

Privacy-aware Synthesizing for Crowdsourced Data

Mengdi Huai¹, Di Wang², Chenglin Miao², Jinhui Xu² and Aidong Zhang¹

¹Department of Computer Science, University of Virginia

²Department of Computer Science and Engineering, State University of New York at Buffalo

¹{mh6ck, aidong}@virginia.edu, ²{dwang45, cmiao, jinhui}@buffalo.edu

Abstract

Although releasing crowdsourced data brings many benefits to the data analyzers to conduct statistical analysis, it may violate crowd users' data privacy. A potential way to address this problem is to employ traditional differential privacy (DP) mechanisms and perturb the data with some noise before releasing them. However, considering that there usually exist conflicts among the crowdsourced data and these data are usually large in volume, directly using these mechanisms can not guarantee good utility in the setting of releasing crowdsourced data. To address this challenge, in this paper, we propose a novel privacy-aware synthesizing method (i.e., PrsCrowd) for crowdsourced data, based on which the data collector can release users' data with strong privacy protection for their private information, while at the same time, the data analyzer can achieve good utility from the released data. Both theoretical analysis and extensive experiments on real-world datasets demonstrate the desired performance of the proposed method.

1 Introduction

In recent years, crowdsourcing has emerged as a popular and fast paradigm to solve many challenging data analysis tasks. Through the power of the crowd, the data collectors (e.g., hospitals, foundations and government agencies) can easily obtain large amounts of useful information. At the same time, the proliferation of new information techniques enables these data collectors to easily share their data that are collected from a crowd of users (e.g., patients, customers) with researchers or data analyzers. From such a wealth of shared data, researchers or data analyzers can discover useful knowledge or patterns to improve the quality of products, the management of public health, and so on. For example, in healthcare applications, the adverse events about a new drug can be easily collected by the hospitals from different patients. If the hospitals are willing to share these medical data, it would be very useful for the drug makers or medical research institutions to understand the efficacy of the drug.

Although the sharing of crowdsourced data brings many benefits, it may introduce privacy issues [Miao *et al.*, 2015;

Shi and Wu, 2017; Miao *et al.*, 2017; Feng *et al.*, 2017]. Considering the above example, the hospital aims to collect the adverse events about a new drug from different patients. The patients usually trust the hospital and are willing to provide all the requested information. But if the hospital directly releases the patients' medical data to the drug makers, the private information of patients would be disclosed. Without effective privacy-preserving mechanisms, the patients may not allow their data to be released. Thus, it is essential to address how to enable the data collectors to release the crowdsourced data without disclosing users' private information.

Among existing privacy-preserving techniques, differential privacy (DP) has drawn significant attention as it can provide very rigorous privacy and utility guarantee [Dwork *et al.*, 2006]. However, this technique has several practical limitations when it is applied in the setting of releasing crowdsourced data. First of all, since the crowdsourced data on an object (e.g., the new drug) are usually collected from multiple users or sources, there inevitably exist conflicts among these data. The reasons include incomplete views of observations, environment noise, different knowledge bases and even the intent to deceive, etc. Directly applying DP on these data can not eliminate the conflicts, and this will certainly degrade the accuracy of the data analysis results. Additionally, DP is usually achieved by adding noise following the Laplace or exponential mechanisms [Dwork *et al.*, 2006]. The noise scale introduced by the Laplace mechanism is proportional to the number of data records, and such noise may make the data useless considering that crowdsourced dataset usually contains large amounts of data records. Although the noise introduced by the exponential mechanism does not depend on the number of data records, it depends on the domains of the input data [Dwork *et al.*, 2014], which may also make the crowdsourced data useless because these data usually have large domains.

To address the above challenges, in this paper, we propose a novel sampling-based **privacy-aware synthesizing** method for **crowdsourced** data (**PrsCrowd**). In this method, the data collector first learns the underlying patterns (i.e., densities) of the data for the objects through assigning each user a fine grained weight (or reliability degree) on each object. Then, for each object, the data collector samples a set of candidate synthetic data from the learned density. Finally, these synthetic data are subjected to our proposed privacy test, and the data collector only releases the synthetics that can pass the privacy

test. The proposed method can not only extract high quality crowdsourced data via differentiating each user’s fine grained reliability degrees on different objects but also achieve DP without injecting noise to the data. Both theoretical analysis and extensive experiments on real-world datasets are provided to verify the desirable performance of the proposed method.

2 Problem Setting

This paper considers a data releasing scenario, where a crowd of users and a data collector are involved. The users (or data sources) are the individuals (e.g., patients, customers) who can observe some objects (e.g., drugs, commodities) and provide claims for them. The data collector is an individual or institution (e.g., a hospital, an online store) who can collect the claims for these objects from a crowd of users and then release these claims to the public either voluntarily or for financial incentives. Here, we assume that the collector is trusted and the security threats mainly come from the public.

Problem formulation. Suppose there are N objects $\mathcal{O} = \{o_i\}_{i=1}^N$ which are observed by M users $\mathcal{U} = \{1, 2, \dots, M\}$. For each object o_i , the claims of users are denoted as $\mathcal{X}_i = \{x_{i,s}\}_{s \in \mathcal{U}_i}$, where $x_{i,s}$ represents the claim provided by user s for object o_i and \mathcal{U}_i represents the set of users who provide claims for this object. The claims collected by the data collector from all users are denoted as $\mathcal{X} = \{\mathcal{X}_i\}_{i=1}^N$, which need to be released to the public. Our goal in this paper is to design a mechanism based on which the data collector can release users’ claims with strong privacy protection for their private information, while at the same time, the data analyzer can achieve good utility from the released data.

3 Preliminary

Definition 1 (Differential Privacy [Dwork *et al.*, 2006]). A randomized algorithm \mathcal{A} is (ϵ, δ) -differentially private if for all neighboring datasets $D, D' \in \mathcal{X}^n$ and for all events S in the output space of \mathcal{A} , the following holds: $\Pr(\mathcal{A}(D) \in S) \leq e^\epsilon \Pr(\mathcal{A}(D') \in S) + \delta$.

The kernel density estimation (KDE) is a statistically-sound method to estimate a continuous distribution. Suppose there are n independent observations $X = \{x_1, \dots, x_n\} \in \mathbb{R}^d$ following an unknown true density $f^*(x)$. The standard KDE $\tilde{f}(x)$ for the estimation of $f^*(x)$ at those points is defined as $\tilde{f}(x) = \frac{1}{n} \sum_{i=1}^n \mathcal{K}_H(x, x_i)$. The following assumption will be used throughout the paper.

Assumption 1. For a vector $x_i \in \mathbb{R}^d$, we assume that the kernel function satisfies $\mathcal{K}_H(x, x_i) = \mathcal{K}_H(x - x_i)$. Furthermore, $\mathcal{K}_H(x - x_i)$ is essentially a bump centered at x_i . More specifically, we take $\mathcal{K}_H(x) = |\mathcal{H}|^{-\frac{1}{2}} \mathcal{K}(\mathcal{H}^{-\frac{1}{2}} z)$, where the kernel \mathcal{K} itself is a probability density with zero mean and identity covariance and satisfying $\lim_{\|x\| \rightarrow \infty} \|x\|^d \mathcal{K}(x) = 0$.

Common choices for \mathcal{K} that satisfy the above assumption include Gaussian and Epanechnikov kernels. As an example, Fig. 6 visualizes the construction of the standard KDE of 5 data points (black circles) using the well-known Gaussian kernel that is defined as $\mathcal{K}_H(x - x_i) = (\frac{1}{\sqrt{2\pi}h})^d \exp(-\frac{\|x-x_i\|^2}{2h^2})$,

where h is the bandwidth. The red curves are the component densities, and each red curve is a scaled version of the normal density curve centered at a datum. The standard KDE is obtained by summing these five scaled components.

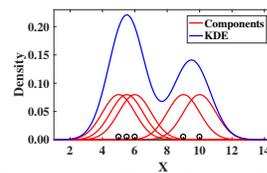


Figure 1: An example for the standard KDE

4 Methodology

4.1 Overview

To achieve the goal described in Section 2, we propose a novel privacy-aware synthesizing method for crowdsourced data (i.e., **PrisCrowd**), which contains two phases. In the first phase, we propose to use the weighted KDE as an intermediate representation of the raw data. This intermediate representation can well capture the statistical properties of the raw data. In the second phase, we first sample a set of candidate synthetic claims from the learned densities in the first phase, then each of these candidate claims is subjected to the proposed privacy test. If the claim passes the privacy test, it will be released, otherwise it will be discarded. The flowchart of the proposed two-phase method is shown in Fig. 2.

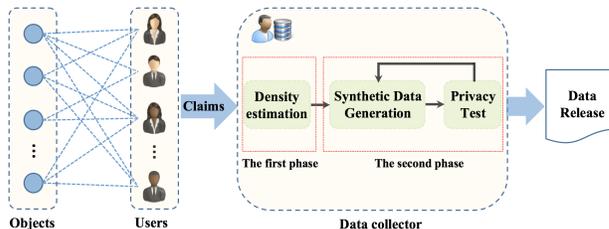


Figure 2: Privacy-aware synthesizing for crowdsourced data

4.2 Weighted KDE-based Data Representation

In order to share “wealth” data with the data analyzer, the data collector first needs to learn the characteristics or the underlying patterns of original data, i.e., the informative density distributions of objects. To estimate the density for each object, the standard kernel density estimation (KDE) can be adopted. Additionally, since different users may provide different claims for the same object, the reliability degrees (or weights) of these users should be taken into account when estimating the densities [Li *et al.*, 2014b; Li *et al.*, 2014a; Li *et al.*, 2016; Miao *et al.*, 2019]. However, the standard KDE cannot differentiate the importance of users (i.e., user reliability degrees). In order to learn users’ reliability and compute the densities of objects simultaneously, we propose a novel method which can estimate users’ global and local weights, and then combine them to learn objects’ informative density distributions. A user’s global weight reflects his capability to provide truthful information for all the objects, and the local weights represent that this user may have different confidence when providing claims for different objects. The advantage of the proposed method is that it can estimate reasonable reliability for each user, and in turn, learn the accurate informative density distributions for objects.

• **Global Weight Estimation.** To evaluate the overall importance of users, the data collector assigns a global weight $g_s \in \mathbb{R}$ to each user s . Meanwhile, we can obtain a global density f_i^g for each object o_i , which should be close to the distribution of claims from reliable users. The distribution of the input claims \mathcal{X}_i can be obtained by $\mathcal{K}_{\mathcal{H}_i}(x, \mathcal{X}_i)$ ($x \in \mathbb{R}$ is a variable), i.e., the kernel function associated with a reproducing kernel Hilbert space \mathcal{H}_i . To minimize the weighted deviation from the estimated density $Q = \{f_i^g(x)\}_{i=1}^N$ to the multi-user input $\mathcal{X} = \{\mathcal{X}_i\}_{i=1}^N$, we propose the following optimization framework

$$\begin{aligned} \min_{G, Q} \sum_{s \in \mathcal{U}} g_s \sum_{i \in \mathcal{E}_s} d_{\mathcal{H}_i}(\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s}), f_i^g(x)) \quad (1) \\ \text{s.t.} \quad \sum_{s \in \mathcal{U}} \exp(g_s) = 1, \end{aligned}$$

where \mathcal{E}_s denotes the set of objects observed by user s , $G = \{g_s\}_{s \in \mathcal{U}}$ and the normalized squared loss $d_{\mathcal{H}_i}(\mathcal{K}_{\mathcal{H}_i}(\cdot, \cdot), f_i^g(x))$ is defined as $d_{\mathcal{H}_i}(\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s}), f_i^g(x)) = \|\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s}) - f_i^g(x)\|_{\mathcal{H}_i}^2$. The global loss function (i.e., Eq. (1)) extends the framework in [Li *et al.*, 2014b] from real space to Hilbert space. We can use an iterative procedure to solve it. Specifically, in the k -th iteration, g_s is updated as

$$g_s^{(k+1)} = -\log \frac{\sum_{i \in \mathcal{E}_s} d_{\mathcal{H}_i}(\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s}), f_i^{g^{(k)}}(x))}{\sum_{s' \in \mathcal{U}} \sum_{i \in \mathcal{E}_{s'}} d_{\mathcal{H}_i}(\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s'}), f_i^{g^{(k)}}(x))}, \quad (2)$$

where $f_i^{g^{(k)}}(x) = \sum_{t \in \mathcal{U}_i} g_t^{(k)} \mathcal{K}_{\mathcal{H}_i}(x, x_{i,t}) / (\sum_{t \in \mathcal{U}_i} g_t^{(k)})$. Eq. (2) shows that a user's global weight is inversely proportional to the distance between its claims and the estimated global densities at the log scale. Users whose claims are close to the derived global densities will have higher global weights.

• **Local Weight Estimation.** As described above, each user may have different confidence when providing claims for different objects. Thus, we need to model the local weight of each user on every object, which will in turn help to infer the accurate density estimations. A potential way to achieve this is to establish a square loss function. However, it leads to a problem that each user would receive the same local weight, and the trustworthiness of the claims provided by different users would be equal. In order to address this problem, we use Hampel loss function [Hampel, 2011]

$$\zeta_{q_1, q_2, q_3}(y) = \begin{cases} y^2/2, & 0 \leq y < q_1 \\ q_1 y - q_1^2/2, & q_1 \leq y < q_2 \\ \frac{q_1(y-q_3)^2}{2(q_2-q_3)} + \frac{q_1(q_2+q_3-q_1)}{2}, & q_2 \leq y < q_3 \\ q_1(q_2+q_3-q_1)/2, & q_3 \leq y, \end{cases}$$

where $q_1 < q_2 < q_3$ are predefined nonnegative parameters. These parameters allow us to decrease the trustworthiness of "bad" claims and increase that of "good" ones for each object, so the importance of users can be well distinguished.

Since we incorporate users' reliability into estimating the local densities, the local kernel density of object o_i can be defined as $f_i^l(x) = \sum_{s \in \mathcal{U}_i} l_{i,s} \mathcal{K}_{\mathcal{H}_i}(x, x_{i,s})$, where $l_{i,s}$ is the

local weight of the user s on object o_i . Thus, the objective function for estimating $l_i = \{l_{i,s}\}_{s \in \mathcal{U}_i}$ is

$$J(l_i) = \min_{l_i} \sum_{s' \in \mathcal{U}_i} \zeta(\|\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s'}) - \sum_{s \in \mathcal{U}_i} l_{i,s} \mathcal{K}_{\mathcal{H}_i}(x, x_{i,s})\|), \quad (3)$$

where $\|\cdot\|$ denotes the difference between users' claims and the estimated local density $f_i^l(x)$. This objective function is not convex, i.e., Eq. (3) does not have a closed form solution. Fortunately, it is possible to approximate $l_i = \{l_{i,s}\}_{s \in \mathcal{U}_i}$ with a standard iteratively re-weighted least squares (IRWLS) algorithm. The iterative procedure for computing $\{l_{i,s}\}_{s \in \mathcal{U}_i}$ is

$$l_{i,s}^{(k+1)} = \frac{\zeta(\|\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s}) - \sum_{t \in \mathcal{U}_i} l_{i,t}^{(k)} \mathcal{K}_{\mathcal{H}_i}(x, x_{i,t})\|)}{\sum_{s' \in \mathcal{U}_i} \zeta(\|\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s'}) - \sum_{t \in \mathcal{U}_i} l_{i,t}^{(k)} \mathcal{K}_{\mathcal{H}_i}(x, x_{i,t})\|)}, \quad (4)$$

where k denotes the number of iterations. Eq. (4) shows that users would receive lower weights when they provide "bad" claims which deviate largely from the center $f_i^l(x) = \sum_{t \in \mathcal{U}_i} l_{i,t} \mathcal{K}_{\mathcal{H}_i}(x, x_{i,t})$.

• **Combined Weight Estimation.** For each user s , to measure the consistency degree of the global and local weights (i.e., g_s and $l_{i,s}$), we define a mixture weight, named combined weight $c_{i,s}$. To learn the combined weight, the relative entropy is employed, which minimizes the information loss between user's global weight and local weight. The smaller the relative entropy value of those weights, the higher the degree of their consistency. The objective of the combined model is

$$\begin{aligned} \min_{\{c_{i,s}\}_{s \in \mathcal{U}_i}} \sum_{s \in \mathcal{U}_i} c_{i,s} \log \frac{c_{i,s}}{l_{i,s}} + \sum_{s \in \mathcal{U}_i} c_{i,s} \log \frac{c_{i,s}}{g_s} \quad (5) \\ \text{s.t.} \quad \sum_{s \in \mathcal{U}_i} c_{i,s} = 1, c_{i,s} \geq 0. \end{aligned}$$

By solving Eq. (5), we can obtain the combined weight $c_{i,s}$ of user s on the object o_i as $c_{i,s} = \sqrt{l_{i,s} g_s} / (\sum_{t \in \mathcal{U}_i} \sqrt{l_{i,t} g_t})$. Based on the learned combined weights, we can obtain the density of object o_i which is the weighted sum of claims in Hilbert space and is given as

$$f_i(x) = \sum_{s \in \mathcal{U}_i} \frac{\sqrt{l_{i,s} g_s}}{\sum_{t \in \mathcal{U}_i} \sqrt{l_{i,t} g_t}} \mathcal{K}_{\mathcal{H}_i}(x, x_{i,s}). \quad (6)$$

4.3 Privacy Test-based Synthetics Release

To provide strong privacy protection for users' private information, in this section, we propose a privacy test-based synthetics release method, which contains two steps: *Candidate synthetics generation* and *Privacy test for candidate synthetics*. In the first step, we sample a set of synthetic claims from the learned density in Eq. (6) as the candidate data to release. Then, in the second step, these sampled synthetics are subjected to a privacy test. If a synthetic claim passes the test, it will be released, otherwise it will be discarded.

Candidate Synthetics Generation. We first discuss how to generate the synthetic claims $\tilde{\mathcal{X}}_i$ for each object o_i . Specifically, we generate each element in $\tilde{\mathcal{X}}_i$ as follows:

1. Select a random integer $s \in \mathcal{U}_i$ with probability $c_{i,s}$;

2. Generate a synthetic claim $\tilde{x}_{i,s}$ through sampling from the probability distribution $\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s})$.

Here $c_{i,s}$ can be treated as the sampling probability that determines whether $x_{i,s}$ is selected or not. In this step, we aim to select some seed data (e.g., $x_{i,s}$) and then probabilistically transform them into the synthetic data. The sampling mechanism used here can increase the uncertainty of the adversary about whether a user's data is in the released dataset or not, and thus it can help to protect users' privacy to some extent. However, it is not enough to only use the sampling mechanism, directly releasing the sampled data can still violate users' privacy [Gehrke *et al.*, 2012]. To tackle this problem, we design the following privacy test mechanism to further prevent users' private information from being disclosed.

Privacy Test for Candidate Synthetics. To prevent an adversary from deducing that a particular claim in \mathcal{X}_i is more responsible for generating the released synthetic data than other claims, the following randomized privacy test mechanism is proposed. Each candidate synthetic data in $\tilde{\mathcal{X}}_i$ is subjected to the randomized privacy test, and it is released only when it passes this test.

Suppose $k \geq 1$ and $\gamma > 1$ are the privacy parameters, and ϵ_0 is the randomness parameter. Let $\mathcal{M}(\cdot)$ denote the above synthetic data generation procedure, which samples a candidate synthetic based on a seed data. Given $x_{i,s} \in \mathcal{X}_i$, we use $\Pr\{\tilde{x}_{i,s} = \mathcal{M}(x_{i,s})\}$ to denote the probability that a synthetic data $\tilde{x}_{i,s}$ is generated based on $\mathcal{M}(\cdot)$. Then the privacy test procedure for $\tilde{x}_{i,s}$ is described as follows:

1. Randomize k by adding a noise: $\tilde{k} = k + z$, where $z \sim \text{Lap}(1/\epsilon_0)$ is sampled from the Laplace Distribution.
2. Let $a \geq 0$ be the integer that satisfies the inequalities $\gamma^{-a-1} < \Pr\{\tilde{x}_{i,s} = \mathcal{M}(x_{i,s})\} \leq \gamma^{-a}$.
3. Let k' be the number of records $x_{i,s'} \in \mathcal{X}_i$ that satisfies $\gamma^{-a-1} < \Pr\{\tilde{x}_{i,s} = \mathcal{M}(x_{i,s'})\} \leq \gamma^{-a}$.
4. If $k' \geq \tilde{k}$, $\tilde{x}_{i,s}$ passes the test, otherwise it fails.

Note that k' denotes the number of possible data seeds that can generate $\tilde{x}_{i,s}$ with a probability value falling into a very stringent interval $[\gamma^{-a-1}, \gamma^{-a}]$. The threshold parameter \tilde{k} prevents releasing sensitive synthetic data. Under this randomized privacy test, a candidate synthetic data is released only when there are at least \tilde{k} possible data seeds that can generate $\tilde{x}_{i,s}$. Intuitively, the larger the value of k , the larger the number of the possible seed data that are indistinguishable from $x_{i,s}$. Also, the less the value of γ , the more difficult to distinguish $x_{i,s}$ from other possible seed data. Algorithm 1 summarizes the proposed privacy test-based synthetics release procedure, in which m denotes the number of synthetic claims that need to be released for object o_i .

4.4 Theoretical Analysis

Consistency Analysis. In Section 4.3, we generate the synthetic claims for object o_i by sampling from the mixture distribution $f_i(x)$, i.e., $\tilde{x}_{i,s} \sim f_i(x)$. After obtaining the dataset $\tilde{\mathcal{X}}_i = \{\tilde{x}_{i,s}\}_{s=1}^m$, a basic question here is that how well the generated dataset can reflect the original density function

Algorithm 1 Private test-based synthetics release for o_i

Input: $\{c_{i,s}\}_{s \in \mathcal{U}_i}$, $\mathcal{X}_i = \{x_{i,s}\}_{s \in \mathcal{U}_i}$, k , γ , ϵ_0 , and m

Output: The synthetic dataset $\tilde{\mathcal{X}}_i$ that can be released

- 1: $\tilde{\mathcal{X}}_i = \emptyset$
 - 2: **while** $|\tilde{\mathcal{X}}_i| < m$ **do**
 - 3: Select a random integer $s \in \mathcal{U}_i$ with probability $c_{i,s}$;
 - 4: Generate a synthetic claim $\tilde{x}_{i,s}$ based on the probability distribution $\mathcal{K}_{\mathcal{H}_i}(x, x_{i,s})$;
 - 5: Conduct randomized privacy test for $\tilde{x}_{i,s}$;
 - 6: **if** $\tilde{x}_{i,s}$ passes the privacy test **then**
 - 7: $\tilde{\mathcal{X}}_i = \tilde{\mathcal{X}}_i \cup \{\tilde{x}_{i,s}\}$;
 - 8: **end if**
 - 9: **end while**
 - 10: **return** $\tilde{\mathcal{X}}_i$;
-

$f_i(x)$. Since each $\tilde{x}_{i,s}$ is sampled from $f_i(x)$ independently, the density function over $\{\tilde{x}_{i,s}\}_{s=1}^m$ can be denoted as $\tilde{f}_i(x) = \frac{1}{m} \sum_{s=1}^m \mathcal{K}_{\mathcal{H}_i}(x, \tilde{x}_{i,s})$. In Theorem 1, we provide the expected squared L_2 -norm distance between $f_i(x)$ and $\tilde{f}_i(x)$.

Theorem 1. Under Assumption 1 for $\mathcal{K}_{\mathcal{H}_i}$ with the diagonal bandwidth matrix $\mathcal{H}_i = \hat{h}^2 I_d$, we further assume that the support of $\mathcal{K}(z)$ satisfies $\|z\| \leq 1$. Then, the expected squared L_2 -norm distance between $f_i(x)$ and $\tilde{f}_i(x)$, i.e., $J = \mathbb{E}[\int (f_i(x) - \tilde{f}_i(x))^2 dx]$, satisfies

$$J \leq 4A\hat{h} + A^2\hat{h}^2V + \frac{B}{m\hat{h}^d} + \frac{ABV}{m\hat{h}^{d-1}}, \quad (7)$$

where $A = \sup_{x \in \mathbb{R}^d} \|\nabla f_i(x)\|$, $B = \int (\mathcal{K}(z))^2 dz$ and V is the volume of the support of $f_i(x)$. The expectation is respected to $\{\tilde{x}_{i,s}\}_{s=1}^m \sim f_i(x)$. This theorem is a general result for d dimensional case, in this paper, the value of d is 1.

Privacy Analysis. Next, we conduct privacy analysis for Algorithm 1. Based on Theorem 2, we know that the proposed algorithm is differentially private.

Theorem 2. Note that the input parameters of Algorithm 1 include $\{c_{i,s}\}_{s \in \mathcal{U}_i}$, $k \geq 1$, $\gamma > 0$, and ϵ_0 . For any neighboring datasets \mathcal{X}_i and \mathcal{X}'_i such that $|\mathcal{X}_i|, |\mathcal{X}'_i| \geq k$ and any integer $1 \leq t < k$, we have that Algorithm 1 is (ϵ, δ) -differentially private, where $\epsilon = \epsilon_0 + \log(1 + \frac{\gamma}{t} \frac{\max_{s \in \mathcal{U}_i} c_{i,s}}{\min_{s \in \mathcal{U}_i} c_{i,s}})$, $\delta = |\mathcal{U}_i| \max_{s \in \mathcal{U}_i} c_{i,s} e^{-\epsilon(k-t)}$.

Remark 1. Note that the proposed Algorithm 1 is different from the mechanism in [Bindschaedler *et al.*, 2017]. The probability of choosing the seed $x_{i,s}$ is non-uniform in Algorithm 1 while that is uniform in [Bindschaedler *et al.*, 2017]. The non-uniform property may generate different parameters of differential privacy. When $\max_{s \in \mathcal{U}_i} c_{i,s} = \min_{s \in \mathcal{U}_i} c_{i,s} = 1/|\mathcal{U}_i|$ (i.e., we uniformly sample the seed $x_{i,s}$), the above Theorem 2 is actually Theorem 1 in [Bindschaedler *et al.*, 2017]. Thus, Theorem 2 in our paper is a generalization of Theorem 1 in [Bindschaedler *et al.*, 2017]. Although the main idea of the proof for Theorem 2 is similar to that in [Bindschaedler *et al.*, 2017], the details are quite different: in [Bindschaedler *et al.*, 2017] the proof consider $\mathcal{X}'_i = \mathcal{X}_i \cup \{x_{i,s'}\}$ as the neighborhood dataset while ours consider $\mathcal{X}'_i = \{\mathcal{X}_i - \{x_{i,s}\}\} \cup \{x_{i,s'}\}$ as the neighborhood dataset. That is because if we add one data record, the probability of sampling seeds,

i.e., $\{c_{i,s}\}$, will be totally changed. So the proof in [Bind-schadler *et al.*, 2017] cannot satisfy our case.

5 Experiments

Performance measure. To evaluate the performance of our method, we adopt the following two measure metrics.

- *ISE*: The integrated squared error (*ISE*) is defined as: $\sum_{i=1}^N \int_{-\infty}^{+\infty} (f_i - \tilde{f}_i)^2 dx$, where f_i and \tilde{f}_i are respectively the original density and the density derived from the synthetic data for object o_i .
- *SISE*: The squared integrated squared error (*SISE*) is defined as: $\sum_{i=1}^N (\int_{-\infty}^{+\infty} (f_i - \tilde{f}_i)^2 dx)^2$. Compared with *ISE*, *SISE* tends to penalize more on the large distance and less on the small distance.

Since the goal of the collector is to release the data whose pattern is similar to the true underlying pattern for the objects, the lower the *ISE* or *SISE*, the better the method.

Datasets. We adopt the following three real-world datasets to evaluate the performance of the proposed method.

- *Population Dataset* [Pasternack and Roth, 2010; Wan *et al.*, 2016]. It is about the population information of some cities at different years, and it contains 2,344 users and 1,124 objects.
- *Stock Dataset* [Li *et al.*, 2012]. This dataset consists of 1000 stock symbols and 16 properties. In this experiment, we only adopt the properties whose data type is continuous. Totally, there are 55 users and 5,521 objects in this dataset.
- *Indoor Floorplan Dataset* [Li *et al.*, 2014a]. It is collected when constructing the indoor floorplans, which is a representative example of social sensing applications. The objects are the hallway segments of a building, the task here is to measure the distances of these segments with the inertial sensors built in the smartphone. It consists of 247 users and 129 objects.

Baselines. Here, we adopt two baselines, i.e. **Basic** and **Uniform**. In the **Basic** method, the data collector adds three level noise to the original data: $\epsilon = 0.1$ (Strong), $\epsilon = 1$ (normal) and $\epsilon = 10$ (Weak). In the **Uniform** method, the collector treats all users equally and the entities' densities are learned with the uniformly weighted kernel density estimation. Here, the synthetic data generation and the privacy tests procedures are the same with those in our proposed method.

Case study. In our proposed method, we take into consideration users' fine grained weights (i.e., the combined weights) when estimating objects' densities. In order to investigate the advantages of the users' combined weights, we conduct case studies on the three real-world datasets. For each dataset, we randomly select two objects as the cases, and then estimate their densities. The estimated densities are shown in Fig. 3. The red line in each subfigure represents the density estimated based on users' combined weights. The black line represents the result estimated only based on the global weight of each user. We also conduct estimations without considering user quality, i.e., treating all users equally, and the estimated density for each object is represented with the green line. The

results in Fig.3 show that the densities estimated based on users' combined weights are the closest to the true densities which are represented with the blue lines. Additionally, we show the claims of each object in this figure with magenta circles and crosses. We can see that some claims (i.e., the magenta crosses) are far away from others (i.e., the magenta circles). These claims are usually provided by the users with low weights, and they can be treated as outliers when estimating each object's density. The results in this figure show that the density estimation method based on the combined weight is more robust to outliers than the methods which only adopt users' global weights or treat all users equally. In other words, the estimated density for each object based Priscrowd can well reflect the underlying true density of this object.

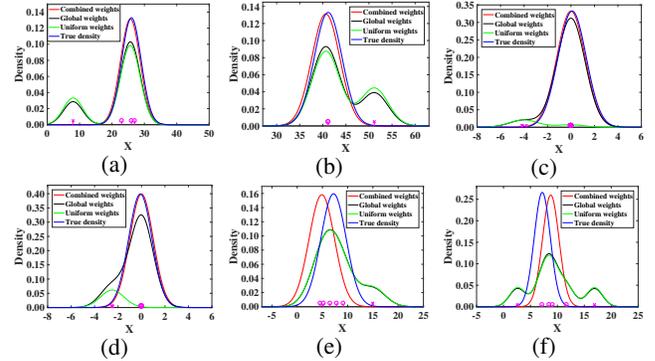


Figure 3: Case study on real-world datasets. (a) and (b): the two cases for Population dataset. (c) and (d): the two cases for Stock dataset. (e) and (f): the two cases for Indoor Floorplan dataset.

Accuracy comparison. As described in Section 4, instead of publishing the raw data collected from the users, the data collector releases the synthetic claims which are sampled from the estimated densities to the public. In this experiment, we evaluate the accuracy (or quality) of the published synthetic data and explore whether these data can well reflect the underlying true densities of the objects. Here we assume that the data collector releases 30 synthetic claims for each object to the public. The parameters γ and k are set as 4 and 5 respectively. In order to evaluate the accuracy of the synthetic claims, we first derive the density (i.e., \tilde{f}_i) of each object from the synthetic data, and then calculate *ISE* and *SISE* for each dataset. The results are shown in Table 1, from which we can see the proposed approach performs much better than the baseline methods on all real-world datasets. That is to say, the synthetic data generated based on our proposed method could well preserve the characteristics of the underlying pattern for the objects, so the data analyzers could achieve much better utility from the published data. Additionally, the results in Table 1 also show that the advantages of our proposed approach on the Stock dataset is larger than that on the Population and Indoor Floorplan datasets. The reason is that there are more outlying data points in the Stock dataset, and our proposed approach is robust to these outliers while the baseline methods are very sensitive to them.

The effect of the number of sampled claims. In our proposed method, the data collector needs to release the synthetic claims which are sampled from the estimated densities. Thus, the number of sampled claims for each object plays an impor-

Measure	Method	Population	Stock	Indoor
ISE	PrisCrowd	0.479	1.699	1.051
	Uniform	0.628	17.628	1.220
	Basic(Strong)	1.183	12.799	1.943
	Basic(Normal)	1.119	9.867	1.937
	Basic(Weak)	0.866	2.125	1.882
SISE	PrisCrowd	6.209	11.430	8.405
	Uniform	8.502	47.217	11.112
	Basic(Strong)	12.013	40.420	15.111
	Basic(Normal)	11.768	35.391	15.046
	Basic(Weak)	10.149	15.412	14.723

Table 1: Accuracy comparison on the real-world datasets

tant role during the data releasing procedure. In this experiment, we evaluate the effect of the number of sampled claims for each object on the performance of the proposed method. Here we vary the number of the sampled claims for each object from 1 to 30 and then calculate the *ISE* and *SISE* on the three real-world datasets. The results are shown in Figure 4, from which we can see the *ISE* and *SISE* gradually get flattened with the increase of the number of the sampled claims for each object. Take the population dataset as an example, when the number of sampled claims is larger than 10, the values of *ISE* and *SISE* are almost the same. That is to say, the released data generated based on our proposed method could well reflect the underlying patterns of the objects even only a few claims are sampled for each object.

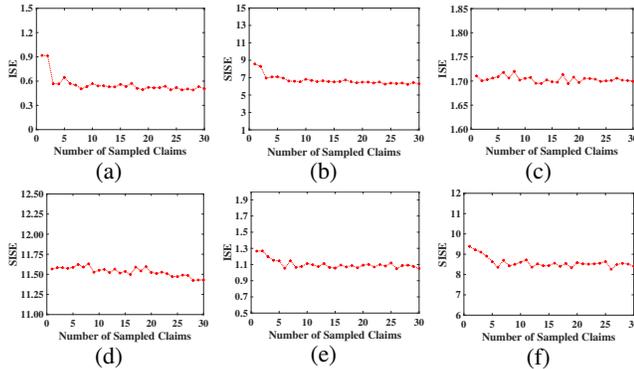


Figure 4: Accuracy w.r.t. Number of Sampled Claims on Real-world Datasets. (a) and (b): Population. (c) and (d): Stock. (e) and (f): Indoor Floorplan.

Computational Cost. Next we evaluate the computational cost of the synthetic claims generation procedure, i.e., the second phase in our proposed method. In this experiment, we only generate synthetic claims for the objects whose ground truths can be achieved from the original datasets, and consider two scenarios, i.e., with privacy tests and without privacy tests. Then we vary the number of the sampled claims for each object from 1 to 30. The running time of the synthetic claims generation procedure for the three datasets is shown in Fig. 5, from which we can see the running time in the two scenarios is approximately linear with respect to the number of sampled objects for each object. Additionally, the results also show that the privacy test step introduce extra computational cost during the released data generation procedure. This is because each candidate synthetic data record needs to be tested in the privacy test step such that strong privacy protection can

be guaranteed for users’ raw data. Since good utility can be achieved based on our proposed method even only a few synthetic claims are generated for each object, the computational cost is tolerable in practice.

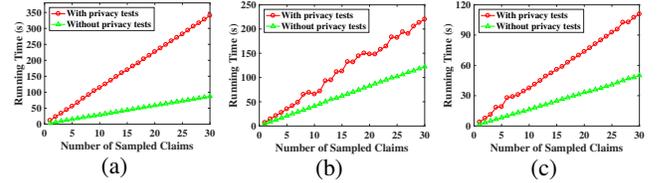


Figure 5: Running time vs. number of sampled claims for each object. (a): Population. (b): Stock. (c): Indoor Floorplan.

6 Related Work

Recently, various differential private data release approaches have been proposed. Those methods can be roughly partitioned into two categories: the interactive ones and the non-interactive ones. In an interactive method [Li *et al.*, 2010; Hardt and Rothblum, 2010; Roth and Roughgarden, 2010], a data analyzer can pose queries via a private mechanism, and a dataset owner answers these queries in response. In the non-interactive framework [Nissim *et al.*, 2007; Bindschaedler and Shokri, 2016; Blum *et al.*, 2013; Wang *et al.*, 2018; Wang *et al.*, 2019], a data owner releases the private version of the original data. Once data are published, the owner has no further control over the published data.

The method in our paper is non-interactive. The typical approach to protect data privacy in the non-interactive context is to directly add noise, which is taken by [Bindschaedler and Shokri, 2016; Blum *et al.*, 2013]. These works are either computationally infeasible on high-dimensional data, or practically ineffective because of their large utility costs. There are also some other works [Bindschaedler and Shokri, 2016] which release private data without adding noise, but they are unsuitable to be used in the newly appearing crowdsourcing setting considered in this paper where multi-sources provide multi-observations for multi-objects.

7 Conclusions

In this paper, we propose a novel privacy-aware synthesizing method for crowdsourced data. Based on this method, the data collector can release the crowdsourced data with strong privacy protection for users’ private information, while at the same time, the data analyzer can achieve good utility from the released data. Both theoretical analysis and extensive experiments on real-world datasets verify the effectiveness of the proposed method.

Acknowledgement

This work was supported in part by the US National Science Foundation (NSF) under grants IIS-1514204 and CCF-1716400. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the NSF.

References

- [Bindschaedler and Shokri, 2016] Vincent Bindschaedler and Reza Shokri. Synthesizing plausible privacy-preserving location traces. In *2016 IEEE Symposium on Security and Privacy (SP)*, pages 546–563. IEEE, 2016.
- [Bindschaedler et al., 2017] Vincent Bindschaedler, Reza Shokri, and Carl A Gunter. Plausible deniability for privacy-preserving data synthesis. *Proceedings of the VLDB Endowment*, 10(5):481–492, 2017.
- [Blum et al., 2013] Avrim Blum, Katrina Ligett, and Aaron Roth. A learning theory approach to noninteractive database privacy. *Journal of the ACM (JACM)*, 60(2):12, 2013.
- [Dwork et al., 2006] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [Dwork et al., 2014] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [Feng et al., 2017] Wei Feng, Zheng Yan, Hengrun Zhang, Kai Zeng, Yu Xiao, and Y Thomas Hou. A survey on security, privacy, and trust in mobile crowdsourcing. *IEEE Internet of Things Journal*, 5(4):2971–2992, 2017.
- [Gehrke et al., 2012] Johannes Gehrke, Michael Hay, Edward Lui, and Rafael Pass. Crowd-blending privacy. In *Advances in Cryptology—CRYPTO 2012*, pages 479–496. Springer, 2012.
- [Hampel, 2011] et al. Hampel, Frank R. *Robust statistics: the approach based on influence functions*. John Wiley & Sons, 2011.
- [Hardt and Rothblum, 2010] Moritz Hardt and Guy N Rothblum. A multiplicative weights mechanism for privacy-preserving data analysis. In *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, pages 61–70. IEEE, 2010.
- [Li et al., 2010] Chao Li, Michael Hay, Vibhor Rastogi, Gerome Miklau, and Andrew McGregor. Optimizing linear counting queries under differential privacy. In *Proceedings of the twenty-ninth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 123–134. ACM, 2010.
- [Li et al., 2012] Xian Li, Xin Luna Dong, Kenneth Lyons, Weiyi Meng, and Divesh Srivastava. Truth finding on the deep web: Is the problem solved? *Proceedings of the VLDB Endowment*, 6(2):97–108, 2012.
- [Li et al., 2014a] Qi Li, Yaliang Li, Jing Gao, Lu Su, Bo Zhao, Murat Demirbas, Wei Fan, and Jiawei Han. A confidence-aware approach for truth discovery on long-tail data. *PVLDB*, 2014.
- [Li et al., 2014b] Qi Li, Yaliang Li, Jing Gao, Bo Zhao, Wei Fan, and Jiawei Han. Resolving conflicts in heterogeneous data by truth discovery and source reliability estimation. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 1187–1198. ACM, 2014.
- [Li et al., 2016] Yaliang Li, Qi Li, Jing Gao, Lu Su, Bo Zhao, Wei Fan, and Jiawei Han. Conflicts to harmony: A framework for resolving conflicts in heterogeneous data by truth discovery. *TKDE*, 2016.
- [Miao et al., 2015] Chenglin Miao, Wenjun Jiang, Lu Su, Yaliang Li, Suxin Guo, Zhan Qin, Houping Xiao, Jing Gao, and Kui Ren. Cloud-enabled privacy-preserving truth discovery in crowd sensing systems. In *Proceedings of the 13th ACM Conference on Embedded Networked Sensor Systems*, pages 183–196. ACM, 2015.
- [Miao et al., 2017] Chenglin Miao, Lu Su, Wenjun Jiang, Yaliang Li, and Miaomiao Tian. A lightweight privacy-preserving truth discovery framework for mobile crowd sensing systems. In *IEEE INFOCOM 2017-IEEE Conference on Computer Communications*, pages 1–9. IEEE, 2017.
- [Miao et al., 2019] Chenglin Miao, Wenjun Jiang, Lu Su, Yaliang Li, Suxin Guo, Zhan Qin, Houping Xiao, Jing Gao, and Kui Ren. Privacy-preserving truth discovery in crowd sensing systems. *ACM Transactions on Sensor Networks (TOSN)*, 15(1):9, 2019.
- [Nissim et al., 2007] Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. Smooth sensitivity and sampling in private data analysis. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 75–84. ACM, 2007.
- [Pasternack and Roth, 2010] Jeff Pasternack and Dan Roth. Knowing what to believe (when you already know something). In *Proc. of Coling*, 2010.
- [Roth and Roughgarden, 2010] Aaron Roth and Tim Roughgarden. Interactive privacy via the median mechanism. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages 765–774. ACM, 2010.
- [Shi and Wu, 2017] Xinghua Shi and Xintao Wu. An overview of human genetic privacy. *Annals of the New York Academy of Sciences*, 1387(1):61–72, 2017.
- [Wan et al., 2016] Mengting Wan, Xiangyu Chen, Lance Kaplan, Jiawei Han, Jing Gao, and Bo Zhao. From truth discovery to trustworthy opinion discovery: An uncertainty-aware quantitative modeling approach. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1885–1894. ACM, 2016.
- [Wang et al., 2018] Di Wang, Marco Gaboardi, and Jinhui Xu. Empirical risk minimization in non-interactive local differential privacy revisited. In *Advances in Neural Information Processing Systems*, pages 965–974, 2018.
- [Wang et al., 2019] Di Wang, Adam Smith, and Jinhui Xu. Noninteractive locally private learning of linear models via polynomial approximations. In *Algorithmic Learning Theory*, pages 897–902, 2019.

A An example for KDE

As an example, Fig. 6 visualizes the construction of the standard KDE of 5 data points (black circles) using the well-known Gaussian kernel that is defined as $\mathcal{K}_{\mathcal{H}}(x - x_i) = (\frac{1}{\sqrt{2\pi}h})^d \exp(-\frac{\|x-x_i\|^2}{2h^2})$, where h is the bandwidth. The red curves are the component densities, and each red curve is a scaled version of the normal density curve centered at a datum. The standard KDE is obtained by summing these five scaled components.

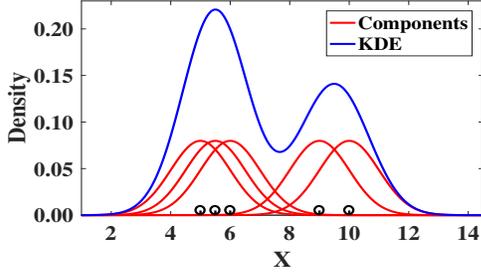


Figure 6: An example for the standard KDE

B Proof of Theorem 1

This section will prove Theorem 1, of which the proof procedure will utilize the results of Lemma 1 and 2. Let $T_1 = \int (\tilde{f}_i(x) - f_i(x))f_i(x)dx$, and $T_2 = \int (E[\tilde{f}_i^2(x)] - f_i^2(x))dx$. Then, we have

$$\begin{aligned} \Delta_{f_i} &= E[(f_i(x) - \tilde{f}_i(x))^2 dx] \\ &= \int f_i^2(x)dx - 2 \int E[\tilde{f}_i(x)]f_i(x)dx + \int \tilde{f}_i^2(x)dx \\ &= \int f_i^2(x)dx - 2 \int E([\tilde{f}_i(x)] - f_i(x))f_i(x)dx \\ &\quad - 2 \int f_i^2(x)dx + \int (\tilde{f}_i^2(x) - f_i^2(x))dx + \int f_i^2(x)dx \\ &= -2 \int E([\tilde{f}_i(x)] - f_i(x))f_i(x)dx + \int (E[\tilde{f}_i^2(x)] - f_i^2(x)) \\ &= -2T_1 + T_2, \end{aligned}$$

Since $\Delta_{f_i} = E[(f_i(x) - \tilde{f}_i(x))^2 dx] \geq 0$, we have $\Delta_{f_i} = -2T_1 + T_2 = |-2T_1 + T_2| \leq |-2T_1| + |T_2| = |2T_1| + |T_2|$.

Combining the results of Lemmas 1 and 2 which provide the upper bounds for T_1 and T_2 , gives us the followings result.

$$\begin{aligned} \Delta_{f_i} &\leq |2T_1| + |T_2| \leq 2A_i\tilde{h}_i + 2A_i\tilde{h}_i + A_i^2\tilde{h}_i^2V_i + \frac{\rho B}{|\tilde{S}_i|\tilde{h}_i^d} + \frac{\rho A_i B V_i}{|\tilde{S}_i|\tilde{h}_i^{d-1}} \\ &= 4A_i\tilde{h}_i + A_i^2\tilde{h}_i^2V_i + \frac{\rho B}{|\tilde{S}_i|\tilde{h}_i^d} + \frac{\rho A_i B V_i}{|\tilde{S}_i|\tilde{h}_i^{d-1}}. \end{aligned}$$

Thus, it is clear that f_i and \tilde{f}_i will be close if $|\tilde{S}_i|$ is large enough and if \tilde{h}_i is chosen properly as a function of $|\tilde{S}_i|$.

Below we will denote $\tilde{x}_{i,s} = \frac{1}{|\tilde{S}_i|}$.

Lemma 1. Let T_1 be $T_1 = \int E([\tilde{f}_i(x)] - f_i(x))f_i(x)dx$. Then, $|T_1| \leq A_i\tilde{h}_i$.

$$\begin{aligned} \text{Proof. } T_1 &= \int E([\tilde{f}_i(x)] - f_i(x))f_i(x)dx \\ &= \int (E[\sum_{s \in \tilde{S}_i} \tilde{l}_{i,s} \mathcal{K}_{\tilde{H}_i}(x - \tilde{x}_{i,s})] - f_i(x))f_i(x)dx \\ &= \int (\sum_{s \in \tilde{S}_i} \tilde{l}_{i,s} \int \mathcal{K}_{\tilde{H}_i}(x - \tilde{x}_{i,s})f_i(\tilde{x}_{i,s})d\tilde{x}_{i,s} - f_i(x))f_i(x)dx \end{aligned}$$

$$\begin{aligned} &\leq \int (\sum_{s \in \tilde{S}_i} \tilde{l}_{i,s} |\int \mathcal{K}_{\tilde{H}_i}(x - \tilde{x}_{i,s})f_i(\tilde{x}_{i,s})d\tilde{x}_{i,s} - f_i(x)|)f_i(x)dx \\ &\leq \int (\sum_{s \in \tilde{S}_i} \tilde{l}_{i,s} A_i \tilde{h}_i) f_i(x)dx = \int (A_i \tilde{h}_i) f_i(x)dx = A_i \tilde{h}_i. \end{aligned}$$

Lemma B.3 and that $\int f_i(x)dx = 1$ are used in the last line. \square

Lemma 2. Let T_2 be as in above, and let V_i be the volume of the support of f_i . Then, $T_2 \leq 2A_i\tilde{h}_i + A_i^2\tilde{h}_i^2V_i + \frac{\rho B}{|\tilde{S}_i|\tilde{h}_i^d} + \frac{\rho A_i B V_i}{|\tilde{S}_i|\tilde{h}_i^{d-1}}$.

Proof. Since $\tilde{f}_i(x) = \sum_{s \in \tilde{S}_i} \tilde{l}_{i,s} \mathcal{K}_{\tilde{H}_i}(x - \tilde{x}_{i,s})$, then $\tilde{f}_i^2(x) = \sum_{s \neq t} \tilde{l}_{i,s} \tilde{l}_{i,t} \mathcal{K}_{\tilde{H}_i}(x - \tilde{x}_{i,s}) \mathcal{K}_{\tilde{H}_i}(x - \tilde{x}_{i,t}) + \sum_s \tilde{l}_{i,s}^2 \mathcal{K}_{\tilde{H}_i}^2(x - \tilde{x}_{i,s})$. For convenience, we drop the \tilde{H}_i subscript and denote $\mathcal{K}_{\tilde{H}_i}(x - \tilde{x}_{i,s})$ as $\mathcal{K}(\tilde{x}_{i,s})$. Then, T_2 can be simplified as follows:

$$\begin{aligned} T_2 &= \int (E[\tilde{f}_i^2(x)] - f_i^2(x))dx \\ &= \int (E[\sum_{s \neq t} \tilde{l}_{i,s} \tilde{l}_{i,t} \mathcal{K}(\tilde{x}_{i,s}) \mathcal{K}(\tilde{x}_{i,t}) + \sum_s \tilde{l}_{i,s}^2 \mathcal{K}^2(\tilde{x}_{i,s})] - f_i^2(x))dx \\ &= \int (\sum_{s \neq t} \tilde{l}_{i,s} \tilde{l}_{i,t} \int \int \mathcal{K}(\tilde{x}_{i,s}) \mathcal{K}(\tilde{x}_{i,t}) f(\tilde{x}_{i,s}) f(\tilde{x}_{i,t}) d\tilde{x}_{i,s} d\tilde{x}_{i,t} \\ &\quad + \sum_s \tilde{l}_{i,s}^2 \int \mathcal{K}^2(\tilde{x}_{i,s}) f_i(\tilde{x}_{i,s}) d\tilde{x}_{i,s} - \sum_{s \neq t} \tilde{l}_{i,s} \tilde{l}_{i,t} f_i^2(x) - \sum_s \tilde{l}_{i,s}^2 f_i^2(x))dx \\ &= \int (\sum_{s \neq t} \tilde{l}_{i,s} \tilde{l}_{i,t} [\int \int \mathcal{K}(\tilde{x}_{i,s}) \mathcal{K}(\tilde{x}_{i,t}) f(\tilde{x}_{i,s}) f(\tilde{x}_{i,t}) d\tilde{x}_{i,s} d\tilde{x}_{i,t} \\ &\quad - f_i^2(x)] + \sum_s \tilde{l}_{i,s}^2 \int \mathcal{K}^2(\tilde{x}_{i,s}) f_i(\tilde{x}_{i,s}) d\tilde{x}_{i,s} - \sum_s \tilde{l}_{i,s}^2 f_i^2(x))dx \\ &\leq \int (\sum_{s \neq t} \tilde{l}_{i,s} \tilde{l}_{i,t} (2A_i \tilde{h}_i f_i(x) + A_i^2 \tilde{h}_i^2) + \sum_s \tilde{l}_{i,s}^2 B(f_i(x) + A_i \tilde{h}_i) \tilde{h}_i^{-d}) dx \\ &= 2A_i \tilde{h}_i + A_i^2 \tilde{h}_i^2 V_i + \frac{\rho B}{|\tilde{S}_i|\tilde{h}_i^d} + \frac{\rho A_i B V_i}{|\tilde{S}_i|\tilde{h}_i^{d-1}}. \end{aligned}$$

In the above, we have made use of Lemmas B.4 and the fact that $\sum_{s \neq t} \tilde{l}_{i,s} \tilde{l}_{i,t} \leq 1$. And, $0 \leq \rho \leq |\tilde{S}_i|$. \square

Lemma 3. For \mathcal{O}_i , let $A_i = \sup_{x \in \mathbb{R}} \|\nabla f_i(x)\|$. Then $\|\int \mathcal{K}_{\tilde{H}_i}(x - y)f_i(y)dy - f_i(x)\| \leq A_i \tilde{h}_i$.

Proof. The proof of this lemma contains the following two steps. Firstly, let $x, y \in \mathbb{R}$ such that $\|x - y\| \leq \tilde{h}_i$. And, $A_i = \sup_{x \in \mathbb{R}} \|\nabla f_i(x)\|$. Define a function $\beta : [0, 1] \rightarrow \mathbb{R}$, $\beta(t) = (1 - t)x + ty$. Then, we have

$$\begin{aligned} |f_i(y) - f_i(x)| &= |f_i(\beta(1)) - f_i(\beta(0))| = |\int_0^1 \frac{d}{dt} f_i(\beta(t)) dt| \\ &= |\int_0^1 \nabla f_i(\beta(t)) \cdot \beta'(t) dt| = |\int_0^1 \nabla f_i(\beta(t)) \cdot (y - x) dt| \\ &\leq |\int_0^1 \nabla f_i(\beta(t)) \cdot (y - x) dt| \leq |\int_0^1 \nabla f_i(\beta(t)) \cdot (y - x) dt| \\ &\leq \int_0^1 A_i \tilde{h}_i dt = A_i \tilde{h}_i, \end{aligned}$$

In the above, we have made use of the Holder and Cauchy-Schwarz inequalities. Based on this, we then have

$$\begin{aligned} |\int \mathcal{K}_{\tilde{H}_i}(x - y)f_i(y)dy - f_i(x)| &= |\int \mathcal{K}_{\tilde{H}_i}(x - y)[f_i(y) - f_i(x)]dy| \\ &\leq \int \mathcal{K}_{\tilde{H}_i}(x - y)|f_i(y) - f_i(x)|dy \leq \int \mathcal{K}_{\tilde{H}_i}(x - y)A_i \tilde{h}_i dy = A_i \tilde{h}_i. \end{aligned}$$

In the above, the fact that $\int \mathcal{K}_{\tilde{H}_i}(\cdot) = 1$ is used. We have also made use of the fact that in the support of $\mathcal{K}_{\tilde{H}_i}(\cdot)$, we have

that $\|y - x\| \leq \tilde{h}_i$ and therefore that $|f_i(y) - f_i(x)| \leq A_i \tilde{h}_i$ applies. \square

Lemma 4. We have

$$\left| \int \int \mathcal{K}_{\tilde{h}_i}(x-y) \mathcal{K}_{\tilde{h}_i}(x-z) f(y) f(z) dy dz - f_i^2(x) \right| \leq 2A_i \tilde{h}_i f_i(x) + A_i \tilde{h}_i^2, \quad (8)$$

$$\int \mathcal{K}_{\tilde{h}_i}^2(x-y) f_i(y) dy \leq B(f_i(x) + A_i \tilde{h}_i) \tilde{h}_i^{-d}. \quad (9)$$

Proof. Firstly, the proof procedure of (8) is as follows.

$$\begin{aligned} &= \left| \left(\int \mathcal{K}_{\tilde{h}_i}(x-y) f_i(y) dy \right) \left(\int \mathcal{K}_{\tilde{h}_i}(x-z) f_i(z) dz \right) - f_i^2(x) \right| \\ &= \left| \left(\int \mathcal{K}_{\tilde{h}_i}(x-y) f_i(y) dy \right)^2 - f_i^2(x) \right| \\ &= \left| \left(\int \mathcal{K}_{\tilde{h}_i}(x-y) f_i(y) dy - f_i(x) \right) \left(\int \mathcal{K}_{\tilde{h}_i}(x-y) f_i(y) dy + f_i(x) \right) \right| \\ &\leq A_i \tilde{h}_i (2f_i(x) + A_i \tilde{h}_i) = 2A_i \tilde{h}_i f_i(x) + A_i \tilde{h}_i^2. \end{aligned}$$

In the above product, the upper bound of the first term is $A_i \tilde{h}_i$. The second term is upper bounded by $2f_i(x) + A_i \tilde{h}_i$ since we can get $\int \mathcal{K}_{\tilde{h}_i}(x-y) f_i(y) dy \leq f_i(x) + A_i \tilde{h}_i$ from Lemma 3.

Secondly, we will prove (9). Based on Lemma 3, we know that within the integrand's support, $|f_i(y) - f_i(x)| \leq A_i \tilde{h}_i$, so that $f_i(y) \leq f_i(x) + A_i \tilde{h}_i$. Also, $\mathcal{K}_{\tilde{h}_i}(z) = \tilde{h}_i^{-1} \mathcal{K}(\tilde{h}_i^{-1}z)$.

Then, we have

$$\begin{aligned} \int \mathcal{K}_{\tilde{h}_i}^2(x-y) f_i(y) dy &\leq \int \mathcal{K}_{\tilde{h}_i}^2(x-y) (f_i(x) + A_i \tilde{h}_i) dy \\ &= (f_i(x) + A_i \tilde{h}_i) \int \mathcal{K}^2(z) dz = (f_i(x) + A_i \tilde{h}_i) \tilde{h}_i^{-1} \int \tilde{h}_i^{-1} \mathcal{K}^2(\tilde{h}_i^{-1}z) dz \\ &= (f_i(x) + A_i \tilde{h}_i) \tilde{h}_i^{-1} \int \mathcal{K}^2(w) dw = B(f_i(x) + A_i \tilde{h}_i) \tilde{h}_i^{-d}. \end{aligned}$$

In the above, a change of variables is used, i.e., $w = \tilde{h}_i^{-1} I_d z$. \square

C Proof of Theorem 2

We consider two neighbor datasets D and D' with size n , we assume $D = \{x_1, x_2, \dots, x_n\}$ and $D' = \{x_1, x_2, \dots, x_{n-1}, x'_n\}$, that is $D' = D - \{x_n\} \cup \{x'_n\}$. We denote $P_d(y) = \Pr\{y = \mathcal{M}(d)\}$ for data record d and the data universe as \mathcal{Z} . Also we denote the Algorithm 1 as \mathcal{F} . Our will show the differential privacy for \mathcal{F} following [Bindschaedler *et al.*, 2017].

Now we fix y , where $y \in \mathcal{Z}$. The records in a dataset D^* can be partitioned by $I_d(y)$, where $I_d(y)$ is the unique integer which satisfies $\gamma^{-I_d(y)-1} < P_d(y) \leq \gamma^{-I_d(y)}$. If $P_d(y) = 0$, then we denote $I_d(y)$ as \emptyset , thus we can define the partition set for i as $C_i(D^*, y) = \{d : d \in D^*, I_d(y) = i\}$.

Lemma 5. For any dataset D^* , if the seed is in partition set for i . The probability of passing the privacy test is giving by: $pt(D^*, i, y) = \Pr\{L \geq k - |C_i(D^*, y)|\}$ where $L \sim \text{Lap}(\frac{1}{\epsilon_0})$.

Now we assume x_n falls in partition set of j while x'_n falls in partition set of k (if either of them falls into \emptyset the proof is the same), W.L.O.G we assume $j \neq k$ (the same proof for $j = k$). We now denote $p(s)$ is the probability we choose the i -th data record, $s \in D$ as the seed of generating the synthetic record.

Thus, it is only depend on the position of the data in D and not dependent on the data record it self, in Algorithm 1, it is just $p(x_{i,j}) = c_{i,j}$

Lemma 6. For any dataset D^* and data record $y \in \mathcal{Z}$, we have

$$\Pr\{\mathcal{F} = y\} = \sum_{i \geq 0} \left(\sum_{s \in C_i(D^*, y)} p(s) P_s(y) \right) pt(D^*, i, y). \quad (10)$$

Proof.

$$\Pr\{\mathcal{F} = y\} = \sum_{s \in D^*} \Pr\{s \text{ is the seed}, \mathcal{F}(D^*) = y\} \quad (11)$$

$$= \sum_s \Pr\{\text{select } s \text{ as the seed}\} \Pr\{(D^*, s, y) \text{ passes the test}\} \quad (12)$$

$$= \sum_s p(s) P_s(y) \Pr\{y \text{ passes the test}\} \quad (13)$$

$$= \sum_{i \geq 0} \left(\sum_{s \in C_i(D^*, y)} p(s) P_s(y) \right) pt(D^*, i, y) \quad (14)$$

\square

We now denote

$$q(D^*, i, y) = \left(\sum_{s \in C_i(D^*, y)} p(s) P_s(y) \right) pt(D^*, i, y).$$

Lemma 7. For D', D as above,

$$e^{-\epsilon_0} pt(D, i, y) \leq pt(D', i, y) \leq e^{\epsilon_0} pt(D, i, y). \quad (15)$$

Proof. If $i \neq j, k$, that means $C_i(D, y) = C_i(D', y)$, so we have $pt(D, i, y) = pt(D', i, y)$.

If $i = j$, that means $|C_i(D', y)| = |C_i(D, y)| + 1$ (since x'_n falls in), by the property of Laplace distribution we have:

$$\begin{aligned} pt(D, i, y) &= \Pr\{L \geq k - |C_i(D, y)|\} \\ &\leq \Pr\{L \geq k - |C_i(D, y)| - 1\} \\ &\leq e^{\epsilon_0} \Pr\{L \geq k - |C_i(D, y)|\} \end{aligned}$$

That is $pt(D, i, y) \leq pt(D', i, y) \leq e^{\epsilon_0} pt(D, i, y)$.

If $i = k$, we have the same $e^{-\epsilon_0} pt(D, i, y) \leq pt(D', i, y) \leq pt(D, i, y)$. Thus in total we have

$$e^{-\epsilon_0} pt(D, i, y) \leq pt(D', i, y) \leq e^{\epsilon_0} pt(D, i, y). \quad \square$$

Lemma 8. For all $i \neq j, k$

$$q(D, i, y) = q(D', i, y) \quad (16)$$

For $i = j$, we have

$$q(D, i, y) \leq q(D', i, y) \quad (17)$$

Furthermore, if $|C_j(D, y)| \leq t$,

$$q(D', i, y) \leq t e^{-\epsilon_0(k-t)} \max p(s) \quad (18)$$

Otherwise,

$$q(D', i, y) \leq e^{\epsilon_0} \left[1 + \frac{r \max p(s)}{t \min p(s)} \right] q(D, i, y) \quad (19)$$

For $i = k$, we have when $|C_k(D, y)| > t$

$$q(D, i, y) \leq e^{\epsilon_0} \left[1 + \frac{r \max p(s)}{t \min p(s)} \right] q(D', i, y) \quad (20)$$

Otherwise $q(D, k, y) \leq t e^{-\epsilon_0(k-t)} \max p(s)$ and also

$$q(D', i, y) \leq q(D, i, y) \quad (21)$$

Proof. (16) is oblivious since for $i \neq j, k$ we have $C_i(D, y) = C_i(D', y)$, also by the proof of Lemma 7 we have $ptt(D, i, y) = pt(D', i, y)$. When $i = j$, we have

$$\begin{aligned} q(D, i, y) &= pt(D, i, y) \sum_{s \in C_i(D, y)} p(s) P_s(y) \\ &\leq pt(D', i, y) \left(\sum_{s \in C_i(D, y)} p(s) P_s(y) + p(x'_n) P_{x'_n}(y) \right) \\ &= q(D', i, y) \end{aligned}$$

If $|C_i(D, y)| < t$, since $pt(D', i, y) = Pr\{L \geq k - |C_i(D', y)|\} \leq Pr\{L \geq k - t\} = \frac{1}{2}e^{-\epsilon_0(k-t)}$, we have $q(D', i, y) \leq e^{-\epsilon_0(k-t)} \sum_{s \in C_i(D', y)} P_s(y) p(s) \leq te^{-\epsilon_0(k-t)} \max p(s)$, which is (18). If $|C_i(D, y)| \geq t$, then we have

$$\begin{aligned} q(D', i, y) &= pt(D', i, y) \sum_{s \in C_i(D', y)} p(s) P_s(y) \\ &\leq e^{\epsilon_0} pt(D, i, y) \left(\sum_{s \in C_i(D, y)} p(s) P_s(y) + p(x'_n) P_{x'_n}(y) \right) \quad (22) \end{aligned}$$

By definition, we know $P_{x'_n}(y) \leq \gamma P_s(y)$ for every $s \in C_i(D, y)$, so we have

$$\begin{aligned} P_{x'_n}(y) &\leq \frac{r}{|C_i(D, y)|} \sum_{s \in C_i(D, y)} \gamma P_s(y) \\ &\leq \frac{r}{t} \sum_{s \in C_i(D, y)} \gamma P_s(y) \end{aligned}$$

Also, by the definition of $p(s)$, we have $\frac{p(x'_n)}{p(s)} = \frac{p(x_n)}{p(s)} \leq \frac{\max p(s)}{\min p(s)}$. Thus, we have the following

$$p(x'_n) P_{x'_n} \leq \frac{\gamma \max p(s)}{t \min p(s)} \sum_{s \in C_i(D, y)} p(s) P_s(y) \quad (23)$$

Take it into (22), we have when $q(D', i, y) \leq e^{\epsilon_0} [1 + \frac{\gamma \max p(s)}{t \min p(s)}] pt(D, i, y) (\sum_{s \in C_i(D, y)} p(s) P_s(y)) = e^{\epsilon_0} [1 + \frac{\gamma \max p(s)}{t \min p(s)}] q(D, i, y)$, which is (19).

When $i = k$, it is the same as $i = j$ as we just change the role of D, D' , by assumption we have $C_k(D', y) = C_k(D, y) - \{x_n\}$, so we have when $|C_k(D, y)| > t$,

$$\begin{aligned} q(D, k, y) &= pt(D, k, y) \sum_{s \in C_k(D, y)} p(s) P_s(y) \\ &= pt(D, k, y) \left(\sum_{s \in C_k(D', y)} p(s) P_s(y) + p(x_n) P_{x_n}(y) \right) \\ &\leq e^{\epsilon_0} pt(D', k, y) \left[1 + \frac{\gamma \max p(s)}{t \min p(s)} \right] \sum_{s \in C_k(D', y)} p(s) P_s(y) \\ &= e^{\epsilon_0} pt(D', k, y) \left[1 + \frac{\gamma \max p(s)}{t \min p(s)} \right] q(D', k, y) \quad (24) \end{aligned}$$

When $|C_k(D, y)| \leq t$,

$$\begin{aligned} q(D, k, y) &= pt(D, k, y) \sum_{s \in C_k(D, y)} p(s) P_s(y) \\ &\leq te^{-\epsilon_0(k-t)} \max p(s) \end{aligned} \quad (25)$$

Others are the same as $i = j$, we omit the proof. \square

The following is the key lemma similar with Lemma 5 in [Bindschaedler *et al.*, 2017].

Lemma 9. For parameters $k \geq 1, \gamma > 1$ and $\epsilon_0 > 0$. Take dataset D, D' as above. Then for ant integer $1 \leq t < k$ and synthetic record $y \in \mathcal{Z}$, we have: $Pr\{\mathcal{F}(D)\} \leq e^{\epsilon_0} [1 + \frac{r \max p(s)}{t \min p(s)}] Pr\{\mathcal{F}(D')\} + te^{-\epsilon_0(k-t)} \max p(s)$, and $Pr\{\mathcal{F}(D')\} \leq e^{\epsilon_0} [1 + \frac{r \max p(s)}{t \min p(s)}] Pr\{\mathcal{F}(D)\} + te^{-\epsilon_0(k-t)} \max p(s)$.

Proof. By definition we have: $Pr\{\mathcal{F}(D)\} = \sum_{i \neq j, k} q(D, i, y) + q(D, j, y) + q(D, k, y)$, and $Pr\{\mathcal{F}(D')\} = \sum_{i \neq j, k} q(D', i, y) + q(D', j, y) + q(D', k, y)$. Since $\sum_{i \neq j, k} q(D, i, y) = \sum_{i \neq j, k} q(D', i, y)$ and also by Lemma 8, we have $q(D, j, y) \leq q(D', i, y)$ and $q(D, k, y) \leq e^{\epsilon_0} [1 + \frac{r \max p(s)}{t \min p(s)}] q(D', i, y) + te^{-\epsilon_0(k-t)} \max p(s)$. Thus, we have the following

$$\begin{aligned} Pr\{\mathcal{F}(D)\} &\leq e^{\epsilon_0} [1 + \frac{r \max p(s)}{t \min p(s)}] Pr\{\mathcal{F}(D')\} \\ &\quad + te^{-\epsilon_0(k-t)} \max p(s). \end{aligned} \quad \square$$

The later proof of Theorem 2 is the same as in [Bindschaedler *et al.*, 2017], we omit them.

D Experimental Setup

In this section, we give the detailed description of the adopted dataset and performance measure.

Performance Measure. To evaluate the performance of the proposed method, we adopt the following two measure metrics.

- *ISE*: The integrated squared error (*ISE*) is a widely adopted measure metric which can measure the distance between the densities derived from the synthetic data and the objects' true densities. It is defined as: $ISE = \sum_{i=1}^{\mathcal{N}} \int_{-\infty}^{+\infty} (f_i^* - \tilde{f}_i)^2 dx$, where f_i^* and \tilde{f}_i are respectively the true density and density derived from the synthetic data for entity o_i . *ISE* tends to penalize more on the small distances.
- *SISE*: The squared integrated squared error (*SISE*) is another metric which can be adopted to measure the distance between the derived densities and the objects' true densities, and it is defined as: $SISE = \sum_{i=1}^{\mathcal{N}} (\int_{-\infty}^{+\infty} (f_i^* - \tilde{f}_i)^2 dx)^2$. Compared with *ISE*, *SISE* tends to penalize more on the large distance and less on the small distance.

Since the goal of the collector is to release the data whose pattern is similar to the true underlying pattern for the objects, the less the *ISE* or *SISE*, the better the method.

Datasets. In this experiment, we adopt the following three real-world datasets to evaluate the performance of the proposed mechanism.

- *Population Dataset* [Pasternack and Roth, 2010; Wan *et al.*, 2016]. It is about the population information of some cities at different years. We process this dataset with the method adopted in [Wan *et al.*, 2016]. We first remove the objects whose claims are all the same and keep only the latest claim for the same user and the same object. Then we remove the obviously-wrong objects whose claims are larger than 10^8 . This dataset contains 2,344 users and 1,124 objects.
- *Stock Dataset* [Li *et al.*, 2012]. This dataset is collected from 55 users by the authors of [Li *et al.*, 2012] on each weekday in July 2011. It consists of 1000 stock symbols and 16 properties. In this experiment, we only adopt the properties whose data type is continuous. Totally, there are 55 users and 5,521 objects in this dataset.
- *Indoor Floorplan Dataset* [Li *et al.*, 2014a]. It is collected when constructing the indoor floorplans, which is a representative example of social sensing applications. The objects are the hallway segments of a building, the task here is to measure the distances of these segments with the inertial sensors built in the smartphone. It consists of 247 users and 129 object.

E Experiments on Simulated Datasets

In this section, we evaluate the performance of **PrisCrowd** on the simulated datasets. We first introduce the data generation procedure, and then report the experimental results on these datasets.

Data Generation. We first generate 30 users and 50 objects. We assume that $30 * p$ users are marked as “unreliable” and the rest $30 * (1 - p)$ users are marked as “reliable”, where p is the percentage of the unreliable users. For each object, the claims provided by reliable users are generated from the Gaussian distribution $N(\mu_1 = 5, \sigma_1 = 0.5)$, and the claims provided by unreliable users are generated from the uniform distribution $U(\mu_2 = 0, \mu_3 = 10)$. In this way, the reliable users always provide high quality claims while the unreliable users may provide many extreme claims (i.e., outliers) for the objects.

Accuracy Comparison. In this experiment, we generate five simulated datasets via varying p from 0.1 to 0.5, and then evaluate the accuracy of the released claims generated based on our proposed method. Here we still assume that the data collector samples 30 claims for each object in the synthetic data generation procedure. The calculated *ISE* and *SISE* for these datasets are shown in Table 2. These results show that our proposed method performs much better than the baseline methods in all cases. The reason is that we take into consideration the fine grained weights of the users in the proposed method so that it can be more robust to the outlying claims

Table 2: Results on the Simulated Datasets

p	Methods	ISE	SISE
p=0.1	PrisCrowd	4.758	14.991
	Uniform	11.302	23.576
	Basic(Strong)	20.297	31.857
	Basic(Normal)	20.282	31.845
	Basic(Low)	20.268	31.834
p=0.2	PrisCrowd	4.842	15.327
	Uniform	12.393	24.668
	Basic(Strong)	20.679	32.155
	Basic(Normal)	20.669	32.155
	Basic(Low)	20.594	32.137
p=0.3	PrisCrowd	4.860	15.325
	Uniform	14.613	26.929
	Basic(Strong)	21.348	32.671
	Basic(Normal)	21.342	32.666
	Basic(Low)	21.298	32.633
p=0.4	PrisCrowd	4.866	15.380
	Uniform	16.703	28.742
	Basic(Strong)	21.650	32.901
	Basic(Normal)	21.615	32.850
	Basic(Low)	21.571	32.817
p=0.5	PrisCrowd	7.475	18.203
	Uniform	18.833	30.579
	Basic(Strong)	22.236	33.320
	Basic(Normal)	22.226	33.243
	Basic(Low)	22.205	33.236

than the baseline methods. Additionally, the results in this table also show that the increase of the value of p (i.e., the percentage of the unreliable users) degrades the performance of all the methods. This is mainly because more outlying claims are involved in the collected data when the percentage of the unreliable users increases.