

Improving Data Utility Through Game Theory in Personalized Differential Privacy

Lei Cui^{1,2,Δ}, *Student Member, IEEE*, Youyang Qu^{2,Δ}, *Student Member, IEEE*
Mohammad Reza Nosouhi³, *Student Member, IEEE*, Shui Yu³, *Senior Member, IEEE*
Jian-Wei Niu⁴, *Senior Member, IEEE*, and Gang Xie^{1,5,*}

¹College of Information and Computer, Taiyuan University of Technology, Taiyuan 030024, China

²School of Information Technology, Deakin University, Melbourne, VIC 3125, Australia

³School of Software, University of Technology Sydney, Sydney, NSW 2007, Australia

⁴School of Computer Science and Engineering, Beihang University, Beijing 100191, China

⁵Shanxi Key Laboratory of Advanced Control and Intelligent Information System, Taiyuan University of Science and Technology, Taiyuan 030024, China

E-mail: {cuil, quyo}@deakin.edu.au; {mohammad.r.nosouhi, shui.yu}@uts.edu.au; niujianwei@buaa.edu.cn
xiegang@tyut.edu.cn

Received June 28, 2018; revised January 27, 2019.

Abstract Due to dramatically increasing information published in social networks, privacy issues have given rise to public concerns. Although the presence of differential privacy provides privacy protection with theoretical foundations, the trade-off between privacy and data utility still demands further improvement. However, most existing studies do not consider the quantitative impact of the adversary when measuring data utility. In this paper, we firstly propose a personalized differential privacy method based on social distance. Then, we analyze the maximum data utility when users and adversaries are blind to the strategy sets of each other. We formalize all the payoff functions in the differential privacy sense, which is followed by the establishment of a static Bayesian game. The trade-off is calculated by deriving the Bayesian Nash equilibrium with a modified reinforcement learning algorithm. The proposed method achieves fast convergence by reducing the cardinality from n to 2. In addition, the in-place trade-off can maximize the user's data utility if the action sets of the user and the adversary are public while the strategy sets are unrevealed. Our extensive experiments on the real-world dataset prove the proposed model is effective and feasible.

Keywords personalized privacy protection, game theory, trade-off, reinforcement learning

1 Introduction

In this big data era, the proliferation of mobile devices results in massive information being published over various social networks. Users tend to publish or share their information, including sensitive data, over social networks^[1]. The utility and the attraction of the social network services origin from data sharing among users^[2]. For example, a follower can browse the tweets of the other users in Twitter to learn recent

developments^[3]. However, the improper collection and abuse of the published data are negatively influencing the service quality of social networks.

Privacy may leak in various perspectives, for instance, identity privacy, social tie privacy, location privacy, interest privacy, and medial content privacy^[4]. Therefore, privacy issues are critical while adversaries keep tracing users' sensitive information, especially in social networks^[5,6]. The adversaries continuously launch structural-based attacks, background knowledge

Regular Paper

Special Section of NSFC Joint Research Fund for Overseas Chinese Scholars and Scholars in Hong Kong and Macao 2014–2017

This work was supported by the Shanxi International Cooperation Project under Grant No. 201803D421039, the China Scholarship Council (CSC) under Grant No. 201708240007, and the China Scholarship Council (CSC) under Grant No. 201808240004.

^ΔLei Cui and Youyang Qu contributed to this work equally.

*Corresponding Author

©2019 Springer Science + Business Media, LLC & Science Press, China

attacks, and collusion attacks to steal sensitive information from users^[7]. The adversaries are never tired of devising new attacks or combining a variety of attacks, which put privacy protection under great threats.

Data utility is equivalently necessary as privacy protection in privacy-preserving models. To achieve better data utility, statistical properties are being discussed in different aspects^[8]. In accordance with differential privacy, global sensitivity^[9] is firstly proposed to maximize the data utility. Then, sampling Laplace noise^[10] and machine learning based method come into existence successively. Moreover, the employment of game theory based method^[11] captures the characters of the desired trade-off in a more appropriate way. However, data utility is still not satisfying because of insufficient optimization.

The trade-off between personalized privacy and improved data utility arises extensive concerns^[12]. Lack of protection results in privacy leakage while over-protection brings about data utility degradation^[13]. Normally, the trade-off is being demonstrated by simply splitting privacy level and data utility, which lacks measurement and optimization. This is not practical and feasible due to that trade-off can be optimized from different perspectives^[14]. With an optimized trade-off, the proposed model can improve its performance based on the pre-set constraints. Therefore, deriving an optimized trade-off is a necessity in a personalized privacy protection model.

Existing privacy protection models have two main branches including clustering-based methods and differential privacy. Clustering-based methods are first put into use, such as K -anonymity^[15], L -diversity^[16], T -closeness, and their extensions. Clustering-based methods are feasible under the scenario of data publishing in the sense of dataset. It requires enough data to satisfy the amount, diversity, and distribution requirements. But the data shared in social networks are relative sparse and therefore the clustering-based methods are not practical. Nowadays, differential privacy has been implemented in more and more scenarios, for example, dataset correlations^[9], location-based services^[17], and so on. Differential privacy^[18] and its extensions provide the privacy preservation with solid theoretical foundations^[19]. Most existing differential privacy-based mechanisms are used in a data publishing scenario as well. The employment in social networks still requires for modification^[20]. In addition, uniform privacy is another big issue. Nevertheless, uniform privacy protection cannot satisfy the ever increas-

ing demands^[21]. Built upon the new requirements, personalized differential privacy-preserving models are presented for various circumstances.

We have an observation that existing studies consider the privacy level to be uniform while they do not consider the improved data utility problem. Normally, the adversary is assumed to be static and both the adversary and the attack are measured qualitatively rather than quantitatively^[22]. Base on this assumption, the data utility is measured in a way where the impact of the adversary is not taken into consideration. However, this is not practical, especially in a personalized differential private model.

To address the trade-off problem, we propose a game theoretical model aiming to maximize the data utility which takes the adversary effect into account. A personalized differential privacy protection model is proposed based on the social distance in social networks, that is, we measure the distance between the users in social networks with the relationship rather than the locational distance. Although this idea is discussed under the scenario of social networks, it can be extended to multiple situations such as IoT and machine learning. Firstly, we observe that users and adversaries are not sure of each other's payoff. Thus, we employ a static Bayesian game to model the confront between users and adversaries. Secondly, we define the two strategy sets in differential privacy sense, which measures the effect of both sides. As an extension of our conference paper^[23] published on IEEE ICC2018, we further define payoff functions of each side and employ a modified reinforcement learning method to fast derive the Nash equilibrium, which accurately denotes the maximum data utility the users can achieve. The fast convergence is guaranteed by cardinality reduction from n to 2. Finally, the evaluation results verify the effectiveness and feasibility of the proposed model.

The contributions are summarized as follows.

- We propose a personalized differential privacy protection model based on social distance. By personalizing the privacy level, we solve the problem of protecting privacy uniformly. Therefore, the personalized differential privacy leads to less overall privacy budget and higher data utility.
- We establish a static Bayesian game to capture the real-world confront characters and take the effect of both the user and the adversary into consideration. We formalize all the parameters in differential privacy sense and thereby measure the adversary and the attack quantitatively. The proposed model eliminates the

uncertainty of data utility measurement. Furthermore, we obtain the improved data utility by deriving the Bayesian Nash equilibrium.

- We derive the Bayesian Nash equilibrium rapidly with a modified Q-learning algorithm. We reduce the cardinality of the data and thereby reduce the updating rules, which brings about fast convergence. We also derive the best static strategy for users to achieve maximum data utility based on static Bayesian game.

- We evaluate the proposed model with extensive experiments on the real-world dataset. The outcomes prove the effectiveness and feasibility of the proposed model.

The rest of this paper is organized as below. In Section 2, we present related work and the preliminaries. We then present the framework of the proposed model in Section 3. Section 4 depicts the statistic Bayesian game while Section 5 demonstrates the Nash equilibrium derivation with reinforcement learning. We illustrate the system analysis in Section 6, which is followed by the performance and evaluation in Section 7. Finally, this paper is concluded in Section 8.

2 Related Work

Despite the advantages social networks have brought to us, the problem of privacy leakage arises extensive concerns^[24], for example, identity privacy^[25], authenticated data redaction^[26], and so on. Clustering-based methods were proposed to deal with this, including K -anonymity^[15], sensitive attributes grouping^[27] and so on. These methods are practical with a small database size^[12]. However, the privacy concern is highly reinforced in big data era^[13]. Motivated by this, differential privacy^[18] arises with higher performance and solid theoretical foundations^[28]. Differential privacy functions well in statistic query and a lot of extensions have been presented to serve different scenarios^[19]. In [29], Du *et al.* proposed a cryptography method which sheds light on privacy protection. Multiple social network privacy has been discussed in recent years^[30].

Laplace mechanism provides random noisy responses which realize differential privacy in real-valued data sense^[28,31]. After privacy is protected, researchers pay more attention to the optimal trade-off between privacy protection and data utility^[14]. Wang and Zhang^[11] used game theory and machine learning to obtain the optimal trade-off.

Jorgensen *et al.*^[32] argued that not all users require

the same level of privacy and introduced a new privacy framework called personalized differential privacy, in which the privacy requirements are specified at the user level, rather than by a single privacy parameter. In the scenario of crowdsourcing data aggregation, Wang *et al.*^[33] proposed an aggregation scheme for histogram estimation, which enables participants to publish data at personalized differential-privacy levels. In [34], He *et al.* proposed a latent data privacy preserving model which can achieve personalized data utility in social networks. The authors tried to optimize the trade-off between latent data privacy protection and personalized data utility. Both prediction data utility loss and structure data utility loss are taken into consideration in this paper. Moreover, attribute sanitization and link sanitization are collaboratively solved rather than discussed separately. They proposed an attribute-link sanitization strategy which can provide satisfying the quality of service while protecting the sensitive latent information with personalized data utility. In addition, a powerful adversary with maximum inference attack ability is also employed to test the model which demonstrates its superior performance. Recently, Nie *et al.*^[35] personalized the traditional definition of local differential privacy and first proposed a utility optimization framework for histogram estimation with personalized multi-level privacy.

3 System Modeling

In this section, the system mode is presented in detail. We propose a personalized differential privacy model based on social distance, which is followed by the establishment of a static Bayesian game.

We use a graph $G = \{n_i, e_i, m_i | n \in N, e \in E, m \in M\}$ to represent a social network in the proposed model. In the graph G , $n_i \in N$ denotes the set of nodes and $e_{ij} \in E$ is the set of edges, which are user set and relationship set between users respectively. M indicates the set of sensitive information $m_i \in M$ that the users in N may publish. For a pair of users (i, j) , if there is at least one edge $\{e_{ik_1}, e_{k_1k_n}, e_{k_nj}\} \in E$ between them, we conclude that user i and user j have a relationship. Based on the relationship, we also have $d_{ij} \in D$ to denote the social distance between the two users.

To better clarify, we assume the graph G to be an undirected graph. However, this assumption can be removed. In addition, we also assume that there is a trusted central authority which processes the data with

ϵ -differential privacy and transmits the data with secure communication. In a trusted central authority, the privacy budget is a constant B . For all the published data, the sum of all privacy budget equals B .

3.1 Personalized Differential Privacy

Using ϵ -differential privacy, we implement the personalized privacy protection system based on social distance. In social networks, the users tend to share their messages to people with close relationships and do not care about the strangers. When a user i diffuses a sensitive message m over the social networks, it is intuitive that recipients with different social distances should obtain m with different privacy protection.

To personalize differential privacy, we use social distance $d_{ij} \in D$ to decide the privacy budget ϵ . If the social distance is relatively short, user i wants the other users to see the message with high-level privacy protection. On the contrary, user i does not care about people with relative long social distances and spends less privacy budget ϵ on them. Moreover, we use $\epsilon(\frac{1}{d_{ij}})$ -differential privacy to personalize the privacy level as below.

Definition 1 ($\epsilon(\frac{1}{d_{ij}})$ -Differential Privacy (DP)). Given $\epsilon(\frac{1}{d_{ij}})$ to be a positive privacy budget function decided by social distance d_{ij} , D_1 and D_2 are two datasets with an adjacent relationship, and \mathcal{M} to be a randomized algorithm that sanitizes the dataset, the algorithm \mathcal{M} is $\epsilon(\frac{1}{d_{ij}})$ -differential private on D_1 and D_2 if

$$\Pr[\mathcal{M}(D_1) \in \Omega] \leq \text{Exp}(\epsilon(\frac{1}{d_{ij}})) \times \Pr[\mathcal{M}(D_2) \in \Omega],$$

where the probability space Ω is taken over the randomness used by \mathcal{M} .

In the case of social distance, we measure the distance between the users in social networks rather than the locational distance. Normally, the social distance is calculated by hop counts based on the graph theory. Therefore, we take hop as the social distance as an example. This model can be further extended to other social distances, such as the shortest social distance, effective social distance, and so on. We use a logarithm function to describe the mapping function between ϵ and d_{ij} where $\epsilon = -\ln \frac{1}{d_{ij}}$.

In this personalized privacy protection model, users with long social distance are allocated with more privacy budget $\epsilon(\frac{1}{d_{ij}})$ while users with short social distance are allocated with less one. $\epsilon(\frac{1}{d_{ij}})$ denotes the

mapping function that decides the privacy budget allocation based on social distance. Moreover, the sum of all privacy levels' budgets equals the fixed budget B .

$$B = \sum \epsilon(\frac{1}{d_{ij}}) \{(i, j) | 1 \leq i, j \leq n, i \neq j\}.$$

3.2 Adversary Model

We consider a static adversary \mathcal{A} who actions after considering the possible payoff functions of the user. Usually, the adversary \mathcal{A} has a prior belief of the sensitive information. The data type of sensitive information in our study is shown as Table 1. The prior belief can be formalized as ϵ_p -differential privacy. If $\epsilon_p \leq \epsilon(\frac{1}{d_{ij}})$, the adversary already has enough background knowledge of this sensitive information and he/she can successfully launch an inference attack. If $\epsilon_p > \epsilon(\frac{1}{d_{ij}})$, the privacy protection is strong enough and the adversary has to get more information to launch attacks, for example, collusion. Given the adversary's prior belief ϵ_p and the obtained information ϵ_b , we formalize the adversary's behavior as

$$DP(\epsilon_p) + DP(\epsilon_b) = DP(\epsilon_p + \epsilon_b),$$

where $DP()$ is short for differential privacy. This equation holds because of the composition character of differential privacy scheme.

$$\begin{cases} DP(\epsilon_p + \epsilon_b) > DP(B), \\ DP(\epsilon_p + \epsilon_b) > DP(\sum_{i,j}^{n,n} \epsilon(\frac{1}{d_{ij}})), \\ DP(\epsilon_p + \epsilon_b) \leq DP(\sum_{i,j}^{n,n} \epsilon(\frac{1}{d_{ij}})). \end{cases} \quad (1)$$

Table 1. Sensitive Data Type

Data Type	Explanation
Binary state	$s_i \in \{0, 1\}$ indicates a binary status like on-line or not
Location	$(longitude_i, latitude_i, height_i) \in R_3$ is the GPS coordinates of a certain entity
Time stamp	$t_i \in R_+$ denotes positive real number such as time or date
Text data	$d_i \in ABC$ represents alphabetic instance like light-weighted encryption cipher
Media data	Media data includes videos, images, sound, and so forth

Based on (1), we have the following observations. Firstly, the sum of ϵ_p and ϵ_b is smaller than $\epsilon(\frac{1}{d_{ij}})$, and the adversary breaches the privacy of a certain level $\frac{1}{d_{ij}}$. Secondly, the sum of ϵ_p and ϵ_b is smaller than B ,

the adversary breaches the privacy of the whole system. Thirdly, the sum of ϵ_p and ϵ_b is larger than $\epsilon(\frac{1}{d_{i1}})$, the adversary fails the attack.

3.3 Static Bayesian Game

In the established personalized privacy protection model, the confrontation between the user i and an adversary captures the features of the game theory. Before successfully breaching the privacy, the adversary is not sure about how much more information he/she can obtain, and the user has no idea about how much background the adversary holds. Therefore, the confrontation between the user and the adversary can be formalized as a static Bayesian game.

This static Bayesian game includes action space A , type space T , and their inference P , and payoff function U . The type of user i is regarded as the personal information, which decides user i 's payoff function $u_i(a_1, \dots, a_n; t_i)$. Meanwhile, i is an element in the possible type set T_i . The user i 's inference $p_i = (t_{-i}|t_i)$ describes the uncertainty of the other $(n-1)$ users' possible type t_{-i} when t_i is the type of user i . Given all the above conditions, we define the static Bayesian game as (2).

$$G = \{A_1, \dots, A_n; T_1, \dots, T_n; P_1, \dots, P_n; U_1, \dots, U_n\}. \quad (2)$$

In Algorithm 1, we can express an incomplete information game as an imperfect information game by the virtue of step 8, which assigns t_i "naturally". The involved players are not aware of the previous game process.

Algorithm 1. Game-Based Optimal Data Utility Derivation Algorithm

Input: static Bayesian game G ;

Output: optimal data utility U_o ;

- 1: Implement $\epsilon(\frac{1}{d_{ij}})$ -differential privacy;
 - 2: Formulate the adversary behavior as $(\epsilon_p + \epsilon_b)$ -differential privacy;
 - 3: Initialize the action space A ;
 - 4: Initialize the type space T ;
 - 5: Define the inference space P based on (A, T) ;
 - 6: Define the payoff function set U based on (A, T, P) ;
 - 7: **while** until convergence **do**
 - 8: Players are assigned with t_i but aware of t_{-i} naturally;
 - 9: Players choose their actions $\exists a_i \in A$ at the same time;
 - 10: Payoff generation $u_i(a_1, \dots, a_n; t_i)$;
 - 11: Update action a_{i+1} and t_{i+1} based on p_i ;
 - 12: Update payoff u_{i+1} ;
 - 13: **end while**
 - 14: Nash equilibrium NE obtained;
 - 15: Derive the optimal data utility U_o based on NE
-

We further discuss how to calculate the reference $p(t_{-i}|t_i)$. We first assume the transcendental probability distribution $p(t)$ and the type vector $t = (t_1, \dots, t_n)$ are common knowledge. After the player i is assigned with t_i , he/she can calculate the reference $p(t_{-i}|t_i)$ according to all the other players' conditional probability with Bayes rules.

$$p(t_{-i}|t_i) = \frac{p(t_{-i}, t_i)}{p(t_i)} = \frac{p(t_{-i}, t_i)}{\sum_{t_{-i} \in T_{-i}} p(t_{-i}, t_i)}.$$

In addition, the other players can obtain i 's various inferences based on t_i , that is, $p(t_{-i}|t_i)$ can be derived for $\forall t_i \in T_i$. Normally, the types of players are assumed to be independent, which means $p(t_{-i})$ does not depend on t_i . But $p(t_{-i})$ is still derived from the transcendental probability distribution $p(t)$. In this situation, the other players realize i 's inferences on their types.

In this algorithm, as we have iteration process to update and optimize the final output, the total iteration times would be n for each cardinality. Therefore, the computation complexity of each cardinality would be $O(n)$. Since the cardinality of the proposed model is reduced from n to 2 as analyzed above, we conclude the final computation complexity is $O(2n)$, which enjoys a significant reduction compared with the traditional models of $O(n^2)$.

4 Bayesian Nash Equilibrium

Built upon the proposed model, we present the Bayesian Nash equilibrium in this part. We first derive the Bayesian Nash equilibrium of the static Bayesian game. The Bayesian Nash equilibrium denotes the optimal data utility of the personalized privacy protection system. We then analyze the confrontation between two players, including one user and one adversary. In addition, the results can be further extended to multiple adversaries. On the one hand, if the multiple adversaries do not collude, then the course of offense-defense can be formulated into multiple independent static Bayesian games. On the other hand, if the multiple adversaries collude with one another, we can formulate them into one collusion adversary according to the composition theorem of differential privacy.

4.1 Bayesian Nash Equilibrium

In order to acquire the Bayesian Nash equilibrium, we must define the strategy space of the players. The strategy is a package plan of actions, including every

possible action corresponding to all the possible situations. Given the time sequence of static Bayesian game and the pre-assigned type $t_i \in T_i$, a strategy of player i has to contain a feasible action of every possible type $t_i \in T_i$.

Therefore, in a static Bayesian game G , a strategy of player i is a function $s_i(t_i)$. With regard to $\forall t_i \in T_i$, $s_i(t_i)$ contains the possible actions $\exists a_i \in A_i$ when t_i is the type of user i .

Unlike complete information game, no strategy space is defined in (2). As a substitute, we can construct the strategy space from the action space A_i and the type space T_i . Player i 's feasible strategy set S_i is the function set whose definition domain is T_i and whose value domain is A_i . For example, every $t_i \in T_i$ chooses a different action $a_i \in A_i$ in a separating strategy while all $t_i \in T_i$ choose the same $a_i \in A_i$ in a pooling strategy.

Based on the analysis, we conclude every player's strategy must be other players' optimal responses of strategy, that is, the following definition of Bayesian Nash equilibrium is the Nash equilibrium of static Bayesian game.

Definition 2 (Bayesian Nash Equilibrium). *Given n players and a static Bayesian game $G = \{A_1, \dots, A_n; T_1, \dots, T_n; P_1, \dots, P_n; U_1, \dots, U_n\}$, a strategy set is a pure strategy Bayesian Nash equilibrium if for $\forall i$ and the type set $\forall t_i \in T_i$, $s_i^*(t_i)$ satisfies*

$$\max_{a_i \in A_i} \sum_{t_{-i} \in T_{-i}} u_i(s_i^*(t_1), \dots, s_{i-1}^*(t_{i-1}), a_i, s_{i+1}^*(t_{i+1}), \dots, s_n^*(t_n); t) \times p_i(t_{-i}|t_i).$$

Namely, no player wants to change his/her strategy even if this change only involves a single action of a specific type.

4.2 Two-Player Game Analysis

In this subsection, we consider the static Bayesian game between one adversary and user i . This can be further extended to multiple adversaries scenario. Firstly, we analyze the action-type based strategy space $S_i(a_i, t_i)$. Secondly, the payoff function U_i is presented based on strategy space. Thirdly, we derive the Bayesian Nash equilibrium which could be seen as the optimal trade-off with maximum data utility.

The player i has two actions, which is actually performed by the trusted central authority. The first action a_i^1 is to publish the data without privacy protection. Without privacy protection, we use $\epsilon \rightarrow \infty$ to

denote the action in differential privacy sense. The second action a_i^2 is to publish the data with certain privacy protection, which is described as $\epsilon(\frac{1}{d_{ij}})$ based on social distance.

The type of specification action is the personal information that the other player does not know. Therefore, a_i^1 's type t_i^1 depends on whether player i adopts any protection measures. In the case of a_i^2 , its type t_i^2 can be regarded as how the mapping function $\epsilon(\frac{1}{d_{ij}})$ functions. Namely, the adversary does not know how the function maps social distance d_{ij} into privacy budget ϵ .

Talking about the adversary j , there are also two actions. The first action a_j^1 means that the adversary launches the attack with background knowledge only, which can be formalized as ϵ_p . The second action a_j^2 is that the adversary launches the attack with background knowledge and other supplemental information. We use $\epsilon_p + \epsilon_b$ to denote a_j^2 .

In the case of types corresponding to the actions, we also consider the uncertainty. Thus, a_j^1 's type t_j^1 is that user i has no idea of how much background knowledge the adversary has, while a_j^2 's type t_j^2 is that user i has no idea of how much more information the adversary can obtain.

As we formulate all the strategies in the differential privacy sense, the payoff function U_i can be easily obtained by minus operation. Moreover, the payoffs of the user i and the adversary j are reciprocal numbers. The detailed U_i is shown as Table 2.

Table 2. Two-Player Static Bayesian Game

A	$P(i)$	
	$\epsilon \rightarrow \infty$	$\epsilon(\frac{1}{d_{ij}})$
ϵ_p	$\pm(\epsilon_p - \epsilon)$	$\pm(\epsilon_p - \epsilon(\frac{1}{d_{ij}}))$
$\epsilon_p + \epsilon_b$	$\pm(\epsilon_p + \epsilon_b - \epsilon)$	$\pm(\epsilon_p + \epsilon_b - \epsilon(\frac{1}{d_{ij}}))$

5 Bayesian Nash Equilibrium Derivation with Reinforcement Learning

In this section, we derive the Bayesian Nash equilibrium built upon Markov decision process with reinforcement learning. We first model the actions of users and adversaries and state transmission. Based on these, we formulate users payoff function and thereby the fast convergence learning algorithm. At last, we present the best strategy of users with optimized data utility.

5.1 Actions of Users and Adversaries

In social networks, users tend to publish data (including sensitive information) on their homepage and

share the data with the other users. Therefore, we formulate the action of users AC_u^i as the granularity of the published data.

We have an assumption that adversaries have limited computing resources, which is feasible and practical in real-world application. Therefore, the adversaries can only access and process n pieces of data published by a certain user.

$$\sum_{i \in n} AC_u^i \leq \Upsilon, \quad 0 \leq AC_u^i \leq 1. \quad (3)$$

In (3), Υ denotes the maximum computing power of adversaries while the granularity of AC_u^i is unified into the range $[0, 1]$. When Υ is larger than n , the computing power of the adversaries can be regarded as unlimited and they can post-process all the data published by users. In opposite, if Υ is smaller than n , we know the adversaries have limited computing power.

From the perspective of adversaries, we formulate the action AC_{ad} as the probability of an adversary re-identifying a specific user. It is intuitive that $0 \leq AC_{ad} \leq 1$. We can learn the probability by observing the attacking results from adversaries, for example, spam emails, personalized advertisements, and so on.

5.2 State Transmission

To model a Markov decision process, we need to investigate the state transmission of the attack and defense confrontation.

From the perspective of the system state, there are two components including the current piece of data and the last attack response. As mentioned above, users can observe the attack results by receiving spam emails with personalized advertisements. This is a good way for a user to measure the degree of sensitive information disclosure. Therefore, the user may change his/her strategy and publish the next piece of data with another granularity. For better data utility, the granularity will be more fine-grained while coarse-grained data can result in higher privacy protection level.

In terms of attack response AR , we make $AR_i = 1$ if the adversary successfully obtains the sensitive data M_{i-1} while $AR_i = 0$ if the adversary fails to get any useful information from M_{i-1} . Therefore, the system status can be formulated as $S = \{D_i, AR_{i-1}\}$.

To model the state transmission of the confrontation between users and adversaries, we define the system states as $S = \{S_u^i, S_{ad}^i\}$. The state is determined by sensitive data and attack response while they are

further decided by actions of users and adversaries, respectively. Therefore, we define the system state transmission as

$$\begin{aligned} & \Pr[(S_u^i, S_{ad}^i) | (S_1^{i-1}, S_2^{i-1}), A_u^{i-1}, A_{ad}^{i-1}] \\ &= \Pr[D_i | D_{i-1}] \Pr[AR^i | AR^{i-1}, A_u^{i-1}, A_{ad}^{i-1}] \\ &= \Pr[D_i | D_{i-1}] \Pr[AR^i | A_u^{i-1}, A_{ad}^{i-1}]. \end{aligned}$$

We derive the second equation because the attack response AR_i at time slot i only depends on the actions of the user and adversaries at time slot $i - 1$. That is the reason why AR_I is not impacted by AR_{i-1} .

5.3 Payoff Function and Nash Equilibrium

We are aiming to maximize data utility in this article. Therefore, we establish a service quality-based data utility measurement as an index. The measurement is as below.

$$R((S_u^i, S_{ad}^i), A_u^i, A_{ad}^i) = Q(A_u^i) - \omega \times L(S_u^i, S_{ad}^i),$$

where $Q(A_u^i)$ is the quality of service and $L(S_u^i, S_{ad}^i)$ denotes the privacy loss.

We use χ to denote the strategy. Thus, we use $\chi^u : S_u \mapsto \delta(A_u)$ to denote the strategy of the user, and $\chi^{ad} : S_{ad} \mapsto \delta(A_{ad})$ to denote the strategy of adversaries, where (S_u, S_{ad}) denotes the state space, and $\delta(A_u)$ and $\delta(A_{ad})$ the probability distribution over A_u and A_{ad} , which represents the action spaces of users and adversaries.

Given the initial time slot $i = 0$ and the initial state $s_0 \in S$, we re-formulate the payoff function as

$$\begin{aligned} R^X(s) &= \sum_{i=0} E[R((S_u^i, S_{ad}^i), A_u^i, A_{ad}^i) | \chi_u, \tau_{ad}, s_0 = s] \\ &= R(s, A_u^i, A_{ad}^i) + \sum_{\hat{s}} \Pr[\hat{s} | s, A_s^i, A_{ad}^i] R^T(\hat{s}). \end{aligned}$$

As both of the adversary and the user want to follow the best strategy of their own, there is a confrontation between them. To model the confrontation, we first define the best strategy as χ_u^* and χ_{ad}^* , and the best strategy pair is $\chi^* = (\chi_u^*, \chi_{ad}^*)$. Given a multi-slot game and system state $s \in S$, the Nash equilibrium of χ^* is

$$\begin{cases} R^{\chi^*}(s) \geq R^{\chi_{ad}^*}(s), \\ R^{\chi^*}(s) \leq R^{\chi_u^*}(s), \end{cases}$$

where $\chi_{ad} = \{\chi_u, \chi_{ad}^*\}$, and $\chi_u = \{\chi_u^*, \chi_{ad}\}$, for all χ_{ad} and χ_u .

5.4 Fast Convergence Learning Algorithm

For the sake of fast derivation of Nash equilibrium, we employ a modified Q-learning method to achieve fast convergence. We define the equivalent $\hat{R}_u^{\chi^*}(AR)$ as the expected value $R_u^{\chi^*}(s)$ where $s = \{AR, D\}$, through which we can get rid of D . We express $\hat{R}_u^{\chi^*}(AR)$ as

$$\begin{aligned} & \hat{R}_u^{\chi^*}(AR) \\ &= E[R_u(s, A^{\chi^*}) + \sum_{AR'} (\Pr[AR'|A^{\chi^*}] \hat{P}U^{\chi^*}(AR'))], \end{aligned}$$

where $A^{\chi^*} = \{A_u^{\chi^*}, A_{ad}^{\chi^*}\}$ is the best action following the best strategy χ^* .

Built upon this, we can derive χ^* with an equivalent problem as below. The equivalent problem can reduce the cardinality from n to 2 and thereby accelerate the derivation process.

$$\begin{aligned} \chi^* = \max_{\chi_u} \min_{\chi_{ad}} & E[R_u(s, A^{\chi^*}) + \\ & \sum_{AR'} (\Pr[AR'|A^{\chi^*}] \hat{P}U^{\chi^*}(AR'))]. \end{aligned} \quad (4)$$

According to (4), we further leverage \hat{R}^{χ^*} to accomplish the following updating rule.

$$\begin{aligned} & \hat{R}^{t+1}(AR) \\ &= (1 - \alpha_{i+1}) \hat{R}^i(AR) + \alpha_{i+1} E[R(s, A_u^i, A_{ad}^i) + \\ & \quad \hat{R}^i(AR')], \end{aligned}$$

where $\alpha_i \in [0, 1]$ denotes the learning rate of the algorithm. We set α_i to decrease with time in order to get a deterministic convergence, which is $\alpha_i = 1/i$. In this update step, $\hat{R}^{i+1}(AR)$ is regarded as the approximate value of $\hat{R}^{\chi^*}(AR)$ and it finally converges to $\hat{R}^{\chi^*}(AR)$ after finite rounds of updates.

5.5 Best Strategy Generation with Optimized Utility

In our model, we need to derive the Nash equilibrium of the multi-slot confrontation to obtain $\hat{R}^i(AR)$. Built upon the above equations, we can reshape the Nash equilibrium modelling as

$$\begin{aligned} & \min_{\chi_{ad}} \max_{\chi_u} \{ \frac{2\text{Exp}(-\rho(A_u^i - \sigma))}{1 + \text{Exp}(-\rho(A_u^i - \sigma))} - D(D)A_{ad}^iA_u^i - 1 \}, \\ & \text{s.t.} \\ & \sum_i A_{ad}^i \leq \Psi, \\ & 0 \leq A_{ad}^i \leq 1, \forall i, \\ & 0 \leq A_u^i \leq 1, \forall i, \end{aligned} \quad (5)$$

where $D(d)$ is the function of the piece of data d that $D(d) = \omega D\text{Sen}(d) + (\hat{R}^{\chi^*}(AR' = 0) - \hat{P}U^{\chi^*}(AR' = 1))$. As $\chi^*(AR' = 0)$ and $\hat{R}^{\chi^*}(AR' = 1)$ maintain the same, $D(d)$ simply depends on d . In above analysis, we have the observation that $\hat{R}^{\chi^*}(AR' = 0) > \hat{R}^{\chi^*}(AR' = 1)$. As $D\text{Sen} \geq 0$, we conclude that the function of message $D(d) > 0$.

To solve (5), we firstly eliminate the effects of the adversary. As we focus on stationary strategy in this paper, χ_u is fixed to a constant value. As we have proved $D(d) > 0$, we assume the adversary to eavesdrop Ψ messages to launch the attack to minimize the value in (5). Therefore, we re-formulate the problem as

$$\begin{aligned} & \max_{\chi_u, \Theta, T'} \{ \frac{2\text{Exp}(-\rho(a_u^i - \sigma))}{1 + \text{Exp}(-\rho(A_u^i - \sigma))} - D(d)A_u^i - 1 \}, \\ & \text{s.t.} \\ & 0 \leq A_u^i \leq 1, \forall i, \\ & A_{ad}^i \leq \Theta, \forall i \in T', \\ & A_u^i \geq \Theta, \forall i \in \{T/T'\}, \end{aligned}$$

where T' is one subset of T consisting of Θ messages. Given a certain T' , we can easily derive the closed form of the best strategy χ_u .

6 System Analysis

In this section, we demonstrate the system analysis in terms of collusion attack-proof and optimized data utility.

6.1 Optimized Data Utility Analysis

Among all the differentially private mechanisms, the most widely used one is the Laplace mechanism. Laplace mechanism adds controllable noise to the outputs, where the noise generation process complies with differential privacy.

Given the mechanism $\mathcal{M} : R^n \rightarrow \Delta(R^n)$ which adds Laplace distributed noise \mathcal{N} as

$$\begin{aligned} & \mathcal{M}(\mathcal{D}) = \mathcal{D} + \mathcal{N}, \\ & \text{s.t.} \\ & N \sim \text{Lap}(\frac{\delta}{\epsilon}), \\ & \text{Lap}(b) \sim d\Pr[N = n] = \text{Exp}(-\frac{\|n\|_2}{b}), \end{aligned} \quad (6)$$

where $d\Pr[N = n]$ is the density of $\text{Lap}(b)$. Formulated by this, we regard \mathcal{M} as a ϵ -differential private mechanism under adjacency relation.

However, the Laplace mechanism cannot be optimal in terms of minimum mean-squared error. Therefore, we target on achieving optimum Laplace mechanism for both minimum entropy and minimum mean-squared error through designing the noise properly.

Theorem 1 (Optimum Laplace Mechanism). *Given the ϵ -differential private mechanism $A : R^n \rightarrow \Delta(R^n)$, A satisfies $y_{ij}^K = d_{ij} + N$, where $N \sim \rho(N) \in \Delta(R^n)$. The mean-squared error is minimized when the noise density f complies with*

$$f_1^n(v) = \left(\frac{\epsilon}{2}\right) \text{Exp}(-\epsilon \|v\|_1),$$

where $f_1^n(v)$ denotes the density of noise in the sense of v . Thus, we have

$$\mathbb{E} \|y_{ij}^t - d_{ij}\|_2^2 = \mathbb{E}_{V \sim \rho} \|V\|^2 \geq \mathbb{E}_{V \sim f_1^n} \|V\|_2^2 = \frac{2n}{\epsilon^2}.$$

The optimum Laplace mechanism provides the solution to achieve optimized data utility when the privacy level is fixed. We further prove the proposed method can satisfy optimum Laplace mechanism, which makes the proposed model more feasible and practical.

6.2 Collusion Attack-Proof Analysis

We present an example of two-fold customized privacy levels for clarity. Given two privacy levels ϵ_i and ϵ_{i+1} , where $\epsilon_{i+1} > \epsilon_i$, there is a mechanism $\mathcal{M}_{\epsilon_i \rightarrow \epsilon_{i+1}} : \mathcal{D} \rightarrow \Delta(\mathcal{Y}^2)$ which releases the data in two different social networks. At first, u_i publishes a noisy outcome y_{ij}^1 to u_j . y_{ij}^1 satisfies ϵ_i -DP. Afterwards, the privacy level is relaxed to ϵ_{i+1} -DP. In case that users collude with others to obtain more precise output, the proposed mechanism should at least satisfy

$$\mathcal{M}_{DP}(\epsilon_i + \epsilon'_{i+1}) = \mathcal{M}_{DP}(\epsilon_{i+1}),$$

where ϵ'_{i+1} is the privacy level of the second noisy response and \mathcal{M}_{DP} denotes the differentially private mechanism. As the upper bound of composition mechanism indicates, we have

$$\mathcal{M}_{DP}(\epsilon'_{i+1}) = \mathcal{M}_{DP}(\epsilon_{i+1} - \epsilon_i),$$

where we can conclude $\epsilon'_{i+1} < \epsilon_{i+1}$. This conclusion implies the second noisy response cannot relax the privacy level to a satisfying degree. Especially in the case that $\epsilon(1) < \epsilon_{i+1} \ll 1$, the privacy level may even be upgraded. The data utility degrades so that it is not suitable for practical applications.

Theorem 2 (Collusion Attack-Proof Mechanism). *Two privacy levels ϵ_1, ϵ_2 , which are short for*

$\epsilon_1(\frac{1}{d_{ij}}, t), \epsilon_2(\frac{1}{d_{ij}}, t)$ and $0 < \epsilon_1(\frac{1}{d_{ij}}, t) < \epsilon_2(\frac{1}{d_{ij}}, t)$, are given. Then, the form of the mechanism can be formulated as

$$y_{i1}^t = d + V_1, \quad y_{i2}^t = d + V_2, \quad (V_1, V_2) \sim \rho \Delta(R^2).$$

Moreover, the density $f_{\epsilon_1(\frac{1}{d_{ij}}, t), \epsilon_2(\frac{1}{d_{ij}}, t)}$ is

$$f_{\epsilon_1, \epsilon_2}(x, y) = \frac{\epsilon_1^2}{2\epsilon_2} \text{Exp}(-\epsilon_2 |y|) \delta(x - y) + \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \text{Exp}(-\epsilon_1 |x - y| - \epsilon_2 |y|).$$

The rationale behind noises complying with Markov stochastic process is that the Markov process requires the current state is only related to the last state. That means the current state is not impacted by the other states before the last state. In this case, the current noise is only decided by the last noise. In the proposed model, the privacy level increases with the trust distance and the noise increases as well. Therefore, the current user has no incentive to collude with the next user who has an inaccurate output with more noises. The proof is presented in Appendix.

7 Performance Evaluation

In addition to theoretical analysis, we testify the proposed model with a real-world dataset. The evaluation results are satisfying and confirm the significance of this work.

We evaluate our model on the “Google+” dataset collected by McAuley and Leskovec^[36]. This dataset contains 107 614 nodes and 13 673 453 edges. We also use the shortest distance to be the distance parameter. This can be extended to other distance metrics.

We randomly capture a piece of the dataset with 4 000 nodes and 56 352 edges. Speaking of privacy budget B , current studies leveraged a method to set the upper boundaries of ϵ since the lower boundary should be 0. The method is formulated as $B \leq \frac{\Delta q}{\Delta v} \ln \left\{ \frac{(n-1)p}{1-p} \right\}$, where $B = \sum_i^n$ is the privacy budget, Δq is global sensitivity, n is the size of dataset, and p is the probability of being successfully attacked. In this paper, we use $n = 5 000$ nodes, $p = 1\%$, $\Delta q / \Delta v = 1$. Therefore, the value range of B would be $[0, 4]$. Therefore, we use the largest ϵ as the parameter to test the protection level of the proposed model. Firstly, we illustrate the privacy-level advantages of personalized differential privacy. Secondly, we testify the improved trade-off with maximum data utility based on various parameters.

7.1 Privacy Level

We consider an 8-level privacy protection model. For a fixed privacy budget $B = 4$, the classic differential privacy divides the budget equally $\{\forall d_{ij} \in D | \epsilon(\frac{1}{d_{ij}}) = 0.5\}$ while personalized differential privacy customizes the privacy levels according to the distance. In the personalized privacy model, only social distance $d_{ij} \leq 6$ is protected. We use the simplest logarithm function to denote this mapping function.

Then, given another 8-level privacy protection model with changing privacy budget B , we still assume social distance $d_{ij} \leq 6$ is protected in personalized privacy model.

In Fig.1, we compare the privacy protection differentiation between classic differential privacy and personalized privacy. The classic differential privacy is also known as uniform differential privacy. It equally divides the privacy budget into n parts and assigns them to each piece of data. However, in the proposed personalized privacy, the privacy level is personalized based on social distance. Social distance is the distance between the users in social networks with the relationship rather than the locational distance. The details are as follows.

From Fig.1(a), the first thing we can tell is that personalized differential privacy is more flexible than the classic differential privacy. Then second one is that personalized differential privacy can protect privacy in a more accurate way over social networks. The privacy level keeps decreasing with the increment of social distance. When a certain social distance is reached, the

users beyond the social distance are not protected anymore, which is feasible and practical in social networks.

Fig.1(b) illustrates how privacy levels increase with the increase of the fixed privacy budget B . Fig.1(b) shows that privacy levels increase almost lineally which is easy to handle. In the proposed personalized privacy protection model, privacy levels are derived by the privacy requirements. The flexibility confirms that there is no over-protection or data utility degradation.

7.2 Data Utility

We further consider an 8-level privacy protection model. For a fixed privacy budget $B = 4$, the classic differential privacy divides the budget equally $\{\forall d_{ij} \in D | \epsilon(\frac{1}{d_{ij}}) = 0.5\}$ while personalized differential privacy customizes the privacy levels according to the distance. In the personalized privacy model, only social distance $d_{ij} \leq 6$ is protected. We use the simplest logarithm function to denote this mapping function and the data utility can be measured with RMSE.

In Fig.2, we compare the data utility tendencies in terms of all three methods, including differential privacy (Classic DP), normal personalized differential privacy (Personalized DP), and the proposed model (OTO). From the overall trends, we can tell that the proposed model has a superior performance from the perspective of data utility improvement.

In Fig.2(a), the data utility is shown for the proposed model (OTO), common personalized differential privacy, and classic differential privacy. The data utility remains a constant for classic differential privacy, which

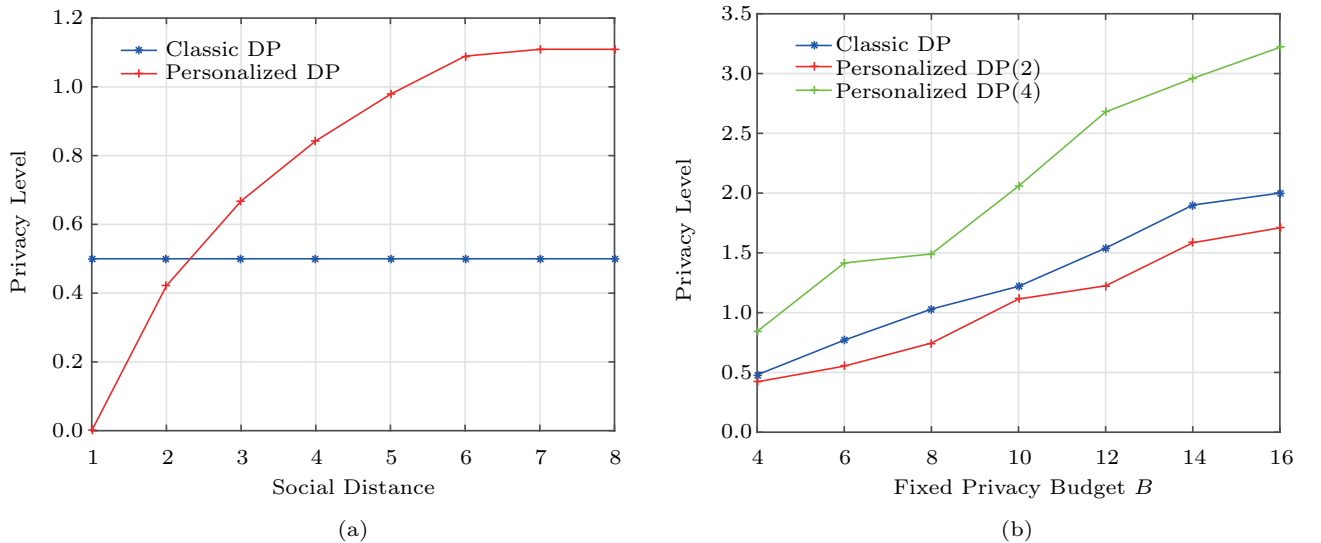


Fig.1. Privacy level comparison in single and multiple social networks. DP(x) means the social distance of DP is x .

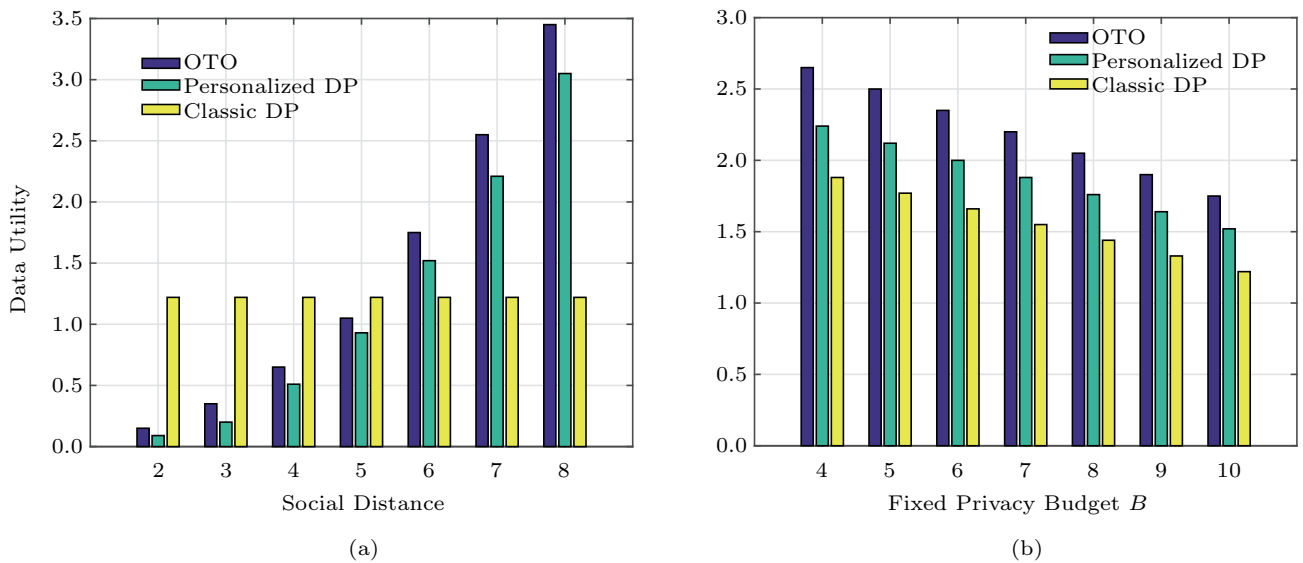


Fig.2. Data utility comparison in single and multiple social networks.

is decided by the fixed privacy-level feature. In the case of the proposed model and common personalized differential privacy, the utility increases with the increment of social distance. But in the proposed model, the starting point and the increasing rate are higher than those of common personalized differential privacy respectively, which proves that the proposed model is more feasible.

Fig.2(b) illustrates the trends of how data utility decreases with the increase of the fixed privacy budget. In all three models, data utility decreases linearly, which is decided by the trade-off feature. However, the proposed model has the maximum data utility from the beginning to the end.

All in all, the data utility of the proposed model is improved compared with classic differential privacy and personalized differential privacy without the game model, which complies with the above analysis.

7.3 Efficiency Evaluation

We introduce a modified Q-learning algorithm to achieve fast convergence and thereby obtain the Nash equilibrium. The reason why we can limit iteration times is that we reduce the cardinality from n to 2. Built upon cardinality deduction, we obtain faster convergence as shown in Fig.3.

As shown in Fig.3(a) and Fig.3(b), we demonstrate that the efficiency of the proposed method outperforms the classic one. The proposed method converges between 10^2 and 10^3 iteration times while the classic method converges between 10^5 and 10^6 iteration times.

Therefore, the magnitude order of the proposed idea is roughly two times less than that of the classic one, which means the proposed idea cost around 1% time the classic one.

Due to the limited iteration times, the proposed idea can derive the Nash equilibrium in a very short time. Furthermore, it can ensure that the system can implement the differentially private algorithm in real time.

8 Conclusions

In this work, we first established a personalized differential privacy model based on social distance. The privacy budgets ϵ increase with the increment of social distance. Building upon this privacy-preserving model, we further modeled the confrontation between the user and the adversary as a static Bayesian game because both sides are not aware of each other's strategy sets. We then formalized the action sets, typesets, strategy sets, and payoff functions in the differential privacy sense, which takes the effect of the adversary into consideration. In addition, the Bayesian Nash equilibrium is derived based on the payoff functions. Finally, we obtained the optimal trade-off with maximum data utility with the theoretical analysis. Extensive experiments were implemented to compare existing work and the proposed model. Evaluation results demonstrated the superior performance of our model.

In the case of future work, we plan to extend this model to multiple adversaries and consider the situation of collusion attack. Besides, we also consider tak-

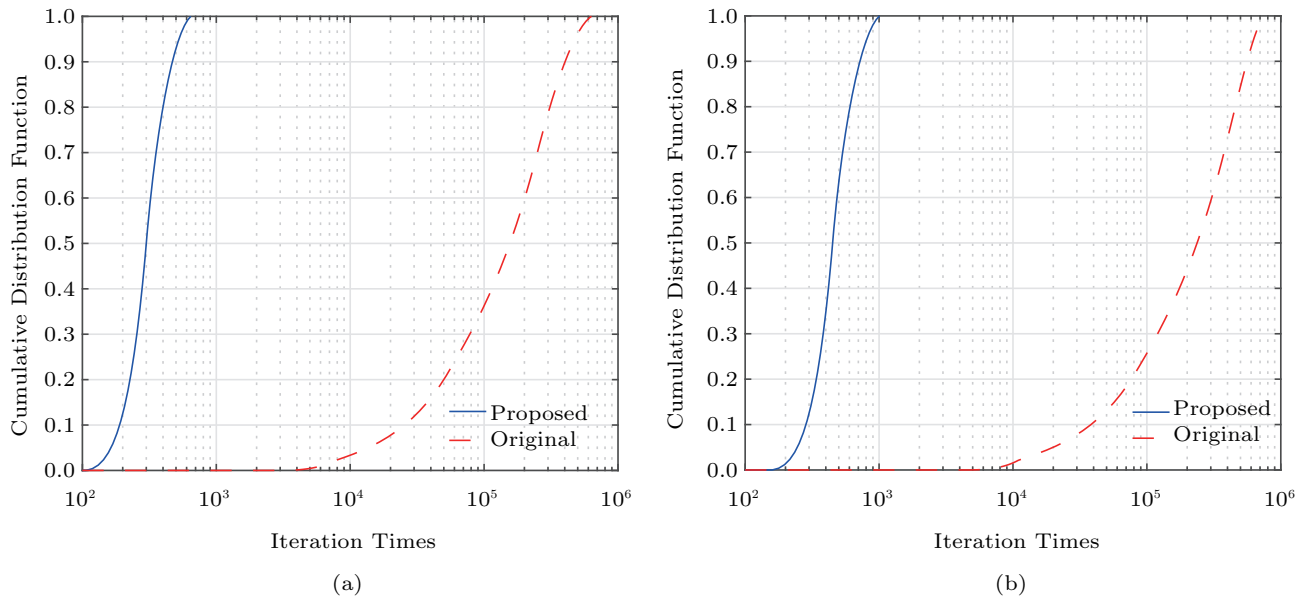


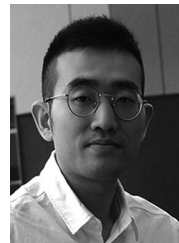
Fig.3. Efficiency comparison in single and multiple social networks.

ing the dynamic continuous data publishing into account to make the model more practical. Furthermore, changing the fixed privacy budget into a varying one could adapt the model into other scenarios.

References

- [1] Garcia D. Leaking privacy and shadow profiles in online social networks. *Science Advances*, 2017, 3(8): Article No. e1701172.
- [2] He Z, Cai Z, Yu J. Latent-data privacy preserving with customized data utility for social network data. *IEEE Transactions on Vehicular Technology*, 2018, 67(1): 665-673.
- [3] Cristofaro E D, Soriente C, Tsudik G, Williams A. Hummingbird: Privacy at the time of twitter. In *Proc. the 2012 IEEE Symposium on Security and Privacy*, May 2012, pp.285-299.
- [4] Abawajy J H, Ninggal M I H, Herawan T. Privacy preserving social network data publication. *IEEE Communications Surveys and Tutorials*, 2016, 18(3): 1974-1997.
- [5] Yu S, Zhou W, Guo S, Guo M. A feasible IP traceback framework through dynamic deterministic packet marking. *IEEE Transactions on Computers*, 2016, 65(5): 1418-1427.
- [6] Qu Y, Yu S, Gao L, Zhou W, Peng S. A hybrid privacy protection scheme in cyber-physical social networks. *IEEE Transactions on Computational Social Systems*, 2018, 5(3): 773-784.
- [7] Qu Y, Yu S, Zhou W, Peng S, Wang G, Xiao K. Privacy of things: Emerging challenges and opportunities in wireless Internet of Things. *IEEE Wireless Communications*, 2018, 25(6): 91-97.
- [8] Yu S, Liu M, Dou W, Liu X, Zhou S. Networking for big data: A survey. *IEEE Communications Surveys and Tutorials*, 2017, 19(1): 531-549.
- [9] Zhu T, Xiong P, Li G, Zhou W. Correlated differential privacy: Hiding information in non-IID data set. *IEEE Transactions on Information Forensics and Security*, 2015, 10(2): 229-242.
- [10] Koufogiannis F, Pappas G J. Diffusing private data over networks. *IEEE Transactions on Control of Network Systems*, 2016, 5(3): 1027-1037.
- [11] Wang W, Zhang Q. Privacy preservation for context sensing on smartphone. *IEEE/ACM Transactions on Networking*, 2016, 24(6): 3235-3247.
- [12] Yu S. Big privacy: Challenges and opportunities of privacy study in the age of big data. *IEEE Access*, 2016, 4: 2751-2763.
- [13] Mohassel P, Zhang Y. SecureML: A system for scalable privacy-preserving machine learning. In *Proc. the 2017 IEEE Symposium on Security and Privacy*, May 2017, pp.19-38.
- [14] Costantino G, Martinelli F, Santi P. Investigating the privacy versus forwarding accuracy tradeoff in opportunistic interest-casting. *IEEE Transactions on Mobile Computing*, 2014, 13(4): 824-837.
- [15] Pierangela S, Latanya S. Protecting privacy when disclosing information: k -anonymity and its enforcement through generalization and suppression. <https://dataprivacylab.org/dataprivacy/projects/kanonymity/paper3.pdf>, May 2018.
- [16] Machanavajjhala A, Kifer D, Gehrke J, Venkitasubramaniam M. L -diversity: Privacy beyond k -anonymity. *ACM Transactions on Knowledge Discovery from Data*, 2007, 1(1): Article No. 3.
- [17] Gong X, Chen X, Xing K, Shin D, Zhang M, Zhang J. Personalized location privacy in mobile networks: A social group utility approach. In *Proc. the 2005 IEEE Conference on Computer Communications*, April 2015, pp.1008-1016.

- [18] Dwork C. Differential privacy. In *Proc. the 33rd International Colloquium on Automata, Languages and Programming*, July 2006, pp.1-12.
- [19] Zhu T, Li G, Zhou W, Yu P S. Differentially private data publishing and analysis: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 2017, 29(8): 1619-1638.
- [20] Wang Q, Zhang Y, Lu X, Wang Z, Qin Z, Ren K. Real-time and spatio-temporal crowd-sourced social network data publishing with differential privacy. *IEEE Transactions on Dependable and Secure Computing*, 2016, 15(4): 591-606.
- [21] Zhang K, Liang X, Lu R, Shen X. PIF: A personalized fine-grained spam filtering scheme with privacy preservation in mobile social networks. *IEEE Transactions on Computational Social Systems*, 2015, 2(3): 41-52.
- [22] Yu S, Guo S, Stojmenovic I. Fool me if you can: Mimicking attacks and anti-attacks in cyberspace. *IEEE Transactions on Computers*, 2015, 64(1): 139-151.
- [23] Qu Y, Cui L, Yu S, Zhou W, Wu J. Improving data utility through game theory in personalized differential privacy. In *Proc. the 2018 IEEE International Conference on Communications*, May 2018, Article No. 656.
- [24] Wu D, Yang B, Wang R. Scalable privacy-preserving big data aggregation mechanism. *Digital Communications and Networks*, 2016, 2(3): 122-129.
- [25] Wang Q, Hu S, Ren K, Wang J, Wang Z, Du M. Catch me in the dark: Effective privacy-preserving outsourcing of feature extractions over image data. In *Proc. the 35th Annual IEEE International Conference on Computer Communications*, April 2016, Article No. 131.
- [26] Ma J, Liu J, Huang X, Xiang Y, Wu W. Authenticated data redaction with fine-grained control. *IEEE Transactions on Emerging Topics in Computing*. doi:10.1109/TETC.2017.2754646.
- [27] Qu Y, Yu S, Gao L, Niu J. Big data set privacy preserving through sensitive attribute-based grouping. In *Proc. the 2017 IEEE International Conference on Communications*, May 2017, Article No. 792.
- [28] Dwork C, McSherry F, Nissim K, Smith A D. Calibrating noise to sensitivity in private data analysis. In *Proc. the 3rd Theory of Cryptography Conference*, March 2006, pp.265-284.
- [29] Du X, Guizani M, Xiao Y, Chen H. Secure and efficient time synchronization in heterogeneous sensor networks. *IEEE Transactions on Vehicular Technology*, 2008, 57(4): 2387-2394.
- [30] Aghasian E, Garg S, Gao L, Yu S, Montgomery J. Scoring users' privacy disclosure across multiple online social networks. *IEEE Access*, 2017, 5: 13118-13130.
- [31] Wasserman L, Zhou S. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 2010, 105(489): 375-389.
- [32] Jorgensen Z, Yu T, Cormode G. Conservative or liberal? Personalized differential privacy. In *Proc. the 31st IEEE International Conference on Data Engineering*, April 2015, pp.1023-1034.
- [33] Wang S, Huang L, Tian M, Yang W, Xu H, Guo H. Personalized privacy-preserving data aggregation for histogram estimation. In *Proc. the 2015 IEEE Global Communications Conference*, December 2015, Article No. 423.
- [34] He Z, Cai Z, Yu J. Latent-data privacy preserving with customized data utility for social network data. *IEEE Transactions on Vehicular Technology*, 2018, 67(1): 665-673.
- [35] Nie Y, Yang W, Huang L, Xie X, Zhao Z, Wang S. A utility-optimized framework for personalized private histogram estimation. *IEEE Transactions on Knowledge and Data Engineering*. doi:10.1109/TKDE.2018.2841360.
- [36] McAuley J, Leskovec J. Social circles: Google+. <https://snap.stanford.edu/data/egonets-Gplus.html>, Nov. 2018.



Lei Cui received his B.S. degree in electrical engineering from Taiyuan University of Technology, Taiyuan, in 2010. He is currently pursuing his Ph.D. degree at the School of Information and Computer, Taiyuan University of Technology, Taiyuan. His research interests include security and privacy issues in the IoT, social networks, and big data.



Youyang Qu received his B.S. and M.S degrees in mechanical engineering from Beijing Institute of Technology, Beijing, in 2012 and 2015, respectively. He is currently pursuing his Ph.D. degree at the School of Information Technology, Deakin University, Melbourne. His research interests focus on addressing security and privacy issues in social networks, cloud computing, IoT, and big data.



Mohammad Reza Nosouhi received his M.S. degree in telecommunications engineering from Isfahan University of Technology, Iran, in 2014. He is currently pursuing his Ph.D. degree in information technology with the School of Software, University of Technology Sydney, Sydney, Australia. He joined the Commonwealth Scientific and Industrial Research Organisation (Data61 unit) in 2018 as a visiting scholar. He has won Australian Post Graduate Award Scholarship, funded by the Australian Federal Government in January 2017. He has published a couple of papers in the area of data security and privacy. His research interests include data security and privacy, applied cryptography and AI applications in security.



Shui Yu is currently a full professor of School of Software, University of Technology Sydney, Sydney, Australia. Dr. Yu's research interest includes security and privacy, networking, big data, and mathematical modelling. He has published two monographs and edited two books, more than 200 technical papers, including top journals and top conferences. Dr. Yu initiated the research field of networking for big data in 2013. His *h*-index is 32. He is currently serving the editorial boards of IEEE Communications Surveys and Tutorials (Area Editor), IEEE Communications Magazine, IEEE Internet of Things Journal, IEEE Communications Letters, IEEE Access, and IEEE Transactions on Computational Social Systems. He is a senior member of IEEE, a member of AAAS and ACM, and a distinguished lecturer of IEEE Communication Society.



Jian-Wei Niu received his Ph.D. degree in computer science from Beihang University, Beijing, in 2002. He was a visiting scholar in the School of Computer Science, Carnegie Mellon University, from Jan. 2010 to Feb. 2011. He is a professor in the School of Computer Science and Engineering, Beihang University, Beijing. He received the New Century Excellent Researcher Award from Ministry of Education of China 2009, and the First Prize of Technical Invention of the Ministry of Education of China 2012. His current research interests include mobile and pervasive computing and mobile video analysis. He is a senior member of IEEE.



Gang Xie received his B.S. degree in control theory and his Ph.D. degree in circuits and systems from Taiyuan University of Technology, Taiyuan, in 1994 and 2006, respectively. He is currently the vice president of Taiyuan University of Science and Technology, Taiyuan, and has also been a professor of Taiyuan University of Technology, Taiyuan, since 2008. His main research interests cover intelligent information processing, complex networks and big data. He has received six provincial science and technology awards, authored more than 100 papers and held five invention patents.

Appendix

Proof of Theorem 2. The noise of mechanism $\mathcal{M} =$

$(\mathcal{M}_1, \mathcal{M}_2)$ is defined by (6). According to this, we prove the propose mechanism satisfies all the required properties.

1) The first coordinate is Laplacian-distributed with parameter $\frac{1}{\epsilon_1}$. When $x > 0$, we can derive (A1).

$$\begin{aligned}
 & \Pr(V_1 = x) \\
 &= \int_R \rho(x, y) dy \\
 &= \frac{\epsilon_1^2}{2\epsilon_2} \text{Exp}(-\epsilon_2 x) + \\
 & \quad \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \int_R \text{Exp}(-\epsilon_1|x - y| - \epsilon_2|y|) dy \\
 &= \frac{\epsilon_1^2}{2\epsilon_2} \text{Exp}(-\epsilon_2 x) + \\
 & \quad \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \left(\int_{-\infty}^0 \text{Exp}(-\epsilon_1 x + (\epsilon_1 + \epsilon_2)y) dy + \right. \\
 & \quad \left. \int_0^x \text{Exp}(-\epsilon_1 x + (\epsilon_2 - \epsilon_1)y) dy + \right. \\
 & \quad \left. \int_0^{+\infty} \text{Exp}(-\epsilon_1 x + (\epsilon_1 + \epsilon_2)y) dy \right) \\
 &= \frac{\epsilon_1^2}{2\epsilon_2} \text{Exp}(-\epsilon_2 x) + \\
 & \quad \frac{\epsilon_1(\epsilon_2 - \epsilon_1)}{4\epsilon_2} \text{Exp}(-\epsilon_1 x) \text{Exp}((\epsilon_1 + \epsilon_2)y)|_{-\infty}^0 + \\
 & \quad \frac{\epsilon_1(\epsilon_2 + \epsilon_1)}{4\epsilon_2} \text{Exp}(-\epsilon_1 x) \text{Exp}((\epsilon_2 - \epsilon_1)y)|_0^x + \\
 & \quad \frac{\epsilon_1(\epsilon_2 - \epsilon_1)}{4\epsilon_2} \text{Exp}((\epsilon_1 + \epsilon_2)y)|_x^{+\infty} \\
 &= \frac{\epsilon_1}{2} \text{Exp}(-\epsilon_1 x). \tag{A1}
 \end{aligned}$$

When $x < 0$, the equation follows the symmetry $(x, y) \rightarrow (-x, -y)$. Thus we can conclude \mathcal{M}_∞ is ϵ_1 -differential private and obtains the best data utility.

2) The second coordinate is Laplacian-distributed with parameter $\frac{1}{\epsilon_2}$. We can derive (A2).

$$\begin{aligned}
 & \Pr(V_2 = y) \\
 &= \int_R \rho(x, y) dx = \frac{\epsilon_1^2}{2\epsilon_2} \text{Exp}(-\epsilon_2|y|) + \\
 & \quad \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \text{Exp}(-\epsilon_2|y|) \int_R \text{Exp}(-\epsilon_1|x - y|) dx \\
 &= \frac{\epsilon_1^2}{2\epsilon_2} \text{Exp}(-\epsilon_2|y|) + \\
 & \quad \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \text{Exp}(-\epsilon_2|y|) \int_R \text{Exp}(-\epsilon_1|x|) dx \\
 &= \frac{\epsilon_1^2}{2\epsilon_2} \text{Exp}(-\epsilon_2|y|) + \frac{\epsilon_2^2 - \epsilon_1^2}{4\epsilon_2} \text{Exp}(-\epsilon_2|y|) \\
 &= \frac{\epsilon_2}{2} \text{Exp}(-\epsilon_2|y|). \tag{A2}
 \end{aligned}$$

Therefore, we prove \mathcal{M}_ϵ is ϵ_2 -differential private and obtains the best data utility.

3) At last, we need to prove the composition mechanism maintains ϵ_2 -differential private. The delta part is separately handled by defining $L = \{x : (x, x) \in \Omega\}$. The probability of landing in Ω is represented by

$$\begin{aligned} & \Pr(\mathcal{M} \in S) \\ &= \frac{\epsilon_1^2}{2\epsilon_2} \int_D \text{Exp}(-\epsilon_2|x-d|)dx + \\ & \quad \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \iint_S \text{Exp}(-\epsilon_1|(x-d) - (y-d)| - \\ & \quad \epsilon_2|y-d|)dxdy. \end{aligned}$$

We take the derivative and use Fubini's theorem to exchange the derivative with integral as (A3).

$$\begin{aligned} & \frac{d}{du} \Pr(\mathcal{M} \in S) \\ &= \frac{\epsilon_1^2}{2\epsilon_2} \int_D \epsilon_2 \text{sgn}(x-d) \text{Exp}(-\epsilon_2|x-d|)dx + \\ & \quad \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \iint_S \epsilon_2 \text{sgn}(x-d) \times \\ & \quad \text{Exp}(-\epsilon_1|x-y| - \epsilon_2|y-d|)dxdy \\ &\Rightarrow \left| \frac{d}{du} \Pr(\mathcal{M} \in S) \right| \leq \frac{\epsilon_1^2}{2\epsilon_2} \int_D \epsilon_2 \text{Exp}(-\epsilon_2|x-d|)dx + \\ & \quad \frac{\epsilon_1(\epsilon_2^2 - \epsilon_1^2)}{4\epsilon_2} \iint_S \epsilon_2 \text{Exp}(-\epsilon_1|x-y| - \epsilon_2|y-d|)dxdy \\ &\Rightarrow \left| \frac{d}{du} \Pr(\mathcal{M} \in S) \right| \leq \epsilon_2 \Pr(\mathcal{M} \in S) \\ &\Rightarrow \left| \frac{d}{du} \ln \Pr(\mathcal{M} \in S) \right| \leq \epsilon_2. \end{aligned} \tag{A3}$$