

Quantitative Models for Privacy Protection

Tsan-sheng Hsu Churn-Jung Liao* Da-Wei Wang
Institute of Information Science
Academia Sinica, Taipei, Taiwan
E-mail: {tshsu, liaucj, wdw}@iis.sinica.edu.tw

Jeremy K.-P. Chen
Wireless Networking and Communications Group (WNCG)
Dept of ECE, 406 ENS Bldg.,
University of Texas, Austin TX 78712-1084
E-mail: kchen@ece.utexas.edu

Abstract

We assume a database consists of records of individuals with private or sensitive fields. Queries on the distribution of a sensitive field within a selected population in the database can be submitted to the data center. The answers to the queries leak private information of individuals though no identification information is provided. Inspired by decision theory, we present two quantitative models for the privacy protection problem in such a database query or linkage environment in this paper. One is for modeling the value of information from the viewpoint of the querier and the other is for the damage and compensation of privacy leakage.

In both models, we define the information state by a class of probability distributions on the set of possible confidential values. These states can be modified and refined by the user's knowledge acquisition actions. In the first model, the value of information is defined as the expected gain of the privacy receiver and the privacy is protected by imposing costs on the answers of the queries for balancing the gain. In the second one, the safety is guaranteed by enforcing that anyone misusing the private information must pay more compensation than the possible gain.

Key words: Privacy, Data table, Decision logic, Quantitative models, Information value.

1 Introduction

With the rapid development of computer and communication technology, it has become much easier to store massive amounts of data in a central location and distribute them to the end

*Corresponding author.

users via Internet. In appropriate use, a data may be valuable information sources for scientists, analysts, or policy and decision makers. However, there may be a high risk of privacy invasion if it is accessed without any restriction. As in [1], “in the past, most individuals lacked the time and resources to conduct surveillance required to invade an individual’s privacy, as well as the means to disseminate the information uncovered, so privacy violations were restricted to those who did, mainly the government and the press”. However, the development of Internet technology has changed the situation radically. Nowadays, any individual Internet user may easily spread information worldwide within seconds. In such a situation, the revelation of private information to unauthorized users, even if not intentionally, may cause a serious invasion of human rights.

There are many technical problems to be addressed for privacy protection. The most basic one is to avoid unauthorized users access to the confidential information. In [7], an epistemic logic is proposed to model the privacy protection problem in a database linking context. In [3], a prototype system is designed to implement this logical model in the context of querying a medical database. The safety criteria of the data is defined rigorously in the logic and data to be disclosed must be generalized to meet this requirement. The safety criteria defined in the requirement are purely qualitative so we can only identify the situation in which the exact confidential information came to the users’ knowledge. However, in many cases, even if the private information is only known with some imprecision, there is still a risk of privacy leakage. Thus it is very important to have the capability of risk assessment in the model of privacy protection.

Someone may benefit from the privacy leakage, but it may also be harmful for others. For example, the health information of a customer would be valuable in the decision-making of an insurance company. However, the dissemination of an individual’s health information without his consent in advance is definitely an invasion of his privacy. Thus the value of confidential information would be an incentive towards invasion of privacy. The information brokers may try to collect and sell personal information for their own interest. On the other hand, it is usually difficult to estimate the damage caused by privacy leakage. However, to discourage the invasion of privacy, the damage of the victim must be appropriately compensated by the one disseminating the information. Therefore, the evaluation of gain and loss of privacy leakage is a crucial problem in privacy protection.

In this paper, we try to tackle the problem from the aspects of information value and the damage caused by privacy leakage. We focus on the following database query environment. In a data center, private information about individuals are collected. There are private or sensitive fields as well as identification fields in each record. Queries on the distribution of a sensitive field within a selected population in the database can be submitted to the data center. The answers to the queries leak private information of individuals though no identification information is provided.

We study two quantitative models for the privacy protection problem in such a database query environment. One is for modeling the value of information from the viewpoint of the querier and the other is for the damage and compensation of privacy leakage. For the former, we will model the value of information as the expected gain of knowledge of the information. For the latter, the safety of data is guaranteed by enforcing that anyone disseminating the private information must

pay more compensation than his gain from such behavior. In both models, we need to represent the knowledge states of an user receiving some kind of information. The knowledge states can be changed or refined by receiving some answer to the user’s query. Thus we also need a formalism to represent the data to be protected and a language to describe which kinds of queries are allowed. The data table and decision logic proposed in [12] will be employed as the data representation formalism and the query language respectively.

In the rest of the paper, we first review the data table formalism and the decision logic in our application context. The basic components of our models—the information states and knowledge acquisition actions—is defined in section 3. In section 4 and 5, the two models for information value and the damage and compensation of privacy leakage are presented. Finally, the results are summarized in the concluding section .

2 Data Representation and Query Language

To state the privacy protection problem, we must first fix the data representation. The most popular data representation is by data table([12]). The data in many application domains, for example, medical records, financial transactions, employee data, etc., can be represented as data tables. A data table can be seen as a simplification of a relational database, since the latter in general consists of a number of data tables. A formal definition of data table is given in [12].

Definition 1 A data table¹ is a pair $T = (U, A)$ such that

- U is a nonempty finite set of individuals, called the population or the universe,
- A is a nonempty finite set of primitive attributes, and
- every primitive attribute $a \in A$ is a total function $a : U \rightarrow V_a$, where V_a is the set of values of a , called the domain of a .

The attributes of a data table can be divided into three sets. The first contains the *key attributes*, which can be used to identify to whom a data record belongs, therefore these attributes are always masked off in response to a query. Since the key attributes uniquely determine the individuals, we can assume that they are associated with elements in the universe U and omit them from this point. Second, we have a set of *easy-to-know attributes*, the values of which are easily discovered by the public. For example, in [16], it is pointed out that some attributes like birth-date, gender, ethnicity, etc., are included in some public databases such as census data or voter registration lists. The last kind of attributes is the *confidential type*, the values of which are mainly the goals we have to protect. Sometimes, there is an asymmetry between the values of a confidential attribute. For example, if the attribute is a HIV test result, the revelation of a ‘+’

¹Also called knowledge representation system, information system, or attribute-value system

value may cause a serious privacy invasion, whereas it does not matter to know that an individual has a '–' value. For simplification, we assume there is exactly one confidential attribute in a data table. Thus a data table is usually written as $T = (U, A \cup \{c\})$ where A is the set of easy-to-know attributes and c is the confidential one.

Let $V_c = \{s_0, s_1, \dots, s_{t-1}\}$ be the set of possible values for the confidential attribute c . It is assumed that the *a priori* information of the user is the probability distribution of the population on V_c . In other words, we assume that the user knows the value $\frac{|\{u \in U \mid c(u) = s_i\}|}{|U|}$ for all $0 \leq i \leq t-1$. Then the user can improve his knowledge by investigating some sampled individuals of the population or querying the data center that stores the data table. By investigation, the user can discover the exact value of the confidential attribute of the chosen individuals. However, much effort is necessary to do the investigation. On the other hand, a query may ask for the probability distribution of sensitive fields in a specific subset of the population. Once the query is correctly answered, the user not only knows the probability distribution of the specific sub-population, but also that of its complement on V_c . Thus we need a language to specify a subset of individuals. To achieve this purpose, we suggest to use the decision logic(DL) proposed in [12]. The DL is originally designed for the representation of rules induced from a data table by data mining techniques. However, it is also perfectly suitable for the query of a data table since each formula of the logic is satisfied by some individuals in the data table.

The atomic formula of a data table $T = (U, A \cup \{c\})$ is of the form (a, v) , where $a \in A$ is an easy-to-know attribute and $v \in V_a$ is a possible value of the attribute a . The well-formed formulas (wff) of the logic is then formed by the Boolean connectives negation(\neg), conjunction(\wedge), disjunction(\vee), and implication(\rightarrow): :

- Each atomic formula is a wff.
- If φ is a wff, so is $\neg\varphi$.
- If φ and ψ are wffs, so are $\varphi \wedge \psi$, $\varphi \vee \psi$, and $\varphi \rightarrow \psi$.

The satisfaction relation \models_T between U and the wffs is defined recursively by the following clauses:

1. $u \models_T (a, v)$ iff $a(u) = v$
2. $u \models_T \neg\varphi$ iff $u \not\models_T \varphi$
3. $u \models_T \varphi \wedge \psi$ iff $u \models_T \varphi$ and $u \models_T \psi$
4. $u \models_T \varphi \vee \psi$ iff $u \models_T \varphi$ or $u \models_T \psi$
5. $u \models_T \varphi \rightarrow \psi$ iff $u \not\models_T \varphi$ or $u \models_T \psi$

It can be seen that intuitively, any individual satisfying (a, v) has v as the value of his attribute a .

From the semantics of decision logic, we define the truth set of a wff φ with respect to the data table T , denoted by $|\varphi|_T$, as $\{u \in U \mid u \models_T \varphi\}$. Thus each wff φ specifies a subset of individuals $|\varphi|_T$ in the data table. When a query φ is submitted by an user to the data center, this means this user wants to know the distribution of the sub-population $|\varphi|_T$ on V_c . If the query is correctly answered, the user would also simultaneously know the distribution of the sub-population $U - |\varphi|_T$ by the axioms of probability. In other words, a correctly answered query would partition the population into two sub-populations and the distributions thereof on the confidential attribute values are known respectively. In this way, the user can subsequently query the data center to refine his knowledge regarding the distributions of the different sub-populations on the confidential attribute values. To model the evolution of the user's information after different queries, we need a formal representation of user's information states. The next section will be devoted to the definitions of such representation.

3 The Information States

From here on, let us fix a data table $T = (U, A \cup \{c\})$. Let $V_c = \{s_0, s_1, \dots, s_{t-1}\}$ be the set of possible values for the confidential attribute and let $U = \{u_1, u_2, \dots, u_n\}$ be the set of individuals. A *logical partition* of U is a subset of DL wffs $\Pi = \{\varphi_1, \varphi_2, \dots, \varphi_m\}$ such that $|\varphi_1|_T \cup \dots \cup |\varphi_m|_T = U$ and $|\varphi_i|_T \cap |\varphi_j|_T = \emptyset$ if $i \neq j$. Each $|\varphi_i|_T$ is called an equivalence class of Π . A piece of information (or knowledge) known to the user is given by a logical partition of U , a set of probability distributions indexed by the wffs of the partition, and the number of investigated individuals. In the following, we use $|\varphi|$ to denote the cardinality of $|\varphi|_T$.

Definition 2 *An information state (or a knowledge state) \mathcal{I} for the set of possible private attribute values V_c and the set of individuals U is a triple $(\Pi, (\mu_i)_{0 \leq i \leq t-1}, (\kappa_i)_{0 \leq i \leq t-1})$, where Π is a logical partition on U and for all $0 \leq i \leq t-1$, $\mu_i : \Pi \rightarrow [0, 1]$ and $\kappa_i : \Pi \rightarrow \mathcal{N}$ (\mathcal{N} denotes the set of natural number) are functions satisfying the following constraints for any $\varphi \in \Pi$,*

$$(i) \sum_{i=0}^{t-1} \mu_i(\varphi) = 1,$$

(ii) $|\varphi| \cdot \mu_i(\varphi)$ is a natural number, and

$$(iii) \kappa_i(\varphi) \leq |\varphi| \cdot \mu_i(\varphi)$$

For ease of description, we use the vector notations in denoting μ_i 's and κ_i 's. Thus $\boldsymbol{\mu} = (\mu_0, \dots, \mu_{t-1})$ and $\boldsymbol{\kappa} = (\kappa_0, \dots, \kappa_{t-1})$ denotes vector mappings which can be applied to elements of Π and the result of such application is a vector consisting of the results of applying its component functions to the element. The dimension of each vector will be self-evident from the context and not explicitly specified. By the vector notation, an information state defined above can be denoted by $(\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa})$.

Let \mathcal{I} be such an information state, then $(\Pi, \boldsymbol{\mu})$ is called a *partial knowledge state* compatible with \mathcal{I} . Note that a partial knowledge state may be compatible with various information states.

Within an information state, the user partitions the population into a number of subpopulations. He knows the probability distribution of each subpopulation on the confidential attribute values. Intuitively, $\mu_i(\varphi)$ is the proportion of the individuals in sub-population $|\varphi|_T$ which have confidential attribute value s_i , whereas $\kappa_i(\varphi)$ is the number of investigated individuals in sub-population $|\varphi|_T$ which have confidential attribute value s_i . Since each DL wff φ is composed from atomic formulas with easy-to-know attributes only, it can be assumed that it takes little effort for the user to verify whether a given individual satisfies φ . Furthermore, it can also be assumed that the cardinality of the truth set of each φ is known to the public. However, note that it may sometimes be very difficult for the user to locate an individual satisfying a specific φ from the whole population U .

The information states of an user can be subsequently changed by his investigation of some individuals in a specific sup-population and by his queries posed to and the answers obtained from the data center. This is a process of knowledge refinement and can be modeled by the knowledge acquisition actions as follows.

A logical partition Π_2 is a refinement of another logical partition Π_1 , denoted by $\Pi_2 \sqsubseteq \Pi_1$, if for all $\varphi_2 \in \Pi_2$, there exists $\varphi_1 \in \Pi_1$ such that $|\varphi_2|_T \subseteq |\varphi_1|_T$. It is clear that if $\Pi_2 \sqsubseteq \Pi_1$, then each $|\varphi_1|_T$ such that $\varphi_1 \in \Pi_1$ can be written as a union of the truth sets of some wffs in Π_2 .

Definition 3 Let $\mathcal{I}_1 = (\Pi_1, \boldsymbol{\mu}_1, \boldsymbol{\kappa}_1)$ and $\mathcal{I}_2 = (\Pi_2, \boldsymbol{\mu}_2, \boldsymbol{\kappa}_2)$ be two information states. \mathcal{I}_2 is a refinement of \mathcal{I}_1 , also denoted by $\mathcal{I}_2 \sqsubseteq \mathcal{I}_1$, if both of the following conditions are satisfied:

1. $\Pi_2 \sqsubseteq \Pi_1$.
2. For each $\varphi \in \Pi_1$, if $|\varphi|_T = \bigcup_{1 \leq i \leq l} |\varphi_i|_T$ for some set $\{\varphi_1, \dots, \varphi_l\} \subseteq \Pi_2$, then

$$|\varphi| \cdot \boldsymbol{\mu}_1(\varphi) = \sum_{i=1}^l |\varphi_i| \cdot \boldsymbol{\mu}_2(\varphi_i),$$

and

$$\boldsymbol{\kappa}_1(\varphi) \leq \sum_{i=1}^l \boldsymbol{\kappa}_2(\varphi_i).$$

Note that the arithmetics (addition and multiplication) and comparison between vectors (and scalars) are defined as usual. For example, the addition of two vectors is carried out point-wise and results in a vector of the same dimension.

In our framework, there are two kinds of knowledge acquisition actions which can refine the user's information states. The first one is the query action. Each query action is represented by

a wff φ in DL. The intended answer of the query is the distribution of the confidential values within the selected population $|\varphi|_T$ in the database. The other is the investigation action, which is specified by a wff φ and a positive integer number k . This means that the user have investigated k individuals from the set $|\varphi|_T$ in this action. For the uniformity of the representation, each knowledge acquisition action is written as $\alpha = (\varphi, k)$ for some DL wff φ and $k \geq 0$. When $k > 0$, it is an investigation action, whereas it is a query one if $k = 0$.

Definition 4 1. A knowledge acquisition action $(\varphi, 0)$ is applicable under the information state $\mathcal{I}_1 = (\Pi_1, \mu_1, \kappa_1)$ and results in a state $\mathcal{I}_2 = (\Pi_2, \mu_2, \kappa_2)$ if

- (a) there exists $\varphi' \in \Pi_1$ such that $|\varphi|_T \subseteq |\varphi'|_T$,
- (b) $\Pi_2 = \Pi_1 - \{\varphi'\} \cup \{\varphi, \varphi' \wedge \neg\varphi\}$,
- (c) \mathcal{I}_2 is a refinement of \mathcal{I}_1 ,
- (d) $\kappa_2(\psi) = \kappa_1(\psi)$ for any $\psi \in \Pi_1 - \{\varphi'\}$, and
- (e) $\kappa_2(\varphi) + \kappa_2(\varphi' \wedge \neg\varphi) = \kappa_1(\varphi')$.

2. A knowledge acquisition action (φ, k) where $k > 0$ is applicable under the information state $\mathcal{I}_1 = (\Pi_1, \mu_1, \kappa_1)$, and $\mathcal{I}_2 = (\Pi_2, \mu_2, \kappa_2)$ is a resultant state of the application if

- (a) $\varphi \in \Pi_1$ and $k \leq |\varphi| - \sum_{i=0}^{t-1} \kappa_i(\varphi)$
- (b) $\Pi_1 = \Pi_2$,
- (c) $\mu_1 = \mu_2$,
- (d) $\kappa_2(\psi) = \kappa_1(\psi)$ for any $\psi \neq \varphi$, and
- (e) $\sum_{i=0}^{t-1} \kappa_{2i}(\varphi) = \sum_{i=0}^{t-1} \kappa_{1i}(\varphi) + k$.

Since the goal of the user is to refine his knowledge by the queries, a rational user would pose his queries so that his knowledge would be improved by the answers of the queries. Thus if the user's information state is (Π_1, μ_1, κ_1) , then he poses a query about a subset of an equivalence class in Π_1 . This is the requirement of Condition 1a in Definition 4. Then, after the query is answered, the corresponding equivalence class is partitioned into two parts — one satisfying φ and the other not, so we have the Condition 1b in Definition 4. Condition 1c in Definition 4 further requires that the answer is correct so that the resultant information state is a refinement of the original one. Furthermore, since the query action does not cause any new individuals being investigated, the κ_2 function agrees with κ_1 in the part of the population which is not split by the query, while for the split part, the number of investigated individuals is not changed in total. This is reflected respectively in Conditions 1d and 1e of the definition.

In the case of investigation action, we assume the user will only investigate the individuals in a sub-population represented by a wff in Π_1 . The assumption is inessential, since, if the investigated individuals are across some different sub-populations, the corresponding investigation action can be decomposed into a sequence of actions satisfying the applicability condition. Since it is assumed that

the user knows the total number of individuals in $|\varphi|_T$ and those which have been investigated by him so far is equal to $\sum_{i=0}^{t-1} \kappa_i(\varphi)$, he would not try to investigate more individuals than all remaining ones. This is exactly required by the applicability condition of Definition 4.2a. Conditions 2b to 2d are obvious since these values are not affected by the investigation. What the investigation can affect is the total number of the investigated individuals in $|\varphi|_T$ and this is reflected in Condition 2e.

4 The Value of Information

To quantitatively determine the value of information, we must have a user model. Let us consider the case where the user is an agent trying to use the private information to aid his decision in a game. The game is played between the agent and individuals in the population U . The agent can decide the rate he wants to charge an individual for playing the game (i.e., the admission fee). The rate is decided on a personalized basis so that each individual may be charged with different rates. However, once an individual agrees to play the game with the agent and pay the fee asked by the agent, he will have a chance to get back some reward which will be the loss of the agent. The reward of an individual is determined by his confidential attribute value. Let r_i denote the reward of an individual with the confidential attribute value s_i for $0 \leq i \leq t-1$, then $\boldsymbol{\rho} = (r_0, r_1, \dots, r_{t-1}) \in \mathfrak{R}^t$ is called the loss vector of the agent.

Let $\mathcal{I}_0 = (\{\top\}, \boldsymbol{\mu}_0, \boldsymbol{\kappa}_0)$ be the initial information state of the user, where \top denotes any tautology in the DL and $\boldsymbol{\kappa}_0(\varphi) = (0, \dots, 0)$. Let $\boldsymbol{\rho}$ be a given loss vector. The agent first decides the *base rate* of the game on the expected loss according to his initial information state, i.e., $R_0 = \boldsymbol{\rho} \cdot \boldsymbol{\mu}_0(\top)$. Thus, in the initial state, the expected payoff of the agent for playing the game is zero. However, once he acquires pieces of information and reaches a new information state, he can utilize the acquired information for making some profit.

We further assume that each individual will go into the game if he is charged with the base rate. However, he can refuse to do so if the agent charges him with a rate higher than the base one. The higher the rate, the more likely the individual refuses to play the game. If the information state is $\mathcal{I} = (\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa})$, where $\Pi = \{\varphi_1, \dots, \varphi_m\}$, a reasonable decision of the agent for the rate of an individual u satisfying φ is as follows:

1. if u has been investigated and it is known that the confidential attribute value of u is s_i , then the most profitable decision of the agent would be to charge the individual with $\max(R_0, r_i)$ so that the agent's payoff is $\max(R_0 - r_i, 0)$;
2. if the individual has not been investigated, the agent knows the probability of the confidential attribute value of u being s_i to be

$$p_i(\varphi) = \frac{|\varphi| \cdot \mu_i(\varphi) - \kappa_i(\varphi)}{|\varphi| - \sum_{i=0}^{t-1} \kappa_i(\varphi)}. \quad (1)$$

In this case, the most reasonable decision of the agent would be to charge the individual with

$\max(R_o, \sum_{i=0}^{t-1} p_i(\varphi) \cdot r_i)$ so that the agent's expected payoff would be $\max(R_o - \sum_{i=0}^{t-1} p_i(\varphi) \cdot r_i, 0)$

Thus, in average, the agent can have the following expected payoff B_φ in playing the game with an individual satisfying φ :

$$B_\varphi = \max(R_o - \sum_{i=0}^{t-1} (p_i(\varphi) \cdot r_i), 0) \cdot \frac{|\varphi| - \sum_{i=0}^{t-1} \kappa_i(\varphi)}{|\varphi|} + \sum_{i=0}^{t-1} \max(R_o - r_i, 0) \cdot \frac{\kappa_i(\varphi)}{|\varphi|} \quad (2)$$

Thus, by using the knowledge about the individuals' confidential attributes, the agent can raise the rates of those who may incur a greater loss to him in order to avoid the possible loss. The value of the information is then dependent on how much he can benefit from obtaining the information. The expected gain of the agent with regard to each individual is computed by

$$B_{\mathcal{I}} = \sum_{\varphi \in \Pi} B_\varphi \cdot \frac{|\varphi|}{|U|},$$

if he decides the rates according to the two principles above.

Example 1 The scenario described above usually occurs between an insurance company and its customers. The base rate is applied to a typical customer if the company does not have any further information about his health condition. However, for the customers of high risk, the company would raise their rates. Thus the health information of the customers would be valuable to the insurance company. To avoid the leakage of privacy, the data center may correspondingly raise the cost of answering a query so that the information value for the company is counter-balanced. The company would not have the incentive to obtain the private information. ■

The notions of the value of information have been extensively studied in decision theory[4, 10]. In our model above, if investigation actions are not allowed, all information states are of the form $(\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa}_0)$, so $\kappa_{0i}(\varphi) = 0$ and $p_i(\varphi) = \mu_i(\varphi)$ for all $0 \leq i \leq t-1$ and $\varphi \in \Pi$. Consequently, $B_{\mathcal{I}}$ would be simplified into

$$\sum_{\varphi \in \Pi} \max(R_o - \boldsymbol{\mu}(\varphi) \cdot \boldsymbol{\rho}, 0) \cdot \frac{|\varphi|}{|U|}$$

which is the value of partial information defined in [10] if our user model is appropriately formulated as a decision problem of the agent. While in our case the partial information is obtained by querying the data center, another approach for obtaining partial information by sampling is suggested in [10]. Though sampling is similar to investigation, the information obtained from these two kinds of actions are quite different. For the sampling actions, even though the chosen individuals may be thoroughly investigated, only the statistical information of these investigated individuals would be kept. In fact, it is the statistical information which would be used in the prediction of the status of the whole population. However, for the investigative actions, the user would indeed keep the

personal information of each investigated individual and not do the statistical inference from the investigated individuals to the whole population.

On the other hand, if no query actions are possible, the information states are always of the form $(\{\top\}, \boldsymbol{\mu}_0, \boldsymbol{\kappa})$. Once all individuals have been fully investigated (though this is hardly possible in any practical case) the information state becomes a perfect state $\mathcal{I} = (\{\top\}, \boldsymbol{\mu}_0, \boldsymbol{\kappa})$, where $\kappa_i(\top) = \mu_{0i}(\top) \cdot |U|$, so $p_i(\top) = 0$ for all $0 \leq i \leq t-1$. Consequently, $B_{\mathcal{I}}$ would be simplified into

$$\sum_{i=0}^{t-1} \max(R_o - r_i, 0) \cdot \mu_{0i}(\top)$$

which is precisely the value of perfect information defined in [10]. Thus we have modeled the value of hybrid information in the above-defined framework.

4.1 Privacy protection by pricing mechanism

According to the user model above, the user can improve his payoff from 0 to $B_{\mathcal{I}}$ when his information state is evolved from the initial state to \mathcal{I} . If the information is free of charge, the user would gladly receive it and consequently, the privacy of the individuals may be invaded. Thus, one approach to privacy protection is to impose costs on the answers of the queries so that the user cannot make a profit from obtaining the private information. This can be achieved by including a pricing mechanism in the data center. However, since the answer to a query may have different effects under different information states, the pricing mechanism must be adaptive according to the query history of the user. In general, it is very difficult to design an adaptive pricing mechanism since the users may have to pay different prices for the same queries under different situations. Therefore, instead of charging each query separately, we shall consider a more restricted setting. Assume that each user is allowed to ask a batch of queries only once. Afterward, he can do any investigative actions he wants. However, the data center would not answer his queries afterwards. Thus the pricing mechanism of the data center is to decide the cost of each batch of queries so that the user cannot benefit from receiving the answers of the queries.

Let $(\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa})$ be the information state of the user after a sequence of queries and follow-up investigative actions, where $\Pi = \{\varphi_1, \varphi_2, \dots, \varphi_m\}$. Since the data center has no control on how the user will carry out his investigation after receiving the answers, it can only guarantee that the cost is high enough so that the user cannot make a profit from the answers of the queries, no matter what investigation be done. Thus, based only on the partial knowledge state $\mathcal{P} = (\Pi, \boldsymbol{\mu})$, the data center must estimate the maximum payoff the agent can have under different information states compatible with \mathcal{P} . Let $\mathbf{k} = (k_1, \dots, k_m)$ be an m -tuple of non-negative integers and define

$$F_{\mathbf{k}} = \{\boldsymbol{\kappa} \mid \sum_{i=0}^{t-1} \kappa_i(\varphi_j) = k_j, \forall 1 \leq j \leq m\}$$

as the set of $\boldsymbol{\kappa}$ functions which denote the possible investigation results when a specific number of individuals has been investigated. The set of information states compatible with \mathcal{P} and \mathbf{k} is defined

as

$$\mathcal{IS}(\mathcal{P}, \mathbf{k}) = \{(\mathcal{P}, \boldsymbol{\kappa}) \mid \boldsymbol{\kappa} \in F_{\mathbf{k}}\}$$

and the maximal value of information of the agent under \mathcal{P} and \mathbf{k} is defined as

$$B(\mathcal{P}, \mathbf{k}) = \max\{B_{\mathcal{I}} \mid \mathcal{I} \in \mathcal{IS}(\mathcal{P}, \mathbf{k})\}.$$

We now further assume that a cost function $\gamma_{inv} : \Phi \times \mathcal{Z}^+ \rightarrow \Re^+$ is available to both the user and the data center, where Φ is the set of DL wffs and \mathcal{Z}^+ and \Re^+ are respectively the set of positive integer and real numbers. The intended meaning of $\gamma_{inv}(\varphi, k)$ is the cost of the investigation of k individuals satisfying φ . It can be assumed that γ_{inv} is a super-linear function in its second argument. Thus, when the user poses a batch of queries Q , the data center can know what the resultant partial knowledge state \mathcal{P} would be once the answer is released. Therefore, the price of Q must be decided before releasing the information. The price $price(Q)$ of the answers to the batch of queries should be decided such that

$$|U| \cdot B(\mathcal{P}, \mathbf{k}) - \sum_{i=1}^m \gamma_{inv}(\varphi_i, k_i) \leq price(Q) \quad (3)$$

holds for any \mathbf{k} . The lowest solution of $price(Q)$ for (3) is

$$\max_{\mathbf{k}} |U| \cdot \max\{B_{\mathcal{I}} \mid \mathcal{I} \in \mathcal{IS}(\mathcal{P}, \mathbf{k})\} - \sum_{i=1}^m \gamma_{inv}(\varphi_i, k_i) \quad (4)$$

where the domain of \mathbf{k} is finite since $0 \leq k_i \leq |\varphi_i|$.

4.2 Usefulness of Information

In our pricing mechanism, the data center assumes that the user can play the above-mentioned game with all individuals in U and charge them based on the total gain he can achieve. However, this may be an over-estimation since the user cannot play the game with all individuals when the population is large. To circumvent the problem, we may assume that the user must spend some resources for playing the game with the individuals. Let $\gamma_{ply} : \Phi \times \mathcal{Z}^+ \rightarrow \Re^+$ be another cost function such that $\gamma_{ply}(\varphi, l)$ denotes the cost of the user playing the game with l individuals satisfying φ . Given an m -tuple of non-negative integers $\mathbf{l} = (l_1, \dots, l_m)$ and an information state \mathcal{I} , define

$$B_{\mathcal{I}}^{\mathbf{l}} = \sum_{i=1}^m B_{\varphi_i} \cdot l_i.$$

The price in (4) can be replaced by

$$\max_{\mathbf{k}, \mathbf{l}} (\max\{B_{\mathcal{I}}^{\mathbf{l}} \mid \mathcal{I} \in \mathcal{IS}(\mathcal{P}, \mathbf{k})\} - \sum_{i=1}^m \gamma_{inv}(\varphi_i, k_i) - \sum_{i=1}^m \gamma_{ply}(\varphi_i, l_i)) \quad (5)$$

where both the domains of \mathbf{k} and \mathbf{l} are restricted to $0 \leq k_i, l_i \leq |\varphi_i|$.

Intuitively, each l_i and k_j represent the *usefulness* of information. Given two equivalent classes in a logical partition, it may be easier to find potential members in one equivalence class than in the other depending on the conditions each equivalence class satisfied. It may also be true that it is easier, and thus cost-effective, to investigate members in one equivalence class than in the other. These two may be closely related, but not necessarily the same.

Example 2 Assume we again use the insurance company model mentioned in Example 1. Assume the world population is represented by all adults in the country. An equivalence class may be characterized as being the people living in the same county while another equivalence class is described as the people with weight between 60 to 65 kilograms. It is easy for the first group of people to be investigated and then to be added as customers, while it is relatively difficult for the second group of people. ■

Thus the data center can decide the price of the answers to the batch of queries Q by a two-level maximization procedure in (4) or (5). The outer level maximization would depend on the form of the cost functions γ_{inv} and/or γ_{ply} , so it is unlikely to find an analytic solution for it. However, the inner maximization can be reduced to a set of m maximization of B_φ for each $\varphi \in \Pi$. More specifically, given φ and $0 \leq k \leq |\varphi|$, it is to find $\kappa(\varphi)$ which maximizes B_φ among all κ satisfying $\sum_{i=0}^{t-1} \kappa_i(\varphi) = k$ and $\kappa_i(\varphi) \leq |\varphi| \cdot \mu_i(\varphi)$ for all $0 \leq i \leq t-1$. This is in turn equivalent to the following constraint optimization problem in the integer domain:

$$\begin{aligned} \text{Maximize} \quad & \max(R_0 - \sum_{i=0}^{t-1} \frac{n_i - x_i}{N - k} \cdot r_i, 0) \cdot \frac{N - k}{N} + \sum_{i=0}^{t-1} \max(R_0 - r_i, 0) \cdot \frac{x_i}{N} \\ & \text{s.t.} \\ & x_0 + x_1 + \dots + x_{t-1} = k \\ & 0 \leq x_i \leq n_i \quad (0 \leq i \leq t-1) \end{aligned} \tag{6}$$

where N and n_i 's correspond to $|\varphi|$ and $|\varphi| \cdot \mu_i(\varphi)$'s respectively. The solution of Equation (6) can be given by the following proposition for $k \leq N$. Without loss of generality, we assume $r_0 \geq r_1 \geq \dots \geq r_{t-1}$ for the loss vector in the proposition.

Proposition 1 Assume $N = \sum_{i=0}^{t-1} n_i$

1. if $k = N$, then the solution of Equation (6) is $x_i = n_i$ for $0 \leq i \leq t-1$ and its maximum value is

$$\sum_{i=0}^{t-1} \max(R_0 - r_i, 0) \cdot \frac{n_i}{N};$$

2. if $k < N$ and l is the smallest natural number such that $\sum_{i=0}^l n_i > k$, then the solution of Equation (6) is

$$x_i = \begin{cases} n_i & \text{if } i < l, \\ k - \sum_{i=0}^{l-1} n_i & \text{if } i = l, \\ 0 & \text{if } i > l, \end{cases}$$

and its maximum value is

$$\max(R_0 - \sum_{i=l+1}^{t-1} \frac{n_i}{N-k} \cdot r_i + \frac{\sum_{i=0}^l n_i - k}{N-k} \cdot r_l, 0) \cdot \frac{N-k}{N} + \sum_{i=0}^{l-1} \max(R_0 - r_i, 0) \cdot \frac{n_i}{N} + \max(R_0 - r_l, 0) \cdot \frac{k - \sum_{i=0}^{l-1} n_i}{N}.$$

The individuals who will incur more loss to the agent are high risk ones. For the low risk individuals, the investigation cannot improve the payoff for the agent. However, for the high risk ones, the investigation can indeed decrease the loss for the agent by raising their admission fees appropriately. The more high risk individuals have been investigated, the more loss the agent can avoid, so the maximum payoff occurs when the investigation is carried out from the most risky individuals to the least risky ones. This intuition is verified by the preceding proposition.

5 The Damage and Compensation of Privacy Leakage

To model the damage and compensation of privacy leakage, we consider another game. The game is played between an agent, called the accuser, and the individuals in U . The accuser try to disseminate the private information of the individuals. Assume $(d_0, d_1, \dots, d_{t-1}) \in \mathbb{R}^t$ and $(c_0, c_1, \dots, c_{t-1}) \in \mathbb{R}^t$ are respectively the damage and compensation vectors of the game. The rule of the game is as follows: If an individual is accused of s_i and he actually has the attribute value s_i , then his damage is d_i which is also the reward of the accuser. However, if he is accused of s_i and his private attribute value is not s_i , then he can receive the compensation c_i from the accuser. Thus if $\mathcal{I} = (\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa})$ is an information state, then the agent who wants to accuse an un-investigated individual satisfying $\varphi \in \Pi$ of s_i would have the risk of losing

$$L_i(\varphi) = (1 - p_i(\varphi))c_i - p_i(\varphi)d_i \quad (7)$$

where $p_i(\varphi)$ is defined by (1). The task of privacy protection then is to make the dissemination of the individuals' privacy unprofitable for the accuser. This is done by raising his expected loss to a threshold level. The threshold level should be high enough so that any rational agent would not risk such a loss. However, for the ease of presentation, we assume the threshold is zero. Thus, an information state $\mathcal{I} = (\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa})$ is said to be *safe* if $L_i(\varphi) \geq 0$ for all $\varphi \in \Pi$ and $0 \leq i \leq t-1$.

Example 3 Assume a person being tested for faulty genes. These faulty genes increase the chance of acquiring some rare disease. This person's job application may be rejected since the company may feel he will be more likely to become ill and have a poor performance. By using this information, the damage can be caused to this person. ■

5.1 The basic model

The query of the agent can be answered only if his information state will be safe after receiving the answer. It can be seen that an information state $\mathcal{I} = (\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa})$ is safe if $p_i(\varphi) \leq \frac{c_i}{c_i + d_i}$ for any

$\varphi \in \Pi$ and $0 \leq i \leq t - 1$. However, since $p_i(\varphi)$ not only depends on $\mu_i(\varphi)$, but also on how many individuals have been investigated by the user, the data center cannot decide whether answering a query will lead to a safe state or not. To guarantee the safety of an information state, the data center can use the worst-case analysis. Assume for each wff φ , the user can investigate at most K_φ individuals in $|\varphi|$ at affordable cost. Then, given a partial knowledge state $\mathcal{P} = (\Pi, \boldsymbol{\mu})$ which results from the answer of a query, the data center can guarantee the safety, no matter which (affordable) investigation is made by the user, if the following condition holds for all $\varphi \in \Pi$ and $0 \leq i \leq t - 1$:

$$\frac{|\varphi| \cdot \mu_i(\varphi)}{|\varphi| - K_\varphi} \leq \frac{c_i}{c_i + d_i}, \quad (8)$$

since by (1),

$$p_i(\varphi) \leq \frac{|\varphi| \cdot \mu_i(\varphi)}{|\varphi| - K_\varphi}.$$

The condition (8) can be rewritten as

$$\mu_i(\varphi) \leq \frac{c_i}{c_i + d_i} \cdot \left(1 - \frac{K_\varphi}{|\varphi|}\right). \quad (9)$$

Let us consider some cases where Equation (9) can be satisfied.

1. If no investigative actions are possible, i.e., $K_\varphi = 0$, then (9) is satisfied if $\mu_i(\varphi) \leq \frac{c_i}{c_i + d_i}$. In such case, if $d_i = 0$ or $c_i \gg d_i$, then the information state is still safe even though $\mu_i(\varphi)$ is approximately equal to 1. This means that knowing that an individual is s_i will not harm that individual or the compensation will be sufficiently large to cover his damage. Hence it does not matter if the user almost certainly discover that a class of individuals has s_i value. On the other hand, if $c_i \approx d_i > 0$, then the information state is safe only when $\mu_i(\varphi)$ is less than 0.5. In other words, if the compensation cannot cover the damage sufficiently, then the user should not know the confidential value with a certainty above 0.5.
2. If investigation is allowed for at most K_φ individuals, then the upper bound of $\mu_i(\varphi)$ is discounted with the ratio $1 - \frac{K_\varphi}{|\varphi|}$ to maintain safety. The discount effect is alleviated when $|\varphi| \gg K_\varphi$. Thus, the larger the size $|\varphi|$, the more the possibility of achieving the safety requirement. This corresponds to the k -anonymity requirement for privacy protection in [14].

Based on the safety criterion, the data center can decide whether the user's query should be answered or refused. However, note that (9) is only a sufficient condition for the safety of data release, so we may not have to test it for every i and φ . In particular, if $d_i = 0$, then $L_i(\varphi) \geq 0$ holds no matter how the investigative actions are carried out, so we only have to test (9) for those i 's such that $d_i > 0$.

An alternative approach is to use the pricing mechanism to discourage the user. To formulate the approach, we need another cost function $\gamma_{acc} : \Phi \times \mathcal{Z}^+ \rightarrow \mathfrak{R}^+$ where $\gamma_{acc}(\varphi, k)$ denotes the cost

the user has to spend for accusing k individuals satisfying φ . The minimum loss the user may incur under the partial knowledge state $\mathcal{P} = (\Pi, \boldsymbol{\mu})$ should then be

$$L^*(\varphi) = \min_{\mathbf{k}, \mathbf{l}} \left[\sum_{i=0}^{t-1} L_i(\varphi, \mathbf{k}) \cdot l_i + \gamma_{inv}(\varphi, \sum_{i=0}^{t-1} k_i) + \gamma_{acc}(\varphi, \sum_{i=0}^{t-1} l_i) \right]$$

where

$$L_i(\varphi, \mathbf{k}) = c_i - (c_i + d_i) \cdot \frac{|\varphi| \cdot \mu_i(\varphi) - k_i}{|\varphi| - \sum_{i=0}^{t-1} k_i}$$

is the result of substituting (1) into (7) when $\kappa_i(\varphi) = k_i$ for $0 \leq i \leq t-1$ and the minimization has taken over all $k_i \leq |\varphi| \cdot \mu_i(\varphi)$ for $0 \leq i \leq t-1$ and l_i such that $\sum_{i=0}^{t-1} l_i \leq |\varphi|$. Then, if the answer to a batch of queries results in a partial knowledge state $(\Pi, \boldsymbol{\mu})$, its price should be determined by $\sum_{\varphi \in \Pi} price(\varphi)$, where price of each φ is in accordance with the following equation:

$$price(\varphi) = \begin{cases} -L^*(\varphi), & \text{if } L^*(\varphi) < 0 \\ 0, & \text{otherwise.} \end{cases}$$

5.2 The extended model

In the basic model for the damage and compensation of privacy leakage, we assume the damage vector is associated with each specific value of the confidential attribute. This means that if an individual is known to have the attribute value s_i , then he will have damage d_i . However, sometimes it is also harmful if an individual is known to have his attribute value in some specific subset of V_c even if the subset is not a singleton.

Example 4 Assume a fatal disease can be diagnosed and classified as a stage 0–5, where 0 means no disease, 1–3 can be cured, and 4–5 are deadly. Then knowing that a person was diagnosed in stage 4 or 5 can be harmful to that person. ■

The basic model should be extended to also cover such case. Since it is reasonable that the compensation should proportionally depend on the damage, we can simplify the model by assuming that there is a function $\alpha : \mathfrak{R} \rightarrow \mathfrak{R}$ mapping each damage value to its corresponding compensation. For example, it may be the case that $\alpha(x) = r \cdot x$ for some positive number r . Thus we can concentrate on the estimate of damage in the extended model. In the model, we assume there is a damage function $\delta : (2^{\{0, \dots, t-1\}} - \{\emptyset\}) \rightarrow \mathfrak{R}$. For any $S \subseteq \{0, \dots, t-1\}$, $\delta(S)$ means the damage of an individual when it is known that his confidential attribute value belongs to $\{s_i \mid i \in S\}$. By using the same game rule as in the basic model, the expected loss of the agent accusing an individual in $\varphi \in \Pi$ of $\{s_i \mid i \in S\}$ would be

$$L_S(\varphi) = (1 - \sum_{i \in S} p_i(\varphi)) \alpha(\delta(S)) - (\sum_{i \in S} p_i(\varphi)) \delta(S)$$

where $p_i(\varphi)$ is defined by (1). Then the safety criterion for an information $\mathcal{I} = (\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa})$ is extended to

$$L_S(\varphi) \geq 0$$

for all $\varphi \in \Pi$ and $S \subseteq \{0, \dots, t-1\}$. This is equivalent to

$$\sum_{i \in S} p_i(\varphi) \leq \frac{\alpha(\delta(S))}{\alpha(\delta(S)) + \delta(S)}. \quad (10)$$

By using the same worst-case analysis as in the basic model, the following must be satisfied for all $\varphi \in \Pi$ and $S \subseteq \{0, \dots, t-1\}$,

$$\sum_{i \in S} \frac{|\varphi| \cdot \mu_i(\varphi)}{|\varphi| - K_\varphi} \leq \frac{\alpha(\delta(S))}{\alpha(\delta(S)) + \delta(S)},$$

or alternatively in the following form:

$$\sum_{i \in S} \mu_i(\varphi) \leq \left(\frac{\alpha(\delta(S))}{\alpha(\delta(S)) + \delta(S)} \right) \cdot \left(1 - \frac{K_\varphi}{|\varphi|} \right). \quad (11)$$

So far, the model does not say anything on estimating the damage function δ . In fact, the damage vector in the basic model should be determined by some external mechanism, such as the legal system or social convention, so we can assume that $\delta(\{i\}) = d_i$ for each $0 \leq i \leq t-1$. However, for other values of the damage function, it should be possible to impose some reasonable constraints so that the damage of a subset S can be (partially) estimated by those of its elements.

Example 5 We list below some possible conditions that a damage function $\delta : (2^{\{0, \dots, t-1\}} - \{\emptyset\}) \rightarrow R$ should satisfy.

1. $\delta(\{0, \dots, t-1\}) = 0$.
2. $\delta(S_1) \leq \delta(S_2)$ if $S_2 \subseteq S_1$.
3. $\delta(S) = 0$ if $|S| > 1$.
4. $\delta(S) = \min_{i \in S} \delta(\{i\})$.

Condition 1 means that if no privacy leakage, no damage. Since it is known that all possible values of the confidential attribute are in V_c , the index set $\{0, \dots, t-1\}$ corresponds to the situation of no privacy leakage. Condition 2 means the more specific information is known, the more damage is caused. Condition 3 corresponds to the basic model in which only the damage value of the singleton is considered. Condition 4 is due to the principle of least commitment. The principle implies that if an individual is accused of a set of possible faults disjunctively, it can only be sure that he has the least harmful fault, so that the damage to him caused by such accusation would be equivalent

to the minimal one in accusing him of a specific fault in the set. Note that Conditions 3 and 4 are not compatible if there are at least two indices i and j such that $\delta(\{i\}) > 0$ and $\delta(\{j\}) > 0$. However, both Conditions 1 and 2 are implied by Conditions 3 and Condition 4 implies Condition 2. Furthermore, Condition 4 also implies Condition 1 provided that i exists such that $\delta(\{i\}) = 0$. ■

An alternative way to estimate the damage value of a subset is by the information-theoretic approach. In other words, if the *a priori* probability function on the possible values of the confidential attribute is given by $\mu(\top)$, then we can compute the *a posteriori* probability for any $S \subseteq \{0, 1, \dots, t-1\}$ as

$$Pr(s_i|S) = \begin{cases} \frac{\mu_i(\top)}{\sum_{j \in S} \mu_j(\top)}, & \text{if } i \in S; \\ 0, & \text{otherwise.} \end{cases}$$

Then a possible constraint on the damage function is

$$\delta(S) = \sum_{i \in S} m(\mu_i(\top), Pr(s_i|S)) \cdot \delta(\{i\}),$$

where $m : [0, 1] \times [0, 1] \rightarrow [0, 1]$ is called an information distance function. The information distance function estimates how the user's information on some specific s_i is increased by knowing its index being in S . Typically, the information distance function can be defined as the relative difference between the entropy values of the two probabilities, i.e.,

$$m(p, q) = \frac{\log p - \log q}{\log p}.$$

6 Related Works

To quantify the value of information is by no means a novel problem. However, the quantitative models for privacy protection provides a new angle to look at the problem. As shown in section 4, our model for the value of information has generalized a standard notion in decision theory[10, 4]. While the decision-theoretic analysis [10] emphasizes the value of information from the decision maker's viewpoint, our model is mainly concerned with privacy protection by the information provider. For the former, a decision maker can decide if he will purchase a piece of information according to the value of the information. For the latter, the information provider can charge the user of the information with appropriate rates.

An alternative model for the value of information in the privacy protection context is proposed in [2, 3]. In their model, the value of information is estimated by the expected cost the user must pay for achieving the perfect knowledge state from the given information. More specifically, let a piece of information be a pair $(\varphi, (n_i)_{0 \leq i \leq t-1})$ such that φ is a DL wff and $|\varphi| = \sum_{i=0}^{t-1} n_i$. The meaning of the information $(\varphi, (n_i)_{0 \leq i \leq t-1})$ is that in the subpopulation $|\varphi|_T$, there are exactly n_i individuals with confidential attribute value s_i for $0 \leq i \leq t-1$. To know the confidential values of all individuals satisfying φ , the user can investigate them one by one. Assume the

investigation cost for each individual is fixed, then the total cost is proportional to the number of individuals he must investigate to know all individuals' confidential values. Though the maximum number of investigative actions the user must carry out may be equal to $|\varphi|$, a piece of information $(\varphi, (n_i)_{0 \leq i \leq t-1})$ may help him in reducing the investigation cost. An extreme case is when $n_i = |\varphi|$ and $n_j = 0$ for all $1 \leq j \neq i \leq t-1$. In such case, the user does not have to carry out any investigation since he knows that all individuals have the confidential value s_i . Though the model of [2, 3] is restricted to the case where the domain of values for the confidential attribute has exactly two elements, the result can be generalized according to the following proposition.

Proposition 2 *Given a piece of information $(\varphi, (n_i)_{0 \leq i \leq t-1})$, the expected number of investigative actions for knowing the confidential values of all individuals satisfying φ is*

$$\sum_{i=0}^{t-1} \frac{n_i \cdot (|\varphi| - n_i)}{|\varphi| - n_i + 1} \quad (12)$$

The value of information by equation (12) is based on the rationale that the more investigation efforts a piece of information can reduce, the more valuable it is. Thus the value of information can be defined as a monotonically decreasing function of (12). However, without regarding the user model, the value of information based on (12) may not reflect the real situation. For example, equation (12) is obviously invariant under the permutation of n_i 's. In other words, if $\sigma : \{0, \dots, t-1\} \rightarrow \{0, \dots, t-1\}$ is a permutation on the index set, then

$$\sum_{i=0}^{t-1} \frac{n_i \cdot (|\varphi| - n_i)}{|\varphi| - n_i + 1} = \sum_{i=0}^{t-1} \frac{n_{\sigma(i)} \cdot (|\varphi| - n_{\sigma(i)})}{|\varphi| - n_{\sigma(i)} + 1}.$$

This means that two pieces of information $(\varphi_1, (n_i)_{0 \leq i \leq t-1})$ and $(\varphi_2, (n_{\sigma(i)})_{0 \leq i \leq t-1})$ have the same value if $|\varphi_1| = |\varphi_2|$. However, in practice, knowing that most individuals have confidential value s_0 may be of a different value than knowing that the same number of individuals have confidential value s_1 .

Example 6 Let the confidential attribute denote the HIV test result and $V_c = \{+, -\}$ (i.e. $t = 2$), then, given two sub-populations of size 10000 characterized by φ_1 and φ_2 respectively, two pieces of information $(\varphi_1, (9999, 1))$ and $(\varphi_2, (1, 9999))$ should have different values from the viewpoint of privacy protection since the former says that most individuals in the subpopulation are infected whereas the latter tells the contrary. On the other hand, our model can deal with the problem by imposing different loss values r_0 and r_1 on the two test results. ■

Besides the decision theoretic analysis, the value of information can also be estimated by some information theoretic measures. The central notion of such measures is the entropy introduced by Shannon[15]. In the machine learning literatures, it is used to define the information gain of an

attribute for a classification problem[11]. By reformulating the notions in our context, the entropy of a partial knowledge state $\mathcal{P} = (\Pi, \boldsymbol{\mu})$ is defined by

$$H(\mathcal{P}) = \sum_{\varphi \in \Pi} \frac{|\varphi|}{|U|} \cdot \sum_{i=0}^{t-1} \mu_i(\varphi) \cdot \frac{1}{\log \mu_i(\varphi)}$$

The information gain of a partial knowledge state is defined as the difference between its entropy and that of the initial information state, i.e.

$$Gain(\mathcal{P}) = H(\mathcal{P}_0) - H(\mathcal{P}),$$

where $\mathcal{P}_0 = (\{\top\}, \boldsymbol{\mu}_0)$ is the partial knowledge state compatible with the initial information state. Though the information gain is an useful index in selecting the most informative features for the classification problem, it still suffers the same problem as the value of information based on (12) since it does not take into account the fact that some confidential attribute values are more sensitive than others.

In contrast with the quantitative approach of this paper, some qualitative criteria for privacy protection have been proposed in [7, 8, 13, 14, 16]. These criteria are designed to protect personal sensitive information in the release of a microdata set, i.e. a set of records containing information on individuals. The main objective is to avoid the re-identification of individuals or in other words, to prevent the possibility of deducing which record corresponds to a particular individual even though the explicit identifier of the individual is not contained in the released information. On the other hand, our models are concerned with the release of statistical information which is less specific than microdata in general. However, microdata release can also be handled in our framework when the queries are specific enough. Let us define a complete specification formula (CSF) as a DL wff of the form $\bigwedge_{a \in A} (a, v_a)$, where A is the set of all easy-to-know attributes and v_a is a value in the domain of A . The answer to the batch of queries Q consisting of all CSF's is equivalent to the microdata release of the whole data table T . Therefore, our models are applicable in a more general context.

No matter the difference of application contexts, our models are still comparable with the qualitative ones. In the description of the μ -ARGUS system[8], it is emphasized that re-identification of an individual can occur when the individual is rare in the population with respect to a certain easy-to-know attribute value. In our notation, this means that when a query φ is posed, if $|\varphi|$ is small, then it is rather unsafe to answer the query. In particular, if $|\varphi| = 1$, then the answer of the query necessarily results in the re-identification of the individual satisfying φ . The intuition is formulated as the *bin size* criterion in the Datafly system [16]. A bin is defined as an equivalence class of individuals who have exactly the same easy-to-know attribute values. The bin size criterion is that the size of each bin must be greater than some threshold level. To achieve the criterion, it may be necessary to generalize the data to a more imprecise level. These data modification techniques, mainly generalization and suppression, are further formally investigated in [13, 14]. In their framework, a formal requirement, called k -anonymity, is defined and the generalization and suppression techniques are employed to ensure that the requirement is satisfied. Roughly speaking, k -anonymity requires that each bin must contain at least k individuals, so it is not far removed from

the bin size criterion. Both bin size criterion and k -anonymity requirement can be easily enforced in our model if it is restricted to the effect that a query φ cannot be answered if the size $|\varphi|$ is less than some threshold. However, instead of generalizing or suppressing the data to, we try to assess the value or the damage of the release of such data and discourage the misuse of the information by some pricing or penalty mechanism.

We have also seen in section 5.1 that when the size $|\varphi|$ is large enough, the safety requirement for answering the query φ in our model can be achieved more easily. This also provides an indirect justification why k -anonymity is an useful requirement for privacy protection. Indeed, the diversity of confidential attribute values tends to be higher in a larger size of (sub)population. However, theoretically speaking, it does not fully exclude the possibility of privacy leakage. Imagine the case when all individuals in the bin have the same confidential value,. The sensitive information of all individuals in the bin would be simultaneously leaked if the data is released even though the bin size criterion is satisfied. To circumvent the problem, a logical criterion is proposed in [7]. The criterion is based on the possible world semantics of epistemic logic[5], so it is possible to rigourously define what an user(or intruder) knows. The release of data is then said to be (logically) safe if it does not result in the confidential information being known by the user. The quantitative models defined in this paper can be seen as a generalization of the logical one. Let us temporarily leave aside the investigative actions and consider an information state $(\Pi, \boldsymbol{\mu}, \boldsymbol{\kappa}_0)$ such that $\kappa_{0i}(\varphi) = 0$ for all $\varphi \in \Pi$ and $0 \leq i \leq t - 1$. If for some $\varphi \in \Pi$, we have $\mu_i(\varphi) = 1$ and $\mu_j(\varphi) = 0$ for $j \neq i$, then according to (7), the information state is unsafe if $d_i > 0$. This corresponds to the case where the user knows that all individuals satisfying φ have the confidential value s_i in the logical model and since $d_i > 0$, the information is sensitive, so the information state is also unsafe according to the logical criterion in [7]. However, even though the user cannot know any individual's privacy with certainty, the knowledge of the distribution of the confidential attribute values within a group of individuals may still result in the risk of privacy invasion. The quantitative models can improve the logical criterion exactly in this aspect. Moreover, the k -anonymity requirement can still be considered as complementary to the logical criterion since the latter does not take into account the size of the population, which is another aspect handled in our models by the investigative actions.

Example 7 Let us still consider as the HIV test result the confidential attribute, where $V_c = \{+, -\}$. Assume the damage and compensation of the “+” result being known are respectively r and $5r$ for some $r > 0$ and that the user can investigate at most one individual at affordable cost. By substituting them into equation (9), the safety condition is $\mu_0(\varphi) \leq \frac{5}{6} \cdot (1 - \frac{1}{|\varphi|})$. Now, for a partial knowledge state $\mathcal{P} = (\{\varphi, \neg\varphi\}, \boldsymbol{\mu})$, if $|\varphi| = |\neg\varphi| = 2$ and $\mu_0(\varphi) = \mu_0(\neg\varphi) = \frac{1}{2}$, then by the logical criterion, \mathcal{P} is a safe knowledge state, though it is unsafe according to the 3-anonymity requirement. Since $\mu_0(\varphi) = \frac{1}{2} > \frac{5}{12} = \frac{5}{6} \cdot (1 - \frac{1}{|\varphi|})$, the unsafe situation can be detected by the quantitative criterion proposed in this paper.

On the other hand, if $|\varphi| = |\neg\varphi| = 10$ but $\mu_0(\varphi) = 1$, then it satisfies the k -anonymity requirement up to $k \leq 10$, though the logical criterion of safety is obviously violated. The violation of logical safety can still be detected by the quantitative criterion since $\mu_0(\varphi) = 1 > \frac{3}{4} = \frac{5}{6} \cdot (1 - \frac{1}{|\varphi|})$. ■

7 Conclusion

In this paper, we present two quantitative models for privacy protection. In both models, a formal representation of the user's information states is given. In the first model, we estimate the value of information for the user by considering a specific user model. Under the user model, the privacy protection task is to ensure that the user cannot profit from obtaining the private information. It must be emphasized that the value of information is defined with respect to the particular user model. When other user models are considered, the value of information may be different. Some examples can be seen in [9].

In the second model, we assume that the damage and compensation of revealing each specific confidential value is known. An information state is safe when the user can discover a specific confidential value only with a sufficiently small probability if the damage of revealing the value is large.

References

- [1] L.J. Camp. *Trust and Risk in Internet Commerce*. The MIT Press, 2000.
- [2] Y.-C. Chiang. Protecting privacy in public database (in Chinese). Master's thesis, Graduate Institute of Information Management, National Taiwan University, 2000.
- [3] Y.-C. Chiang, T.-s. Hsu, S. Kuo, and D.-W. Wang. Preserving confidentiality when sharing medical data. In *Proceedings of Asia Pacific Medical Informatics Conference*, 2000.
- [4] G.D. Eppen and F.J. Gould. *Quantitative Concepts for Management*. Prentice Hall, 1985.
- [5] R. Fagin, J.Y. Halpern, Y. Moses, and M.Y. Vardi. *Reasoning about Knowledge*. MIT Press, 1996.
- [6] R.L. Graham, D.E. Knuth, and O. Patashnik. *Concrete Mathematics :A Foundation for Computer Science*. Addison-Wesley, 1994.
- [7] T.-s. Hsu, C.-J. Liao, and D.-W. Wang. A logical model for privacy protection. In *Proceedings of the 4th International Conference on Information Security*, LNCS 2200, pages 110–124. Springer-Verlag, 2001.
- [8] A.J. Hundepool and L.C.R.J. Willenborg. “ μ - and τ -ARGUS: Software for statistical disclosure control”. In *Proceedings of the 3rd International Seminar on Statistical Confidentiality*, 1996.
- [9] J. Kleinberg, C.H. Papadimitriou, and P. Raghavan. “On the value of private information”. In *Proc. 8th Conf. on Theoretical Aspects of Rationality and Knowledge*, 2001.
- [10] D.V. Lindley. *Making Decisions*. John Wiley & Sons, 1985.

- [11] T. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- [12] Z. Pawlak. *Rough Sets—Theoretical Aspects of Reasoning about Data*. Kluwer Academic Publishers, 1991.
- [13] P. Samarati. “Protecting respondents’ identities in microdata release”. *IEEE Transactions on Knowledge and Data Engineering*, 13(6):1010–1027, 2001.
- [14] P. Samarati and L. Sweeney. Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and suppression. Technical report SRI-CSL-98-04, Computer Science Laboratory, SRI International, 1998.
- [15] C.E. Shannon. “The mathematical theory of communication”. *The Bell System Technical Journal*, 27(3&4):379–423,623–656, 1948.
- [16] L. Sweeney. “Guaranteeing anonymity when sharing medical data, the Datafly system”. In *Proceedings of American Medical Informatics Association*, 1997.

A Proof of Proposition 2

Let us define a $(\varphi, (n_i)_{0 \leq i \leq t-1})$ -trace (or simply trace) as a string of symbols in V_c , the domain of values for the confidential attribute, where s_i occurs n_i times for $0 \leq i \leq t-1$ and set $N = |\varphi|$. A trace is a possible result when the user investigates the individuals in $|\varphi|_T$ one by one. If the last k symbols of a trace are all s_i for some i and the symbol before the last k symbols is s_j for some $j \neq i$, then the number of investigative actions the user has to carry out is in fact $N - k$ since once the $(N - k)$ -th individual is investigated, he knows the remaining individuals have the confidential value s_i according to the information. Let us call such trace a (k, s_i) -trace. For each i and $0 < k \leq n_i$, there are in total

$$\sum_{j \neq i} \frac{(N - k - 1)!}{\prod_{l \neq i, j} n_l! \cdot (n_j - 1)! \cdot (n_i - k)!} \quad (13)$$

(k, s_i) -traces. Since $N = \sum_{i=0}^{t-1} n_i$, the expression (13) can be rewritten into

$$\sum_{j \neq i} \frac{(N - k - 1)! \cdot n_j}{\prod_{l \neq i} n_l! \cdot (n_i - k)!} = \frac{(N - k - 1)! \cdot (N - n_i)}{\prod_{l \neq i} n_l! \cdot (n_i - k)!} \quad (14)$$

Since the total number of possible traces is

$$\frac{N!}{\prod_{l=0}^{t-1} n_l!}, \quad (15)$$

the probability of a trace being a (k, s_i) -trace is the division of (14) by (15), i.e.

$$\frac{(N - k - 1)! \cdot n_i! \cdot (N - n_i)}{N! \cdot (n_i - k)!} \quad (16)$$

Thus the expected number of investigative actions the user has to carry out for discovering all individuals' confidential values is equal to

$$\sum_{i=0}^{t-1} \sum_{k=1}^{n_i} \frac{(N-k-1)! \cdot n_i! \cdot (N-n_i)}{N! \cdot (n_i-k)!} \cdot (N-k) \quad (17)$$

which can be further rewritten into

$$\sum_{i=0}^{t-1} \left(\sum_{k=1}^{n_i} \frac{(N-k)! \cdot n_i!}{N! \cdot (n_i-k)!} \right) \cdot (N-n_i) \quad (18)$$

According to [6](pp. 173-174, Problem 1),

$$\sum_{k=0}^{n_i} \frac{\binom{n_i}{k}}{\binom{N}{k}} = \frac{N+1}{N-n_i+1}$$

so

$$\sum_{k=1}^{n_i} \frac{(N-k)! \cdot n_i!}{N! \cdot (n_i-k)!} = \frac{n_i}{N-n_i+1},$$

so (18) can be simplified into

$$\sum_{i=0}^{t-1} \frac{n_i \cdot (N-n_i)}{N-n_i+1} \quad (19)$$

which is exactly (12) since $|\varphi| = N$. ■