

# $\ell$ -Diversity: Privacy Beyond $k$ -Anonymity

Ashwin Machanavajjhala      Johannes Gehrke      Daniel Kifer  
Muthuramakrishnan Venkitasubramaniam  
Department of Computer Science, Cornell University  
{mvnak, johannes, dkifer, vmuthu}@cs.cornell.edu

## Abstract

*Publishing data about individuals without revealing sensitive information about them is an important problem. In recent years, a new definition of privacy called  $k$ -anonymity has gained popularity. In a  $k$ -anonymized dataset, each record is indistinguishable from at least  $k - 1$  other records with respect to certain “identifying” attributes.*

*In this paper we show with two simple attacks that a  $k$ -anonymized dataset has some subtle, but severe privacy problems. First, we show that an attacker can discover the values of sensitive attributes when there is little diversity in those sensitive attributes. Second, attackers often have background knowledge, and we show that  $k$ -anonymity does not guarantee privacy against attackers using background knowledge. We give a detailed analysis of these two attacks and we propose a novel and powerful privacy definition called  $\ell$ -diversity. In addition to building a formal foundation for  $\ell$ -diversity, we show in an experimental evaluation that  $\ell$ -diversity is practical and can be implemented efficiently.*

## 1. Introduction

Many organizations are increasingly publishing microdata – tables that contain unaggregated information about individuals. These tables can include medical, voter registration, census, and customer data. Microdata is a valuable source of information for the allocation of public funds, medical research, and trend analysis. However, if individuals can be uniquely identified in the microdata then their private information (such as their medical condition) would be disclosed, and this is unacceptable.

To avoid the identification of records in microdata, uniquely identifying information like names and social security numbers are removed from the table. However, this first sanitization still does not ensure the privacy of individuals in the data. A recent study estimated that 87% of the population of the United States can be uniquely identi-

fied using the seemingly innocuous attributes gender, date of birth, and 5-digit zip code [23]. In fact, those three attributes were used to link Massachusetts voter registration records (which included the name, gender, zip code, and date of birth) to supposedly anonymized medical data from GIC<sup>1</sup> (which included gender, zip code, date of birth and diagnosis). This “linking attack” managed to uniquely identify the medical records of the governor of Massachusetts in the medical data [24].

Sets of attributes (like gender, date of birth, and zip code in the example above) that can be linked with external data to uniquely identify individuals in the population are called *quasi-identifiers*. To counter linking attacks using quasi-identifiers, Samarati and Sweeney proposed a definition of privacy called  $k$ -anonymity [21, 24]. A table satisfies  $k$ -anonymity if every record in the table is indistinguishable from at least  $k - 1$  other records with respect to every set of quasi-identifier attributes; such a table is called a  $k$ -anonymous table. Hence, for every combination of values of the quasi-identifiers in the  $k$ -anonymous table, there are at least  $k$  records that share those values. This ensures that individuals cannot be uniquely identified by linking attacks.

**An Example.** Figure 1 shows medical records from a fictitious hospital located in upstate New York. Note that the table contains no uniquely identifying attributes like name, social security number, etc. In this example, we divide the attributes into two groups: the *sensitive* attributes (consisting only of medical condition) and the *non-sensitive* attributes (zip code, age, and nationality). An attribute is marked sensitive if an adversary must not be allowed to discover the value of that attribute for any individual in the dataset. Attributes not marked sensitive are non-sensitive. Furthermore, let the collection of attributes {zip code, age, nationality} be the quasi-identifier for this dataset. Figure 2 shows a 4-anonymous table derived from the table in Figure 1 (here “\*” denotes a suppressed value so, for example, “zip code = 1485\*” means that the zip code is in the range [14850 – 14859] and “age=3\*” means the age is in the range

<sup>1</sup>Group Insurance Company (GIC) is responsible for purchasing health insurance for Massachusetts state employees.

	Non-Sensitive			Sensitive
	Zip Code	Age	Nationality	Condition
1	13053	28	Russian	Heart Disease
2	13068	29	American	Heart Disease
3	13068	21	Japanese	Viral Infection
4	13053	23	American	Viral Infection
5	14853	50	Indian	Cancer
6	14853	55	Russian	Heart Disease
7	14850	47	American	Viral Infection
8	14850	49	American	Viral Infection
9	13053	31	American	Cancer
10	13053	37	Indian	Cancer
11	13068	36	Japanese	Cancer
12	13068	35	American	Cancer

**Figure 1. Inpatient Microdata**

	Non-Sensitive			Sensitive
	Zip Code	Age	Nationality	Condition
1	130**	< 30	*	Heart Disease
2	130**	< 30	*	Heart Disease
3	130**	< 30	*	Viral Infection
4	130**	< 30	*	Viral Infection
5	1485*	≥ 40	*	Cancer
6	1485*	≥ 40	*	Heart Disease
7	1485*	≥ 40	*	Viral Infection
8	1485*	≥ 40	*	Viral Infection
9	130**	3*	*	Cancer
10	130**	3*	*	Cancer
11	130**	3*	*	Cancer
12	130**	3*	*	Cancer

**Figure 2. 4-anonymous Inpatient Microdata**

[30 – 39]). Note that in the 4-anonymous table, each tuple has the same values for the quasi-identifier as at least three other tuples in the table.

Because of its conceptual simplicity,  $k$ -anonymity has been widely discussed as a viable definition of privacy in data publishing, and due to algorithmic advances in creating  $k$ -anonymous versions of a dataset [3, 6, 16, 18, 21, 24, 25],  $k$ -anonymity has grown in popularity. However, does  $k$ -anonymity really guarantee privacy? In the next section, we will show that the answer to this question is interestingly *no*. We give examples of two simple, yet subtle attacks on a  $k$ -anonymous dataset that allow an attacker to identify individual records. Defending against these attacks requires a stronger notion of privacy that we call  $\ell$ -diversity, the focus of this paper. But we are jumping ahead in our story. Let us first show the two attacks to give the intuition behind the problems with  $k$ -anonymity.

### 1.1. Attacks On $k$ -Anonymity

In this section we present two attacks, the *homogeneity attack* and the *background knowledge attack*, and we show how they can be used to compromise a  $k$ -anonymous dataset.

**Homogeneity Attack:** Alice and Bob are antagonistic neighbors. One day Bob falls ill and is taken by ambulance to the hospital. Having seen the ambulance, Alice sets out to discover what disease Bob is suffering from. Alice discovers the 4-anonymous table of current inpatient records published by the hospital (Figure 2), and so she knows that one of the records in this table contains Bob’s data. Since Alice is Bob’s neighbor, she knows that Bob is a 31-year-old American male who lives in the zip code 13053. Therefore, Alice knows that Bob’s record number is 9,10,11, or 12. Now, all of those patients have the same medical condition (cancer), and so Alice concludes that Bob has cancer.

**Observation 1**  $k$ -Anonymity can create groups that leak information due to lack of diversity in the sensitive attribute.

Note that such a situation is not uncommon. As a back-of-the-envelope calculation, suppose we have a dataset containing 60,000 distinct tuples where the sensitive attribute can take 3 distinct values and is not correlated with the non-sensitive attributes. A 5-anonymization of this table will have around 12,000 groups<sup>2</sup> and, on average, 1 out of every 81 groups will have no diversity (the values for the sensitive attribute will all be the same). Thus we should expect about 148 groups with no diversity. Therefore, information about 740 people would be compromised by a homogeneity attack. This suggests that in addition to  $k$ -anonymity, the sanitized table should also ensure “diversity” – all tuples that share the same values of their quasi-identifiers should have diverse values for their sensitive attributes.

Our next observation is that an adversary could use “background” knowledge to discover sensitive information.

**Background Knowledge Attack:** Alice has a pen-friend named Umeko who is admitted to the same hospital as Bob, and whose patient records also appear in the table shown in Figure 2. Alice knows that Umeko is a 21 year-old Japanese female who currently lives in zip code 13068. Based on this information, Alice learns that Umeko’s information is contained in record number 1,2,3, or 4. Without additional information, Alice is not sure whether Umeko caught a virus or has heart disease. However, it is well-known that Japanese have an extremely low incidence of heart disease. Therefore Alice concludes with near certainty that Umeko has a viral infection.

**Observation 2**  $k$ -Anonymity does not protect against attacks based on background knowledge.

<sup>2</sup>Our experiments on real data sets show that data is often very skewed and a 5-anonymous table might not have so many groups

We have demonstrated (using the homogeneity and background knowledge attacks) that a  $k$ -anonymous table may disclose sensitive information. Since both of these attacks are plausible in real life, we need a stronger definition of privacy that takes into account diversity and background knowledge. This paper addresses this very issue.

## 1.2. Contributions and Paper Outline

In the previous section, we showed that  $k$ -anonymity is susceptible to homogeneity and background knowledge attacks; thus a stronger definition of privacy is needed. In the remainder of the paper, we derive our solution. We start by introducing an ideal notion of privacy called *Bayes-optimal* for the case that both data publisher and the adversary have full (and identical) background knowledge (Section 3). Unfortunately in practice, the data publisher is unlikely to possess all this information, and in addition, the adversary may have more specific background knowledge than the data publisher. Hence, while Bayes-optimal privacy sounds great in theory, it is unlikely that it can be guaranteed in practice. To address this problem, we show that the notion of Bayes-optimal privacy naturally leads to a novel *practical* definition that we call  $\ell$ -diversity.  $\ell$ -Diversity provides privacy even when the data publisher does not know what kind of knowledge is possessed by the adversary. The main idea behind  $\ell$ -diversity is the requirement that the values of the sensitive attributes are well-represented in each group (Section 4).

We show that existing algorithms for  $k$ -anonymity can be adapted to compute  $\ell$ -diverse tables (Section 5), and in an experimental evaluation we show that  $\ell$ -diversity is practical and can be implemented efficiently (Section 6). We discuss related work in Section 7, and we conclude in Section 8. Before jumping into the contributions of this paper, we introduce the notation needed to formally discuss data privacy in the next section.

## 2. Model and Notation

In this section we will introduce some basic notation that will be used in the remainder of the paper. We will also discuss how a table can be anonymized and what kind of background knowledge an adversary may possess.

**Basic Notation.** Let  $T = \{t_1, t_2, \dots, t_n\}$  be a table with attributes  $A_1, \dots, A_m$ . We assume that  $T$  is a subset of some larger population  $\Omega$  where each tuple represents an individual from the population. For example, if  $T$  is a medical dataset then  $\Omega$  could be the population of the United States. Let  $\mathcal{A}$  denote the set of all attributes  $\{A_1, A_2, \dots, A_m\}$  and  $t[A_i]$  denote the value of attribute  $A_i$  for tuple  $t$ . If  $\mathcal{C} = \{C_1, C_2, \dots, C_p\} \subseteq \mathcal{A}$  then we

use the notation  $t[\mathcal{C}]$  to denote the tuple  $(t[C_1], \dots, t[C_p])$ , which is the projection of  $t$  onto the attributes in  $\mathcal{C}$ .

In privacy-preserving data publishing, there exist several important subsets of  $\mathcal{A}$ . A *sensitive attribute* is an attribute whose value for any particular individual must be kept secret from people who have no direct access to the original data. Let  $\mathcal{S}$  denote the set of all sensitive attributes. An example of a sensitive attribute is *Medical Condition* from Figure 1. The association between individuals and *Medical Condition* should be kept secret; thus we should not disclose which particular patients have cancer, but it is permissible to disclose the information that there exist cancer patients in the hospital. We assume that the data publisher knows which attributes are sensitive. All attributes that are not sensitive are called *nonsensitive* attributes. Let  $\mathcal{N}$  denote the set of nonsensitive attributes. We are now ready to formally define the notion of a quasi-identifier.

**Definition 2.1 (Quasi-identifier)** A set of nonsensitive attributes  $\{Q_1, \dots, Q_w\}$  of a table is called a quasi-identifier if these attributes can be linked with external data to uniquely identify at least one individual in the general population  $\Omega$ .

One example of a quasi-identifier is a primary key like social security number. Another example is the set  $\{\text{Gender}, \text{Age}, \text{Zip Code}\}$  in the GIC dataset that was used to identify the governor of Massachusetts as described in the introduction. Let us denote the set of all quasi-identifiers by  $\mathcal{QI}$ . We are now ready to formally define  $k$ -anonymity.

**Definition 2.2 ( $k$ -Anonymity)** A table  $T$  satisfies  $k$ -anonymity if for every tuple  $t \in T$  there exist  $k - 1$  other tuples  $t_{i_1}, t_{i_2}, \dots, t_{i_{k-1}} \in T$  such that  $t[\mathcal{C}] = t_{i_1}[\mathcal{C}] = t_{i_2}[\mathcal{C}] = \dots = t_{i_{k-1}}[\mathcal{C}]$  for all  $\mathcal{C} \in \mathcal{QI}$ .

**The Anonymized Table  $T^*$ .** Since the quasi-identifiers might uniquely identify tuples in  $T$ , the table  $T$  is not published; it is subjected to an *anonymization procedure* and the resulting table  $T^*$  is published instead.

There has been a lot of research on techniques for anonymization (see Section 7 for a discussion of related work). These techniques can be broadly classified into *generalization* techniques [3, 16], *generalization with tuple suppression* techniques [6, 22], and *data swapping and randomization* techniques [1, 13]. In this paper we limit our discussion only to generalization techniques.

**Definition 2.3 (Domain Generalization)** A domain  $D^* = \{P_1, P_2, \dots\}$  is a generalization (partition) of a domain  $D$  if  $\bigcup P_i = D$  and  $P_i \cap P_j = \emptyset$  whenever  $i \neq j$ . For  $x \in D$  we let  $\phi_{D^*}(x)$  denote the element  $P \in D^*$  that contains  $x$ .

Note that we can create a partial order  $\prec_G$  on domains by requiring  $D \prec_G D^*$  if and only if  $D^*$  is a generalization of

D. Given a table  $T = \{t_1, \dots, t_n\}$  with the set of nonsensitive attributes  $\mathcal{N}$  and a generalization  $D_N^*$  of  $\text{domain}(\mathcal{N})$ , we can construct a table  $T^* = \{t_1^*, \dots, t_n^*\}$  by replacing the value of  $t_i[\mathcal{N}]$  with the generalized value  $\phi_{D_N^*}(t_i[\mathcal{N}])$  to get a new tuple  $t_i^*$ . The tuple  $t_i^*$  is called a *generalization* of the tuple  $t_i$  and we use the notation  $t_i \xrightarrow{*} t_i^*$  to mean “ $t_i^*$  generalizes  $t_i$ ”. Extending the notation to tables,  $T \xrightarrow{*} T^*$  means “ $T^*$  is a generalization of  $T$ ”. Typically, ordered attributes are partitioned into intervals, and categorical attributes are partitioned according to a user-defined hierarchy (for example, cities are generalized to counties, counties to states, and states to regions).

**Example 1 (Continued).** The table in Figure 2 is a generalization of the table in Figure 1. We generalized on the *Zip Code* attribute by partitioning it into two sets: “1485\*” (representing all zip codes whose first four digits are 1485) and “130\*” (representing all zip codes whose first three digits are 130). Then we partitioned *Age* into three groups: “< 30”, “3\*” (representing all ages between 30 and 39), and “≥ 40”. Finally, we partitioned *Nationality* into just one set “\*” representing all nationalities.

**The Adversary’s Background Knowledge.** Since the background knowledge attack was due to the adversary’s additional knowledge about the table, let us briefly discuss the type of background knowledge that we are modeling.

First, the adversary has access to the published table  $T^*$  and she knows that  $T^*$  is a generalization of some base table  $T$ . The adversary also knows the domain of each attribute of  $T$ .

Second, the adversary may know that some individuals are in the table. This knowledge is often easy to acquire. For example, GIC published medical data about Massachusetts state employees. If the adversary Alice knows that her neighbor Bob is a Massachusetts state employee then Alice is almost certain that Bob’s information is contained in that table. In this case, we assume that Alice knows all of Bob’s nonsensitive attributes. In addition, the adversary could have knowledge about the sensitive attributes of specific individuals in the population and/or the table. For example, the adversary Alice might know that neighbor Bob does not have pneumonia since Bob does not show any of the symptoms of pneumonia. We call such knowledge “instance-level background knowledge,” since it is associated with specific instances in the table.

Third, the adversary could have partial knowledge about the distribution of sensitive and nonsensitive attributes in the population. We call this “demographic background knowledge.” For example, the adversary may know  $P(t[\text{Condition}] = \text{“cancer”} \mid t[\text{Age}] \geq 40)$ , and may use it to make additional inferences about records in the table.

Now armed with the right notation, let us start looking into principles and definitions of privacy that leak little information.

### 3. Bayes-Optimal Privacy

In this section we analyze an ideal notion of privacy called *Bayes-Optimal Privacy* since it involves modeling background knowledge as a probability distribution over the attributes and uses Bayesian inference techniques to reason about privacy. We introduce tools for reasoning about privacy (Section 3.1), we use them to discuss theoretical principles of privacy (Section 3.2), and then we point out the difficulties that need to be overcome to arrive at a practical definition of privacy (Section 3.3).

#### 3.1. Changes in Belief Due to Data Publishing

For simplicity of discussion, we will combine all the nonsensitive attributes into a single, multi-dimensional quasi-identifier attribute  $Q$  whose values are generalized to create the anonymized table  $T^*$  from the base table  $T$ . Since Bayes-optimal privacy is only used to motivate a practical definition, we make the following two simplifying assumptions: first, we assume that  $T$  is a simple random sample from some larger population  $\Omega$  (a sample of size  $n$  drawn without replacement is called a *simple random sample* if every sample of size  $n$  is equally likely); second, we assume that there is a single sensitive attribute. We would like to emphasize that both these assumptions will be dropped in Section 4 when we introduce a practical definition of privacy.

Recall that in our attack model, the adversary Alice has partial knowledge of the distribution of the sensitive and non-sensitive attributes. Let us assume a worst case scenario where Alice knows the complete joint distribution  $f$  of  $Q$  and  $S$  (i.e. she knows their frequency in the population  $\Omega$ ). She knows that Bob corresponds to a record  $t \in T$  that has been generalized to a record  $t^*$  in the published table  $T^*$ , and she also knows the value of Bob’s non-sensitive attributes (i.e., she knows that  $t[Q] = q$ ). Alice’s goal is to use her background knowledge to discover Bob’s sensitive information — the value of  $t[S]$ . We gauge her success using two quantities: Alice’s *prior belief*, and her *posterior belief*.

Alice’s *prior belief*,  $\alpha_{(q,s)}$ , that Bob’s sensitive attribute is  $s$  given that his nonsensitive attribute is  $q$ , is just her background knowledge:

$$\alpha_{(q,s)} = P_f(t[S] = s \mid t[Q] = q)$$

After Alice observes the table  $T^*$ , her belief about Bob’s sensitive attribute changes. This new belief,  $\beta_{(q,s,T^*)}$ , is her *posterior belief*:

$$\beta_{(q,s,T^*)} = P_f(t[S] = s \mid t[Q] = q \wedge \exists t^* \in T^*, t \xrightarrow{*} t^*)$$

Given  $f$  and  $T^*$ , we can derive a formula for  $\beta_{(q,s,T^*)}$  which will help us formulate our new privacy definition in Section 4. The main idea behind the derivation is to find a set of equally likely disjoint random worlds (like in [5]) such that the conditional probability  $P(A|B)$  is the number of worlds satisfying the condition  $A \wedge B$  divided by the number of worlds satisfying the condition  $B$ . We avoid double-counting because the random worlds are disjoint. In our case, a random world is any permutation of a simple random sample of size  $n$  that is drawn from the population  $\Omega$  and which is *compatible* with the published table  $T^*$ .<sup>3</sup>

**Theorem 3.1** *Let  $q$  be a value of the nonsensitive attribute  $Q$  in the base table  $T$ ; let  $q^*$  be the generalized value of  $q$  in the published table  $T^*$ ; let  $s$  be a possible value of the sensitive attribute; let  $n_{(q^*,s')}$  be the number of tuples  $t \in T^*$  where  $t^*[Q] = q^*$  and  $t^*[S] = s'$ ; and let  $f(s' | q^*)$  be the conditional probability of the sensitive attribute conditioned on the fact that the nonsensitive attribute  $Q$  can be generalized to  $q^*$ . Then the following relationship holds:*

$$\beta_{(q,s,T^*)} = \frac{n_{(q^*,s)} \frac{f(s|q)}{f(s|q^*)}}{\sum_{s' \in S} n_{(q^*,s')} \frac{f(s'|q^*)}{f(s'|q^*)}} \quad (1)$$

Armed with a way of calculating Alice's belief about Bob's private data after she has seen  $T^*$ , let us now examine some principles for building definitions of privacy.

### 3.2. Privacy Principles

Given the adversary's background knowledge, a published table  $T^*$  might disclose information in two important ways: *positive disclosure* and *negative disclosure*.

**Definition 3.1 (Positive disclosure)** *Publishing the table  $T^*$  that was derived from  $T$  results in a positive disclosure if the adversary can correctly identify the value of a sensitive attribute with high probability; i.e., given a  $\delta > 0$ , there is a positive disclosure if  $\beta_{(q,s,T^*)} > 1 - \delta$  and there exists  $t \in T$  such that  $t[Q] = q$  and  $t[S] = s$ .*

**Definition 3.2 (Negative disclosure)** *Publishing the table  $T^*$  that was derived from  $T$  results in a negative disclosure if the adversary can correctly eliminate some possible values of the sensitive attribute (with high probability); i.e., given an  $\epsilon > 0$ , there is a negative disclosure if  $\beta_{(q,s,T^*)} < \epsilon$  and there exists a  $t \in T$  such that  $t[Q] = q$  but  $t[S] \neq s$ .*

The homogeneity attack in Section 1.1 where Alice determined that Bob has cancer is an example of a positive

disclosure. Similarly, in the example from Section 1.1, even without background knowledge Alice can deduce that Umeko does not have cancer. This is an example of a negative disclosure.

Note that not all positive disclosures are disastrous. If the prior belief was that  $\alpha_{(q,s)} > 1 - \delta$ , the adversary would not have learned anything new. Similarly, negative disclosures are not always bad: discovering that Bob does not have Ebola might not be very serious because the prior belief of this event was small. Hence, the ideal definition of privacy can be based on the following principle:

**Principle 1 (Uninformative Principle)** *The published table should provide the adversary with little additional information beyond the background knowledge. In other words, there should not be a large difference between the prior and posterior beliefs.*

The uninformative principle can be instantiated in several ways, for example with the  $(\rho_1, \rho_2)$ -privacy breach definition [14]. Under this definition, privacy is breached either when  $\alpha_{(q,s)} < \rho_1 \wedge \beta_{(q,s,T^*)} > \rho_2$  or when  $\alpha_{(q,s)} > 1 - \rho_1 \wedge \beta_{(q,s,T^*)} < 1 - \rho_2$ . An alternative privacy definition based on the uninformative principle would bound the maximum difference between  $\alpha_{(q,s)}$  and  $\beta_{(q,s,T^*)}$  using any of the functions commonly used to measure the difference between probability distributions. Any privacy definition that is based on the uninformative principle, and instantiated either by a  $(\rho_1, \rho_2)$ -privacy breach definition or by bounding the difference between  $\alpha_{(q,s)}$  and  $\beta_{(q,s,T^*)}$  is a Bayes-optimal privacy definition. The specific choice of definition depends on the application.

Note that any Bayes-optimal privacy definition captures diversity as well as background knowledge. To see how it captures diversity, suppose that all the tuples whose nonsensitive attribute  $Q$  have been generalized to  $q^*$  have the same value  $s$  for their sensitive attribute. Then  $n_{(q^*,s')} = 0$  for all  $s' \neq s$  and hence the value of the observed belief  $\beta_{(q,s,T^*)}$  becomes 1 in Equation 1. This will be flagged as a breach whenever the prior belief is not close to 1.

### 3.3. Limitations of the Bayes-Optimal Privacy

For the purposes of our discussion, we are more interested in the properties of Bayes-optimal privacy rather than its exact instantiation. In particular, Bayes-optimal privacy has several drawbacks that make it hard to use in practice.

**Insufficient Knowledge.** The data publisher is unlikely to know the full distribution  $f$  of sensitive and nonsensitive attributes over the general population  $\Omega$  from which  $T$  is a sample.

**The Adversary's Knowledge is Unknown.** It is also unlikely that the adversary has knowledge of the complete

<sup>3</sup>Due to space constraints we had to omit the proof of the following theorem; see [17] for the derivation of Equation 1.

joint distribution between the non-sensitive and sensitive attributes. However, the data publisher does not know how much the adversary knows. For example, in the background knowledge attack in Section 1.1, Alice knew that Japanese have a low incidence of heart disease, but the data publisher did not know that Alice knew this piece of information.

**Instance-Level Knowledge.** The theoretical definition does not protect against knowledge that cannot be modeled probabilistically. For example, suppose Bob's son tells Alice that Bob does not have diabetes. The theoretical definition of privacy will not be able to protect against such adversaries.

**Multiple Adversaries.** There will likely be multiple adversaries with different levels of knowledge, each of which is consistent with the full joint distribution. Suppose Bob has a disease that is (a) very likely among people in the age group [30-50], but (b) is very rare for people of that age group who are doctors. An adversary who only knows the interaction of age and illness will think that it is very likely for Bob to have that disease. However, an adversary who also knows that Bob is a doctor is more likely to think that Bob does not have that disease. Thus, although additional knowledge can yield better inferences on average, there are specific instances where it does not. Thus the data publisher must take into account all possible levels of background knowledge.

In the next section, we present a definition that eliminates these drawbacks.

## 4. $\ell$ -Diversity: A Practical Privacy Definition

In this section we discuss how to overcome the difficulties outlined at the end of the previous section. We derive the  $\ell$ -diversity principle (Section 4.1), show how to instantiate it with specific definitions of privacy (Section 4.2), outline how to handle multiple sensitive attributes (Section 4.3), and discuss how  $\ell$ -diversity addresses the issues raised in the previous section (Section 4.4).

### 4.1. The $\ell$ -Diversity Principle

Recall that Theorem 3.1 allows us to calculate the observed belief of the adversary. Let us define a  $q^*$ -block to be the set of tuples in  $T^*$  whose nonsensitive attribute values generalize to  $q^*$ . Consider the case of positive disclosures; i.e., Alice wants to determine that Bob has  $t[S] = s$  with very high probability. From Theorem 3.1, this can happen only when:

$$\exists s, \forall s' \neq s, \quad n_{(q^*, s')} \frac{f(s'|q)}{f(s'|q^*)} \ll n_{(q^*, s)} \frac{f(s|q)}{f(s|q^*)} \quad (2)$$

The condition in Equation (2) could occur due to a combination of two factors: (i) a lack of diversity in the sensi-

tive attributes in the  $q^*$ -block, and/or (ii) strong background knowledge. Let us discuss these in turn.

**Lack of Diversity.** Lack of diversity in the sensitive attribute manifests itself as follows:

$$\forall s' \neq s, \quad n_{(q^*, s')} \ll n_{(q^*, s)} \quad (3)$$

In this case, almost all tuples have the same value  $s$  for the sensitive attribute  $S$ , and thus  $\beta_{(q, s, T^*)} \approx 1$ . Note that this condition can be easily checked since it only involves counting the values of  $S$  in the published table  $T^*$ . We can ensure diversity by requiring that *all* the possible values  $s' \in \text{domain}(S)$  occur in the  $q^*$ -block with roughly equal proportions. This, however, is likely to cause significant loss of information: if  $\text{domain}(S)$  is large then the  $q^*$ -blocks will necessarily be large and so the data will be partitioned into a small number of  $q^*$ -blocks. Another way to ensure diversity and to guard against Equation 3 is to require that a  $q^*$ -block has at least  $\ell \geq 2$  different sensitive values such that the  $\ell$  most frequent values (in the  $q^*$ -block) have roughly the same frequency. We say that such a  $q^*$ -block is *well-represented by  $\ell$  sensitive values*.

**Strong Background Knowledge.** The other factor that could lead to a positive disclosure (Equation 2) is strong background knowledge. Even though a  $q^*$ -block may have  $\ell$  “well-represented” sensitive values, Alice may still be able to use her background knowledge to eliminate sensitive values when the following is true:

$$\exists s', \quad \frac{f(s'|q)}{f(s'|q^*)} \approx 0 \quad (4)$$

This equation states that Bob with quasi-identifier  $t[Q] = q$  is much less likely to have sensitive value  $s'$  than any other individual in the  $q^*$ -block. For example, Alice may know that Bob never travels, and thus he is extremely unlikely to have Ebola. It is not possible for a data publisher to guard against attacks employing arbitrary amounts of background knowledge. However, the data publisher can still guard against many attacks even without having access to Alice's background knowledge. In our model, Alice might know the distribution  $f(q, s)$  over the sensitive and non-sensitive attributes, in addition to the conditional distribution  $f(s|q)$ . The most damaging type of such information has the form  $f(s|q) \approx 0$ , e.g., “men do not have breast cancer”, or the form of Equation 4, e.g., “among Asians, Japanese have a very low incidence of heart disease”. Note that *a priori* information of the form  $f(s|q) = 1$  is not as harmful since this positive disclosure is independent of the published table  $T^*$ . Alice can also eliminate sensitive values with instance-level knowledge such as “Bob does not have diabetes”.

In spite of such background knowledge, if there are  $\ell$  “well represented” sensitive values in a  $q^*$ -block, then Alice needs  $\ell - 1$  damaging pieces of background knowledge

	Non-Sensitive			Sensitive
	Zip Code	Age	Nationality	Condition
1	1305*	≤ 40	*	Heart Disease
4	1305*	≤ 40	*	Viral Infection
9	1305*	≤ 40	*	Cancer
10	1305*	≤ 40	*	Cancer
5	1485*	> 40	*	Cancer
6	1485*	> 40	*	Heart Disease
7	1485*	> 40	*	Viral Infection
8	1485*	> 40	*	Viral Infection
2	1306*	≤ 40	*	Heart Disease
3	1306*	≤ 40	*	Viral Infection
11	1306*	≤ 40	*	Cancer
12	1306*	≤ 40	*	Cancer

**Figure 3. 3-Diverse Inpatient Microdata**

to eliminate  $\ell - 1$  possible sensitive values and infer a positive disclosure! Thus, by setting the parameter  $\ell$ , the data publisher can determine how much protection is provided against background knowledge — even if this background knowledge is unknown to the publisher.

Putting these two arguments together, we arrive at the following principle.

**Principle 2 ( $\ell$ -Diversity Principle)** A  $q^*$ -block is  $\ell$ -diverse if contains at least  $\ell$  “well-represented” values for the sensitive attribute  $S$ . A table is  $\ell$ -diverse if every  $q^*$ -block is  $\ell$ -diverse.

Returning to our example, consider the inpatient records shown in Figure 1. We present a 3-diverse version of the table in Figure 3. Comparing it with the 4-anonymous table in Figure 2 we see that the attacks against the 4-anonymous table are prevented by the 3-diverse table. For example, Alice cannot infer from the 3-diverse table that Bob (a 31 year old American from zip code 13053) has cancer. Even though Umeko (a 21 year old Japanese from zip code 13068) is extremely unlikely to have heart disease, Alice is still unsure whether Umeko has a viral infection or cancer.

The  $\ell$ -diversity principle advocates ensuring  $\ell$  “well represented” values for the sensitive attribute in every  $q^*$ -block, but does not clearly state what “well represented” means. Note that we called it a “principle” instead of a theorem — we will use it to give two concrete instantiations of the  $\ell$ -diversity principle and discuss their relative trade-offs.

#### 4.2. $\ell$ -Diversity: Instantiations

Our first instantiation of the  $\ell$ -diversity principle uses the information-theoretic notion of entropy:

**Definition 4.1 (Entropy  $\ell$ -Diversity)** A table is Entropy  $\ell$ -

Diverse if for every  $q^*$ -block

$$-\sum_{s \in S} p_{(q^*, s)} \log(p_{(q^*, s)}) \geq \log(\ell)$$

where  $p_{(q^*, s)} = \frac{n_{(q^*, s)}}{\sum_{s' \in S} n_{(q^*, s')}} is the fraction of tuples in the  $q^*$ -block with sensitive attribute value equal to  $s$ .$

As a consequence of this condition, every  $q^*$ -block has at least  $\ell$  distinct values for the sensitive attribute. Using this definition, Figure 3 is actually 2.8-diverse.

Since  $-x \log(x)$  is a concave function, it can be shown that if we split a  $q^*$ -block into two sub-blocks  $q_a^*$  and  $q_b^*$  then  $\text{entropy}(q^*) \geq \min(\text{entropy}(q_a^*), \text{entropy}(q_b^*))$ . This implies that in order for entropy  $\ell$ -diversity to be possible, the entropy of the entire table must be at least  $\log(\ell)$ . This might not be the case, especially if one value of the sensitive attribute is very common — for example, if 90% of the patients have “heart problems” as the value for the “Medical Condition” attribute.

Thus entropy  $\ell$ -diversity may sometimes be too restrictive. If some positive disclosures are acceptable (for example, a clinic is allowed to disclose that a patient has a “heart problem” because it is well known that most patients who visit the clinic have heart problems) then we can do better. This reasoning allows us to develop a less conservative instantiation of the  $\ell$ -diversity principle called *recursive  $\ell$ -diversity*.

Let  $s_1, \dots, s_m$  be the possible values of the sensitive attribute  $S$  in a  $q^*$ -block. Assume that we sort the counts  $n_{(q^*, s_1)}, \dots, n_{(q^*, s_m)}$  in descending order and name the elements of the resulting sequence  $r_1, \dots, r_m$ . One way to think about  $\ell$ -diversity is the following: the adversary needs to eliminate at least  $\ell - 1$  possible values of  $S$  in order to infer a positive disclosure. This means that, for example, in a 2-diverse table, none of the sensitive values should appear too frequently. We say that a  $q^*$ -block is  $(c, 2)$ -diverse if  $r_1 < c(r_2 + \dots + r_m)$  for some user-specified constant  $c$ . For  $\ell > 2$ , we say that a  $q^*$ -block satisfies *recursive  $(c, \ell)$ -diversity* if we can eliminate one possible sensitive value in the  $q^*$ -block and still have a  $(c, \ell - 1)$ -diverse block. This recursive definition can be succinctly stated as follows:

**Definition 4.2 (Recursive  $(c, \ell)$ -Diversity)** In a given  $q^*$ -block, let  $r_i$  denote the number of times the  $i^{\text{th}}$  most frequent sensitive value appears in that  $q^*$ -block. Given a constant  $c$ , the  $q^*$ -block satisfies recursive  $(c, \ell)$ -diversity if  $r_1 < c(r_\ell + r_{\ell+1} + \dots + r_m)$ . A table  $T^*$  satisfies recursive  $(c, \ell)$ -diversity if every  $q^*$ -block satisfies recursive  $\ell$ -diversity. We say that 1-diversity is always satisfied.

Now suppose that  $Y$  is the set of sensitive values for which positive disclosure is allowed (for example, because they are extremely frequent, or because they may not be an

invasion of privacy – like “Medical Condition”=“Healthy”). Since we are not worried about those values being too frequent, let  $s_y$  be the most frequent sensitive value in the  $q^*$ -block that is *not* in  $Y$  and let  $r_y$  be the associated frequency. Then the  $q^*$ -block satisfies  $\ell$ -diversity if we can eliminate the  $\ell - 2$  most frequent values of  $S$  *not including*  $r_y$  without making  $s_y$  too frequent in the resulting set. This is the same as saying that after we remove the sensitive values with counts  $r_1, \dots, r_{y-1}$ , then the result is  $(\ell - y + 1)$ -diverse. This brings us to the following definition.

**Definition 4.3 (Positive Disclosure-Recursive  $(c, \ell)$ -Diversity).** Let  $Y$  denote the set of sensitive values for which positive disclosure is allowed. In a given  $q^*$ -block, let the most frequent sensitive value not in  $Y$  be the  $y^{\text{th}}$  most frequent sensitive value. Let  $r_i$  denote the frequency of the  $i^{\text{th}}$  most frequent sensitive value in the  $q^*$ -block. Such a  $q^*$ -block satisfies pd-recursive  $(c, \ell)$ -diversity if one of the following hold:

- $y \leq \ell - 1$  and  $r_y < c \sum_{j=\ell}^m r_j$
- $y > \ell - 1$  and  $r_y < c \sum_{j=\ell-1}^{y-1} r_j + c \sum_{j=y+1}^m r_j$

We denote the summations on the right hand side of the both conditions by  $\text{tail}_{q^*}(s_y)$ .

Now, note that if  $r_y = 0$  then the  $q^*$ -block only has sensitive values that can be disclosed and so both conditions in Definition 4.3 are trivially satisfied. Second, note that if  $c > 1$  then the second condition clearly reduces to just the condition  $y > \ell - 1$  because  $r_y \leq r_{\ell-1}$ . The second condition states that even though the  $\ell - 1$  most frequent values can be disclosed, we still do not want  $r_y$  to be too frequent if  $\ell - 2$  of them have been eliminated (i.e., we want the result to be 2-diverse).

Until now we have treated negative disclosure as relatively unimportant compared to positive disclosure. However, negative disclosure may also be important. If  $W$  is the set of values for the sensitive attribute for which negative disclosure is not allowed then, given a user-specified constant  $c_2 < 100$ , we require that each  $s \in W$  appear in at least  $c_2$ -percent of the tuples in every  $q^*$ -block, resulting in the following definition.

**Definition 4.4 (Negative/Positive Disclosure-Recursive  $(c_1, c_2, \ell)$ -Diversity).** Let  $W$  be the set of sensitive values for which negative disclosure is not allowed. A table satisfies npd-recursive  $(c_1, c_2, \ell)$ -diversity if it satisfies pd-recursive  $(c_1, \ell)$ -diversity and if every  $s \in W$  occurs in at least  $c_2$  percent of the tuples in every  $q^*$ -block.

### 4.3. Multiple Sensitive Attributes

Multiple sensitive attributes present some additional challenges. Suppose  $S$  and  $V$  are two sensitive attributes, and consider the  $q^*$ -block with the following tuples:  $\{(q^*, s_1, v_1), (q^*, s_1, v_2), (q^*, s_2, v_3), (q^*, s_3, v_3)\}$ . This  $q^*$ -block is 3-diverse (actually recursive (2,3)-diverse) with respect to  $S$  (ignoring  $V$ ) and 3-diverse with respect to  $V$  (ignoring  $S$ ). However, if we know that Bob is in this block and his value for  $S$  is not  $s_1$  then his value for attribute  $V$  cannot be  $v_1$  or  $v_2$ , and therefore must be  $v_3$ . One piece of information destroyed his privacy. Thus we see that a  $q^*$ -block that is  $\ell$ -diverse in each sensitive attribute separately may still violate the principle of  $\ell$ -diversity.

Intuitively, the problem occurred because within the  $q^*$ -block,  $V$  was not well-represented for each value of  $S$ . Had we treated  $S$  as part of the quasi-identifier when checking for diversity in  $V$  (and vice versa), we would have ensured that the  $\ell$ -diversity principle held for the entire table. Formally,

**Definition 4.5 (Multi-Attribute  $\ell$ -Diversity)** Let  $T$  be a table with nonsensitive attributes  $Q_1, \dots, Q_{m_1}$  and sensitive attributes  $S_1, \dots, S_{m_2}$ . We say that  $T$  is  $\ell$ -diverse if for all  $i = 1 \dots m_2$ , the table  $T$  is  $\ell$ -diverse when  $S_i$  is treated as the sole sensitive attribute and  $\{Q_1, \dots, Q_{m_1}, S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_{m_2}\}$  is treated as the quasi-identifier.

As the number of sensitive attributes grows, it is not hard to see that we will necessarily need larger and larger  $q^*$ -blocks to ensure diversity. This problem may be ameliorated through tuple suppression and generalization on the sensitive attributes, and is a subject for future work.

### 4.4. Discussion

Recall that we started our journey into Section 4 motivated by the weaknesses of Bayes-optimal privacy. Let us now revisit these issues one by one.

- $\ell$ -Diversity no longer requires knowledge of the full distribution of the sensitive and nonsensitive attributes.
- $\ell$ -Diversity does not even require the data publisher to have as much information as the adversary. The parameter  $\ell$  protects against more knowledgeable adversaries; the larger the value of  $\ell$ , the more information is needed to rule out possible values of the sensitive attribute.
- Instance-level knowledge (Bob’s son tells Alice that Bob does not have diabetes) is automatically covered. It is treated as just another way of ruling out possible values of the sensitive attribute.



- Different adversaries can have different background knowledge leading to different inferences.  $\ell$ -Diversity simultaneously protects against all of them without the need for checking which inferences can be made with which levels of background knowledge.

Overall, we believe that  $\ell$ -diversity is practical, easy to understand, and addresses the shortcomings of  $k$ -anonymity with respect to the background knowledge and homogeneity attacks. Let us now see whether we can give efficient algorithms to implement  $\ell$ -diversity. We will see that, unlike Bayes-optimal privacy,  $\ell$ -diversity possesses a property called *monotonicity*. We will define this concept in Section 5, and we show how this property can be used to efficiently generate  $\ell$ -diverse tables.

## 5. Implementing Privacy Preserving Data Publishing

In this section we discuss how to build algorithms for privacy-preserving data publishing using domain generalization. Let us first review the search space for privacy-preserving data publishing using domain generalization [6, 16]. For ease of explanation, we will combine all the nonsensitive attributes into a single multi-dimensional attribute  $Q$ . For attribute  $Q$ , there is a user-defined generalization lattice. Formally, we define a generalization lattice to be a set of domains partially ordered by a generalization relation  $\prec_G$  (as described in Section 2). The bottom element of this lattice is  $\text{domain}(Q)$  and the top element is the domain where each dimension of  $Q$  is generalized to a single value. Given a base table  $T$ , each domain  $D_Q^*$  in the lattice defines an anonymized table  $T^*$  which is constructed by replacing each tuple  $t \in T$  by the tuple  $t^*$ , such that the value  $t^*[Q] \in D_Q^*$  is the generalization of the value  $t[Q] \in \text{domain}(Q)$ . An algorithm for data publishing should find a point on the lattice such that the corresponding generalized table  $T^*$  preserves privacy and retains as much utility as possible. In the literature, the utility of a generalized table is usually defined as a distance metric on the lattice – the closer the lattice point is to the bottom, the larger the utility of the corresponding table  $T^*$ . Hence, finding a suitable anonymized table  $T^*$  is essentially a lattice search problem. There has been work on search strategies for  $k$ -anonymous tables that explore the lattice top-down [6] or bottom-up [16].

In general, searching the entire lattice is computationally intractable. However, lattice searches can be made efficient if there is a stopping condition of the form: if  $T^*$  preserves privacy then every generalization of  $T^*$  also preserves privacy [16, 22]. This is called the *monotonicity property*, and it has been used extensively in frequent itemset mining algorithms [4].  $k$ -Anonymity satisfies the monotonicity prop-

erty, and it is this property which guarantees the correctness of all efficient algorithms [6, 16]. Thus, if we show that  $\ell$ -diversity also possesses the monotonicity property, then we can re-use these efficient lattice search algorithms to find the  $\ell$ -diverse table with optimal utility. Although more of theoretical interest, we can prove the following theorem that gives a computational reason why Bayes-optimal privacy does not lend itself to efficient algorithmic implementations.

**Theorem 5.1** *Bayes-optimal privacy does not satisfy the monotonicity property.*

However, we can prove that all variants of  $\ell$ -diversity satisfy monotonicity.

**Theorem 5.2 (Monotonicity of Entropy  $\ell$ -diversity)**

*Entropy  $\ell$ -diversity satisfies the monotonicity property: if a table  $T^*$  satisfies entropy  $\ell$ -diversity, then any generalization  $T^{**}$  of  $T^*$  also satisfies entropy  $\ell$ -diversity.*

**Theorem 5.3 (Monotonicity of NPD Recursive**

**$\ell$ -diversity)** *npd recursive  $(c_1, c_2, \ell)$ -diversity satisfies the monotonicity property: if a table  $T^*$  satisfies npd recursive  $(c_1, c_2, \ell)$ -diversity, then any generalization  $T^{**}$  of  $T^*$  also satisfies npd recursive  $(c_1, c_2, \ell)$ -diversity.*

Thus to create an algorithm for  $\ell$ -diversity, we simply take any algorithm for  $k$ -anonymity and make the following change: every time a table  $T^*$  is tested for  $k$ -anonymity, we check for  $\ell$ -diversity instead. Since  $\ell$ -diversity is a property that is local to each  $q^*$ -block and since all  $\ell$ -diversity tests are solely based on the counts of the sensitive values, this test can be performed very efficiently.

## 6. Experiments

In our experiments, we used an implementation of Incognito, as described in [16], for generating  $k$ -anonymous tables. We modified this implementation so that it produces  $\ell$ -diverse tables as well. Incognito is implemented in Java and uses the database manager IBM DB2 v8.1 to store its data. All experiments were run under Linux (Fedora Core 3) on a machine with a 3 GHz Intel Pentium 4 processor and 1 GB RAM.

We ran our experiments on the Adult Database from the UCI Machine Learning Repository [20] and the Lands End Database. The Adult Database contains 45,222 tuples from US Census data and the Lands End Database contains 4,591,581 tuples of point-of-sale information. We removed tuples with missing values and adopted the same domain generalizations as [16]. Figure 4 provides a brief description of the data including the attributes we used, the number of distinct values for each attribute, the type of generalization that was used (for non-sensitive attributes), and the height of the generalization hierarchy for each attribute.

Adults

	Attribute	Domain size	Generalizations type	Ht.
1	Age	74	ranges-5,10,20	4
2	Gender	2	Suppression	1
3	Race	5	Suppression	1
4	Marital Status	7	Taxonomy tree	2
5	Education	16	Taxonomy tree	3
6	Native Country	41	Taxonomy tree	2
7	Work Class	7	Taxonomy tree	2
8	Salary class	2	<i>Sensitive att.</i>	
9	Occupation	41	<i>Sensitive att.</i>	

Lands End

	Attribute	Domain size	Generalizations type	Ht.
1	Zipcode	31953	Round each digit	5
2	Order date	320	Taxonomy tree	3
3	Gender	2	Suppression	1
4	Style	1509	Suppression	1
5	Price	346	Round each digit	4
6	Quantity	1	Suppression	1
7	Shipment	2	Suppression	1
8	Cost	147	<i>Sensitive att.</i>	

Figure 4. Description of Adults and Lands End Databases

Due to space restrictions, we report only a small subset of our experiments. An exhaustive set of experimental results can be found in our technical report [17]; those results are qualitatively similar to the ones we present here.

**Homogeneity Attack.** We illustrate the *homogeneity attack* on a  $k$ -anonymized dataset with the Lands End and Adult databases. For the Lands End Database, we treated the first 5 attributes in Figure 4 as the quasi-identifier. We partitioned the Cost attribute into 147 buckets of size 100 and used this as the sensitive attribute. We then generated all 3-anonymous tables that were minimal with respect to the generalization lattice (i.e. no table at a lower level of generalization was 3-anonymous). There were 3 minimal tables, and 2 of them were vulnerable to the homogeneity attack. In fact, more than 1,000 tuples had their sensitive value revealed. Surprisingly, in each of the vulnerable tables, the average size of a homogeneous group was larger than 100. The table that was not vulnerable to the homogeneity attack was entropy 2.61-diverse.

For the Adult Database, we treated the first 5 attributes in Figure 4 as the quasi-identifier. When we used Occupation as the sensitive attribute, there were a total of 12 minimal 6-anonymous tables, and one of them was vulnerable to the homogeneity attack. On the other hand, when we used Salary Class as the sensitive attribute, there were 9 minimal 6-anonymous tables, and 8 of them were vulnerable. The 9<sup>th</sup> table was recursive (6,2)-diverse. This large value of  $c$  (from the definition of recursive  $(c, \ell)$ -diversity) is due to the distribution of values of the Salary Class attribute: Salary Class is a binary attribute with one value occurring 4 times as frequently as the other.

**Performance.** In our next set of experiments, we compare the running times of entropy  $\ell$ -diversity and  $k$ -anonymity. The results are shown in Figures 5 and 6. For the Adult Database, we used Occupation as the sensitive attribute, and for Lands End we used Cost. We varied the quasi-identifier size from 3 attributes up to 8 attributes; a quasi-identifier of size  $j$  consisted of the first  $j$  attributes of its dataset as listed in Figure 4. We measured the time taken to return all

6-anonymous tables and compared it to the time taken to return all 6-diverse tables. In both datasets, the running times for  $k$ -anonymity and  $\ell$ -diversity were similar. Sometimes the running time for  $\ell$ -diversity was faster, which happened when the algorithm pruned parts of the generalization lattice earlier than it did for  $k$ -anonymity.

**Utility.** The utility of a dataset is a property that is difficult to quantify. As a result, we used three different metrics to gauge the utility of  $\ell$ -diverse and  $k$ -anonymous tables. The first metric, generalization height [16, 21], is the height of an anonymized table in the generalization lattice; intuitively, it is the number of generalization steps that were performed. The second metric is the average size of the  $q^*$ -blocks generated by the anonymization algorithm. The third metric is the *discernibility* metric [6]. The discernibility metric measures the number of tuples that are indistinguishable from each other. Each tuple in a  $q^*$  block  $B_i$  incurs a cost  $|B_i|$  and each tuple that is completely suppressed incurs a cost  $|D|$  (where  $D$  is the original dataset). Since we did not perform any tuple suppressions, the discernibility metric is equivalent to the sum of the squares of the sizes of the  $q^*$ -blocks.

The first graph in Figure 7 shows the minimum generalization height of  $k$ -anonymous and  $\ell$ -diverse tables for  $k, \ell = 2, 4, 6, 8$ . As the graph shows, ensuring diversity in the sensitive attribute does not require many more generalization steps than for  $k$ -anonymity (note that an  $\ell$ -diverse table is automatically  $\ell$ -anonymous); the minimum generalization heights for identical values of  $k$  and  $\ell$  were either the same or differed by one.

Nevertheless, we found that generalization height [21] was not an ideal utility metric because tables with small generalization heights can still have very large group sizes. For example, using full-domain generalization on the Adult Database with 5 quasi-identifiers, we found minimal (with respect to the generalization lattice) 4-anonymous tables that had average group sizes larger than 1,000 tuples. The large groups were caused by data skew. For example, there were only 114 tuples with age between 81 and 90, while

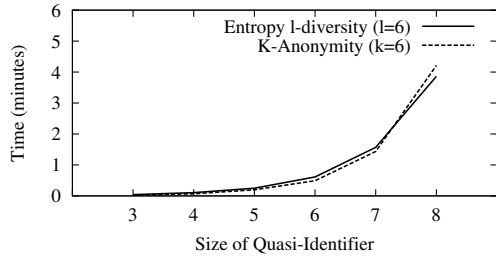


Figure 5. Adults Database

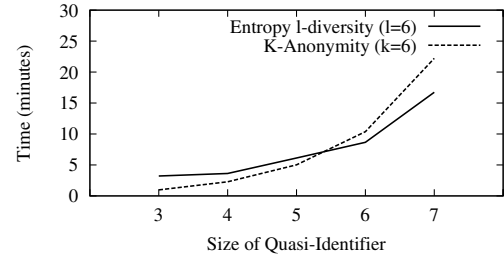


Figure 6. Lands End Database

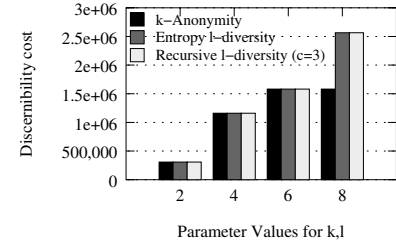
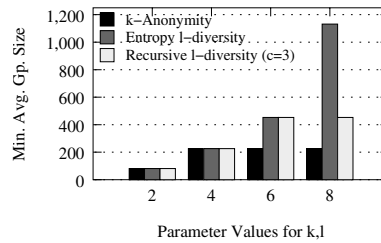
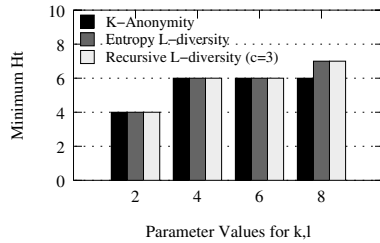


Figure 7. Adults Database.  $Q = \{\text{age, gender, race, marital\_status}\}$

there were 12,291 tuples with age between 31 and 40. So if age groups of length 5 (i.e. [1-5], [6-10], [11-15], etc) were generalized to age groups of length 10 (i.e. [1-10], [11-20], etc), we would end up with very large  $q^*$ -blocks. Generalization hierarchies that are aware of data skew may yield higher quality anonymizations. This is a promising avenue for future work because some recent algorithms [6] can handle certain dynamic generalization hierarchies.

In order to understand the loss of utility due to domain generalization better, we chose to study a subsample of the Adults Database with a lesser data skew in the Age attribute. It turned out that a 5% Bernoulli subsample of the Adult Database suited our requirements – most of the Age values appeared in around 20 tuples each, while only a few values appeared in less than 10 tuples each. The second and third graphs in Figure 7 show the minimum average group size and the discernibility metric cost, respectively, of  $k$ -anonymous and  $\ell$ -diverse tables for  $k, \ell = 2, 4, 6, 8$ . Smaller values for utility metrics represent higher utility. We found that the best  $t$ -anonymous and  $t$ -diverse tables often had comparable utility. We also found that, in some cases,  $\ell$ -diversity had worse utility because some utility must be traded off for privacy. It is interesting to note that recursive  $(3, \ell)$ -diversity permits tables which have better utility than entropy  $\ell$ -diversity. Figure 7 shows that both the instantiations of  $\ell$ -diversity have similar costs for the discernibility metric, but recursive  $\ell$ -diversity permits smaller average group sizes than the entropy definition. Recursive  $(c, \ell)$ -diversity is generally less restrictive than entropy  $\ell$ -diversity, because the extra parameter,  $c$ , allows us to con-

trol how much skew is acceptable in a  $q^*$ -block. Since there is still some residual skew even in our 5% subsample, the entropy definition performs worse than the recursive definition.

## 7. Related Work

The problem of publishing public-use microdata has been extensively studied in both the statistics and computer science communities. The statistics literature, motivated by the need to publish census data, focuses on identifying and protecting the privacy of sensitive entries in contingency tables, or tables of counts which represent the complete cross-classification of the data. Two main approaches have been proposed for protecting the privacy of sensitive cells: *data swapping* and *data suppression*. The data swapping approach involves moving data entries from one cell to another in the contingency table in a manner that is consistent with the set of published marginals [9, 10, 13]. In the data suppression approach [8], cells with low counts are simply deleted, which in turn might lead to the deletion of additional cells. An alternate approach is to determine a *safety range* or *protection interval* for each cell [12], and publish only those marginals which ensure that the feasibility intervals (i.e. upper and lower bounds on the values a cell may take) contain the protection intervals for all the cell entries. The above techniques, however, do not provide a strong theoretical guarantee of the privacy ensured.

Computer science research also has tried to solve the data publishing problem. A technique called  $k$ -anonymity

has been proposed which guarantees that every individual is hidden in a group of size  $k$  with respect to the non-sensitive attributes [24]. It has been shown that the problem of  $k$ -anonymization by suppressing cells in the table is NP-hard [18] and approximation algorithms have been proposed for the same [3]. There has been a lot of study into creating efficient algorithms for  $k$ -anonymity using generalization and tuple suppression techniques [2, 6, 16, 22]. A different formal definition of privacy was proposed for published data based on the notion of *blending in a crowd* in [7]. However, since it is an inter-tuple distance centric measure of privacy, the privacy definition fails to capture scenarios where identification of even a single attribute may constitute a privacy breach.

Query answering techniques are very related to the data publishing approach, where instead of publishing the data, the database answers queries as long as the answers do not breach privacy. There has been work on characterizing the set of views that can be published while keeping some query answer information-theoretically secret [19]. The paper shows that the privacy required is too strong and most interesting queries like aggregates are not allowed to be published. Related techniques in the statistical database literature (see [1] for a survey), especially auditing [15] and output perturbation [11], require maintaining state about the previous queries, while data publishing does not need to maintain any state of the queries asked. The reader is referred to our technical report [17] for a more extensive survey of related work.

## 8. Conclusions and Future Work

In this paper we have shown that a  $k$ -anonymized dataset permits strong attacks due to lack of diversity in the sensitive attributes. We have introduced  $\ell$ -diversity, a framework that gives stronger privacy guarantees.

There are several avenues for future work. First, we want to extend our initial ideas for handling multiple sensitive attributes, and we want to develop methods for continuous sensitive attributes. Second, although privacy and utility are duals of each other, privacy has received much more attention than the utility of a published table. As a result, the concept of utility is not well-understood.

**Acknowledgments.** We thank Joe Halpern for an insightful discussion on the proposed privacy model and Kristen LeFevre for the Incognito source code. This work was partially supported by the National Science Foundation under Grant IIS-0541507, a Sloan Foundation Fellowship, and by a gift from Microsoft. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsors.

## References

- [1] N. R. Adam and J. C. Wortmann. Security-control methods for statistical databases: A comparative study. *ACM Comput. Surv.*, 21(4):515–556, 1989.
- [2] C. C. Aggarwal and P. S. Yu. A condensation approach to privacy preserving data mining. In *EDBT*, pages 183–199, 2004.
- [3] G. Aggarwal, T. Feder, K. Kenthapadi, R. Motwani, R. Panigrahy, D. Thomas, and A. Zhu.  $k$ -anonymity: Algorithms and hardness. Technical report, Stanford University, 2004.
- [4] R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules in Large Databases. In *VLDB*, 1994.
- [5] F. Bacchus, A. J. Grove, J. Y. Halpern, and D. Koller. From statistical knowledge bases to degrees of belief. *A.I.*, 87(1-2), 1996.
- [6] R. J. Bayardo and R. Agrawal. Data privacy through optimal  $k$ -anonymization. In *ICDE-2005*, 2005.
- [7] S. Chawla, C. Dwork, F. McSherry, A. Smith, and H. Wee. Toward privacy in public databases. In *TCC*, 2005.
- [8] L. H. Cox. Suppression, methodology and statistical disclosure control. *Journal of the American Statistical Association*, 75, 1980.
- [9] T. Dalenius and S. Reiss. Data swapping: A technique for disclosure control. *Journal of Statistical Planning and Inference*, 6, 1982.
- [10] P. Diaconis and B. Sturmfels. Algebraic algorithms for sampling from conditional distributions. *Annals of Statistics*, 1:363–397, 1998.
- [11] I. Dinur and K. Nissim. Revealing information while preserving privacy. In *PODS*, pages 202–210, 2003.
- [12] A. Dobra. *Statistical Tools for Disclosure Limitation in Multiway Contingency Tables*. PhD thesis, Carnegie Mellon University, 2002.
- [13] G. T. Duncan and S. E. Feinberg. Obtaining information while preserving privacy: A markov perturbation method for tabular data. In *Joint Statistical Meetings*, Anaheim, CA, 1997.
- [14] A. Evfimievski, J. Gehrke, and R. Srikant. Limiting privacy breaches in privacy preserving data mining. In *PODS*, 2003.
- [15] K. Kenthapadi, N. Mishra, and K. Nissim. Simulatable auditing. In *PODS*, 2005.
- [16] K. LeFevre, D. DeWitt, and R. Ramakrishnan. Incognito: Efficient fulldomain  $k$ -anonymity. In *SIGMOD*, 2005.
- [17] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkitasubramaniam.  $\ell$ -diversity: Privacy beyond  $k$ -anonymity. Available at <http://www.cs.cornell.edu/~mvnarak>, 2005.
- [18] A. Meyerson and R. Williams. On the complexity of optimal  $k$ -anonymity. In *PODS*, 2004.
- [19] G. Miklau and D. Suciu. A formal analysis of information disclosure in data exchange. In *SIGMOD*, 2004.
- [20] U.C.Irvine Machine Learning Repository. <http://www.ics.uci.edu/mllearn/mlrepository.html>.
- [21] P. Samarati. Protecting respondents' identities in microdata release. In *IEEE Transactions on Knowledge and Data Engineering*, 2001.
- [22] P. Samarati and L. Sweeney. Protecting privacy when disclosing information:  $k$ -anonymity and its enforcement through generalization and suppression. Technical report, CMU, SRI, 1998.
- [23] L. Sweeney. Uniqueness of simple demographics in the u.s. population. Technical report, Carnegie Mellon University, 2000.
- [24] L. Sweeney.  $k$ -anonymity: a model for protecting privacy. *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 10(5):557–570, 2002.
- [25] S. Zhong, Z. Yang, and R. N. Wright. Privacy-enhancing  $k$ -anonymization of customer data. In *PODS*, 2005.