

# Collective sanitization methods to prevent and detect the inference attacks in online social networking sites

V.UMA SANKARI<sup>1</sup> AND G.ELAYAROJA<sup>2</sup>

<sup>#</sup> Dept. of CSE, CK College of Engineering & Technology, Cuddalore, Tamil Nadu, India

<sup>\*</sup> Dept. of CSE, CK College of Engineering & Technology, Cuddalore, Tamil Nadu, India

**Abstract**— social networking service is a virtual stage for many users to connect or build social relation among people. Releasing social network data could seriously breach the security of the online site users. Inference attack is a data mining technique performed by analyzing a database. Inference attack is hard to mitigate due to the release of user's data attributes such as sharing, friendship links, nodes and users sensitive information. Here we have used the latent data privacy method to ensure the users privacy in social networking sites form adversary prediction. This paper attempts to reduce the adversary's prediction on sensitive information by using data sanitization method called as "collective inference method" to protect against inference attacks in online social networking sites.

**Index Terms**— Latent data privacy, anonymous, networking, inference attack.

## I. INTRODUCTION

In recent years, worldwide online social network usage has enlarged sharply as these networks have become interlinked the daily lives of people as simulated meeting places that facilitate communication. Facebook users have a total of over 110 billion acquaintance connections. Unfortunately, many OSN users are unaware of the security risks which happen in these types of communications, including privacy risk. The privacy issues are classified as direct attack and indirect attack. One such privacy issue is the inference attack.

Users can control who can access their private information by providing privileges[1] to their acquaintances. To stop private information drip, it is essential to be well informing of the ways by which attacker can attack a social network to learn user's isolated sensitive information. Protecting the privacy of individuals in social networks is done by focusing on the identification of fake nodes in online social networks. Online social networking sites have their business needs to encourage the users to easily find each other and wider their friendship networks. Social media sites should

possess new safety challenges to avoid security fear as users may put themselves in unsecured and members of their social networks at risk to a variety of attacks. The goal of this paper is to show how social network data can be exploited to predict unseen essential evidence in the anonymization process.

Privacy concerns in social networks can be mainly categorized into two types: inherent-data privacy and latent data privacy. Inherent-data privacy[2] is related to sensitive data contained in the data profile submitted by users in order to receive data-related services. For example, age and gender are unavoidable data for health related services yet unwilling to be released by most users. De-anonymization towards sensitive data is an inherent-data privacy instance. For example, two New York Time journalists used to successfully identify personal information from the published search logs involving 650,000 users made available by AOL. The logs include the information of name, age, sex, location, etc., and such information is associated with a specific individual. The latent data privacy[3] is related to unreleased sensitive information, yet such sensitive information can be inferred from released data or users' social relationships. For instance, Jenny does not publish her political opinions online, yet such information could be inferred by mining her friends' data as Jenny's social relationships may be public. In this paper, we focus on latent-data privacy.

We assume third party users may collect anonymous user data from social networks. Some users disclose their sensitive information, while others do not. However, third party users can carry out de-anonymization actions and further infer sensitive information of users. We first collect the de-anonymized dataset of users in online social networking sites to infer sensitive information hidden[4] in the released data. Then, we propose some effective data sanitization strategies to prevent information inference attacks. On the other hand, the sanitized data obtained by these strategies should not reduce the valuable benefit brought by the abundant data resources, so that non-sensitive information can still be inferred and utilized by third party users.

The previous works of preventing inference attack as certain deficiencies in user's friendship links and type of manipulation method such as filtering the dataset, perturbation, data smoothing and anonymization. Therefore the previous works cannot predict accuracy of the adversary prediction methods. In this work we tend to provide balance between the data utility and data privacy. Our strategies will ensure that third party users cannot obtain the sensitive and non-sensitive information from online social networking sites without the users privacy permission.

In the previous works of inference attack we have used the K-means clustering for producing the good clustering of database for predicting attacker or adversary background knowledge of collection sensitive attributes of the users in online social networking sites. The other algorithm naïve bayes classification algorithm is used to cluster the data from one person to another person sensitive attributes collection for analyzing the friendship links. Naïve Bayes classification is used for classification of a dataset. Imagine that dataset are drawn from a number of classes of datasets which can be used independent probability that the  $i$ -th word of a given dataset occurs in a dataset from class  $C$  can be written as

$$P(w_i | C)$$

Simplifying things further by assuming that words are randomly distributed in the dataset - that is, words are not dependent on the length of the dataset, position within the dataset with relation to other words or other dataset-context. Then the probability that a given dataset  $D$  contains all of the words, given a class  $C$ , is

$$P(D|C) \equiv \prod_i P(w_i | C)$$

#### A. Inference attack

Inference attack is a data mining technique by which the attacker gets sensitive data from a user in online social networks from the non-sensitive data. In this process, researchers discover situations in which collective inference technique improves on using a simple local classification technique to identify friendship links between users. Users have strong outlooks of privacy on such data. When social network data is made open in one way or another, it is far from necessary to protect privacy by simply interchanging the identifying attributes.

#### B. Latent data privacy

In this paper we propose a collective inference method called as latent-data privacy to protect the user's sensitive information in online social networking sites. Some users disclose their sensitive information, while others do not. However, third party users can carry out de-anonymization actions and further infer sensitive information of users in social networking sites.

#### C. Collective inference method

The attribute-based Classifier just considers the attribute sets of the users it is classifying. Conversely,

relation-based classifier only considers the friendship information of a user. However, third party users may launch an inference attack by exploiting all the publicly available information. Moreover, a major problem of relation-based classifier is that it requires that at least one of the neighbors of each unlabeled user to be located in the training set (i.e., the set of users with known labels). Collective inference attempts to tackle the above two issues by considering both attribute-based classifier and relation-based classifier in a collaborative manner to improve prediction accuracy.

## II. RELATED WORKS

One promising way to defend against Inference attacks in social networks is to leverage the social network topologies.

Heathery et al., [5] proposed the Bayesian network algorithm to detect the inference attacks in online social networks for protecting privacy. Here techniques are provided that can help with choosing the most real details or links that need to be unconcerned for protecting privacy. Finally, it explores the effect of collective inference techniques in possible inference attacks. Collective inference is a method of sorting social network data using a combination of node points and connecting links in the social graph. After the results obtained from the collective inference method SVMs as a classification technique is implemented.

Sun et al., [6] proposed the K-anonymity model and privacy inference graph, introduces the concept of valid inference path to deduce some privacy information resulting in privacy disclosure. The inference attack in K-anonymous data collections and analyzing their privacy disclosures, builds up a privacy inference graph created on attack graph theory, which is a leeway for attack graph. The privacy inference graph labels widely the inference attack in k-anonymous databases by seeing attacker background information and external factors.

Liu and yang [7] proposed the increasing attractiveness of social networks in various areas has raised privacy concerns for the users involves proposed overall framework for preventing link inference attacks, which accepts a novel lineage tracing mechanism to efficiently cut off the reading paths of sensitive relationships. The increasing attractiveness of social networks in various application areas has raised privacy concerns for the individuals involved.

Friedman et al., [8] proposed the  $k$ -anonymity and use them to prove that a specified data mining model does not disrupt the  $k$ -anonymity of the persons represented in the learning examples. Finally, it contributes new and efficient ways to anonymized data and preserve patterns during anonymization.

Ma et al., [9] proposed the vertex and edge modification algorithm which comprises the finding the optimal target degree of each vertex, deciding the candidates to rise the vertex degree and adding the edges among vertices to satisfy the requirement.

Dewri et al., [10] proposed the U2G matching algorithm is designed to match users with groups with the goal of enhancing compactness. The U2G matching is a worldwide activity distributed on the user and group agents fitting to the agent. network. Each user agent and respective group agent periodically finishes the following user agent task resp. group agent task, where calls each time the task is executed, and denote as T the constant time period passing between two consecutive agents.

Papagelis et al., [11] proposed sampling-based algorithm on improving the performance of information collection to efficiently explore a user's social network regarding its construction and representations of centralized and distributed social networks are considered. Algorithms can be utilized to rank objects in the community.

Wang et al., [12] proposed semantic-based friend recommendation system for social networks, which applauds friends to users based on their lifetime styles instead of social graphs. By taking benefit of sensor-rich smartphones, Friendbook determines life styles of users from user-centric sensor data, measures the resemblance of life styles between users, and recommends friends to users if their lifetime styles have high similarity.

Fire et al., [13] proposed many online social network users are unaware of the abundant security risks that exist in these networks, including privacy violations, identity theft. In addition, presenting an overview of existing explanations that can offer well protection, security and privacy OSN users also offer simple-to-implement references for OSN users, which can increase their security and privacy when using these platforms.

Hiro et al., [14] proposed covert channel analysis models with inference rules and reading and writing operations. The histories of accesses are required to analyze information leakage with inference rules. The histories are used for checking the conditions under which some inference rules are fired. The management system of access histories should not require a trusted third party, and nay participants of the system cannot falsify stored histories.

Gao et al., [15] proposed SocInf that can breach the membership privacy of a given machine learning model's training data set under a demanding scenario, in which SocInf can get nothing except the prediction results of the victim model on input data. The key observation of SocInf is that the machine learning models often have different prediction behaviors on the data that they were trained on versus the data that they "meet" for the first time.

Nasr et al., [16] proposed a comprehensive framework for the privacy analysis of deep neural networks, using white-box membership inference attacks. We take all major scenarios where deep learning is used for training and fine-tuning or updating models, with one or multiple collaborative data holders, when attacker only passively observes the model updates or actively influences the target model in order to extract more information, and for

attackers with different types of prior knowledge. Despite differences in knowledge, observation, and actions of the adversary, their objective is the same: membership inference.

Qian et al., [17] proposed we will model the attacker's prior knowledge using knowledge graphs and use them to express the two attacks stages – de-anonymization and privacy inference. We transform the problem of de-anonymizing a group of people into the maximum weighted bipartite matching problem.

Zhao et al., [18] proposed the first theoretical foundation that gives a non-asymptotic bound on the performance of k-anonymity against inference attacks, taking into consideration of adversaries' background information. The intuition stems from the observation that the deterioration of the performance of k-anonymity is largely resulted from adversaries' background information about the cloaked data.

### III. PROPOSED SYSTEM

The main objective of the proposed work is to detect inference attack in the online social network. In the proposed work, classification and clustering algorithms are used to detect inference attack.

Detection of inference attack is done in two stages. Two techniques included in the detection are as follows:

- Identification of inference attack
- Prevention of inference attack

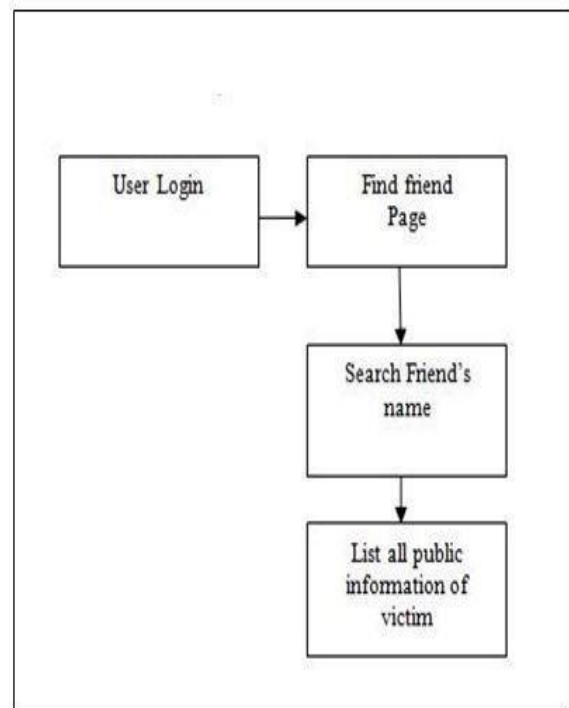
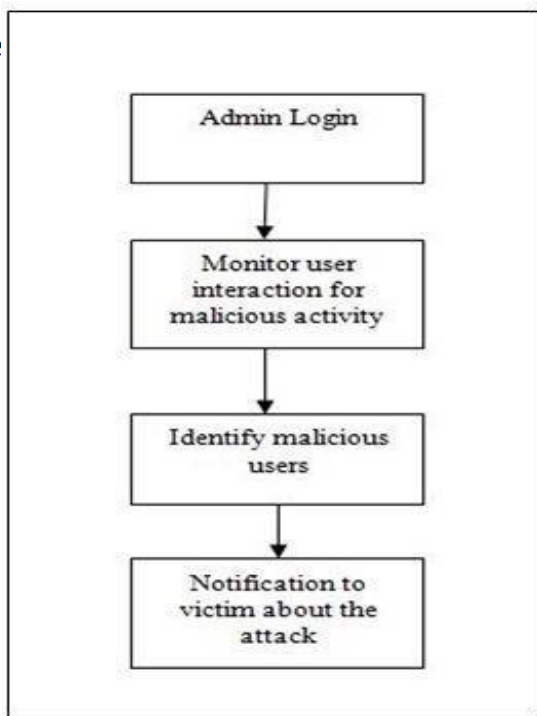
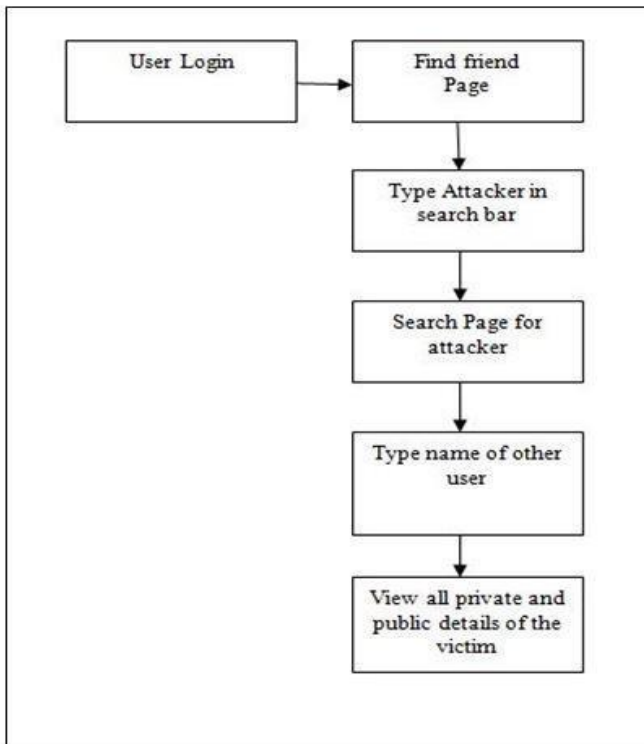


Figure1. Normal user actions in an online social networking site



**Figure 2.** Attacker actions in an online social networking site



**Figure 3.** Administrator actions in an online social networking site

*A. Collection of Attacker Attributes*

In the existing system a web application similar to facebook is built and periodically collects the attacker attributes from the victims of facebook users. These gathered attributes are stored as attributes in the database.

*B. Data Filtering*

In the proposed system, data filtering is done to predict recently updated attributes of attackers in the dataset. In order to compare and detect the attributes obtained recently from the victim users of facebook. K -means algorithm is used to gather the attacker attributes from the victims of facebook. The gathered attacker attributes are compared with the database to detect the attacker. After gathering the attacker attributes, mitigation of attacker from not moving again to the same node, friendship links is done by classification algorithms.

*C. Data Smoothing*

The gathered attributes from the victims of facebook users are leveled by Naïve Bayes classification algorithm. These leveled attributes are blocked by the users of facebook. Naïve Bayes classification is based on the Bayesian theorem. It is mostly suited when the dimensionality of the inputs is high.

*D. Collective inference method*

Collective inference attempts to tackle the above issues by considering both attribute-based classifier and relation-based classifier in a collaborative manner to improve prediction accuracy. Formally, we consider the following network prediction problem.

Collective inference. Given social graph of the users in online social networking sites  $G(V = V^k \cup V^u, E, X, Y = Y^k \cup Y^u, H_s)$  with user set  $V$ , friendship link set  $E$ , the set of user attribute sets  $X$ , and the set of sensitive categories  $H_s$ .  $y^i \in Y$  is a class label of  $u_i$  for an arbitrary category  $h_r \in H_s$ .  $L^k$  is the known labels for users  $u_i \in V^k$ . Collective inference is to predict  $Y^u$  for users  $u_i \in V^u$ , where  $V^u = V - V^k$

This problem is challenging as some of the user labels are unknown. A fundamental idea is to first predict a class label approximately and then refine the predicted result iteratively. Several collective classification algorithms have been proposed to increase accuracy when the network users are interrelated, such as the Iterative Classification Algorithm (ICA) and Gibbs sampling (Gibbs). Many collective classification algorithms and variants, including ICA, use an attribute-based classifier  $M_A$  to predict the approximate class label at the bootstrap stage; then, they use both attribute and link based classifier,  $M_{AR}$ , to refine the results. The algorithms repeat these two operations until the class labels converge. We present an algorithm under the framework of ICA that takes RST as a local classifier (one that uses local information, e.g., attribute sets of users), denoted by ICA-RST. ICA-RST is shown in Algorithm. It first learns an attribute-based classifier  $M_A$  based on the known labels  $Y^k$  (step 1), which is a set of RST decision rules. Then, by  $M_A$ , it predicts the labels of the users with unknown labels,  $V^u$  (steps 2-3). Step 5 stores the known labels  $Y^k$  and the predicted labels  $V_{Ug}$  in set  $Y^R$ . The known labels and the predicted labels are utilized to calculate link features for each user in  $V^u$  (step 7). Step 8 then learns a classifier MAR based on all of the attributes and labels. Step10

utilizes MAR to predict unknown labels. Finally, Step 11 returns the predicted results.

**Algorithm 1: ICA-RST**

Input:  $V = \text{users}$ ,  $E = \text{links}$ ,  $X = \text{attribute set}$ ,  $Y^K = \text{labels of unknown users}$  ( $Y^U = \{y_i | u_i \in V^U; \forall u_i \in V - V^K\}$ );

1.  $M_A = \text{learn RST Rule}(V^K; Y^K)$ ;
2. **for** each user  $u_i \in V^U$  **do**
3.  $y_i = M_A(\sim X_i)$ ;
4. **for**  $t = 1$  to  $n$  **do**
5.  $Y^R = Y^K \cup \{y_i | u_i \in V^U\}$ ;
6. **for** each user  $u_i \in V^U$  **do**
7.  $f_i = \text{calReF eats}(V; E; Y^R)$ ;
8.  $M_{AR} = \text{learn RST Rule}(V; Y^R)$ ;
9. **for** each user  $u_i \in V^U$  **do**
10.  $y_i = M_{AR}(\sim X_i; \sim f_i)$ ;
11. **return**  $Y^U$

Generic Attribute Hierarchy. A Generic Attribute Hierarchy (GAH) is a finite hierarchical ordering. The first layer of the ordering is one of the privacy dependent attributes, and each parent layer is a generic of the sublayer. This generic attribute hierarchy indicates that the ancestor of the GAH is the highest level of generalization of initial attributes. Substituting one privacy-dependent attribute with the ancestor of the GAH would render the highest level of privacy. For example, if one attribute value in core is for category favorite movies, the corresponding GAH can be

Star Wars  $\rightarrow$  Fantasy  $\rightarrow$  American film

This indicates that we can substitute original attribute “Star Wars” with “American film”, in order to get the highest level of generalization. We could also substitute it with “Fantasy” to give more preference to utility for the utility-privacy tradeoff since “Fantasy” is more specific than “American film”. Hence, GAH guarantees that we can programmatically determine which level of generic value should be chosen to optimize the privacy-utility tradeoff.

**Algorithm 2: Generate generic value**

Input: Core = utility threshold

Output: GAH

- 1 while
- 2  $\max_{C \subseteq G} c_2(C, K, X_{\text{non}}) - \max_{C' \subseteq C} c_2(C', K, X_{\text{non}}) \geq \text{Core}$  do
- 3 further generate all the current attributes;
- 3 return Perturbed Core

Information table -Reduct Systems for SNAP, Caltec and MIT datasets.

Decisions attribute No.of condition attributes

Gender in SNAP 19 $\rightarrow$ 13

Flag in Caltech 6 $\rightarrow$ 5

Flag in MIT 6 $\rightarrow$ 5

**IV.EVALUATION**

**4.1 Datasets**

In our experiments, we investigate three different Facebook datasets. The first one is the SNAP Facebook dataset1 which contains user friendships and a number of node attributes such as gender, birthday, position, employer, location, etc. The other two are the Facebook

dataset containing all the Facebook friendships at Caltech and MIT in 2005, as well as a number of node attributes such as student/faculty status flag, gender, graduation year, academic major, etc. 2 For convenience, we denote these three datasets as SNAP, Caltech, and MIT, respectively. In Caltech and MIT, each attribute is specified by a numeric value and each of which indicates a Corresponding attribute. However, in SNAP, each attribute is specified by a 0/1 value and each of which indicates the absence/ presence of the corresponding attribute. For example, attribute “Education Degree: undergraduate; master; PHD” with attribute value 010 means that the attribute value is master. For convenience, we map each attribute in SNAP into a unique numeric value in each attribute category. For example, the above attribute value 010 in Education degree is mapped to 2.

**V. CONCLUSIONS**

We address two issues in this paper: (a) how exactly third party users launch an inference attack to predict sensitive information of users, and (b) are there effective strategies to protect against such an attack to achieve a desired privacy utility tradeoff. For the first issue, we show that collectively utilizing both attribute and link information can significantly increase prediction accuracy for sensitive information. For the second issue, we explore the dependence relationships for utility/public attributes, and privacy/public attributes. Based on these results, we propose a Collective Method that take advantages of various data manipulating methods to guarantee sanitizing user data does not incur a bad impact on data utility. Using Collective Method, we are able to effectively sanitize social network data prior to release. The solutions for the two addressed issues are proven to be effective towards three real social datasets.

**VI. REFERENCES**

- [1] Chandra and Antony Rosewell, “Sanitations To Prevent Inference Attack On Social Network Data”, International Journal of Innovative Research in Science, Engineering and Technology, Volume 3, Special Issue 1, February 2014.
- [2][https://en.wikipedia.org/wiki/Inference\\_attack](https://en.wikipedia.org/wiki/Inference_attack).
- [3][https://en.wikipedia.org/wiki/k\\_means\\_clustering](https://en.wikipedia.org/wiki/k_means_clustering).
- [4][https://en.wikipedia.org/wiki/Naive\\_Bayes\\_classifier](https://en.wikipedia.org/wiki/Naive_Bayes_classifier)
- [5] Raymond Heathery, Murat Kantarcioglu and Bhavani Thuraisingham, “Preventing Private Information Inference Attacks on Social Networks”, IEEE Transactions on Knowledge and Data Engineering, Vol. 25, No. 8, 2013.
- [6] Yan Sun, Lihua Yin, Licai Liu, Shuang Xin, “Toward inference attacks for k-anonymity”, Springer Personal and Ubiquitous Computing, Vol. 18, Issue 8:1871-1880, 2014.
- [7] Xiangyu Liu, Xiaochun Yang, “Protecting Sensitive Relationships against Inference Attacks in Social Networks”, Springer Database Systems for Advanced Applications, Lecture Notes in Computer Science, Vol. 238, Issue 335-350, 2012.
- [8] Arik Friedman, Ran Wolff Assaf, Schuster, “Providing k-anonymity in data mining”, Springer the VLDB Journal, Vol. 17, Issue 789-804, 2008.
- [9] Tinghuai Ma, Yuliang, Zhang, Jie Cao, Jian, Shen, Meili Tang, Yuan Tian, Abdullah Al-Dhelaan, Mznah Al-Rodhaan, “KDDEM: A k-degree anonymity with vertex and edge modification algorithm Computing”, Springer DOI 10.1007/s00607-015-0453, 2015.

- [10] Rinku Dewri, Member, Indrajit Ray, Member, Indrakshi Ray, Darrell Whittle, “k-Anonymization in the Presence of Publisher Preferences”, IEEE Transactions on Knowledge and Data Engineering, Vol. 23, No. 11, 2011.
- [11] Michael Fire, Roy Goldschmidt, and Yuval Elovici, “Online Social Networks: Threats and Solutions”, IEEE Communication Surveys & Tutorials, Vol. 16, No. 4, fourth quarter, 2014.
- [12] Zhibo Wang, Jilong Liao, Qing Cao, Hairong Qi, and Zhi Wang, “Friendbook: A Semantic-Based Friend Recommendation System for Social Networks”, IEEE Transactions on Mobile Computing, Vol. 14, No. 3, 2015.
- [13] Manos Papagelis, Gautam Das, and Nick Koudas, “Sampling Online Social Networks”, IEEE Transactions on Knowledge and Data Engineering, Vol. 25, No. 3, 2013.
- [14] Hiro “Access control model for the inference attacks with access histories”, 2017, Annual Computer Software and Applications Conference
- [15] gaoyang liu, chen wang , kai peng ,Haojun huang, yutong li, and wenqing cheng, “socinf: membership inference attacks on social Media health data with machine learning”, iee transactions on computational social systems, vol. 6, no. 5, october 2019
- [16] Milad Nasr, Reza Shokr, Amir Houmansadr, “Comprehensive Privacy Analysis of Deep Learning”,2019 iee symposium on security and privacy.
- [17] Jianwei Qian, Xiang-Yang Li, Chunhong Zhang, Linlin Chen,Taeho Jung, “Social Network De-Anonymization and Privacy Inference with Knowledge Graph Model”,unpublished.
- [18] Ping Zhao, Hongbo Jiang,Chen Wang, Haojun Huang, Gaoyang Liu, and Yang Yang, “On the Performance of k-Anonymity against Inference Attacks with Background Information”, IEEE internet of things journal.