}$ , in which the objective is to accurately classify their sampled attribute  $t \in [d]$ .

**3.3.3 Partial-Knowledge Plus Synthetic Profiles.** This last setting combines both NK and PK models, in which the attacker has knowledge about the sampled attribute of  $n_{pk} < n$  real users and augments the subset  $\mathcal{U}_{pk} \subset \mathcal{U}$  with  $s$  synthetic profiles. In this **hybrid model (HM)**, the learning set  $\mathbf{D}_l$  is dependent on both the number of synthetic profiles  $s$  the attacker generates and the number of (compromised) profiles  $n_{pk}$  the attacker has access to. Similarly to the PK model, the testing set  $\mathbf{D}_t$  has  $n - n_{pk}$  sanitized tuples  $\mathbf{y}$  of users  $u \in \mathcal{U} \setminus \mathcal{U}_{pk}$ , and the goal is to accurately classify their sampled attribute  $t \in [d]$ .

## 4 EXPERIMENTAL EVALUATION

In this section, we introduce the general setup of our experiments. Next, we present the experimental setting and results on the risks of re-identification of the SMP solution. Afterwards, we describe the setup of experiments carried out to uncover the sampled attribute

<sup>3</sup>If  $\mathcal{U}_{pk} \subseteq \mathcal{U}$ , this will correspond to a full-knowledge model in which the adversary has knowledge of all users’ sampled attribute (*i.e.*, SMP solution).

of the RS+FD solution. Finally, we detail the experimental setting and results on the risks of re-identification of the RS+FD solution.

## 4.1 Experimental Setup

**Environment.** All algorithms were implemented in Python 3. In all experiments, we report the results averaged over 20 runs.

**Datasets.** For ease of reproducibility, we conduct our experiments on two census-based multidimensional and open datasets.

- **ACSEmployment.** This dataset is generated from the Folktables Python package [16] that provides access to datasets derived from the US Census. We have selected the “Montana” state only, which results in  $n = 10,336$  samples with  $d = 18$  discrete attributes (target included) and domain size  $\mathbf{k} = [92, 25, 5, 2, 2, 9, 4, 5, 5, 4, 2, 18, 2, 2, 3, 9, 3, 6]$ .
- **Adult.** This is a classical dataset from the UCI ML repository [18] with  $n = 45,222$  samples after cleaning. We selected  $d = 10$  attributes (“age”, “workclass”, “education”, “marital-status”, “occupation”, “relationship”, “race”, “sex”, “native-country” and “salary”) with domain size  $\mathbf{k} = [74, 7, 16, 7, 14, 6, 5, 2, 41, 2]$ , respectively.

## 4.2 Re-identification Risk of the SMP Solution

**Methods evaluated.** We consider for evaluation all five LDP protocols described in Section 2.2: GRR, OLH,  $\omega$ -SS, SUE and OUE.

**Privacy protection.** We vary the privacy budget in the interval  $\epsilon = [1, 2, \dots, 9, 10]$ , which corresponds to values used by industry nowadays [14, 40] and experiments found in the LDP attacking literature with single [9, 23, 36] and multiple [25] collections.

**Attack performance metric.** We measure the quality of the re-identification attack with the attacker’s re-identification accuracy (RID-ACC) metric, which corresponds to how many times the user is correctly re-identified in the top- $k$  groups, for top- $k \in \{1, 10\}$ .

**Baseline.** For each top- $k$ , the baseline re-identification model follows top- $k$  random guesses (*i.e.*,  $\hat{u}_{id} = \text{Uniform}([n])$ ) without replacement with expected RID-ACC: top- $k/n$ .

**Experimental evaluation.** We set #surveys = 5, in which each survey  $sv \in [\text{\#surveys}]$ , has a different number of attributes  $d_{sv} = \text{Uniform}\left(\frac{d}{2}, \dots, d\right)$  (*i.e.*, with at least  $\frac{d}{2}$  attributes). The attributes are also selected at random per survey. Due to space constraints, we only present here the experiments with the FK-RI model (*cf.* Section 3.2.4), considering the  $d$ -dimensional dataset as background knowledge  $\mathcal{D}_{BK}$ , and with the uniform privacy metric setting from Section 3.2.2. Finally, we measure the attacker’s RID-ACC after #surveys  $\geq 2$ , which results in the inferred profile of each user having respectively 2, 3, 4 or 5 attributes, to be used for the re-identification attack.

**Results.** Fig. 2 illustrates the attacker’s RID-ACC metric on the Adult dataset for top- $k$  re-identification using the SMP solution, the FK-RI model with uniform  $\epsilon$ -LDP privacy metric across users, by varying the LDP protocol and the number of surveys. Additional results with all LDP protocols, Adult and ACSEmployment datasets, FK-RI and PK-RI models, uniform and non-uniform privacy metric settings as well as with the relaxed LDP metric of [36] are presented in Appendix C of [6].

**Analysis.** In general, the experimental results of Fig. 2 match the numerical results of the expected values from Fig. 1. From Fig. 2, one

can observe that our re-identification attacks present significant improvement over a random baseline model that has  $\text{RID-ACC} \ll 1\%$  (*i.e.*, top- $k/n$ ). For instance, with a single shot (*i.e.*, top-1), the attacker’s RID-ACC is already significant for GRR (and  $\omega$ -SS) and SUE after about #surveys  $\geq 4$ , with at most  $\sim 10\%$  of RID-ACC. In comparison, both OUE and OLH protocols have about 10x less RID-ACC, (*i.e.*, at most  $\sim 1\%$  of RID-ACC for top-1). On the other hand, when there is a set of top-10 profiles, the adversary achieves  $\text{RID-ACC} \geq 2.5\%$  for GRR (and  $\omega$ -SS) after only 2 surveys with an upper bound of about 33% of RID-ACC after 5 surveys. Though with slightly smaller RID-ACC, the SUE protocol also achieves about 28% of RID-ACC after 5 surveys, and both OUE and OLH are upper bounded by about 5% of RID-ACC. Although the user is not uniquely re-identified, this still represents a threat due to the possibility of performing, *e.g.*, homogeneity attacks [12, 31, 34].

Overall, these “high” re-identification rates may be explained by many factors. First, the combination of multiple attributes within the Adult dataset leads to several unique people or small groups of people (this is also the case for the ACSEmployment dataset in Fig. 9 of Appendix C of [6]). Additionally, the uniform privacy metric setting require the users to always sample a new attribute, increasing the privacy leakage. In a more realistic scenario, the non-uniform privacy metric setting minimizes the RID-ACC (see Fig. 11 of Appendix C of [6]) as already shown in Fig. 1. Furthermore, the FK-RI model allows the attacker to use the whole background knowledge  $\mathcal{D}_{BK}$  to match the inferred profiles. For instance, the attacker’s RID-ACC metric decreased by almost half when considering the PK-RI model (*cf.* Fig. 10 of Appendix C of [6]) since there are fewer attributes as background information to use for the matching algorithm  $\mathcal{R}$  (see Section 3.2.4). Lastly, we used the same dataset for private data collection and as (partial) background knowledge. A different set of experiments could mix demographic attributes and (synthetic) application usage in each survey, limiting the number of demographic attributes per user to constitute a profile.

## 4.3 Uncovering the Sampled Attribute of the RS+FD Solution ( $\rightarrow$ SMP)

**Classifier.** We use the state-of-the-art XGBoost [10] algorithm to predict the sampled attribute of users in a multiclass classification framework (*i.e.*,  $d$  attributes) with default parameters.

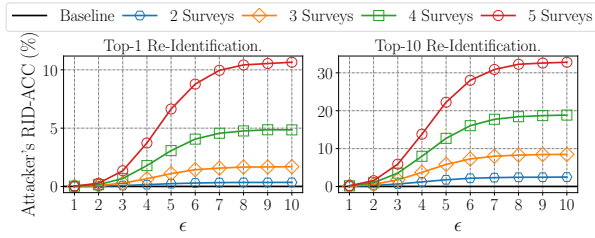
**Methods evaluated.** We consider for evaluation five protocols within the RS+FD solution from Section 2.3.2, namely RS+FD[GRR], RS+FD[SUE-z], RS+FD[SUE-r], RS+FD[OUE-z] and RS+FD[OUE-r].

**Metrics.** Similar to Section 4.2, we vary the privacy budget in the interval  $\epsilon = [1, 2, \dots, 9, 10]$ . Besides, we use the attacker’s attribute inference accuracy (AIF-ACC) metric to measure the quality of the attack, which corresponds to how many times the attacker can correctly predict the users’ sampled attribute.

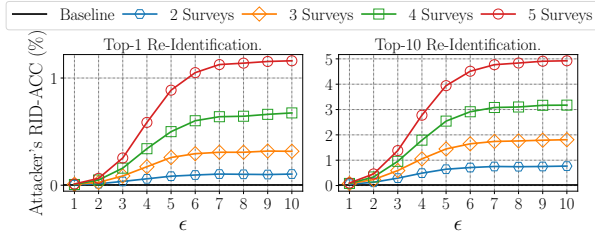
**Baseline.** The baseline classification model is a random guess  $\hat{t} = \text{Uniform}([d])$  with expected AIF-ACC:  $1/d$ .

**Experimental evaluation.** All five protocols are evaluated with the three settings of Section 3.3, namely No Knowledge (NK), Partial-Knowledge (PK) and Hybrid Model (HM). For the NK model, we vary the number of synthetic profiles  $s$  the attacker generates in the interval  $s = [1n, 3n, 5n]$ . For the PK model, we vary the number of compromised profiles  $n_{pk}$  the attacker has access to in the interval

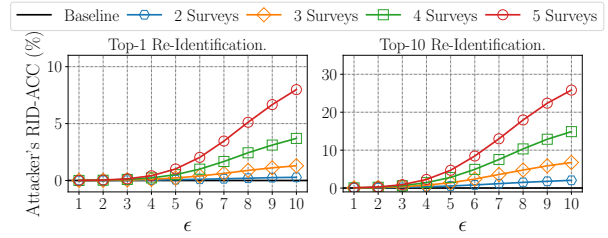




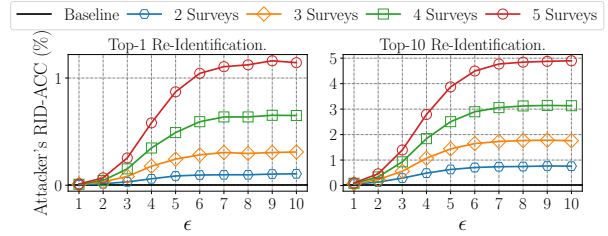
(a) Re-identification risk of the GRR [27, 28] protocol.



(c) Re-identification risk of the OLH [51] protocol.



(b) Re-identification risk of the SUE (a.k.a. RAPPOR) [24] protocol.



(d) Re-identification risk of the OUE [51] protocol.

**Figure 2: Attacker’s re-identification accuracy (RID-ACC) on the Adult dataset for top- $k$  re-identification on using the SMP solution, the full knowledge FK-RI model with uniform  $\epsilon$ -LDP privacy metric across users, and by varying the LDP protocol and the number of surveys (i.e., data collections). Omitted results for the  $\omega$ -SS protocol [49, 56] is due to similarity to plot (a).**

$n_{pk} = [0.1n, 0.3n, 0.5n]$ . Finally, for the HM setting, we combined both intervals, i.e.,  $(s, n_{pk}) = [(1n, 0.1n), (3n, 0.3n), (5n, 0.5n)]$ .

**Results.** Fig. 3 illustrates the attacker’s AIF-ACC metric on the ACSEmployment dataset with the three attack models (i.e., NK, PK and HM) and all five protocols (i.e., RS+FD[GRR], RS+FD[SUE-z], RS+FD[OUE-z], RS+FD[SUE-r] and RS+FD[OUE-r]), varying  $\epsilon$ , the number of synthetic profiles  $s$  and the number of compromised profiles  $n_{pk}$ . Additional results (Adult and Nursery datasets [18]) are presented in Appendix D of [6].

**Analysis.** From Fig. 3, one can notice that the proposed attack models, namely, NK, PK and HM present significant 2-20 fold increments in the attacker’s AIF-ACC over the Baseline model. Surprisingly, even under an NK model in which the attacker has access only to the estimated frequencies satisfying  $\epsilon$ -LDP, generating  $s = [1n, 3n, 5n]$  synthetic profiles to train a classifier provides higher attacker’s AIF-ACC than having compromised  $n_{pk} = 0.5n$  profiles in the PK model. On the other hand, increasing the number of synthetic profiles  $s$  that the attacker generates in the NK model has less impact than increasing the number of compromised profiles  $n_{pk}$  that the attacker has access to in the PK model. Due to this, results for both NK and HM models are quite similar.

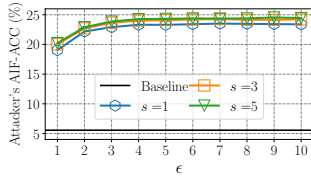
In this adversarial analysis, the attacker’s AIF-ACC now depends on both the LDP protocol and how fake data are generated. For the former (i.e., different LDP protocols), the difference between RS+FD[GRR] and RS+FD[UE-r] protocols lies in the encoding and randomization steps, which directly affects the attacker’s AIF-ACC with a difference of about 5% favoring the RS+FD[GRR] protocol. Since GRR requires no particular encoding, there is less noise compared to a randomized unary encoded vector. Furthermore, with respect to different fake data generation procedures, when fake data are generated with a uniformly random (encoded) value (i.e.,

RS+FD[GRR] and RS+FD[UE-r]), the attacker’s AIF-ACC is upper-bounded by about 25%. On the other hand, generating fake data through applying a UE protocol on zero-vectors led to an attacker’s AIF-ACC of about 50% with RS+FD[OUE-z] and almost 100% with RS+FD[SUE-z] when  $\epsilon = 10$ . This high accuracy with RS+FD[UE-z] protocols is because there is only one parameter to perturb each bit when generating fake data, i.e.,  $\Pr[0 \rightarrow 1] = q$  (cf. Section 2.2.4). When using different UE protocols, the randomization parameters  $p$  and  $q$  (cf. Section 2.2.4) also influence the attacker’s AIF-ACC, which led RS+FD[SUE] protocols to have lower attacker’s AIF-ACC when  $\epsilon$  is small, but higher attacker’s AIF-ACC in low privacy regimes.

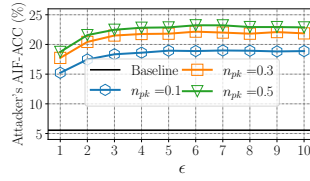
Lastly, we remark that due to the original formulation of RS+FD in [4] to generate fake data uniformly at random, a classifier was able to learn the sampled attribute from the users, as the distribution of the attributes was not always uniform with the ACSEmployment dataset. Nevertheless, when the attributes follow uniform-like distribution, none of the three attack models NK, PK or HM achieves a meaningful increment over the Baseline model (cf. results with the Nursery dataset [18] in Appendix D of [6]).

#### 4.4 Re-identification Risk of the RS+FD Solution

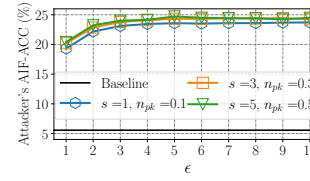
In this section, we experiment with multiple data collections following the RS+FD solution to measure the attacker’s RID-ACC. We follow a similar **experimental evaluation** of Section 4.2 with the addition of the attribute’s inference attack (cf. Section 4.3) in each data collection (i.e., survey). To this end, we use the NK model by generating  $s = 1n$  profiles as accuracy did not substantially increased with higher  $s$  (cf. Fig. 3). We selected the RS+FD[GRR] [4] protocol as it provides an intermediate guarantee between RS+FD[UE-r] (lower bound) and RS+FD[UE-z] (upper bound) protocols. We only evaluated the FK-RI model with  $\mathcal{D}_{BK}$  and uniform  $\epsilon$ -LDP privacy



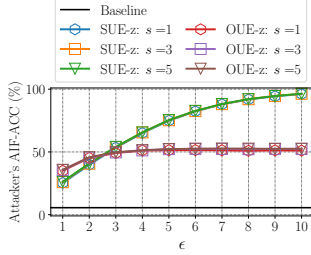
(a) NK model with RS+FD[GRR] protocol.



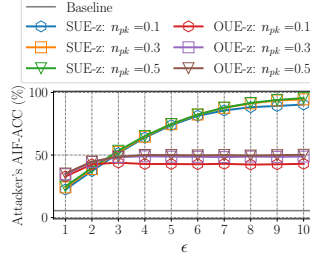
(b) PK model with RS+FD[GRR] protocol.



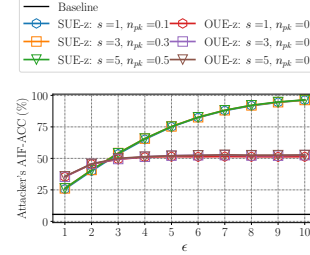
(c) Hybrid model with RS+FD[GRR] protocol.



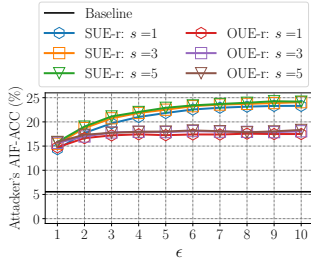
(d) NK model with RS+FD[UE-z] protocols.



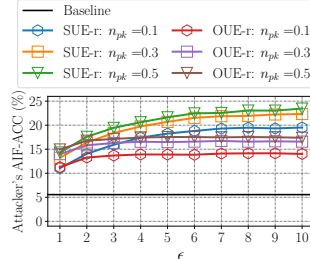
(e) PK model with RS+FD[UE-z] protocols.



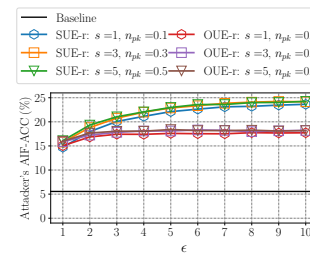
(f) Hybrid model with RS+FD[UE-z] protocols.



(g) NK model with RS+FD[UE-r] protocols.



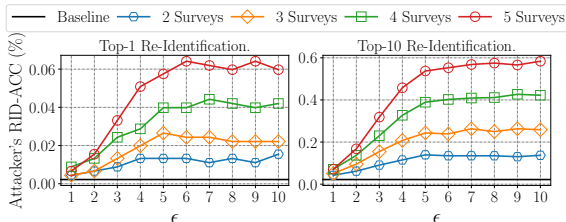
(h) PK model with RS+FD[UE-r] protocols.



(i) Hybrid model with RS+FD[UE-r] protocols.

**Figure 3: Attacker’s AIF-ACC on the ACSEmployment dataset with three attack models (i.e., NK, PK and hybrid) and five protocols (i.e., RS+FD[GRR], RS+FD[SUE-z], RS+FD[OUE-z], RS+FD[SUE-r] and RS+FD[OUE-r]), varying  $\epsilon$ , the number of synthetic profiles  $s$  the attacker generates and the number of compromised profiles  $n_{pk}$  the attacker has access to.**

metric across users (i.e., users select a new attribute for each survey) as they led to higher re-identification rates using the SMP solution. **Results.** Fig. 4 illustrates the attacker’s RID-ACC metric on the Adult dataset for top- $k$  re-identification using the FK-RI model and the RS+FD[GRR] protocol and by varying the uniform  $\epsilon$ -LDP privacy metric and the number of surveys.



**Figure 4: Attacker’s re-identification accuracy (RID-ACC) on the Adult dataset for top- $k$  re-identification using the FK-RI model and the RS+FD[GRR] protocol and by varying the uniform  $\epsilon$ -LDP privacy metric and the number of surveys.**

**Analysis.** From Fig. 4, one can note that the re-identification rates with RS+FD has drastically decreased in comparison with the results of the SMP solution in Fig. 2. Re-identification attacks on the RS+FD solution are not trivial, as the attacker has no guarantee that the predicted attribute is correct. Indeed, from Fig. 14 in Appendix D of [6], the attacker’s AIF-ACC on the Adult dataset with the RS+FD[GRR] protocol is upper bounded in 40%, which leads to chained errors when profiling a target user in multiple collections. For instance, the attacker’s RID-ACC for the top-1 group is nearly equal the random Baseline model. Even for the top-10 group the attacker’s RID-ACC has meaningful improvement over the Baseline model. These results with the RS+FD[GRR] protocol indicates that RS+FD is already (to some extent) a countermeasure to re-identification attacks, except for RS+FD[SUE-z] in which the attacker can predict the attribute with high confidence with high  $\epsilon$ .

## 5 COUNTERMEASURE

As shown in Section 4.4, the RS+FD solution already provides some resistance to re-identification attacks. Thus, we now present an improvement of the RS+FD solution and the experimental results.

## 5.1 Random Sampling Plus Realistic Fake Data

As briefly described in Section 2.3, the client-side of RS+FD [4] is split into two steps (*i.e.*, local randomization and *uniform* fake data generation). We now present an improvement of RS+FD, which we call Random Sampling Plus Realistic Fake Data (RS+RFD) as fake data will follow (potentially prior) *non-uniform* distributions. For instance, several demographic attributes have national statistics released by the Census [1] the previous year. Therefore, more “realistic” profiles can be generated by users to counter the inference of the sampled attribute and consequently the risk of re-identification. **Client-Side.** Alg. 1 displays the pseudocode of our RS+RFD solution at the client-side. The input of RS+RFD is the user’s true tuple of values  $\mathbf{v} = [v_1, v_2, \dots, v_d]$ , the domain size of attributes  $\mathbf{k} = [k_1, k_2, \dots, k_d]$ , the attributes’ prior distributions  $\tilde{\mathbf{f}} = [\tilde{f}_1, \tilde{f}_2, \dots, \tilde{f}_d]$  (transmitted by the server in advance), the privacy parameter  $\epsilon$  and a local randomizer  $\mathcal{M}$ . The output is a tuple  $\mathbf{y} = [y_1, y_2, \dots, y_d]$  of values (LDP and fake). In Alg. 1, line 6, *Sample* means a random sample is generated following prior  $\tilde{f}_i$  of the attribute  $i \in [d] \setminus \{j\}$ .

---

### Algorithm 1 Random Sampling plus Realistic Fake Data (RS+RFD)

---

**Input :** tuple  $\mathbf{v} = [v_1, v_2, \dots, v_d]$ , domain size of attributes  $\mathbf{k} = [k_1, k_2, \dots, k_d]$ , prior distribution of attributes  $\tilde{\mathbf{f}} = [\tilde{f}_1, \tilde{f}_2, \dots, \tilde{f}_d]$ , privacy parameter  $\epsilon$  and local randomizer  $\mathcal{M}$ .  
**Output :** sanitized tuple  $\mathbf{y} = [y_1, y_2, \dots, y_d]$ .

- 1:  $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$  ▷ Amplification by sampling [32]
- 2:  $j \leftarrow \text{Uniform}([d])$  ▷ Selection of attribute to sanitize
- 3:  $B_j \leftarrow \text{Encode}(v_j)$  ▷ Encode (if needed)
- 4:  $y_j \leftarrow \mathcal{M}(B_j, k_j, \epsilon')$  ▷ Sanitize data of the sampled attribute
- 5: **for**  $i \in [d] \setminus \{j\}$  **do** ▷ For each non-sampled attributes
- 6:      $y_i \leftarrow \text{Sample}(\{1, \dots, k_i\}, \tilde{f}_i)$  ▷ Generate one fake data
- 7: **end for**

**return :**  $\mathbf{y} = [y_1, y_2, \dots, y_d]$  ▷ Sanitized tuple

---

**Server-Side.** The aggregator performs multiple frequency estimation on the collected data by removing bias introduced by the local randomizer  $\mathcal{M}$  and fake data. The new estimators of using RS+RFD with GRR or UE-based protocols (*e.g.*, SUE [24] or OUE [51]) as local randomizer  $\mathcal{M}$  in Alg. 1 is presented in the following. For each attribute  $j \in [d]$ , the aggregator estimates  $\hat{f}(v_i)$  for the frequency of each value  $i \in [k_j]$  as:

- **RS+RFD[GRR].** The RS+RFD[GRR] estimator is:

$$\hat{f}_{\text{GRR}}(v_i) = \frac{dC(v_i) - n \left( q + (d-1)\tilde{f}_j(v_i) \right)}{n(p-q)}, \quad (6)$$

in which  $C(v_i)$  is the number of times  $v_i$  has been reported,  $\tilde{f}_j(v_i)$  is the prior distribution of value  $v_i \in A_j$ ,  $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$ ,  $p = \frac{e^{\epsilon'}}{e^{\epsilon'} + k_j - 1}$  and  $q = \frac{1-p}{k_j-1}$ . The probability tree of the RS+RFD[GRR] protocol, the proof that the estimator in Eq. (6) is unbiased and its variance computation are provided in Appendix A of [6].

- **RS+RFD[UE-r].** Similar to the RS+FD[UE-r] protocol in Section 2.3.2, in Line 6 of Alg. 1, for each non-sampled attribute  $i$ , for  $i \in [d] \setminus \{j\}$ , the user generates fake data by

applying a UE protocol to encoded random data following prior distribution  $\tilde{f}_i$ . The RS+RFD[UE-r] estimator is:

$$\hat{f}_{\text{UE-r}}(v_i) = \frac{dC(v_i) - n \left( q + (p-q)(d-1)\tilde{f}_j(v_i) + q(d-1) \right)}{n(p-q)}, \quad (7)$$

in which  $C(v_i)$  is the number of times  $v_i$  has been reported,  $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$  and  $\tilde{f}_j(v_i)$  is the prior distribution of value  $v_i \in A_j$ . Parameters  $p$  and  $q$  can be selected following the SUE [24] protocol ( $p = \frac{e^{\epsilon'/2}}{e^{\epsilon'/2} + 1}$  and  $q = \frac{1}{e^{\epsilon'/2} + 1}$ ) or OUE [51] protocol ( $p = \frac{1}{2}$  and  $q = \frac{1}{e^{\epsilon'} + 1}$ ). The probability tree of the RS+RFD[UE-r] protocol, the proof that the estimator in Eq. (7) is unbiased and its variance calculation is provided in Appendix B of [6].

**Privacy analysis.** Similar to the RS+FD solution [4], let  $\mathcal{M}$  be any existing LDP mechanism, Alg. 1 satisfies  $\epsilon$ -LDP, in a way that  $\epsilon' = \ln(d \cdot (e^\epsilon - 1) + 1)$ , in which  $d$  is the number of attributes.

**Limitations.** Besides known limits of the RS+FD solution [4, 47], RS+RFD adds a limitation on being dependent on the underlying prior distributions  $\tilde{\mathbf{f}}$  to generate realistic fake data. Yet, many demographic attributes have Census data [1] and other attributes’ priors can be defined following domain expert knowledge.

## 5.2 Experimental Results

In this section, we present the general setup of experiments with the RS+RFD solution, which includes: the frequency estimation of multiple attributes and the inference attack of the sampled attribute.

**5.2.1 General Experimental Setup.** We use the ACSEmployment dataset described in Section 4.1.

**Prior distribution.** To simulate “Correct” prior distributions  $\tilde{\mathbf{f}} = [\tilde{f}_1, \tilde{f}_2, \dots, \tilde{f}_d]$  to be used to generate realistic fake data with RS+RFD, we perturb the real frequency of each attribute  $j \in [d]$  with the standard Laplace mechanism [20–22] in centralized DP satisfying  $\epsilon = 0.1/d$  (*i.e.*, split  $\epsilon = 0.1$  by  $d$  attributes). In addition, to simulate an “Incorrect” scenario in which prior distributions are wrongly specified, we use Dirichlet distributions with parameter 1. **Methods evaluated.** We consider for evaluation three protocols within the RS+RFD solution from Section 5.1, namely, RS+RFD[GRR], RS+RFD[SUE-r] and RS+RFD[OUE-r].

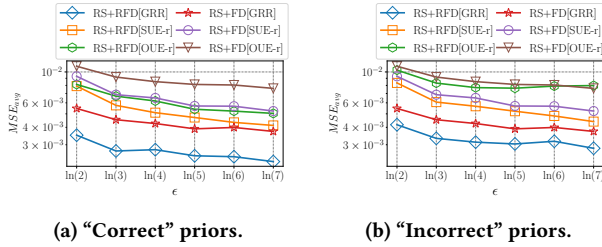
**5.2.2 Frequency Estimation of Multiple Attributes.** We compare the results of our RS+RFD protocols with their respective version within the RS+FD [4] solution, *i.e.*, RS+FD[GRR], RS+FD[SUE-r] and RS+FD[OUE-r] (*cf.* Section 2.3.2).

**Evaluation metrics.** To compare with [4], we vary  $\epsilon$  in the interval  $\epsilon = [\ln(2), \ln(3), \dots, \ln(7)]$  and we measure the quality of the estimated frequencies with the averaged mean squared error metric:  $MSE_{\text{avg}} = \frac{1}{d} \sum_{j \in [d]} \frac{1}{|A_j|} \sum_{v \in A_j} (f(v) - \hat{f}(v))^2$ .

**Results.** Fig. 5 illustrates for all methods the  $MSE_{\text{avg}}$  metric ( $y$ -axis) according to the privacy parameter  $\epsilon$  ( $x$ -axis) for both “Correct” and “Incorrect” priors. Additional empirical and analytical results with the Adult dataset are provided in Appendix E of [6].

**Analysis.** For the “Correct” prior, one can observe that the  $MSE_{\text{avg}}$  metric of our proposed RS+RFD protocols consistently and significantly outperform the utility of their respective version within the

RS+FD solution. The intuition is that since random noise is drawn from realistic prior distributions, the fake data also contributes to the estimation of the attribute. Indeed, even with “Incorrect” priors, our RS+RFD protocols still outperform the RS+FD protocols, with the exception of RS+RFD[OUE-r] with similar utility RS+RD[OUE-r] in low privacy regimes. On the other hand, when random noise follows uniform distributions, as with RS+FD, fake data can only increase the estimation of non-correct items.



**Figure 5: Averaged MSE metric varying  $\epsilon$  for (a) “Correct” and (b) “Incorrect” priors for multidimensional frequency estimation with the RS+RFD and RS+FD solutions.**

5.2.3 *Uncovering the Sampled Attribute of the RS+RFD Solution ( $\rightarrow$  SMP).* This section follows similar parameters (dataset,  $\epsilon$  range and attacker’s AIF-ACC metric) used in the experiments of Section 4.3.

**Results.** Fig. 6 illustrates the attacker’s AIF-ACC metric on the ACSEmployment dataset with three attack models (*i.e.*, NK, PK and hybrid) and our three protocols (*i.e.*, RS+RFD[GRR], RS+RFD[SUE-r] and RS+RFD[OUE-r] with “Correct” priors, varying  $\epsilon$ , the number of synthetic profiles  $s$  and the number of compromised profiles  $n_{pk}$ . Further results with “Incorrect” priors are in Appendix E of [6].

**Analysis.** We highlight that the non-stability in the plots of Fig. 6 is due to different sources of randomness:  $\epsilon = 0.1$ -DP for “Correct” prior distributions  $\hat{f}$ ,  $\epsilon$ -LDP randomization, fake data generation and the XGBoost algorithm. From Fig. 6, one can remark that our RS+RFD protocols considerably decrease the attacker’s AIF-ACC when comparing with their respective RS+FD version in Fig. 3. In contrast with the results of Section 4.3, the results with the PK model has higher attacker’s AIF-ACCs than the NK model. This is intuitive since the attacker gained “real” information of the sampled attribute, increasing the attacker’s AIF-ACC as the number of compromised profiles  $n_{pk}$  gets higher. Nevertheless, for all three NK, PK and HM models, the accuracy gain over a random Baseline model is still minor, highlighting the benefits of our RS+RFD proposal.

## 6 DISCUSSION

In brief, we identified and evaluated empirically two threats to users’ privacy when collecting multidimensional data with the state-of-the-art solutions SMP and RS+FD, namely re-identification attack and inference of the sampled attribute. These threats are generic to any LDP protocol and can be modelled by extending the “plausible deniability” attack analysis of Section 3.2.1. Hereafter, we summarize the key findings that can be used by practitioners and help substantiate the main claims of this paper.

Regarding the SMP solution, in our experiments, the GRR and  $\omega$ -SS protocol had the highest RID-ACC as the probability of accurately

inferring the user’s full profile was higher with relatively small  $k_j$  values (see also Fig. 1). With other protocols, such as OLH and OUE, which are the current state-of-the-art for preserving utility [51], the adversary cannot accurately infer the profile of users when using  $\epsilon$ -LDP as privacy model, which leads to lower re-identification risks (see Fig. 2 (c) and (d)). On the other hand, as shown in Appendix C of [6], when using the relaxed version of LDP from [36], the RID-ACC increases considerably for both OLH and OUE protocols. Though we only experimented with #surveys  $\leq 5$ , we believe that more data collections can lead to higher RID-ACC as long as the profile is accurately inferred. Yet, under standard sequential composition [22], the overall privacy loss is excessive when using high values for  $\epsilon$ , but we have also considered them due to their use in practical deployments [45, 46] and similar experiments found in [25] (though with higher #surveys  $\in \{7, 30, 90, 180\}$ ).

On the other hand, when using the RS+FD to “hide” the sampled attribute, the utility-oriented protocol RS+FD[UE-z] has the highest AIF-ACC due to generating fake data with zero-vectors and we recommend not using it in practice. Even with the RS+FD[GRR] or RS+FD[UE-r] protocols the attacker’s AIF-ACC is considerably greater than a random guess. Yet, since there are chained errors in multiple collections on accurately predicting the sampled attribute and on inferring the user’s value, the RS+FD considerably minimizes the risks of re-identification presented by the SMP solution.

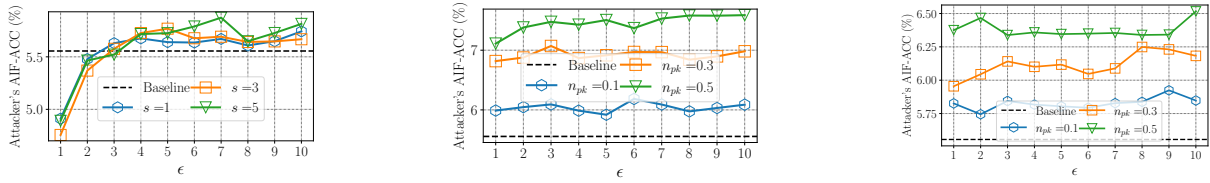
Overall, though some LDP protocols minimized the RID-ACC or AIF-ACC in our experiments (see the main body and Appendices C, D and E of [6]), they did not fully mitigate the risks when increasing  $\epsilon$  as done in practice to get more accurate estimations. This means they still allow a small portion of users to leak more information than others and corroborate with DP consensus of using  $\epsilon \leq 1$ .

Therefore, considering the setting described in Section 3.1, the overall recommendation when using the SMP solution is to select: the standard  $\epsilon$ -LDP as privacy model, the OUE and/or OLH protocols (depending on  $k_j$  due to communication costs [51]), the non-uniform privacy metric setting (*i.e.*, allowing users to sample with replacement and enforce memoization [5, 15, 24]) and to keep  $\epsilon \leq 1$ . On the other hand, when using the RS+FD solution, even when no prior is available, we highly recommend the proposed version in this paper, *i.e.*, RS+RFD with non-uniform fake data.

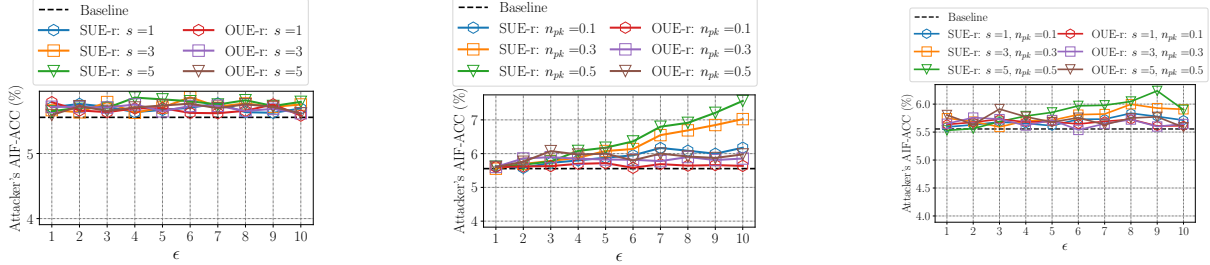
## 7 RELATED WORK

The literature on the local DP model has largely explored the issue of improving the utility of LDP protocols [4, 5, 15, 19, 24, 27, 28, 38, 46–48, 51–53]. Recently, a few works have started to design attacks on LDP protocols. Some authors focused on maliciously modifying the estimated statistic on the server through targeted or untargeted attacks [7, 11, 33, 55]. To counter such kinds of attacks, some works [3, 30] investigated cryptography-based approaches.

These targeted or untargeted attacks raise awareness of potential **security vulnerabilities** of LDP protocols. However, these attacks do not aim to attack **users’ privacy** as initially investigated in [9, 23, 25, 36] and in this work. For instance, Chatzikokolakis *et al.* [9] proposed the Bayes security measure to quantify the expected gain over a random guess of an adversary that observes a report of the RR protocol. Similar to [9], this paper provides a “plausible deniability” attacking interpretation of five state-of-the-art LDP



(a) NK model with RS+RFD[GRR] protocol. (b) PK model with RS+RFD[GRR] protocol. (c) Hybrid model with RS+RFD[GRR] protocol.



(d) NK model with RS+RFD[UE-r] protocols. (e) PK model with RS+RFD[UE-r] protocols. (f) Hybrid model with RS+RFD[UE-r] protocols.

**Figure 6: Attacker’s AIF-ACC on the ACSEmployment dataset with three attack models (i.e., NK, PK and hybrid) and our three protocols (i.e., RS+RFD[GRR], RS+RFD[SUE-r] and RS+RFD[OUE-r] with “Correct” priors), varying  $\epsilon$ , the number of synthetic profiles  $s$  the attacker generates and the number of compromised profiles  $n_{pk}$  the attacker has access to.**

protocols to infer the user’s true value by observing an LDP report. In an independent and concurrent work, Gursoy *et al.* [23] proposed a formalized Bayes adversary for the same attack, which was referred in Section 3.2.1 to give the expected accuracy of our analyses. Besides, we extended our attack to multiple collections of multidimensional data in Sections 3.2.2 and 3.2.3, which were proposed to account for the consequent risks of re-identification [26, 35–37] in Section 3.2.4. Re-identification risks in the LDP model for single-frequency estimation were first investigated by Murakami and Takahashi [36]. However, different from [36] that focused on a single attribute (e.g., location traces), our work considers multiple attributes being collected multiple times. Regarding multiple collections, Gadotti *et al.* [25] introduced pool inference attacks to LDP protocols for *single-frequency estimation* in a way that an adversary can infer the user’s preferred pool (e.g., skin tone used in emojis).

On the other hand, Arcolezzi *et al.* [4] introduced the RS+FD solution focusing only on the utility of the protocols, which was also later studied in [47]. In this work, we are the first to propose three attack models to the RS+FD solution, showing it is possible to distinguish the  $\epsilon$ -LDP report from fake data. Consequently, in multiple collections, we also show that RS+FD is still subject to (reduced) re-identification risks. We thus proposed an improvement of the RS+FD solution that generates non-uniform fake data (i.e., RS+RFD of Section 5.1) and can serve as a countermeasure solution.

## 8 CONCLUSION AND PERSPECTIVES

In this paper, we studied privacy threats against LDP protocols for multidimensional data following two state-of-the-art solutions for frequency estimation of multiple attributes, i.e., SMP and RS+FD [4]. On the one hand, we presented inference attacks based on “plausible deniability” [54] of five widely used LDP protocols (i.e., GRR [27, 28], OLH [51],  $\omega$ -SS [49, 56], RAPPOR [24] and OUE [51]) under

multiple collections following the SMP solution. This analysis also empirically clarifies the risks of re-identification when an attacker is able to build complete and/or partial profiles of users and can correlate them with prior knowledge.

In addition, we introduced three attack models to infer the sampled attribute of the RS+FD [4] solution, which allowed us to still reconstruct complete and/or partial profiles of users and lead to re-identification (although to a much lesser extent than the SMP solution). Finally, we proposed a refinement to the RS+FD solution, called RS+RFD that improves both utility and privacy. That is, in our experiments, RS+RFD minimized the estimation error in comparison with the RS+FD solution, as well as almost fully mitigated the inference of the sampled attribute attack.

Though we identified and investigated two privacy threats for LDP protocols for multidimensional data in single and multiple data collections, these are not unique and we believe that our work opens new avenues of research in this direction. For future work, we suggest and aim to formalize the re-identification risks considering different LDP and  $d$ -privacy [2, 8, 50] protocols, the number of collections, the number of attributes and the “uniqueness” of users in a given dataset. Such a formalization will allow to design other countermeasure solutions beyond RS+FD [4] and our RS+RFD.

## ACKNOWLEDGMENTS

The authors deeply thank the anonymous reviewers for their insightful suggestions. This work was partially supported by the ERC project HYPATIA with grant agreement N° 835294 and by the EIPHI-BFC Graduate School (contract “ANR-17-EURE-0002”). Sébastien Gamsb is supported by the Canada Research Chair program as well as a Discovery Grant from NSERC. All computations were performed on the “Mésocentre de Calcul de Franche-Comté”.

## REFERENCES

- [1] John M. Abowd. 2018. The U.S. Census Bureau Adopts Differential Privacy. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM. <https://doi.org/10.1145/3219819.3226070>
- [2] Mario Alvim, Konstantinos Chatzikokolakis, Catuscia Palamidessi, and Anna Pazi. 2018. Invited Paper: Local Differential Privacy on Metric Spaces: Optimizing the Trade-Off with Utility. In *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*. IEEE. <https://doi.org/10.1109/csf.2018.00026>
- [3] Andris Ambainis, Markus Jakobsson, and Helger Lipmaa. 2004. Cryptographic Randomized Response Techniques. In *Public Key Cryptography – PKC 2004*. Springer Berlin Heidelberg, 425–438. [https://doi.org/10.1007/978-3-540-24632-9\\_31](https://doi.org/10.1007/978-3-540-24632-9_31)
- [4] Héber H. Arcolezzi, Jean-François Couchot, Bechara Al Bouna, and Xiaokui Xiao. 2021. Random Sampling Plus Fake Data: Multidimensional Frequency Estimates With Local Differential Privacy. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (Virtual Event, Queensland, Australia) (CIKM '21)*. Association for Computing Machinery, New York, NY, USA, 47–57. <https://doi.org/10.1145/3459637.3482467>
- [5] Héber H. Arcolezzi, Jean-François Couchot, Bechara Al Bouna, and Xiaokui Xiao. 2022. Improving the utility of locally differentially private protocols for longitudinal and multidimensional frequency estimates. *Digital Communications and Networks* (2022). <https://doi.org/10.1016/j.dcan.2022.07.003>
- [6] Héber H. Arcolezzi, Sébastien Gambs, Jean-François Couchot, and Catuscia Palamidessi. 2022. On the Risks of Collecting Multidimensional Data Under Local Differential Privacy. *arXiv preprint arXiv:2209.01684* (2022).
- [7] Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. 2021. Data Poisoning Attacks to Local Differential Privacy Protocols. In *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, 947–964.
- [8] Konstantinos Chatzikokolakis, Miguel E. Andrés, Nicolás Emilio Bordenabe, and Catuscia Palamidessi. 2013. Broadening the Scope of Differential Privacy Using Metrics. In *Privacy Enhancing Technologies, Emiliano De Cristofaro and Matthew Wright (Eds.)*. Springer Berlin Heidelberg, Berlin, Heidelberg, 82–102. [https://doi.org/10.1007/978-3-642-39077-7\\_5](https://doi.org/10.1007/978-3-642-39077-7_5)
- [9] Konstantinos Chatzikokolakis, Giovanni Cherubini, Catuscia Palamidessi, and Carmela Troncoso. 2023. Bayes Security: A Not So Average Metric. In *2023 IEEE 36th Computer Security Foundations Symposium (CSF)*. IEEE Computer Society, 159–177. <https://doi.org/10.1109/CSF57540.2023.00011>
- [10] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM. <https://doi.org/10.1145/2939672.2939785>
- [11] Albert Cheu, Adam Smith, and Jonathan Ullman. 2021. Manipulation Attacks in Local Differential Privacy. In *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE. <https://doi.org/10.1109/sp40001.2021.00001>
- [12] Aloni Cohen. 2022. Attacks on Deidentification’s Defenses. In *31st USENIX Security Symposium (USENIX Security 22)*. USENIX Association, Boston, MA, 1469–1486.
- [13] Graham Cormode, Samuel Maddock, and Carsten Maple. 2021. Frequency estimation under local differential privacy. *Proceedings of the VLDB Endowment* 14, 11 (July 2021), 2046–2058. <https://doi.org/10.14778/3476249.3476261>
- [14] Damien Desfontaines. 2021. A list of real-world uses of differential privacy. Available online: <https://desfontain.es/privacy/real-world-differential-privacy.html> (accessed on 27 May 2022).
- [15] Bolin Ding, Janardhan Kulkarni, and Sergey Yekhanin. 2017. Collecting Telemetry Data Privately. In *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 3571–3580.
- [16] Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. 2021. Retiring adult: New datasets for fair machine learning. *Advances in Neural Information Processing Systems* 34 (2021).
- [17] Josep Domingo-Ferrer and Jordi Soria-Comas. 2018. Connecting randomized response, post-randomization, differential privacy and t-closeness via deniability and permutation. *arXiv preprint arXiv:1803.02139* (2018).
- [18] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. Available online: <http://archive.ics.uci.edu/ml> (accessed on 12 January 2023).
- [19] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. 2013. Local Privacy and Statistical Minimax Rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*. IEEE. <https://doi.org/10.1109/focs.2013.53>
- [20] Cynthia Dwork. 2006. Differential Privacy. In *Automata, Languages and Programming*, Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 1–12.
- [21] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating Noise to Sensitivity in Private Data Analysis. In *Theory of Cryptography*. Springer Berlin Heidelberg, 265–284. [https://doi.org/10.1007/11681878\\_14](https://doi.org/10.1007/11681878_14)
- [22] Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science* 9, 3–4 (2014), 211–407.
- [23] M. Emre Gursoy, Ling Liu, Ka-Ho Chow, Stacey Truex, and Wenqi Wei. 2022. An Adversarial Approach to Protocol Analysis and Selection in Local Differential Privacy. *IEEE Transactions on Information Forensics and Security* 17 (2022), 1785–1799. <https://doi.org/10.1109/TIFS.2022.3170242>
- [24] Úlfar Erlingsson, Vasyl Pihur, and Aleksandra Korolova. 2014. RAPPOR: Randomized Aggregatable Privacy-Preserving Ordinal Response. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security (Scottsdale, Arizona, USA)*. ACM, New York, NY, USA, 1054–1067. <https://doi.org/10.1145/2660267.2660348>
- [25] Andrea Gadotti, Florimond Houssiau, Meenatchi Sundaram Muthu Selva Annamalai, and Yves-Alexandre de Montjoye. 2022. Pool Inference Attacks on Local Differential Privacy: Quantifying the Privacy Guarantees of Apple’s Count Mean Sketch in Practice. In *31st USENIX Security Symposium (USENIX Security 22)*. USENIX Association, Boston, MA, 501–518.
- [26] Sébastien Gambs, Marc-Olivier Killijian, and Miguel Núñez del Prado Cortez. 2014. De-anonymization attack on geolocated data. *J. Comput. System Sci.* 80, 8 (2014), 1597–1614. <https://doi.org/10.1016/j.jcss.2014.04.024>
- [27] Peter Kairouz, Keith Bonawitz, and Daniel Ramage. 2016. Discrete distribution estimation under local privacy. In *International Conference on Machine Learning*. PMLR, 2436–2444.
- [28] Peter Kairouz, Sewoong Oh, and Pramod Viswanath. 2016. Extremal mechanisms for local differential privacy. *The Journal of Machine Learning Research* 17, 1 (2016), 492–542.
- [29] Shiva Prasad Kasiviswanathan, Homin K. Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. 2008. What Can We Learn Privately?. In *2008 49th Annual IEEE Symposium on Foundations of Computer Science*. 531–540. <https://doi.org/10.1109/FOCS.2008.27>
- [30] Fumiyuki Kato, Yang Cao, and Masatoshi Yoshikawa. 2021. Preventing Manipulation Attack in Local Differential Privacy Using Verifiable Randomization Mechanism. In *Data and Applications Security and Privacy XXXV*. Springer International Publishing, 43–60. [https://doi.org/10.1007/978-3-030-81242-3\\_3](https://doi.org/10.1007/978-3-030-81242-3_3)
- [31] Ninghui Li, Tiancheng Li, and Suresh Venkatasubramanian. 2007. t-Closeness: Privacy Beyond k-Anonymity and l-Diversity. In *2007 IEEE 23rd International Conference on Data Engineering*. IEEE. <https://doi.org/10.1109/icde.2007.367856>
- [32] Ninghui Li, Wahbeh Qardaji, and Dong Su. 2012. On sampling, anonymization, and differential privacy or, k-anonymization meets differential privacy. In *Proceedings of the 7th ACM Symposium on Information, Computer and Communications Security - ASIACCS '12*. ACM Press. <https://doi.org/10.1145/2414456.2414474>
- [33] Xiaoguang Li, Neil Zhenqiang Gong, Ninghui Li, Wenhai Sun, and Hui Li. 2022. Fine-grained Poisoning Attacks to Local Differential Privacy Protocols for Mean and Variance Estimation. *arXiv preprint arXiv:2205.11782* (2022).
- [34] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian. 2006. L-diversity: privacy beyond k-anonymity. In *22nd International Conference on Data Engineering (ICDE'06)*. IEEE. <https://doi.org/10.1109/icde.2006.1>
- [35] Takao Murakami, Atsunori Kanemura, and Hideitsu Hino. 2017. Group Sparsity Tensor Factorization for Re-Identification of Open Mobility Traces. *IEEE Transactions on Information Forensics and Security* 12, 3 (2017), 689–704. <https://doi.org/10.1109/TIFS.2016.2631952>
- [36] Takao Murakami and Kenta Takahashi. 2021. Toward Evaluating Re-identification Risks in the Local Privacy Model. *Transactions on Data Privacy* 14, 3 (2021), 79–116.
- [37] Arvind Narayanan and Vitaly Shmatikov. 2008. Robust De-anonymization of Large Sparse Datasets. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*. 111–125. <https://doi.org/10.1109/SP.2008.33>
- [38] Thông T Nguyễn, Xiaokui Xiao, Yin Yang, Siu Cheung Hui, Hyejin Shin, and Junbum Shin. 2016. Collecting and analyzing data from smart device users with local differential privacy. *arXiv preprint arXiv:1606.05053* (2016).
- [39] Xuebin Ren, Chia-mu Yu, Weiren Yu, Shusen Yang, Senior Member, Xinyu Yang, Julie A Mccann, Philip S Yu, and Life Fellow. 2018. LoPub : High-Dimensional Crowdsourced Data. 13, 9 (2018), 2151–2166. <https://doi.org/10.1109/TIFS.2018.2812146>
- [40] Ryan Rogers, Subbu Subramanian, Sean Peng, David Durfee, Seunghyun Lee, Santosh Kumar Kancha, Shraddha Sahay, and Parvez Ahammad. 2021. LinkedIn’s Audience Engagements API: A Privacy Preserving Data Analytics System at Scale. *Journal of Privacy and Confidentiality* 11, 3 (Dec. 2021). <https://doi.org/10.29012/jpc.782>
- [41] P. Samarati. 2001. Protecting respondents identities in microdata release. *IEEE Transactions on Knowledge and Data Engineering* 13, 6 (2001), 1010–1027. <https://doi.org/10.1109/69.971193>
- [42] Pierangela Samarati and Latanya Sweeney. 1998. Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and suppression. (1998).
- [43] Latanya Sweeney. 2002. k-Anonymity: A Model for Protecting Privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10, 05 (Oct. 2002), 557–570. <https://doi.org/10.1142/s0218488502001648>
- [44] Latanya Sweeney. 2015. Only you, your doctor, and many others may know. *Technology Science* 2015092903, 9 (2015), 29.

- [45] Jun Tang, Aleksandra Korolova, Xiaolong Bai, Xueqiang Wang, and Xiaofeng Wang. 2017. Privacy loss in apple’s implementation of differential privacy on macos 10.12. *arXiv preprint arXiv:1709.02753* (2017).
- [46] Apple Differential Privacy Team. 2017. Learning with privacy at scale. <https://docs-assets.developer.apple.com/ml-research/papers/learning-with-privacy-at-scale.pdf>. Online; accessed 11 December 2021.
- [47] Gatha Varma, Ritu Chauhan, and Dhananjay Singh. 2022. Sarve: synthetic data and local differential privacy for private frequency estimation. *Cybersecurity* 5, 1 (2022), 1–20. <https://doi.org/10.1186/s42400-022-00129-6>
- [48] Ning Wang, Xiaokui Xiao, Yin Yang, Jun Zhao, Siu Cheung Hui, Hyejin Shin, Junbum Shin, and Ge Yu. 2019. Collecting and Analyzing Multidimensional Data with Local Differential Privacy. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. IEEE. <https://doi.org/10.1109/icde.2019.00063>
- [49] Shaowei Wang, Liusheng Huang, Pengzhan Wang, Yiwen Nie, Hongli Xu, Wei Yang, Xiang-Yang Li, and Chunming Qiao. 2016. Mutual information optimally local private discrete distribution estimation. *arXiv preprint arXiv:1607.08025* (2016).
- [50] Shaowei Wang, Yiwen Nie, Pengzhan Wang, Hongli Xu, Wei Yang, and Liusheng Huang. 2017. Local private ordinal data distribution estimation. In *IEEE INFOCOM 2017 - IEEE Conference on Computer Communications*. IEEE. <https://doi.org/10.1109/infocom.2017.8056977>
- [51] Tianhao Wang, Jeremiah Blocki, Ninghui Li, and Somesh Jha. 2017. Locally Differentially Private Protocols for Frequency Estimation. In *26th USENIX Security Symposium (USENIX Security 17)*. USENIX Association, Vancouver, BC, 729–745.
- [52] Tianhao Wang, Ninghui Li, and Somesh Jha. 2018. Locally Differentially Private Frequent Itemset Mining. In *2018 IEEE Symposium on Security and Privacy (SP)*. IEEE. <https://doi.org/10.1109/sp.2018.00035>
- [53] Tianhao Wang, Milan Lopuhaa-Zwakenberg, Zitao Li, Boris Skorin, and Ninghui Li. 2020. Locally Differentially Private Frequency Estimation with Consistency. In *Proceedings 2020 Network and Distributed System Security Symposium*. Internet Society. <https://doi.org/10.14722/ndss.2020.24157>
- [54] Stanley L. Warner. 1965. Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias. *J. Amer. Statist. Assoc.* 60, 309 (March 1965), 63–69. <https://doi.org/10.1080/01621459.1965.10480775>
- [55] Yongji Wu, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. 2022. Poisoning Attacks to Local Differential Privacy Protocols for Key-Value Data. In *31st USENIX Security Symposium (USENIX Security 22)*. USENIX Association, Boston, MA, 519–536.
- [56] Min Ye and Alexander Barg. 2018. Optimal Schemes for Discrete Distribution Estimation Under Locally Differential Privacy. *IEEE Transactions on Information Theory* 64, 8 (2018), 5662–5676. <https://doi.org/10.1109/TIT.2018.2809790>
- [57] Zhikun Zhang, Tianhao Wang, Ninghui Li, Shibo He, and Jiming Chen. 2018. CALM: Consistent adaptive local marginal for marginal release under local differential privacy. *Proceedings of the ACM Conference on Computer and Communications Security* (2018), 212–229. <https://doi.org/10.1145/3243734.3243742>