

LDPGuard: Defenses against Data Poisoning Attacks to Local Differential Privacy Protocols

Kai Huang[§], Gaoya Ouyang[‡], Qingqing Ye[‡], Haibo Hu[‡], Bolong Zheng[†], Xi Zhao[§], Xiaofang Zhou[§]

[§]Department of Computer Science and Engineering, The Hong Kong University of Science and Technology

[‡]Department of Electronic and Information Engineering, Hong Kong Polytechnic University

[†]School of Computer Science and Technology, Huazhong University of Science and Technology

ustkhuang|xizhao|zxf@ust.hk, ouyanggaoya|zblchris@gmail.com, qqing.ye|haibo.hu@polyu.edu.hk

ABSTRACT

The protocols that satisfy Local Differential Privacy (LDP) enable untrusted third parties to collect aggregate information about a population without disclosing each user's privacy. In particular, each user locally encodes and perturbs his private data before the data is sent to the untrusted data collector, who aggregates and estimates the statistics about the population based on the collected perturbed values from individuals. Owing to the growing importance, LDP protocols have been widely studied and deployed in real-world scenarios (e.g., GOOGLE's RAPPOR). However, as data poisoning attacks may be injected by attackers who introduce many fake users, the utility of the statistics is heavily poisoned.

In this paper, we present a generic and extensible framework called LDPGuard to address the problem. LDPGuard provides effective defenses against data poisoning attacks to LDP protocols for frequency estimation, which is to estimate the frequency of the item of interest and is a basic task of other data analytics tasks (e.g., heavy hitter identification). In particular, it first precisely estimates the percentage of fake users and then provides adversarial schemes to defend against particular data poisoning attacks. Experimental study on real-world and synthetic datasets demonstrates the superiority of LDPGuard compared to existing techniques.

ACM Reference Format:

Kai Huang[§], Gaoya Ouyang[‡], Qingqing Ye[‡], Haibo Hu[‡], Bolong Zheng[†], Xi Zhao[§], Xiaofang Zhou[§]. 2023. LDPGuard: Defenses against Data Poisoning Attacks to Local Differential Privacy Protocols. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 16 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Local Differential Privacy (LDP) that enables an untrusted data collector to collect aggregate information about a population has been accepted as a *de facto* standard for data privacy in the local setting. In particular, each user locally encodes and perturbs his private data before sending the data to the untrusted data collector. After receiving the distributed perturbed data from users, the data collector designs detailed methods to aggregate and estimate the statistics

about the population without knowing the true values of each user. As the standard supports not only privacy-preserving data collection for untrusted third parties but also plausible deniability for each individual user, it has received a great amount of attention in both the research and industry community. In recent years, a lot of protocols [2–18] that satisfy Local Differential Privacy (*a.k.a.*, LDP protocols) have been developed. Some of the protocols have been deployed in industries such as Google, Apple, and Microsoft. In particular, Google integrated the RAPPOR [7] into Chrome browser to collect users' answers to questions such as default homepages for Chrome while preserving their privacy. Apple [19] deployed LDP protocols on iOS to gather users' preferences for emojis to identify popular emojis used in the population and recommend them to more users. Microsoft [20] developed LDP protocols to enable Windows 10 systems to aggregate statistics such as application usage. In addition, Samsung [21] also proposed LDP protocols to enable privacy-preserving data (e.g., screen resolution and battery volume) collection. Most these protocols are designed for frequency estimation, which is to estimate the frequency of the item of interest and is a basic task of more advanced data analytics tasks (e.g., heavy hitter identification [11]). Since each user involved in the protocol is required to perturb his value before sending to data collector, it largely sacrifices the utility of analytics results obtained by the data collector. Therefore, almost all existing protocols focused on designing effective strategies to improve the data utility.

However, these protocols are exposed to data poisoning attacks, which are injected by attackers who introduce fake users to LDP protocols to craft the data and manipulate the data analytics results as they desire. In particular, fake users first choose some items (*a.k.a.*, target items) and adopt various attacking methods to increase the estimated frequencies for the attacker-chosen items. Previous studies [22, 57] have already discovered that attackers have access to an amount of compromised accounts of web services (e.g., Hotmail, Twitter, and Google). In addition, they can purchase these compromised accounts from underground markets with very cheap prices (e.g., 0.004 – 0.03 for a Hotmail account and 0.03 – 0.50 for a phone verified Google account). In general, the main poisoning attacks consist of random perturbed-value attack (RPA), random item attack (RIA) and maximal gain attack (MGA) [22]. For RPA, each fake user randomly selects a value from the encoded space of LDP protocols and sends it to the data collector; Unlike RPA, each fake user involved in RIA selects a value from target items and perturbs it before sending it to a data collector; as for MGA, each fake user tries to maximize the frequency gains for the target items by crafting his value via solving a optimization problem (detailed in Section 3). These attacks are posing serious secure problems to

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

existing LDP protocols and threats to LDP-based data analytics. For example, if a attack imposes data poisoning attacks to LDP protocols for default homepages in Chrome, he can promote a phishing webpage as a default homepage of Chrome. Similarly, they can degrade the user experience on emojis by increasing the frequencies of some emojis via the data poisoning attacks. In addition, the popularity of malicious applications in Windows 10 systems could be increased in a similar way. In short, all these attacks can do damage to user experience and commercial interests. In recent years, two main countermeasures [22] including *Normalization* and *Fake users detection* are proposed for defending against these attacks. *Normalization* is to normalize the estimated item frequencies such that each estimated item frequency is non-negative and overall item frequencies sum to 1. *Fake users detection* is to detect possible fake users based on their submitted values and remove their impact on the final results. However, these countermeasures are facing two major challenges. The first is that the countermeasures are ad-hoc in nature, they can take effect on some scenarios but cannot adapt to others. The second is that they have limited effectiveness in reducing the impacts introduced by data poisoning attacks.

In this paper, we present a generic and extensible framework called LDPGuard to address the aforementioned challenges. First, LDPGuard provides such a unified framework that supports different adversarial schemes to offset the adverse effects of various data poisoning attacks. It supports state-of-the-art LDP protocols for frequency estimation such as kRR [10], OUE [13], and OLH [13]. Second, LDPGuard adopts the two-round collection method to precisely estimate the percentage of fake users to facilitate effective countermeasures. In summary, this paper makes the following contributions.

- We introduce a unified framework called LDPGuard for defenses against data poisoning attacks to local differential privacy protocols. LDPGuard works for all existing data poisoning attacks and state-of-the-art LDP protocols.
- We present an effective two-round collection method to precisely estimate the percentage of fake users to facilitate effective countermeasures.
- We conduct an extensive experimental study on real-world and synthetic datasets to demonstrate the superiority of LDPGuard compared to existing techniques.

The rest of the paper is organized as follows. Related research is discussed in Section 2. Section 3 provides the preliminaries. The LDPGuard framework is described in Section 4 and the estimation of the percentage of fake users is detailed in Section 5. Experimental results are presented in Section 6. Section 7 concludes the paper. Formal proofs of theorems and lemmas are in the Appendix [1].

2 RELATED WORK

Most germane to this research is traditional privacy and local differential privacy protocols, data poisoning attacks to LDP protocols, data poisoning attacks to machine learning, and countermeasures against data poisoning attacks to LDP protocols.

Traditional privacy and Local Differential Privacy protocols. Traditional privacy protocols mainly include generalization techniques [23–26, 28] and differential privacy [27]. Generalization techniques aim to generalize several values into a single one (*a.k.a.*,

an equivalence class) to hide the sensitive information of a single value so that the attacker can only see the generalized one. In particular, three classic generalization techniques have been proposed, namely, k -anonymity [23], ℓ -diversity [24] and t -closeness [25]. k -anonymity requires that each equivalence class contains at least k records so that each record is indistinguishable with at least $k - 1$ other records. ℓ -diversity requires that the distribution of a sensitive value in each equivalence class has at least ℓ “well-represented” values. While t -closeness requires the label distribution in each equivalence class is no more than t distance away from that in the whole set of values. Differential privacy [27], the privacy standard that provides semantic, information-theoretic guarantees on individuals’ privacy, is more stringent than generalization techniques as the former is defined regardless of the underlying dataset or apriori knowledge. Formally,

Definition 2.1. (Differential Privacy (DP)) A algorithm \mathcal{A} satisfies Differential Privacy (or ϵ -differential privacy (ϵ -DP)) iff for any two neighbouring datasets D and D' (i.e., D and D' differ in one element),

$$\forall y \in \text{Range}(\mathcal{A}) : \frac{\Pr(\mathcal{A}(D) = y)}{\Pr(\mathcal{A}(D') = y)} \leq e^\epsilon \quad (1)$$

where $\epsilon (\geq 0)$ is the privacy budget and $\text{Range}(\mathcal{A})$ the set of all possible outputs of the protocol \mathcal{A} .

The Laplace mechanism [27] that adds noise to the output of a numerical function is the first and most widely used mechanism for DP. In recent years, many follow-up algorithms for DP have been proposed [29, 30]. However, DP requires a trusted data curator to collect data from individual users and publish it. In real-world life, users do not fully trust the data curator for privacy. Therefore, much attention has been shifted to the local differential privacy (LDP) (the formal definition refers to Section 3), which can be cast as DP in the local setting and enables an untrusted data curator to collect aggregate information.

A long line of research has been introduced for LDP especially frequency estimation and heavy hitter identification under LDP protocols [2–18]. Randomized response is a classical technique for structured survey interview [31] and motivates the first formalized LDP model [32]. Duchi et al. [6] propose a treatment for mean estimation in location family models and convex risk minimization, providing lower and upper bounds for the estimation of population quantities that match up to constant factors. In contrast, Hsu et al. [33] adopt concentration of measure and random projection to aggregate and estimate heavy hitters. Similarly, Bassily et al. [4] provide an efficient protocol called random matrix projection and matching accuracy with lower bounds for succinct histogram estimation. In addition, some protocols extended from the randomized response mechanism have been developed. For example, RAPPOR framework [7] extends randomized response mechanism and adopts Bloom filter to defend against the adversary to infer users’ private information. The follow-up work [34] also extends RAPPOR to estimate more advanced statistics (e.g., joint-distributions, association testing, and categorical attributes). For these LDP protocols or frameworks such as RAPPOR and random matrix projection, the basic task is frequency estimation. The state-of-the-art LDP protocols for frequency estimation consist of kRR [10], OUE [13],

and OLH [13] (details refer to Section 3). For kRR [10], each user keeps his true value unchanged with a probability p and perturbs it to a different random value with probability q . For OUE [13], each user encodes his value into a binary vector and then flips (resp. keeps) each bit with an optimized probability q (resp. p). Tianhao et al. [13] further propose the optimized local hashing approach called OLH. These protocols focused on improving the data utility but may be exposed to data poisoning attacks, which are injected by attackers who introduce fake users to LDP protocols to craft the data and manipulate the data analytics results as they desire.

Data poisoning attacks to LDP protocols. Data poisoning attacks consist of targeted attacks [22] and untargeted attacks [35]. The targeted attack is to increase the estimated frequencies for the attacker-chosen items (*i.e.*, targets). It consists of random perturbed-value attack (RPA), random item attack (RIA), and maximal gain attack (MGA). For RPA, each fake user randomly selects a value from the encoded space of LDP protocols and sends it to the data collector; Unlike RPA, each fake user involved in RIA selects a value from the target item set and perturbs it before sending it to a data collector; as for MGA, each fake user tries to maximize the frequency gains for the target items by crafting his value via solving a optimization problem (detailed in Section 3). These attacks are posing serious secure problems to existing LDP protocols and threats to LDP-based data analytics. In contrast, the untargeted attack is to distort the item frequency distribution by manipulating L_p -norm distance between frequency vectors before and after attacking.

Data poisoning attacks to machine learning. There are a lot of works on data poisoning attacks to machine learning models including traditional machine learning models [36–39], deep networks [40–42] and recommendation systems [43–46]. In particular, [36] presents data poisoning attacks against autoregressive models. [37] develops data poisoning attacks to the classic a classification model (Support Vector Machines). Data poisoning attacks to byzantine robust federated learning are also presented in [38]. [40–42] focus on data poisoning attacks to deep neural networks. Another line of research focus on top-n recommendation systems [43], graph-based recommendation systems [44], factorization-based and collaborative filtering [45]. Observe that most attacks aim to poison the training data to train a machine learning model with bad performance, they are orthogonal to the problem studied in this paper as the latter focuses on developing effective countermeasures against data poisoning attacks to LDP protocols.

Countermeasures against data poisoning attacks to LDP protocols. Although there are some countermeasures against data poisoning attacks to LDP protocols [39, 47] and Sybil attacks [48–55], the former does not focus on LDP protocols for frequency estimation but machine learning models, the latter utilizes various information such as content, behavior, and social graphs to detect fake users in social networks. Recently, Xiaoyu et al. present three countermeasures against data poisoning attacks to LDP protocols for frequency estimation, namely, *Normalization*, *Fake users detection*, and *Conditional probability based Detection*. *Normalization* is to normalize the estimated item frequencies such that each estimated item frequency is non-negative and overall item frequencies sum to *Fake users detection* is to detect possible fake users based on their submitted values and remove their impact on the final results. *Conditional probability based Detection* is to detect possible fake users

Table 1: List of key notations.

Notation	Description
$\epsilon, \epsilon_1, \epsilon_2$	privacy budgets
$d, \{1, 2, \dots, d\}$	number of items, original data space
\mathcal{D}	encoded data space
Encode(\cdot), Perturb(\cdot)	Encode algorithm, Decode algorithm
PE(\cdot)	composition operation of Encode and Perturb
SUPPORT(y)	support set of the reported value y
p	$Pr[PE(v) \in \{y v \in \text{SUPPORT}(y)\}]$
q	$Pr[PE(v') \in \{y v \in \text{SUPPORT}(y)\}]$
p_1, q_1	p and q in the first round
p_2, q_2	p and q in the second round
p', q'	$p' = p_1$ or p_2 , $q' = q_1$ or q_2
$\mathbb{I}_{\text{SUPPORT}(y_i)}(v)$	indicator function
N, M	number of genuine users and fake users
$Y_{\text{true}}, Y_{\text{fake}}$	reported values from genuine users and fake users
y_i, Y	reported values from i -th user and all users
H, \mathcal{H}	a hash function, hash function family
r, τ	number of target items and chosen target items
$T = \{t_1, t_2, \dots, t_r\}$	target items
$\hat{f}_{t,b}, \hat{f}_{t,a}$	estimated frequencies of t before & after attacking
f_t	frequency of the target t over all genuine users
$\beta, \hat{\beta}$	percentage and estimated percentage of fake users

based on conditional probability, but it is applicable when there is only one target item. Moreover, the former two countermeasures are facing two major challenges. The first is that the countermeasures are ad-hoc in nature, they can take effect on some scenarios but cannot adapt to others. The second is that they have limited effectiveness in reducing the impacts introduced by data poisoning attacks.

3 PRELIMINARIES

Table 1 lists the key notations and acronyms used in this paper. Recall that LDP enables an untrusted **data collector** to collect aggregate information of **users** who are willing to help with data collecting but do not fully trust the data collector for privacy concerns. Without loss of generality, we assume that there are N users, each has a value $v \in \{1, 2, 3, \dots, d\}$ where $\{1, 2, 3, \dots, d\}$ is **original data space** and d the **number of items**. This paper focuses on LDP protocols for frequency estimation. That is, the data collector aims to aggregate and estimate the frequencies of values among the N users. To this end, the following three algorithms are performed one by one:

- **Encode:** A user who holds a value $v \in \{1, 2, 3, \dots, d\}$ first encodes v to $\text{Encode}(v) \in \mathcal{D}$ where \mathcal{D} is **encoded data space** (*i.e.*, the space of encoded values).
- **Perturb:** The user then perturbs the encoded value $\text{Encode}(v)$ to $\text{Perturb}(\text{Encode}(v))$ and submits the perturbed value to the data collector. We can also use $\text{PE}(\cdot)$ to denote the composition operation of the Encode and Perturb algorithms for simplification.
- **Aggregate:** The data collector receives the reported values from users and then estimates the statistics of interest.

The Perturb algorithm should satisfy the following ϵ -local differential privacy (ϵ -LDP).

Definition 3.1. (Local Differential Privacy (LDP)) A algorithm \mathcal{A} satisfies Local Differential Privacy (or ϵ -local differential privacy (ϵ -LDP)) iff for any two inputs v and v' users hold,

$$\forall y \in \text{Range}(\mathcal{A}) : \frac{\Pr[\mathcal{A}(v) = y]}{\Pr[\mathcal{A}(v') = y]} \leq e^\epsilon \quad (2)$$

where $\epsilon (\geq 0)$ is the privacy budget and $\text{Range}(\mathcal{A})$ the set of all possible outputs of the protocol \mathcal{A} .

For example, the randomized response [31] reports the true value with probability p and flips it with probability $1 - p$, let $\frac{p}{1-p} \leq e^\epsilon$, randomized response satisfies $\ln \frac{p}{1-p}$ -LDP.

To precisely analyze the accuracy of LDP protocols, the **Pure LDP** that supports simple and fast aggregation is developed in [13]. Formally,

Definition 3.2. (Pure Local Differential Privacy (Pure LDP))

A protocol \mathcal{A} with an Encode and Perturb function (denoted by $\text{PE}(\cdot)$) satisfies Pure Local Differential Privacy (or ϵ -pure local differential privacy (ϵ -Pure LDP)) iff for any input value v a user holds,

$$\forall y \in \text{Range}(\mathcal{A}) : \frac{p}{q} \leq e^\epsilon$$

such that

$$\begin{aligned} \Pr[\text{PE}(v) \in \{y | v \in \text{SUPPORT}(y)\}] &= p \\ \forall v' \neq v : \Pr[\text{PE}(v') \in \{y | v \in \text{SUPPORT}(y)\}] &= q \end{aligned} \quad (3)$$

where $\text{SUPPORT}(y)$ is the set of values that y supports [13], $\epsilon (\geq 0)$ the privacy budget and $\text{Range}(\mathcal{A})$ the set of all possible outputs of the protocol \mathcal{A} .

In particular, $\{y | v \in \text{SUPPORT}(y)\}$ contains all outputs $\{y\}$ that “support” the occurrences of v (the support set of v for short). Note that the definition of the support set depends on a particular LDP protocol. For example, for the basic RAPPOR protocol [7] that encodes a value to a binary vector B whose v -th bit is 1, $\text{SUPPORT}(B) = \{v | B[v] = 1\}$.

When each user follows a pure LDP to report his value y_i ($i \in \{1, 2, \dots, d\}$) to the data collector, the estimated frequency of a value v is

$$\tilde{f}_v = \frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(v) - q}{p - q} \quad (4)$$

where $\mathbb{I}_{\text{SUPPORT}(y_i)}(v)$ is an indicator function. If $v \in \text{SUPPORT}(y_i)$, $\mathbb{I}_{\text{SUPPORT}(y_i)}(v) = 1$, and 0 otherwise. The estimated frequency \tilde{f}_v is an unbiased estimation of the true frequency f_v of v (see Lemma 3.3), thus, $\mathbb{E}[\tilde{f}_v] = f_v$ and

$$\sum_{i=1}^N \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(v)] = N(f_v(p - q) + q). \quad (5)$$

LEMMA 3.3. $\tilde{f}_v = \frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(v) - q}{p - q}$ is the unbiased estimation of the true frequency f_v of the item v .

3.1 Local Differential Privacy Protocols

As discussed in Section 2, the state-of-the-art LDP protocols for frequency estimation consist of kRR [10], OUE [13], and OLH [13]. Their formal definitions and corresponding procedures are given in this section.

3.1.1 kRR.

Encode: Each user encodes his value v to itself, i.e., $\text{Encode}(v) = v$. The encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$, which is the same as the original data space.

Perturb: The user keeps the true value v with probability p and perturbs it to a different value $v' \in \mathcal{D}$ with probability q , then, reports the perturbed value $y = \text{PE}(v)$ to the data collector, formally,

$$\Pr[y = a] = \begin{cases} \frac{e^\epsilon}{e^\epsilon + d - 1} \triangleq p & \text{if } a = v \\ \frac{1}{e^\epsilon + d - 1} \triangleq q & \text{otherwise} \end{cases} \quad (6)$$

Aggregate: The data collector receives the reported values $\{y_i | i \in \{1, 2, \dots, d\}\}$ from users, and estimates the frequency of a particular item v with Equation (4). In particular, the support set $\text{SUPPORT}(y_i)(v)$ is $\{y | y = v\}$.

3.1.2 OUE.

Encode: Each user encodes his value v to a d -dimensional binary vector B where v -th bit is 1, i.e., $\text{Encode}(v) = B$ where $B[v] = 1$ and $B[v'] = 0$ for $v' \neq v$. The encoded data space $\mathcal{D} = \{0, 1\}^d$ since each bit could be 0 or 1.

Perturb: The user keeps each bit with value 1 (resp. 0) in the binary vector B with probability p (resp. $1 - q$) and flips it with probability $1 - p$ (resp. q), then, reports the perturbed value $y = \text{PE}(v)$ (or equally, $y = \text{Perturb}(B)$) to the data collector, formally,

$$\Pr[y[pos] = 1] = \begin{cases} \frac{1}{2} \triangleq p & \text{if } B[pos] = 1 \\ \frac{1}{e^\epsilon + 1} \triangleq q & \text{if } B[pos] = 0 \end{cases} \quad (7)$$

Aggregate: The data collector receives the reported values $\{y_i | i \in \{1, 2, \dots, d\}\}$ from users, and estimates the frequency of a particular item v with Equation (4). The support set $\text{SUPPORT}(y_i)(v)$ for OUE is $\{y | y[v] = 1\}$.

3.1.3 OLH.

Encode: Given a universal hash function family \mathcal{H} such that the range of $H \in \mathcal{H}$ belongs to $\{1, 2, \dots, d'\}$, each user randomly selects a hash function H from \mathcal{H} and encodes his value v to $\langle H, H(v) \rangle$. The encoded data space $\mathcal{D} = \{\langle H, h \rangle | H \in \mathcal{H}, h \in \{1, 2, \dots, d'\}\}$. Note that \mathcal{H} is usually generated by the hash algorithm called xxHash [56], which implements the hash function family with different random seeds. For a better estimation performance, we follow existing work [13, 22] to set $d' = e^\epsilon + 1$.

Perturb: The user keeps the $H(v)$ with probability p and perturbs it to a different value $H(v)' \in \{1, 2, \dots, d'\}$ with probability q , then, reports the perturbed value $y = \text{PE}(v) = \langle H, H(v)' \rangle$ to the data collector, formally,

$$\Pr[y = \langle H, h \rangle] = \begin{cases} \frac{e^\epsilon}{e^\epsilon + d' - 1} \triangleq p^* & \text{if } h = H(v) \\ \frac{1}{e^\epsilon + d' - 1} \triangleq q^* & \text{if } h \neq H(v) \end{cases} \quad (8)$$

Aggregate: The data collector receives the reported values $\{y_i | i \in \{1, 2, \dots, d\}\}$ from users, and estimates the frequency of a particular item v with Equation (4) with the support set $\text{SUPPORT}(y_i)(v) = \{y | y = \langle H, h \rangle \text{ and } H(v) = h\}$. Note that p^* and q^* are for hash function H instead of the pair of $\langle H, H(\cdot) \rangle$. As each hash function is randomly selected from \mathcal{H} , the overall perturbation probability parameters are $p = p^* = \frac{e^\epsilon}{e^\epsilon + d' - 1}$ and $q = \frac{1}{d'} p^* + (1 - \frac{1}{d'}) q^* = \frac{1}{d'}$.

3.2 Threat Model and Data Poisoning Attacks

Data poisoning attack is to increase the estimated frequencies for the attacker-chosen items (*i.e.*, targets). In particular, given a target set with r target items, $T = \{t_1, t_2, \dots, t_r\}$ where $t_j \in \{1, 2, \dots, d\}$, suppose the estimated frequencies of a particular item t before and after attacking are $\tilde{f}_{t,b}$ and $\tilde{f}_{t,a}$ respectively, data poisoning attack is to injects M fake users to increase the **overall frequency gain** over T , *i.e.*, $\sum_{t \in T} \Delta \tilde{f}_t = \tilde{f}_{t,a} - \tilde{f}_{t,b}$ where $\Delta \tilde{f}_t$ is the **frequency gain** of a single target item t .

In general, it consists of **random perturbed-value attack (RPA)**, **random item attack (RIA)**, and **maximal gain attack (MGA)** [22]. For RPA, each fake user randomly selects a value from the encoded space of LDP protocols and sends it to the data collector; Unlike RPA, each fake user involved in RIA selects a value from the target item set and perturbs it before sending it to a data collector; as for MGA, each fake user tries to maximize the frequency gains for the target items by crafting his value via solving the following optimization problem. Let the reported values from M fake users be Y_{fake} , and its overall frequency gain over T be $Gain(Y_{fake})$ MGA is to maximize the optimization problem:

$$\begin{aligned}
 \text{Max}_{Y_{fake}} \text{Gain}(Y_{fake}) &\iff \text{Max}_{Y_{fake}} \sum_{t \in T} \mathbb{E}[\Delta \tilde{f}_t] \\
 &\iff \text{Max}_{Y_{fake}} \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] \iff \text{Max}_{Y_{fake}} \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\
 \text{Equ(4)} &\iff \text{Max}_{Y_{fake}} \sum_{t \in T} \mathbb{E} \left[\frac{\frac{1}{N+M} \sum_{i=1}^{N+M} \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q}{p - q} \right] - \\
 &\quad \mathbb{E} \left[\frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q}{p - q} \right] \quad (\text{according to Equation (4)}) \\
 &\iff \text{Max}_{Y_{fake}} \sum_{t \in T} \mathbb{E} \left[\frac{\sum_{i=N+1}^{N+M} \mathbb{I}_{\text{SUPPORT}(y_i)}(t)}{(N+M)(p-q)} \right] - \mathbb{E} \left[\frac{M \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(t)}{N(N+M)(p-q)} \right] \\
 &\iff \text{Max}_{Y_{fake}} \sum_{t \in T} \frac{\sum_{i=N+1}^{N+M} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)]}{(N+M)(p-q)} - \frac{M \sum_{i=1}^N \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)]}{N(N+M)(p-q)} \\
 \text{Equ(5)} &\iff \text{Max}_{Y_{fake}} \frac{\sum_{i=N+1}^{N+M} \sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)]}{(N+M)(p-q)} - \frac{\sum_{t \in T} MN(f_t(p-q) + q)}{N(N+M)(p-q)}, \quad (9)
 \end{aligned}$$

where f_t is the frequency of the target t over all genuine users. Observe that $\frac{\sum_{t \in T} MN(f_t(p-q) + q)}{N(N+M)(p-q)}$ only depends on genuine users,

MGA is to craft Y_{fake} such that $\frac{\sum_{i=N+1}^{N+M} \sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)]}{(N+M)(p-q)}$ is maximized. To this end, different implementation strategies are adopted for kRR, OUE, and OLH as listed below.

MGA for kRR. For kRR, observe that $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)] \leq 1$ and $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)] = 1$ if and only if y_i belongs to T . Therefore, MGA for kRR is obtained by randomly selecting a target item for each fake user.

MGA for OUE. For OUE, observe that $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)] \leq r$ and $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)] = r$ if and only if $\forall t \in T, y_i[t] = 1$. Intuitively, implementing MGA for OUE only needs to craft each

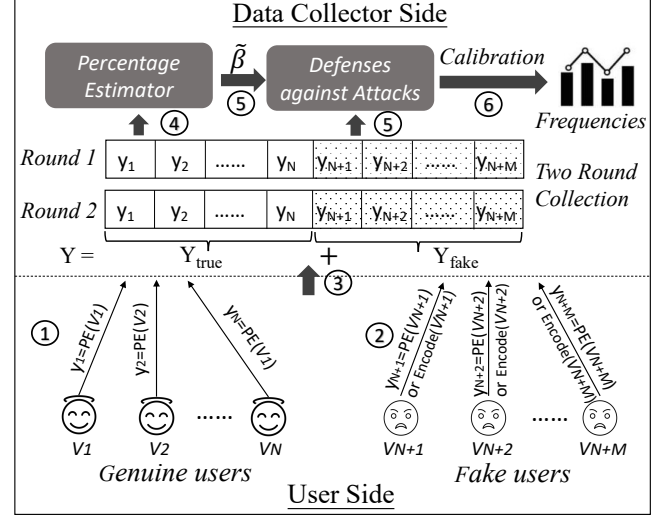


Figure 1: The LDPGuard framework.

fake user's value with $y_i[t] = 1$ for each target t , but this will result in only r bits 1 in the encoded binary vector. In contrast, each genuine user may report y_i with $\lfloor p + (d-1)q \rfloor$ bits 1. As such, MGA for OUE reports the binary vector with $y_i[t] = 1$ for each target t and $l = \lfloor p + (d-1)q - r \rfloor$ bits 1 for non-target bits.

MGA for OLH. For OLH, observe that $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)] \leq r$ and $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUPPORT}(y_i)}(t)] = r$ if and only if the selected hash function hashes all targets to the same value. In our implementation, we follow [22] to sample 1,000 hash functions from xxHash and select the one that results in the most targets to the same value.

4 LDPGUARD: THE FRAMEWORK

In this section, we overview our LDPGuard framework towards effective defenses against data poisoning attacks to local differential privacy protocols. In general, LDPGuard first precisely estimates the percentage of fake users, it then utilizes the estimated percentage of fake users to provide adversarial schemes to defend against particular data poisoning attacks including RPA, RIA, and MGA.

4.1 Overview of LDPGuard

The workflow of LDPGuard is depicted in Figure 1. Suppose there are N **genuine users**, each user holds a value v_i ($i \in \{1, 2, \dots, N\}$) in the **original data space** (*i.e.*, $v_i \in \{1, 2, \dots, d\}$), encodes (*s.t.*, $\text{Encode}(v_i)$ belongs to **encoded data space** \mathcal{D}), perturbs v_i to satisfy LDP with privacy budget ϵ_1 , and output the perturbed value y_i such that $y_i = \text{PE}(v_i)$. The perturbed value y_i is then reported to the data collector (step ①). In addition, there are M **fake users**, each fake user holds a value v_i ($i \in \{N+1, N+2, \dots, N+M\}$) in the original data space (*i.e.*, $v_i \in \{1, 2, \dots, d\}$), and adopts one of the following procedures to obtained the reported value y_i (step ②).

- **Encode only.** The fake user only encodes the input v_i to obtain y_i , *i.e.*, $y_i = \text{Encode}(v_i)$. This applies to both RPA and MGA attacks, which do not utilize any LDP protocols to perturb inputs. In particular, each fake user launched RPA randomly selects a value from the encoded data space of

LDP protocols and sends it to the data collector, similarly, one involved in MGA is to craft his value via solving the optimization problem (see Equation (9)) without perturbing the input value.

- **Encode and Perturb.** The fake user first encodes the input v_i and then perturbs the encoded value to obtain y_i , i.e., $y_i = \text{PE}(v_i)$. This applies to RIA attack where each fake user first selects a value from the target item set and then perturbs it before sending it to the data collector.

Therefore, the reported value y_i ($i \in \{N+1, N+2, \dots, N+M\}$) from a fake user is either $y_i = \text{PE}(v_i)$ or $y_i = \text{Encode}(v_i)$. The values $\{y_i | i \in \{1, 2, \dots, N+M\}\}$ are then reported to the data collector (step ③) for further processing.

Upon receiving the reported values Y consisting of $Y_{\text{true}} = \{y_i | i \in \{1, 2, \dots, N\}\}$ from genuine users and $Y_{\text{fake}} = \{y_i | i \in \{N+1, N+2, \dots, N+M\}\}$ from fake users in the first round, the second round with the same process is launched by the data collector (step ③). Note that the privacy budget in the second round is ϵ_2 such that $\epsilon_1 + \epsilon_2 = \epsilon$ (by default, $\epsilon_1 = \epsilon_2 = \frac{1}{2}\epsilon$). **We call this entire process “Two Round Collection”.** Then, the data collector utilizes **percentage estimator (detailed in Section 5)** to estimate the percentage of fake users (step ④) by comparing the reported values Y in the two rounds (denoted by $Y^1 = Y_{\text{true}}^1 + Y_{\text{fake}}^1$ and $Y^2 = Y_{\text{true}}^2 + Y_{\text{fake}}^2$ where “+” is a concatenation operator). The estimated percentage of fake users (denoted by $\tilde{\beta}$) is then used to design effective **defenses against attacks (detailed in Section 4.2)** (step ⑤). Finally, the data collector adopts calibration strategies to output the estimated item frequencies (step ⑥).

4.2 Defenses against Data Poisoning Attacks

In this section, we assume the estimated percentage $\tilde{\beta}$ of fake users is already derived and will remove this assumption in Section 5). Based on estimated percentage $\tilde{\beta}$, we can provide adversarial schemes to defend against data poisoning attacks including RPA, RIA, and MGA to start-of-the-art LDP protocols such as kRR, OUE, and OLH. The motivation behind the adversarial schemes is offset the adverse effects of various data poisoning attacks with the operations in the opposite direction to the attacks.

4.2.1 Defenses against Data Poisoning Attacks to kRR.

We begin by setting $p_1 = \frac{e^{\epsilon_1}}{e^{\epsilon_1+d}-1}$, $q_1 = \frac{1}{e^{\epsilon_1+d}-1}$, and $p_2 = \frac{e^{\epsilon_2}}{e^{\epsilon_2+d}-1}$, $q_2 = \frac{1}{e^{\epsilon_2+d}-1}$ where $\epsilon_1 + \epsilon_2 = \epsilon$. In addition, let $p = \frac{e^{\epsilon}}{e^{\epsilon+d}-1}$, $q = \frac{1}{e^{\epsilon+d}-1}$.

Defenses against RPA to kRR. First, LDPGuard randomly samples $(N+M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$. For each sampled value, LDPGuard then removes the record $v \in \tilde{Y}_{\text{fake}}$ from all reported Y if $v \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency for item v with

$$\tilde{f}_v = \frac{\frac{1}{|Y \setminus \tilde{Y}_{\text{fake}}|} \sum_{y_i \in Y \setminus \tilde{Y}_{\text{fake}}} \mathbb{I}_{\text{SUPPORT}(y_i)}(v) - q'}{p' - q'} \quad (10)$$

where $p' = p_1$ or p_2 , $q' = q_1$ or q_2 , and $Y \setminus \tilde{Y}_{\text{fake}}$ is the remaining records in Y after removing records in \tilde{Y}_{fake} .

The expectation of \tilde{f}_v is

$$\begin{aligned} \mathbb{E}[\tilde{f}_v] &= \frac{\mathbb{E}\left[\frac{1}{|Y \setminus \tilde{Y}_{\text{fake}}|} \sum_{y_i \in Y \setminus \tilde{Y}_{\text{fake}}} \mathbb{I}_{\text{SUPPORT}(y_i)}(v) - q'\right]}{p' - q'} \\ &= \frac{\frac{N_1(f_v(p'-q')+q')+N_2f_v^*}{N_1+N_2} - q'}{p' - q'} \end{aligned} \quad (11)$$

where $N_1 = |Y_{\text{true}} \cap (Y \setminus \tilde{Y}_{\text{fake}})|$, $N_2 = |Y \setminus \tilde{Y}_{\text{fake}} \setminus Y_{\text{true}}|$, and $f_v^* = \frac{1}{d}$ is the probability that v equals the value sampled from $Y \setminus \tilde{Y}_{\text{fake}} \setminus Y_{\text{true}}$. Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{\text{fake}} = Y_{\text{true}}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$. The overall frequency gain $\text{Gain}(Y_{\text{fake}})$ is :

$$\begin{aligned} \text{Gain}(Y_{\text{fake}}) &= \sum_{t \in T} \mathbb{E}[\tilde{f}_t] - \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\ &\stackrel{\text{Equ(4)}}{=} \sum_{t \in T} \mathbb{E}\left[\frac{\frac{1}{|Y \setminus \tilde{Y}_{\text{fake}}|} \sum_{i \in Y \setminus \tilde{Y}_{\text{fake}}} \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q'}{p' - q'}\right] - \\ &\quad \mathbb{E}\left[\frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q}{p - q}\right] \quad (\text{according to Equation (4)}) \\ &\stackrel{\text{Equ(11)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*1/d}{N_1+N_2} - q'}{p' - q'} - \mathbb{E}\left[\frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q}{p - q}\right] \\ &\stackrel{\text{Equ(5)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*1/d}{N_1+N_2} - q'}{p' - q'} - \frac{(f_t(p-q)+q) - q}{p - q} \\ &= \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*1/d}{N_1+N_2} - f_t(p' - q')}{p' - q'}. \end{aligned} \quad (12)$$

Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{\text{fake}} = Y_{\text{true}}$), $N_2 = 0$, the overall frequency gain $\text{Gain}(Y_{\text{fake}})$ is 0.

Defenses against RIA to kRR. For RIA, LDPGuard first randomly samples $(N+M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from target set $T = \{t_1, t_2, \dots, t_r\}$, and then encodes each value to the encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$. For the encoded data, it keeps the true value with probability $p = p'$ ($p' = p_1$ or p_2) and perturbs it to a different value in \mathcal{D} with probability $q = q'$ ($q' = q_1$ or q_2). Finally, LDPGuard removes the perturbed value from all reported Y if the perturbed value is found in Y , and calibrates the results to obtain the estimated frequency \tilde{f}_v for item v with Equation (10). According to Equation (11), the expectation of \tilde{f}_v is

$$\mathbb{E}[\tilde{f}_v] = \frac{\frac{N_1(f_v(p'-q')+q')+N_2f_v^*}{N_1+N_2} - q'}{p' - q'}, f_v^* = \begin{cases} \frac{1}{r}(p' - q') + q' & \text{if } v \in T \\ q' & \text{if } v \notin T \end{cases} \quad (13)$$

where $N_1 = |Y_{\text{true}} \cap (Y \setminus \tilde{Y}_{\text{fake}})|$, $N_2 = |Y \setminus \tilde{Y}_{\text{fake}} \setminus Y_{\text{true}}|$, and f_v^* is the probability that v equals to the value sampled from $Y \setminus \tilde{Y}_{\text{fake}} \setminus Y_{\text{true}}$. Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{\text{fake}} = Y_{\text{true}}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$. The overall

frequency gain $\text{Gain}(Y_{fake})$ is :

$$\begin{aligned}
 \text{Gain}(Y_{fake}) &= \sum_{t \in T} \mathbb{E}[\Delta \tilde{f}_t] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\
 &\stackrel{\text{Equ(4)}}{=} \sum_{t \in T} \mathbb{E} \left[\frac{\frac{1}{|Y \setminus \tilde{Y}_{fake}|} \sum_{i \in Y \setminus \tilde{Y}_{fake}} \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q'}{p' - q'} \right] - \\
 &\quad \mathbb{E} \left[\frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q}{p - q} \right] \quad (\text{according to Equation (4)}) \\
 &\stackrel{\text{Equ(5)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2(\frac{1}{r}(p' - q') + q')}{N_1 + N_2} - q'}{p' - q'} - \frac{(f_t(p - q) + q) - q}{p - q} \\
 &\stackrel{\text{Equ(13)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2(1/r(p' - q') + q')}{N_1 + N_2} - f_t(p' - q')}{p' - q'}. \quad (14)
 \end{aligned}$$

Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, the overall frequency gain $\text{Gain}(Y_{fake})$ is 0.

Defenses against MGA to kRR. Unlike RPA, LDPGuard first randomly samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from target set $T = \{t_1, t_2, \dots, t_r\}$ instead of \mathcal{D} . For each sampled value, LDPGuard then removes the record $v \in \tilde{Y}_{fake}$ from all reported Y if $v \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency \tilde{f}_v for item v with Equation (10). According to Equation (11), the expectation of \tilde{f}_v is

$$\mathbb{E}[\tilde{f}_v] = \frac{\frac{N_1(f_v(p' - q') + q') + N_2 f_v^*}{N_1 + N_2} - q'}{p' - q'}, f_v^* = \begin{cases} \frac{1}{r} & \text{if } v \in T \\ 0 & \text{if } v \notin T \end{cases} \quad (15)$$

where $N_1 = |Y_{true} \cap (Y \setminus \tilde{Y}_{fake})|$, $N_2 = |Y \setminus \tilde{Y}_{fake} \setminus Y_{true}|$, and f_v^* is the probability that v equals to the value sampled from $Y \setminus \tilde{Y}_{fake} \setminus Y_{true}$. Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$. The overall frequency gain $\text{Gain}(Y_{fake})$ is :

$$\begin{aligned}
 \text{Gain}(Y_{fake}) &= \sum_{t \in T} \mathbb{E}[\Delta \tilde{f}_t] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\
 &\stackrel{\text{Equ(4)}}{=} \sum_{t \in T} \mathbb{E} \left[\frac{\frac{1}{|Y \setminus \tilde{Y}_{fake}|} \sum_{i \in Y \setminus \tilde{Y}_{fake}} \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q}{p - q} \right] - \\
 &\quad \mathbb{E} \left[\frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}(y_i)}(t) - q}{p - q} \right] \quad (\text{according to Equation (4)}) \\
 &\stackrel{\text{Equ(5)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/r}{N_1 + N_2} - q'}{p' - q'} - \frac{(f_t(p - q) + q) - q}{p - q} \\
 &\stackrel{\text{Equ(15)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/r}{N_1 + N_2} - f_t(p' - q')}{p' - q'}. \quad (16)
 \end{aligned}$$

Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, the overall frequency gain $\text{Gain}(Y_{fake})$ is 0.

4.2.2 Defenses against Data Poisoning Attacks to OUE.

For OUE, we set $p_1 = \frac{1}{2}$, $q_1 = \frac{1}{e^{\epsilon_1} + 1}$, and $p_2 = \frac{1}{2}$, $q_2 = \frac{1}{e^{\epsilon_2} + 1}$ where $\epsilon_1 + \epsilon_2 = \epsilon$. In addition, let $p = \frac{1}{2}$, $q = \frac{1}{e^{\epsilon} + 1}$, $p' = p_1$ or p_2 , and $q' = q_1$ or q_2 .

Defenses against RPA to OUE. LDPGuard first randomly samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from encoded data space $\mathcal{D} = \{0, 1\}^d$. For each sampled value, LDPGuard then removes the record $v \in \tilde{Y}_{fake}$ from all reported Y if $v \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency \tilde{f}_v for item v with Equation (10). In average, each value sampled from $\mathcal{D} = \{0, 1\}^d$ supports v with probability $\frac{1}{2}$, according to Equation (11), the expectation of \tilde{f}_v is

$$\mathbb{E}[\tilde{f}_v] = \frac{\frac{N_1(f_v(p' - q') + q') + N_2 * 1/2}{N_1 + N_2} - q'}{p' - q'} \quad (17)$$

where $N_1 = |Y_{true} \cap (Y \setminus \tilde{Y}_{fake})|$ and $N_2 = |Y \setminus \tilde{Y}_{fake} \setminus Y_{true}|$. Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is an unbiased estimate of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$. Following Equation (12), the overall frequency gain $\text{Gain}(Y_{fake})$ is:

$$\text{Gain}(Y_{fake}) \stackrel{\text{Equ(17)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/2}{N_1 + N_2} - f_t(p' - q')}{p' - q'}. \quad (18)$$

Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, the overall frequency gain $\text{Gain}(Y_{fake})$ is 0.

Defenses against RIA to OUE. For RIA, LDPGuard first randomly samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from target set $T = \{t_1, t_2, \dots, t_r\}$, and then encodes each value to the encoded data space $\mathcal{D} = \{0, 1\}^d$. Obviously, each encoded value is a d -bits binary vector, where each bit with value 1 (resp. 0) is kept with probability $p' = p_1$ or p_2 (resp. $1 - q'$) and flipped with probability $1 - p'$ (resp. $q' = q_1$ or q_2). Finally, LDPGuard removes the perturbed value from all reported Y if the perturbed value is found in Y , and calibrates the results to obtain the estimated frequency for item v with Equation (10). Similar to RIA to kRR, the expectation of \tilde{f}_v is $\frac{(N_1(f_v(p' - q') + q') + N_2 * f_v^*) / (N_1 + N_2) - q')}{p' - q'}$ where $f_v^* = (1/r(p' - q') + q')$, if $v \in T$; and $f_v^* = q'$, otherwise. Accordingly, the overall frequency gain $\text{Gain}(Y_{fake})$ is

$$\text{Gain}(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2(1/r(p' - q') + q')}{N_1 + N_2} - f_t(p' - q')}{p' - q'}. \quad (19)$$

Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and the total frequency gain $\text{Gain}(Y_{fake})$ is 0.

Defenses against MGA to OUE. As for MGA to OUE, LDPGuard first initializes $(N + M)\tilde{\beta}$ d -bits zero vectors (denoted by \tilde{Y}_{fake}), for $y_i \in \tilde{Y}_{fake}$, LDPGuard sets $y_i[t] = 1$ for $t \in \mathcal{T}$ and $l = \lfloor p + (d - 1)q - r \rfloor$ randomly selected non-target bits. After that, LDPGuard then removes the record $y_i \in \tilde{Y}_{fake}$ from all reported Y if $y_i \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency for item v with Equation (10). The expectation of \tilde{f}_v is

$\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = 1$, if $v \in T$; and $f_v^* = 0$, otherwise. The overall frequency gain $Gain(Y_{fake})$ is

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2}{N_1+N_2} - f_t(p'-q')}{p'-q'}. \quad (20)$$

\tilde{f}_v is the unbiased estimation of f_v and the overall frequency gain $Gain(Y_{fake})$ is 0 if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$.

4.2.3 Defenses against Data Poisoning Attacks to OLH.

Let $p^* = \frac{e^\epsilon}{e^\epsilon + d' - 1}$, $q^* = \frac{1}{e^\epsilon + d' - 1}$, $p = p^* = \frac{e^\epsilon}{e^\epsilon + d' - 1}$ and $q = \frac{1}{d'} p^* + (1 - \frac{1}{d'}) q^* = \frac{1}{d'}$. Similarly, p_1^* , q_1^* , p_1 and q_1 (resp. p_2^* , q_2^* , p_2 and q_2) are defined w.r.t. ϵ_1 (resp. ϵ_2). In addition, let $p' = p_1$ or p_2 , $q' = q_1$ or q_2 .

Defenses against RPA to OLH. First, LDPGuard crafts $(N + M)\tilde{\beta}$ key-value pairs $\{\langle H, h \rangle | H \in \mathcal{H}, h \in \{1, 2, \dots, d'\}\}$ (denoted by \tilde{Y}_{fake}) where H is randomly selected from \mathcal{H} and h is randomly drawn from $\{1, 2, \dots, d'\}$ ($d' = e^\epsilon + 1$). Each record in \tilde{Y}_{fake} is then removed from Y if it exists in Y . Finally, the data collector calibrates the results to obtain the estimated frequency for item v with Equation (10). The expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = \frac{1}{d'}$, if $v \in T$; and $f_v^* = 0$, otherwise. The overall frequency gain $Gain(Y_{fake})$ is therefore

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*1/d'}{N_1+N_2} - f_t(p'-q')}{p'-q'}. \quad (21)$$

If $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and the overall frequency gain $Gain(Y_{fake})$ is 0.

Defenses against RIA to OLH. To defend against RIA to OLH, LDPGuard first samples $(N + M)\tilde{\beta}$ values from target set \mathcal{T} . For each sampled value v , a hash function H is randomly selected from \mathcal{H} . Then, the key-value pair $\langle H, h \rangle$ is perturbed to $\langle H, h' \rangle$ ($h' \neq h$) with probability q_1^* or q_2^* , and kept unchanged with probability p_1^* or p_2^* . All these perturbed values are removed from Y and the data collector calibrates the results to obtain the estimated frequency for the item v with Equation (10). The expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = \frac{1}{r}(p' - q') + q'$, if $v \in T$; and $f_v^* = q'$, otherwise. The overall frequency gain $Gain(Y_{fake})$ is therefore

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*(\frac{1}{r}(p'-q')+q')}{N_1+N_2} - f_t(p'-q')}{p'-q'}. \quad (22)$$

If $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and the overall frequency gain $Gain(Y_{fake})$ is 0.

Defenses against MGA to OLH. As for MGA to OLH, LDPGuard crafts $(N + M)\tilde{\beta}$ key-value pairs $\{\langle H, h \rangle | H \in \mathcal{H}, h \in \{1, 2, \dots, d'\}\}$ (denoted by \tilde{Y}_{fake}) where H is randomly selected from \mathcal{H} and h is the hash function that results in the most targets to the same value. Each record in \tilde{Y}_{fake} is then removed from Y if it exists in Y . Finally, the data collector calibrates the results to obtain the

estimated frequency for item v with Equation (10). The expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = 1$, if $v \in T$; and $f_v^* = 0$, otherwise. Therefore, the overall frequency gain $Gain(Y_{fake})$ is

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2}{N_1+N_2} - f_t(p'-q')}{p'-q'}. \quad (23)$$

Similarly, if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and the overall frequency gain $Gain(Y_{fake})$ is 0.

4.3 Discussions

Communication and Computing Cost. The communication cost of each user is $O(\log d)$, $O(d)$, and $O(\log d)$ for kRR, OUE and OLH, respectively. The computing cost of the data collector is $O(N + M + d)$, $O((N + M) * d)$, and $O((N + M) * d)$ for kRR, OUE and OLH, respectively.

Privacy Analysis. The first round collection satisfies ϵ_1 -LDP and the second round collection satisfies ϵ_2 -LDP, according to *Sequential Composition* (see Lemma 4.1), LDPGuard satisfies ϵ -LDP where $\epsilon = \epsilon_1 + \epsilon_2$.

LEMMA 4.1. *Given C algorithms $\{\mathcal{A}_i | i \in \{1, 2, \dots, C\}\}$, each \mathcal{A}_i satisfies ϵ_i -LDP, the sequence of algorithms \mathcal{A}_i ($i \in \{1, 2, \dots, C\}$) collectively satisfies $\sum_{i \in \{1, 2, \dots, C\}} \epsilon_i$ -LDP.*

Moreover, observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is the unbiased estimation of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and the overall frequency gain $Gain(Y_{fake})$ is 0. In other words, the effects of poisoning attacks are offset with the proposed defenses if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$). To this end, \tilde{Y}_{fake} should be sampled from the same distribution of that injected for attacks, as we mentioned in this section. In addition, the same number of noise data should be introduced, that is, $(N + M)\tilde{\beta}$ should close to the real number of fake users (i.e., M). We address this problem by precisely estimating of the percentage of fake users (i.e., $\tilde{\beta}$) in the following section.

5 ESTIMATION OF THE PERCENTAGE OF FAKE USERS

Observe that the reported values from fake users may be statistically abnormal compared to those from genuine users. For example, the values from fake users with RPA attacks are too random in each data collection, while these from fake users with MGA attacks are too determinate to be distinguished from the reported values from genuine users. This observation motivates us to precisely estimate the percentage of fake users by comparing the values reported in two round collection.

In particular, let P_1 and P_2 be the probability that the data collector collects the same value from a genuine user and fake user in two different rounds, respectively, and CNT be the number of observed same values collected from all $N + M$ users (i.e., N genuine users and M fake users), the percentage of fake users (i.e., $\beta = \frac{M}{N+M}$) can be estimated by

$$\tilde{\beta} = \frac{(N + M)P_1 - CNT}{(N + M)(P_1 - P_2)}. \quad (24)$$

The estimation $\tilde{\beta}$ is an unbiased estimation of the percentage of fake users.

THEOREM 5.1. $\tilde{\beta} = \frac{(N+M)P_1 - CNT}{(N+M)(P_1 - P_2)}$ is an unbiased estimation of the percentage of fake users.

The problem becomes how to obtain P_1 and P_2 . We answer this question in the remaining parts of this section.

5.1 Percentage Estimation for kRR

Recall that each genuine user reports his value with kRR follows the following three steps: 1) encodes his value v to itself, i.e., $\text{Encode}(v) = v$; 2) keeps the true value v with probability p and perturbs it to a different value $v' \in \mathcal{D} = \{1, 2, \dots, d\}$ with probability q . For two round collection, let $p' = p_1$ or p_2 and $q' = q_1$ or q_2 . 3) and finally reports the perturbed value $y = \text{PE}(v)$ to the data collector. Therefore, the genuine user reports the same value when 1) reporting the true value v in both first and second round; 2) reporting the same value v' ($v' \neq v$). The probability P_1 that the data collector collects the same value from a genuine user is

$$P_1 = (p')^2 + (d-1)(q')^2. \quad (25)$$

where $(p')^2$ and $(q')^2$ are the probability that reports the true value v and the same value v' ($v' \neq v$), respectively.

Computing P_2 for RPA to kRR. Since each fake user with the RPA attack randomly selects a value from the encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$, each item $v \in \mathcal{D}$ is selected in each round with a probability $\frac{1}{d}$. Therefore, both the first and second rounds select the same item v with probability $\frac{1}{d} * \frac{1}{d}$. Furthermore, since v could be any element in \mathcal{D} , P_2 is

$$\begin{aligned} P_2 &= \sum_{v \in \mathcal{D}} \Pr(v \text{ in 1st round}) * \Pr(v \text{ in 2nd round}) \\ &= \sum_{v \in \mathcal{D}} \frac{1}{d} * \frac{1}{d} = d * \left(\frac{1}{d} * \frac{1}{d}\right) = \frac{1}{d}. \end{aligned}$$

Computing P_2 for RIA to kRR. As a fake user is to select a item v from target set $T = \{t_1, t_2, \dots, t_r\}$ and perturbs (resp. keeps) it to a different value in \mathcal{D} (resp. the same value in T) with probability $q' = q_1$ or q_2 (resp. $p' = p_1$ or p_2), computing P_2 for RIA to kRR should be divided into two cases. The first is that the values reported in two rounds are the same and belong to T . In this case, the probability that the fake user reports the same value is $\sum_{v \in T} (\frac{1}{r}p' + (1 - \frac{1}{r})q')^2$ where $\frac{1}{r}p'$ and $(1 - \frac{1}{r})q'$ are the probability that the reported value is not perturbed and perturbed, respectively. The second is that the reported values in two rounds are the same and not T . In this case, the reported value must be perturbed from his true value. Thus, the probability that the fake user reports the same value is $\sum_{v \in \mathcal{D} \setminus T} (q')^2$. Therefore, P_2 is

$$\begin{aligned} P_2 &= \sum_{v \in T} \Pr(v \text{ in 1st round}) * \Pr(v \text{ in 2nd round}) \\ &+ \sum_{v \in \mathcal{D} \setminus T} \Pr(v \text{ in 1st round}) * \Pr(v \text{ in 2nd round}) \\ &= \sum_{v \in T} \left(\frac{1}{r}p' + (1 - \frac{1}{r})q'\right)^2 + \sum_{v \in \mathcal{D} \setminus T} (q')^2 \\ &= r\left(\frac{1}{r}p' + (1 - \frac{1}{r})q'\right)^2 + (d-r)(q')^2. \end{aligned}$$

Computing P_2 for MGA to kRR. Unlike RPA to kRR, each fake user randomly selects a value from target set T instead of encoded data space \mathcal{D} , each item $v \in T$ is selected in each round with a probability $\frac{1}{r}$. Therefore, both the first round and the second round select the same item v with probability $\frac{1}{r} * \frac{1}{r}$. Since v could be any item in T , P_2 is

$$\begin{aligned} P_2 &= \sum_{v \in T} \Pr(v \text{ in 1st round}) * \Pr(v \text{ in 2nd round}) \\ &= \sum_{v \in T} \frac{1}{r} * \frac{1}{r} = r * \left(\frac{1}{r} * \frac{1}{r}\right) = \frac{1}{r}. \end{aligned}$$

5.2 Percentage Estimation for OUE

For OUE, each genuine user encodes his value v to a d -dimensional binary vector B where v -th bit is 1, i.e., $\text{Encode}(v) = B$ where $B[v] = 1$ and $B[v'] = 0$ for $v' \neq v$. Each nonzero entry (resp. zero entry) in B is kept unchanged with probability $p' = p_1$ or p_2 (resp. $1 - q'$) and flipped with probability $1 - p'$ (resp. $q' = q_1$ or q_2). Let B_1 and B_2 be reported binary vectors in first and second round, respectively, $B_1 = B_2$ requires that two corresponding bits in the vectors are the same, and each bit i keeps the same in two rounds when both $B_1[i]$ and $B_2[i]$ are perturbed from $B[i]$ or kept as $B[i]$. For $i = v$, the probability of $B_1[i] = B_2[i]$ is $(p')^2 + (1 - p')^2$ where $(p')^2$ is for $B_1[i] = B_2[i] = B[i]$ and $(1 - p')^2$ is for $B_1[i] = B_2[i] \neq B[i]$. For $i \neq v$, the probability of $B_1[i] = B_2[i]$ is $(q')^2 + (1 - q')^2$. Therefore, the probability P_1 that the data collector collects the same value from a genuine user is

$$\begin{aligned} P_1 &= \sum_{i=v} \Pr(B_1[i] = B_2[i]) * \sum_{i \neq v} \Pr(B_1[i] = B_2[i]) \\ &= \sum_{i=v} (p')^2 + (1 - p')^2 * \sum_{i \neq v} (q')^2 + (1 - q')^2 \\ &= ((p')^2 + (d-1)(q')^2)((q')^2 + (1 - q')^2)^{d-1}. \end{aligned} \quad (26)$$

Observe that the dimension of a binary vector could be very large such that P_1 has a very small value, and scarce observations of the same value are obtained in two rounds. To solve this problem, P_1 for any τ bits of B_1 and B_2 (denoted by $P_1(\tau)$) should be given. Since the τ bits are randomly selected from d bits, they contain v -th bit with probability $\frac{\tau}{d}$. As such, $P_1(\tau)$ is

$$\begin{aligned} P_1(\tau) &= \frac{\tau}{d} * \sum_{i=v} \Pr(B_1[i] = B_2[i]) * \sum_{i \neq v} \Pr(B_1[i] = B_2[i]) \\ &= + (1 - \frac{\tau}{d}) * \sum_{i \neq v} \Pr(B_1[i] = B_2[i]) \\ &= \frac{\tau}{d} * ((p')^2 + (d-1)(q')^2)((q')^2 + (1 - q')^2)^{\tau-1} \\ &= + (1 - \frac{\tau}{d})((q')^2 + (1 - q')^2)^{\tau}. \end{aligned}$$

Computing P_2 for RPA to OUE. Similar to the computation of P_2 for RPA to kRR, the fake user randomly selects a value from the encoded data space, the difference lie in that the encoded data space is now $\mathcal{D} = \{0, 1\}^d$. Thus, the fake user selects the same binary vector B in both the first round and second round with probability

$\frac{1}{2^d} * \frac{1}{2^d}$, P_2 is

$$\begin{aligned} P_2 &= \sum_{B \in \mathcal{D}} Pr(B \text{ in 1st round}) * Pr(B \text{ in 2nd round}) \\ &= \sum_{B \in \mathcal{D}} \frac{1}{2^d} * \frac{1}{2^d} = 2^d * \left(\frac{1}{2^d} * \frac{1}{2^d} \right) = \frac{1}{2^d}. \end{aligned}$$

Accordingly, $P_2(\tau)$ is $2^\tau * (\frac{1}{2^\tau} * \frac{1}{2^\tau}) = \frac{1}{2^\tau}$.

Computing P_2 for RIA to OUE. For RIA to OUE, the fake user first selects a value from target set T and then adopts the same encoding and perturbation procedures used by the genuine users. Let B_1 and B_2 be reported binary vectors in the first and second round, respectively, $B_1 = B_2$ requires that two corresponding bits in the vectors are the same, and each bit i keeps the same in two rounds when both $B_1[i]$ and $B_2[i]$ are perturbed from $B[i]$ or kept as $B[i]$. For $i \in T$, the probability of $B_1[i] = B_2[i]$ should be divided into two cases. First, $B_1[i] = B_2[i] = 1$, the probability is $(\frac{1}{r}p' + (1 - \frac{1}{r})q')^2$. Second, $B_1[i] = B_2[i] = 0$, the probability is $(\frac{1}{r}(1 - p') + (1 - \frac{1}{r})(1 - q'))^2$. For $i \notin T$, the probability of $B_1[i] = B_2[i]$ is $(q')^2 + (1 - q')^2$. Therefore, the probability P_2 that the data collector collects the same value from a fake user is

$$\begin{aligned} P_2 &= \sum_{i \in T} Pr(B_1[i] = B_2[i]) * \sum_{i \notin T} Pr(B_1[i] = B_2[i]) \\ &= \sum_{i \in T} \left(\frac{1}{r}p' + (1 - \frac{1}{r})q' \right)^2 + \left(\frac{1}{r}(1 - p') + (1 - \frac{1}{r})(1 - q') \right)^2 \\ &\quad * \sum_{i \notin T} (q')^2 + (1 - q')^2 \\ &= \left[\left(\frac{1}{r}p' + (1 - \frac{1}{r})q' \right)^2 + \left(\frac{1}{r}(1 - p') + (1 - \frac{1}{r})(1 - q') \right)^2 \right]^r \\ &\quad * \left[(q')^2 + (1 - q')^2 \right]^{d-r} \end{aligned}$$

The computation of P_2 for any chosen τ bits, $P_2(\tau)$, depends on how many bits (denoted by j) are chosen from T . Therefore, P_2 is the expectation over j , formally,

$$\begin{aligned} P_2(\tau) &= \sum_{j=0}^r \prod_{i=0}^d Pr(B_1[i] = B_2[i], j) = \sum_{j=0}^r Pr(j) \prod_{i=0}^d Pr(B_1[i] = B_2[i]|j) \\ &= \sum_{j=0}^r Pr(j) \left[\prod_{i \in T} Pr(B_1[i] = B_2[i]|j) * \prod_{i \notin T} Pr(B_1[i] = B_2[i]|j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[\prod_{i \in T} Pr(B_1[i] = B_2[i]|j) * \prod_{i \notin T} Pr(B_1[i] = B_2[i]|j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[\left(\left(\frac{p'}{r} + \frac{r-1}{r}q' \right)^2 + \left(\frac{1-p'}{r} + \frac{r-1}{r}(1-q') \right)^2 \right)^j \right. \\ &\quad \left. * \left[(q')^2 + (1 - q')^2 \right]^{\tau-j} \right] \end{aligned}$$

where $C_{X_1}^{X_2}$ is the number of X_2 -combinations of the set with X_1 elements.

Computing P_2 for MGA to OUE. Recall that each fake user attacks OUE with MGA by sending a binary vector that has all target bits and randomly chosen $l = \lfloor p + (d-1)q - r \rfloor$ bits with 1, and the

remaining bits with 0. Let B_1 and B_2 be reported binary vectors in the first and second rounds, respectively, whether or not $B_1 = B_2$ depends on the randomly chosen l bits in two different rounds. Since these l bits are randomly selected from all non-target bits, each possible l bits is selected with a probability of $1/C_{d-r}^l$. Let \mathbb{B}_ℓ be all possible binary vectors, each of which has all target bits and additional l bits with 1, P_2 is therefore

$$\begin{aligned} P_2 &= \sum_{B \in \mathbb{B}_\ell} Pr(B \text{ in 1st round}) * Pr(B \text{ in 2nd round}) \\ &= \sum_{B \in \mathbb{B}_\ell} \frac{1}{C_{d-r}^l} * \frac{1}{C_{d-r}^l} = C_{d-r}^l * \left(\frac{1}{C_{d-r}^l} * \frac{1}{C_{d-r}^l} \right) = \frac{1}{C_{d-r}^l} \end{aligned}$$

Observe that $P_2(\tau)$ for any chosen τ bits depends on how many bits (denoted by j) are chosen from T . In particular, if j bits are chosen from T , the values in the corresponding entries in B_1 and B_2 are always the same, thus, $P_2(\tau)$ is computed as follows:

$$\begin{aligned} P_2(\tau) &= \sum_{j=0}^r \prod_{i=0}^d Pr(B_1[i] = B_2[i], j) = \sum_{j=0}^r Pr(j) \prod_{i=0}^d Pr(B_1[i] = B_2[i]|j) \\ &= \sum_{j=0}^r Pr(j) \left[\prod_{i \in T} Pr(B_1[i] = B_2[i]|j) * \prod_{i \notin T} Pr(B_1[i] = B_2[i]|j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[1 * \prod_{i \notin T} Pr(B_1[i] = B_2[i]|j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[\left(\frac{l}{d-r} \right)^2 + \left(\frac{d-r-l}{d-r} \right)^2 \right]^{\tau-j} \end{aligned}$$

where $C_{X_1}^{X_2}$ is the number of X_2 -combinations of the set with X_1 elements, and $\left(\frac{l}{d-r} \right)^2$ (resp. $\left(\frac{d-r-l}{d-r} \right)^2$) is the corresponding probability if the bit i is in (resp. not in) the randomly chosen l bits.

5.3 Percentage Estimation for OLH

We follow Section 4.2.3 to set $p^* = \frac{e^\epsilon}{e^\epsilon + d'^r - 1}$, $q^* = \frac{1}{e^\epsilon + d'^r - 1}$, $p = p^*$ = $\frac{e^\epsilon}{e^\epsilon + d'^r - 1}$ and $q = \frac{1}{d'}p^* + (1 - \frac{1}{d'})q^* = \frac{1}{d'}$. Similarly, p_1^* , q_1^* , p_1 and q_1 (resp. p_2^* , q_2^* , p_2 and q_2) are defined w.r.t. ϵ_1 (resp. ϵ_2). In addition, let $p_{1,2}^* = p_1^*$ or p_2^* , $q_{1,2}^* = q_1^*$ or q_2^* .

Recall that each genuine user reports his value v with OLH first randomly selects a hash function H from \mathcal{H} to obtained $H(v)$, and then reports the key-value pair $\langle H, H(v)' \rangle$ where $H(v)'$ is perturbed (resp. kept) from $H(v)$ with probability $q_{1,2}^*$ (resp. $p_{1,2}^*$). As such, the genuine user reports the same value when 1) reporting the true value $H(v)$ in both first and second round; 2) reporting the same value $H(v)'$ ($H(v)' \neq H(v)$). The probability P_1 that the data collector collects the same value from a genuine user is

$$P_1 = (p_{1,2}^*)^2 + (d' - 1)(q_{1,2}^*)^2. \quad (27)$$

where $p_{1,2}^*$ and $q_{1,2}^*$ are the probability that reports the true value $H(v)$ and the same value $H(v)'$ ($H(v)' \neq H(v)$), respectively.

Computing P_2 for RPA to OLH. For RPA to OLH, each fake user reports a key-value pair $\langle H, h \rangle$ where h is randomly selected from $\mathcal{D} = \{1, 2, \dots, d'\}$ ($d' = e^\epsilon + 1$) without perturbation. As such, each item $h \in \mathcal{D}$ is selected in each round with a probability $\frac{1}{d'}$. Therefore, both first round and second round select the same item

h with probability $\frac{1}{d'} * \frac{1}{d'}$, P_2 is

$$P_2 = \sum_{h \in \mathcal{D}} Pr(h \text{ in 1st round}) * Pr(h \text{ in 2nd round}) \\ = \sum_{h \in \mathcal{D}} \frac{1}{d'} * \frac{1}{d'} = d' * (\frac{1}{d'} * \frac{1}{d'}) = \frac{1}{d'}.$$

Computing P_2 for RIA to OLH. Observe that each fake user randomly selects a value from target set T and perturbs it to a value in $\mathcal{D} = \{1, 2, \dots, d'\}$ where $d' = e^\epsilon + 1$, each value in $\mathcal{D} = \{1, 2, \dots, d'\}$ occurs in the report collected from a particular fake user is with probability $\frac{1}{d'} * p_{1,2}^* + (1 - \frac{1}{d'}) * q_{1,2}^*$, as such,

$$P_2 = \sum_{h \in \mathcal{D}} Pr(h \text{ in 1st round}) * Pr(h \text{ in 2nd round}) \\ = \sum_{h \in \mathcal{D}} (\frac{1}{d'} * p_{1,2}^* + (1 - \frac{1}{d'}) * q_{1,2}^*)^2 \\ = d' * (\frac{1}{d'} * p_{1,2}^* + (1 - \frac{1}{d'}) * q_{1,2}^*)^2.$$

Computing P_2 for MGA to OLH. Since MGA to OLH is to seek for a hash function that maps all target values to the same encoded value, thus, P_2 is always be 1. But in practice, there may has no such hash function, and we find the hash function that results in the most targets to the same value. Suppose there are maximally NUM targets hashed to the same value, $P_2 = (\frac{NUM}{r})^2 + (1 - \frac{NUM}{r})^2$.

6 EXPERIMENTAL EVALUATION

In this section, we investigate the performance of LDPGuard and report the key findings. LDPGuard is implemented with Python 3.10. All experiments are conducted on macOS Monterey 12.2.1 with Quad-Core Intel Core i5 4-Core CPU (2GHz) and 16GB of main memory.

6.1 Experimental Setup

Datasets. LDPGuard is evaluated on three datasets consisting of two real-world datasets (*i.e.*, IPUMS [58] and Fire [59]) and one synthetic dataset. The datasets are summarized in Table 2 and their details are listed as follows:

- **IPUMS [58].** IPUMS dataset consists of harmonized census and American Community Survey data over the years. We select the data of 2017 and the geographic variable “CITY” as the item set. The item set contains 102 items each of which is a city code. The number of users is 389,894.
- **Fire [59].** Fire calls-for-service contains all fire units’ responses to calls in San Francisco. Each record in the includes at least call type, call number, and incident number at all relevant time intervals. We follow [22] to filter out other call types but “Alarms” and select the “Unit ID” as the item field to construct the Fire dataset, which has 262 items and covers 632,543 users.
- **Zipf.** Zipf is a synthetic dataset that follows the Zipf distribution. It includes 128 items and 500,000 users.

Baselines. We compare LDPGuard against two related countermeasures [22] against data poisoning attacks¹: (1) *Normalization*

¹Conditional probability based Detection [22] is to detect possible fake users based on conditional probability, but it is applicable when there is only one target item.

Table 2: Datasets

Dataset	Number of Items	Number of Users
IPUMS [58]	102	389,894
Fire [59]	262	632,543
Zipf	128	500,000

Table 3: Parameter Settings for Experiments

Parameter	Symbol	Default Value
number of targets	r	10
percentage of fake users	β	0.05
privacy budget	ϵ	1.0
	ϵ_1	$\frac{\epsilon}{2}$
	ϵ_2	$\frac{\epsilon}{2}$

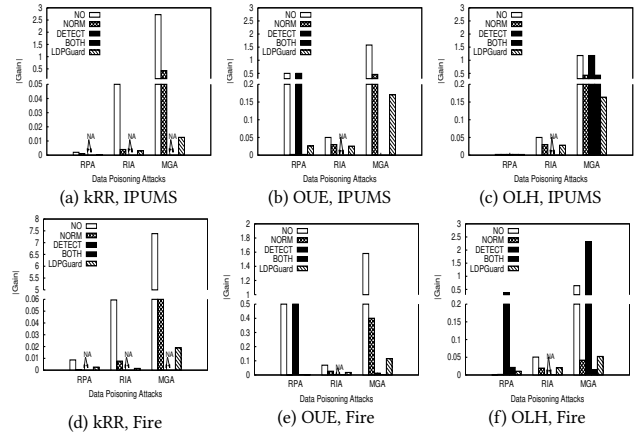


Figure 2: Absolute Gain $|Gain|$ on IPUMS and Fire.

(denoted by **NORM**) that normalizes the estimated item frequencies such that each estimated item frequency is nonnegative and the overall item frequencies sum to 1; (2) *Fake users detection* (denoted by **DETECT**) that detects possible fake users based on their submitted values to remove their impact on the final results. We also compare it with the results where no countermeasure is adopted (denoted by **NO**), and both DETECT and NORM are applied (denoted by **BOTH**). Noted that the resulting frequency of LDPGuard is also normalized to avoid negative results.

Performance Metrics. We use **absolute gain**, *i.e.*, absolute value of the frequency gain (denoted by $|Gain|$), to compare LDPGuard with baselines.² We use **NA** to denote the case that a countermeasure is not applicable for a data poisoning attack. We also compare LDPGuard with DETECT in terms of the estimated percentage (*i.e.*, β) of fake users.

Parameter Settings. Unless otherwise stated, we use the default parameter values in Table 3. In addition, we follow existing work [13, 22] to set $d' = e^\epsilon + 1$ for the OLH protocol.

²One may argue that he can follow the existing work [22] to use frequency gain $Gain$ instead of its absolute value, but he may ignore the fact that the frequency of targets is largely distorted if $Gain$ is a negative number with a large absolute value, which results in unfair experimental comparison for countermeasures.

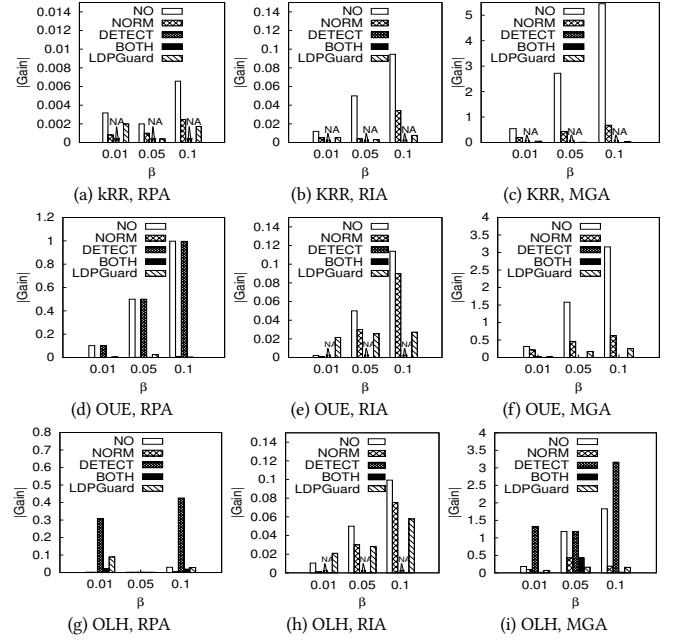
Table 4: Percentage Estimation on IPUMS and Fire, $\beta = 0.05$

Protocols	Attacks	β (IPUMS)		β (Zipf)	
		DETECT	LDPGuard	DETECT	LDPGuard
KRR	RPA	—	0.117	—	0.593
	RIA	—	0.129	—	0.382
	MGA	—	0.050	—	0.049
OUE	RPA	7.0E-05	0.060	7.0E-05	0.042
	RIA	—	0.232	—	0.292
	MGA	0.050	0.044	0.050	0.041
OLH	RPA	0.029	0.013	0.025	0.081
	RIA	—	0.024	—	0.045
	MGA	0.346	0.048	0.267	0.040

6.2 Experimental Results

Exp 1: Performances of Frequency Estimation. We first compare LDPGuard with baselines on three datasets and plot the results of frequency estimation on IPUMS and Fire in Figure 2.³ As depicted in Figure 2(a) and (d), LDPGuard always outperforms baselines for kRR regardless of data poisoning attacks. Compared to NO, NORM reduces $|\text{Gain}|$ to a relatively smaller value, but the impact is limited as NORM does not reduce $|\text{Gain}|$ directly but normalizes each estimated item frequency to a non-negative value. For kRR, both DETECT and BOTH cannot work and their results are denoted by NA. For OUE (Figure 2(b) and (e)) and OLH (Figure 2 (c) and (f)), NORM is ineffective in the scenario where MGA attacks are deployed, and the reason is discussed above. DETECT and BOTH are effective in this scenario, as they can detect abnormal statistical patterns in the reported values generated by MGA. However, DETECT is ineffective for RPA because each fake user under RPA randomly selects perturbed values in the encoded space, and thus the perturbed values do not show meaningful statistical patterns. In addition, both DETECT and BOTH are not applicable for RIA since each fake user under RIA perturbs the value sampled from the target item set before sending it to a data collector, which results in incapability of distinguishing fake users from genuine users. In contrast, LDPGuard can obtain relatively less $|\text{Gain}|$ than NORM for MGA. Compared to DETECT, it also performs better performance for RIA. Moreover, it is applicable regardless of LDP protocols and data poisoning attacks.

In summary, LDPGuard can effectively defend against data poisoning attacks and outperforms baselines in terms of absolute gain. **Exp 2: Performances of Percentage Estimation.** We further compare LDPGuard with DETECT on IPUMS and Fire⁴ and report the estimated percentages of fake users in Table 4. In this experiment, we set the true percentage of fake users to 0.05 and use “—” to denote the case where the countermeasure is not applicable. Observe that DETECT can accurately estimate the percentage of MGA to OUE because each reported value of OUE is a binary vector and the target bits in the array tends to expose abnormal statistical patterns. We also observe that it has limitations in estimating the percentage in the scenarios where there is no deterministic and meaningful statistical patterns. For example, for RPA to OUE, each fake user randomly selects the value from a very large encoding space \mathcal{D} (i.e., $\mathcal{D} = \{0, 1\}^d$) and reports it without perturbation. This

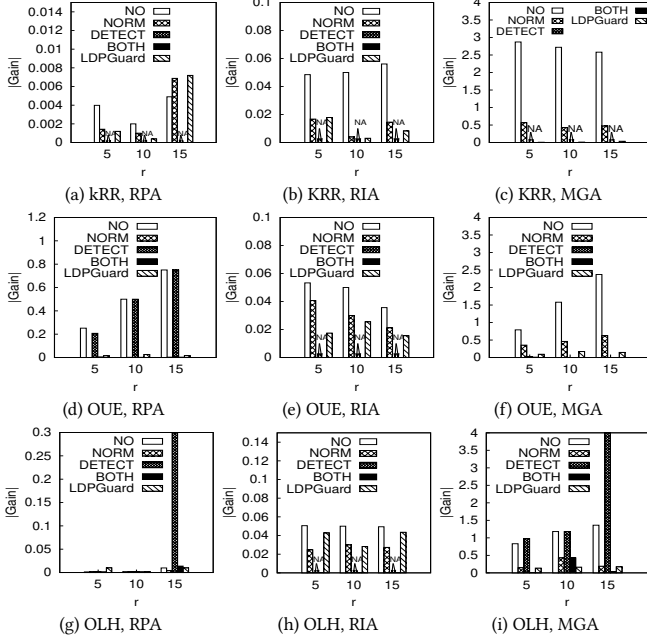
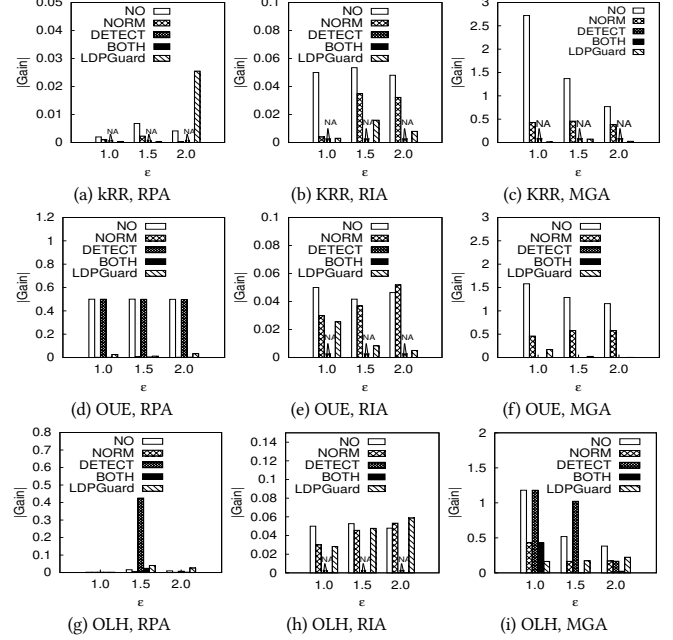
**Figure 3: Effect of β .****Table 5: Percentage Estimation, $\beta=0.01, 0.05, 0.1$**

Protocol, Attack	β (DETECT)			β (LDPGuard)		
	$\beta=0.01$	$\beta=0.05$	$\beta=0.1$	$\beta=0.01$	$\beta=0.05$	$\beta=0.1$
KRR, RPA	—	—	—	0.621	0.117	0.020
KRR, RIA	—	—	—	0.189	0.129	0.022
KRR, MGA	—	—	—	0.009	0.050	0.099
OUE, RPA	1E-05	7E-05	9E-05	0.011	0.060	0.143
OUE, RIA	—	—	—	0.204	0.232	0.379
OUE, MGA	0.010	0.050	0.100	0.014	0.044	0.088
OLH, RPA	0.019	0.029	0.032	0.009	0.013	0.104
OLH, RIA	—	—	—	0.006	0.024	0.109
OLH, MGA	0.122	0.346	0.370	0.008	0.267	0.094

will result in each report vector from fake users with about $d/2$ non-zero bits, and further lead to the difficulty of distinguishing a fake user from a genuine user who reports the value with nearly $1/2 + (d-1)/(e+1)$ non-zero bits when $\epsilon = 1$. In addition, DETECT is not applicable for KRR as well as OUE (and OLH) to RIA. In contrast, LDPGuard is applicable for all LDP protocols and data poisoning attacks, although the estimated percentages under RPA and RIA where randomnesses are here and there are not very precise. In general, the estimated percentages by LDPGuard are close to the true percentages (i.e., 0.05 in this experiment) of fake users. Hence, LDPGuard is able to precisely estimate the percentage of fake users.

Exp 3: Effect of β . Next, we vary the value of β and report the impact of β on IPUMS in Figure 3. Observe that $|\text{Gain}|$ increases with β if no countermeasure is adopted (i.e., NO, Figure 3) because more fake users are introduced to distort the true data distribution. Although NORM can decrease $|\text{Gain}|$, it has little impact on the RIA and MGA (e.g., Figure 3 (b), (c), (e), (f)). DETECT is effective in defending against MGA to OUE (see Figure 3(f)) and has good

³Results of frequency estimation on Zipf are given in the Appendix [1].⁴Results of percentage estimation on Zipf are given in the Appendix [1].

Figure 4: Effect of r .Figure 5: Effect of ϵ .Table 6: Percentage Estimation, $r=5, 10, 15$

Protocol, Attack	β (DETECT)			β (LDPGuard)		
	$r=5$	$r=10$	$r=15$	$r=5$	$r=10$	$r=15$
KRR,RPA	—	—	—	0.055	0.117	0.055
KRR,RIA	—	—	—	0.602	0.129	0.255
KRR,MGA	—	—	—	0.051	0.050	0.050
OUE,RPA	0.003	7.06E-05	4.87E-06	0.134	0.060	0.031
OUE,RIA	—	—	—	0.177	0.232	0.101
OUE,MGA	0.051	0.050	0.050	0.044	0.044	0.046
OLH,RPA	0	0.029	0.176	0.010	0.013	0.068
OLH,RIA	—	—	—	0.021	0.025	0.154
OLH,MGA	0.149	0.346	0.621	0.049	0.048	0.042

performance for MGA to OLH when combined with NORM (*i.e.*, BOTH, Figure 3(i)). But DETECT and BOTH are not applicable for kRR protocols as well as RIA attacks. For MGA to OUE, as shown in Figure 3(f), LDPGuard is comparable to DETECT. Moreover, it outperforms baselines in other scenarios regardless of the value of β . We also report the estimated percentage of fake users (*i.e.*, $\tilde{\beta}$) in Table 5 and observe similar results. In the following experiments, we set $\beta = 0.05$.

Exp 4: Effect of r . We also study the performance of LDPGuard with different numbers of targets (*i.e.*, r) on IPUMS. First, we observe that $|\text{Gain}|$ s of all these countermeasures are close to zero (exactly less than 0.008) (Figure 4(a)), this indicates that RPA has very limited impacts on kRR even if no countermeasure is adopted. For other attacks on kRR, LDPGuard outperforms baselines no matter what r is. Second, observe that DETECT is ineffective in RPA to OUE but very effective in MGA to OUE (Figure 4(d)-(f)). The reason is as discussed above. For OLH, DETECT is incapable of decreasing $|\text{Gain}|$ of MGA attacks when $r = 15$ because MGA can only find a hash function among the 1,000 random seeds that hashes a subset

Table 7: Percentage Estimation, $\epsilon=1, 1.5, 2$

Protocol, Attack	$\tilde{\beta}$ (DETECT)			$\tilde{\beta}$ (LDPGuard)		
	$\epsilon=1$	$\epsilon=1.5$	$\epsilon=2$	$\epsilon=1$	$\epsilon=1.5$	$\epsilon=2$
KRR,RPA	—	—	—	0.117	0.301	0.271
KRR,RIA	—	—	—	0.129	0.054	0.199
KRR,MGA	—	0.050	—	0.050	0.048	0.050
OUE,RPA	7.07E-05	5.36E-05	7.07E-05	0.060	0.035	0.050
OUE,RIA	—	—	—	0.232	0.105	0.065
OUE,MGA	0.050	0.050	0.050	0.044	0.050	0.051
OLH,RPA	0.029	0.049	0	0.013	0.057	0.015
OLH,RIA	—	—	—	0.025	0.025	0.069
OLH,MGA	0.346	0.179	0.058	0.048	0.046	0.038

of the target items to the same value for each fake user. But in general, our LDPGuard can outperform baselines regardless of r . Moreover, it also outperforms baselines in terms of estimated percentages as shown in Table 6. We set $r = 10$ in the following experiments unless specified.

Exp 5: Effect of ϵ . We finally vary privacy budget ϵ on IPUMS and report the results in Figure 5. As depicted in Figure 5(c), (f) and (i), $|\text{Gain}|$ of MGA decreases as the privacy budget ϵ increases. This is because for larger ϵ , genuine users inject less noise to their data, while the values reported by fake users are irrelevant to privacy budget. As such, fake users contribute less to the frequencies of target items and lead to the decrease of $|\text{Gain}|$. But $|\text{Gain}|$ under other attacks are less sensitive to ϵ (*e.g.*, Figure 5(b), (e) and (h)). For example, both genuine users and fake users with RIA need to perturb their data and hence $|\text{Gain}|$ under RIA fluctuates in a small range. We also report the estimated percentages with different privacy budgets in Table 7 and observe that LDPGuard outperforms DETECT in most scenarios and is comparable to DETECT for MGA

to OUE. In short, LDPGuard outperforms baselines regardless of privacy budget.

7 CONCLUSIONS

In this paper, we focus on the problem of defending against data poisoning attacks to local differential privacy protocols. We present a novel framework called LDPGuard to this end. In particular, LDPGuard first adopts the two-round collection method to accurately estimate the percentage of fake users. Then it uses the estimated percentage to compensate for the adverse effects of various data poisoning attacks. LDPGuard supports state-of-the-art LDP protocols (kRR, OUE, and OLH) for frequency estimation and can effectively defend against existing data poisoning attacks. Experimental study on real-world and synthetic datasets demonstrates the superiority of LDPGuard compared to existing techniques. As part of our future work, we will explore the data poisoning attacks and countermeasures to other data types such as graph data.

REFERENCES

- [1] Technical Report. Available at: <https://github.com/TechReport2023/LDPGuard/blob/main/LDPGuard-TR.pdf>.
- [2] Brendan Avent, Aleksandra Korolova, David Zeber, Torgeir Hovden, and Benjamin Livshits. BLENDER: Enabling local search with a hybrid differential privacy model. In *USENIX Security*, 2017.
- [3] Raef Bassily, Kobbi Nissim, Uri Stemmer, and Abhradeep Guha Thakurta. Practical locally private heavy hitters. In *NeurIPS*, 2017.
- [4] Raef Bassily and Adam Smith. Local, private, efficient protocols for succinct histograms. In *STOC*, 2015.
- [5] Graham Cormode, Tejas Kulkarni, and Divesh Srivastava. Marginal release under local differential privacy. In *SIGMOD*, 2018.
- [6] John C Duchi, Michael I Jordan, and Martin J Wainwright. Local privacy and statistical minimax rates. In *FOCS*, 2013.
- [7] Ulfar Erlingsson, Vasyl Pihur, and Aleksandra Korolova. Rappor: Randomized aggregatable privacy-preserving ordinal response. In *CCS*, 2014.
- [8] Jinyuan Jia and Neil Zhenqiang Gong. Calibrate: Frequency estimation and heavy hitter identification with local differential privacy via incorporating prior knowledge. In *INFOCOM*, 2019.
- [9] Peter Kairouz, Keith Bonawitz, and Daniel Ramage. Discrete distribution estimation under local privacy. In *ICML*, 2016.
- [10] Peter Kairouz, Sewoong Oh, and Pramod Viswanath. Extremal mechanisms for local differential privacy. In *NeurIPS*, 2014.
- [11] Zhan Qin, Yin Yang, Ting Yu, Issa Khalil, Xiaokui Xiao, and Kui Ren. Heavy hitter estimation over set-valued data with local differential privacy. In *CCS*, 2016.
- [12] Xuebin Ren, Chia-Mu Yu, Weiren Yu, Shusen Yang, Xinyu Yang, Julie A McCann, and S Yu Philip. LoPub: High-dimensional crowdsourced data publication with local differential privacy. *TIFS*, 2018.
- [13] Tianhao Wang, Jeremiah Blocki, Ninghui Li, and Somesh Jha. Locally differentially private protocols for frequency estimation. In *USENIX Security*, 2017.
- [14] Tianhao Wang, Bolin Ding, Jingren Zhou, Cheng Hong, Zhicong Huang, Ninghui Li, and Somesh Jha. Answering multi-dimensional analytical queries under local differential privacy. In *SIGMOD*, 2019.
- [15] Tianhao Wang, Ninghui Li, and Somesh Jha. Locally differentially private frequent itemset mining. In *S&P*, 2018.
- [16] Tianhao Wang, Ninghui Li, and Somesh Jha. Locally differentially private heavy hitter identification. *TDSC*, 2019.
- [17] Tianhao Wang, Milan Lopuhaä-Zwakenberg, Zitao Li, Boris Skoric, and Ninghui Li. Locally differentially private frequency estimation with consistency. In *NDSS*, 2020.
- [18] Zhikun Zhang, Tianhao Wang, Ninghui Li, Shibo He, and Jiming Chen. Calm: Consistent adaptive local marginal for marginal release under local differential privacy. In *CCS*, 2018.
- [19] Apple Differential Privacy Team. Learning with privacy at scale. *Machine Learning Journal*, 2017.
- [20] Bolin Ding, Janardhan Kulkarni, and Sergey Yekhanin. Collecting telemetry data privately. In *NeurIPS*, 2017.
- [21] Nguyễn, T Thông, Xiaokui Xiao, et al. Collecting and analyzing data from smart device users with local differential privacy. *arXiv preprint arXiv:1606.05053*, 2016.
- [22] Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Data poisoning attacks to local differential privacy protocols. In *USENIX Security*, 2021.
- [23] Sweeney, Latanya. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 2002.
- [24] Machanavajjhala Ashwin, Kifer Daniel, Gehrke Johannes, et al. l-diversity: Privacy beyond k-anonymity. In *ICDE*, 2006.
- [25] Ninghui Li, Tiancheng Li, and Venkatasubramanian Suresh. t-closeness: Privacy beyond k-anonymity and l-diversity. In *ICDE*, 2007.
- [26] Chang Zhao, Lei Zou, and Feifei Li. Privacy preserving subgraph matching on large graphs in cloud. In *SIGMOD*, 2016.
- [27] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, 2006.
- [28] Privacy and efficiency guaranteed social subgraph matching. Kai Huang, Haibo Hu, Shuigeng Zhou, Jihong Guan, Qingqing Ye, and Xiaofang Zhou. *The VLDB Journal*, 2021.
- [29] Cynthia Dwork, Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 2014.
- [30] Ninghui Li, Min Lyu, Dong Su, Weining Yang. Differential Privacy: From Theory to Practice. *Synthesis Lectures on Information Security, Privacy, and Trust*, 2016.
- [31] Stanley L Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 1965.
- [32] Shiva Prasad Kasiviswanathan et al. 2011. What can we learn privately? *SIAM Journal on Computing*, 2011.
- [33] Justin Hsu, Sanjeev Khanna, and Aaron Roth. Distributed private heavy hitters. In *Automata, Languages, and Programming*, 2012.
- [34] Giulia Fanti, Vasyl Pihur, Ulfar Erlingsson. Building a RAPPOR with the Unknown: Privacy-Preserving Learning of Associations and Data Dictionaries. In *PoPETS*,

- 2016.
- [35] Albert Cheu, Adam Smith, and Jonathan Ullman. Manipulation attacks in local differential privacy. In *S&P*, 2021.
- [36] Scott Alfeld, Xiaojin Zhu, and Paul Barford. Data poisoning attacks against autoregressive models. In *AAAI*, 2016.
- [37] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. In *ICML*, 2012.
- [38] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. Local model poisoning attacks to byzantine robust federated learning. In *USENIX Security*, 2020.
- [39] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. In *S&P*, 2018.
- [40] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating backdoor attacks on deep neural networks. In *IEEE Access*, 2019.
- [41] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *NeurIPS*, 2018.
- [42] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In *NDSS*, 2018.
- [43] Minghong Fang, Neil Zhenqiang Gong, and Jia Liu. Influence function based data poisoning attacks to top-n recommender systems. In *WWW*, 2020.
- [44] Minghong Fang, Guolei Yang, Neil Zhenqiang Gong, and Jia Liu. Poisoning attacks to graph-based recommender systems. In *ACSAC*, 2018.
- [45] Bo Li, Yining Wang, Aarti Singh, and Yevgeniy Vorobeychik. Data poisoning attacks on factorization-based collaborative filtering. In *NeurIPS*, 2016.
- [46] Guolei Yang, Neil Zhenqiang Gong, and Ying Cai. Fake co-visitation injection attacks to recommender systems. In *NDSS*, 2017.
- [47] Mehran Mozaffari-Kermani, Susmita Sur-Kolay, Anand Raghunathan, and Niraj K Jha. Systematic poisoning attacks on and defenses for machine learning in healthcare. *IEEE journal of biomedical and health informatics*, 2014.
- [48] Qiang Cao, Xiaowei Yang, Jieqi Yu, and Christopher Palow. Uncovering large groups of active malicious accounts in online social networks. In *CCS*, 2014.
- [49] George Danezis and Prateek Mittal. Sybilinifer: Detecting sybil nodes using social networks. In *NDSS*, 2009.
- [50] Neil Zhenqiang Gong, Mario Frank, and Prateek Mittal. Sybilbelief: A semi-supervised learning approach for structure-based sybil detection. *TIFS*, 2014.
- [51] Gianluca Stringhini, Christopher Kruegel, and Giovanni Vigna. Detecting spammers on social networks. In *ACSAC*, 2010.
- [52] Binghui Wang, Jinyuan Jia, and Neil Zhenqiang Gong. Graph-based security and privacy analytics via collective classification with joint weight learning and propagation. In *NDSS*, 2019.
- [53] Gang Wang, Tristan Konolige, Christo Wilson, Xiao Wang, Haitao Zheng, and Ben Y Zhao. You are how you click: Clickstream analysis for sybil detection. In *USENIX Security*, 2013.
- [54] Haifeng Yu, Haifeng Yu, Michael Kaminsky, Phillip B Gibbons, and Abraham Flaxman. Sybilguard: defending against sybil attacks via social networks. In *SIGCOMM*, 2006.
- [55] Dong Yuan, Yuanli Miao, Neil Zhenqiang Gong, Zheng Yang, Qi Li, Dawn Song, Qian Wang, and Xiao Liang. Detecting fake accounts in online social networks at the time of registrations. In *CCS*, 2019.
- [56] Yann Collet. xxHash: Extremely fast hash algorithm. <https://github.com/Cyan4973/xxHash>, 2016.
- [57] Kurt Thomas, Damon McCoy, Chris Grier, Alek Kolcz, and Vern Paxson. Trafficking fraudulent accounts: The role of the underground market in twitter spam and abuse. In *USENIX Security*, 2013.
- [58] Ruggles Steven, Flood Sarah, Goeken Ronald, Grover Josiah, Meyer Erin, Pacas Jose, and Sobek Matthew. IPUMS USA: Version 9.0 [dataset]. minneapolis, mn: Ipums, 2019. <https://doi.org/10.18128/D010.V9.0>, 2019.
- [59] San francisco fire department calls for service. <http://bit.ly/336sddl>, 2019.

A PROOFS AND ADDITIONAL RESULTS

A.1 Proofs

Proof of Lemma 3.3.

$$\begin{aligned}
 \mathbb{E}[\tilde{f}_v] &= \mathbb{E}\left[\frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}}(y_i)(v) - q}{p - q}\right] \\
 &= \frac{\frac{1}{N} \sum_{i=1}^N \mathbb{E}\left[\mathbb{I}_{\text{SUPPORT}}(y_i)(v)\right] - q}{p - q} \\
 &= \frac{\frac{1}{N} * N(f_v(p - q) + q) - q}{p - q} = f_v
 \end{aligned}$$

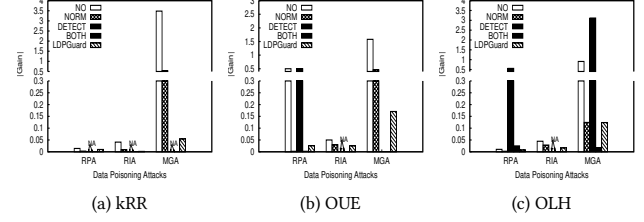


Figure 6: Absolute Gain $|Gain|$ on Zipf.

Table 8: Percentage Estimation on Zipf, $\beta = 0.05$

Protocols	Attacks	$\tilde{\beta}$	
		DETECT	LDPGuard
KRR	RPA	—	0.084
	RIA	—	0.176
	MGA	—	0.049
OUE	RPA	5.7E-05	0.222
	RIA	—	0.442
	MGA	0.050	0.099
OLH	RPA	0.039	0.023
	RIA	—	0.040
	MGA	0.346	0.047

$\mathbb{E}[\tilde{f}_v] = f_v$, $\tilde{f}_v = \frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUPPORT}}(y_i)(v) - q}{p - q}$ is the unbiased estimation of the true frequency f_v of the item v .

Proof of Lemma 4.1. Without loss of generality, we prove it by setting $C = 2$,

$$\begin{aligned}
 \forall y \in \text{Range}(\mathcal{A}_2) : & \frac{\Pr[\mathcal{A}_2(\mathcal{A}_1(v), v) = y]}{\Pr[\mathcal{A}_2(\mathcal{A}_1(v), v') = y]} \\
 &= \frac{\sum_{s \in S} \Pr[\mathcal{A}_1(v) = s] \Pr[\mathcal{A}_2(s, v) = y]}{\sum_{s \in S} \Pr[\mathcal{A}_1(v') = s] \Pr[\mathcal{A}_2(s, v') = y]} \\
 &\leq \sum_{s \in S} \frac{\Pr[\mathcal{A}_1(v) = s] \Pr[\mathcal{A}_2(v) = y]}{\Pr[\mathcal{A}_1(v') = s] \Pr[\mathcal{A}_2(v') = y]} \\
 &\leq e^{\epsilon_1} e^{\epsilon_2} = e^{\epsilon}
 \end{aligned}$$

where $S = \text{Range}(\mathcal{A}_1)$ and $\epsilon_1 + \epsilon_2 = \epsilon$.

Proof of Lemma 5.1.

$$\begin{aligned}
 \mathbb{E}[\tilde{\beta}] &= \mathbb{E}\left[\frac{(N + M)P_1 - CNT}{(N + M)(P_1 - P_2)}\right] \\
 &= \frac{(N + M)P_1 - \mathbb{E}[CNT]}{(N + M)(P_1 - P_2)} \\
 &= \frac{(N + M)P_1 - (NP_1 + MP_2)}{(N + M)(P_1 - P_2)} = \frac{M}{N + M} = \beta
 \end{aligned}$$

$\mathbb{E}[\tilde{\beta}] = \beta$, hence, $\tilde{\beta} = \frac{(N+M)P_1 - CNT}{(N+M)(P_1 - P_2)}$ is the unbiased estimation of the percentage of fake users.

A.2 Additional Experimental Results

Exp 6: Performances of Frequency Estimation on Zipf. We plot results of frequency estimation on Zipf in Figure 6. Similar to the results on IPUMS and Fire, DETECT performs well for MGA to

OLH but is incapable of defending against attacks to kRR and RIA to OUE. In addition, NORM has a limited impact on reducing $|\text{Gain}|$. In contrast, LDPGuard obtains relatively smaller $|\text{Gain}|$ in most scenarios. This observation further demonstrates the effectiveness of LDPGuard in alleviating the effect of data poisoning attacks.

Exp 7: Performances of Percentage Estimation on Zipf. We also compare LDPGuard with DETECT on Zipf and report the estimated percentages of fake users in Table 8. Similar to the results

on IPUMS and Fire, DETECT can accurately estimate the percentage of MGA to OUE where abnormal statistical patterns may exist, but underperforms for RPA to OUE. In contrast, LDPGuard is applicable for all LDP protocols and data poisoning attacks, and its estimated percentages are close to the true percentages of fake users in most scenarios. This observation demonstrates that LDPGuard is able to precisely estimate the percentage of fake users.