

COSC130

Fundamentals of Cybersecurity and Privacy

LECTURE 5: INTRODUCTION TO PRIVACY

Lecture Overview

1. What is Privacy?
 - Meaning of the term “privacy”
 - History of Privacy
 - Psychological Aspects of Privacy
2. Abstract Database Model, Compromise and Basic Privacy Techniques
3. Attack Models
4. k-anonymity
5. Differential Privacy

Much of this lecture is based on

[Fung et al., 2011] B. C. M. Fung, K. Wang, A. W.-C. Fu and P. S. Yu, “Introduction to Privacy-Preserving Data Publishing - Concepts and Techniques”, *CRC Press*, Tylor & Francis Group, 2011.

In-text references to this source are typically omitted for readability.

What is Privacy?

DICTIONARY MEANING

Oxford Dictionary:

“A state in which one is not observed or disturbed by other people.”

“The state of being free from public attention.”

dictionary.com:

“The state of being apart from other people or concealed from their view; solitude; seclusion.”

“The state of being free from unwanted or undue intrusion or disturbance in one's private life or affairs; freedom to be let alone.”

What is Privacy – Legal meaning

Common Law

“The right of people to lead their lives in a manner that is reasonably secluded from public scrutiny, whether such scrutiny comes from a neighbor's prying eyes, an investigator's eavesdropping ears, or a news photographer's intrusive camera; and in statutory law, the right of people to be free from unwarranted drug testing and electronic surveillance.” (*legal-dictionary/thefreedictionary.com*)

Australian privacy law and practice (ALRC Report 108):

Information privacy, which involves the establishment of rules governing the collection and handling of personal data such as credit information, and medical and government records. It is also known as ‘data protection’;

Bodily privacy, which concerns the protection of people’s physical selves against invasive procedures such as genetic tests, drug testing and cavity searches;

Privacy of communications, which covers the security and privacy of mail, telephones, e-mail and other forms of communication; and

Territorial privacy, which concerns the setting of limits on intrusion into the domestic and other environments such as the workplace or public space. This includes searches, video surveillance and ID checks.

What is Privacy?

In this unit we are interested in information privacy.

We adopt the following definition: “Privacy may be defined as the claim of individuals, groups or institutions to determine when, how and to what extent information about them is communicated to others.” (Alan F. Westin, *Privacy and Freedom*, New York: Atheneum, 1967, page 7)

In other words, privacy is right of individuals to control personal information about themselves.

History of Privacy

Privacy as we know it is only 150 years old – often referred to as “luxury goods”.

However, there is evidence of privacy attitudes and behaviours across different cultures including Ancient Rome and Greece, preindustrial Javanese, Balinese and Tuareg society (Acquisti et al., 2015).

Privacy is also mentioned in different religions (Acquisti et al., 2015):

- **Bible (Genesis 3.7)**

“And the eyes of them both were opened, and they knew that they were naked; and they sewed fig leaves together, and made themselves aprons.”

- **Quran (49.12)**

“And do not spy or backbite each other.”

- **Talmud (Bava Batra 60a)** instructs homebuilders to position the windows in such a way that they do not directly face windows of their neighbours.

Psychological Aspect of Privacy

In his paper “Some Psychological Aspects of Privacy”, 1966, Sidney M. Jourard (1926-1974), identifies the following issues and solutions:

- the pressures to conform (punishment or invalidation)
- conformity and health (repression)
- the therapeutic and socially necessary functions of privacy
- the psychological function of self-disclosure and concealment
- disclosure inhibited and privacy denied
 - the social risk of private places
 - institutional life (no privacy implies conformity and no individuality)
 - “hell is other people” (changelessness)
 - opting out: “the beatniks”
- possible solutions
 - Individual stratagems and check-out places
 - Education for private life

In short, privacy is experienced as “room to grow in,” as freedom from interference, and as freedom to explore, to pursue experimental projects in science, art, work, play, and living. In the name of the *status quo* and other, even more attractive goals, privacy may be eroded. But without privacy and its concomitant, freedom, the cost to be paid for the ends achieved—in terms of lost health, weak commitment to the society, and social stagnation—may be too great.

Abstract database Model

Name	City	Age	Gender	Status	Post-traumatic stress disorder	Attempted suicide
White	Sydney	34	F	W	4.1	no
Scarlet	Dubbo	27	F	D	3.9	no
Brown	Sydney	45	M	M	4.3	no
Mustard	Perth	32	M	S	2.1	yes
Green	Ballina	76	NB	M	4.8	no
Green	Darwin	32	F	M	4.6	no
Plum	Hobart	25	M	D	2.9	no
Mustard	Darwin	24	M	W	4.2	no
White	Dubbo	51	F	D	3.8	no
Peacock	Sydney	40	NB	M	4.1	no
Black	Ballina	68	F	W	3.6	no
Violet	Dubbo	33	F	M	2.7	no
Aureate	Sydney	28	F	S	3.5	no

Database Compromise

Example 1:

COUNT(City=Darwin and Sex=M and Age<30) = 1

COUNT(City=Darwin and Sex=M and Age<30 and AS=no) = 1

AVG(City=Darwin and Sex=M and Age<30; PTSD) = 4.2

Example 2:

COUNT(City=Sydney and Age<37) = 2

COUNT(City=Sydney and Age<37 and AS=no) = 2

Basic Privacy Techniques

Restriction

- query set size control
- query set overlap control
- maximum order control
- partitioning
- cell suppression
- auditing

Modification

- data perturbation
- response perturbation
- data swapping (shuffling)
- random sample

Published Data Table

Name *	City	Age	Gen der	Status	Post- traumatic stress disorder	Attempt ed suicide
White	Sydney	34	F	W	4.1	no
Scarlet	Dubbo	27	F	D	3.9	no
Brown	Sydney	45	M	M	4.3	no
Mustard	Perth	32	M	S	2.1	yes
Green	Ballina	76	NB	M	4.8	no
Green	Darwin	32	F	M	4.6	no
Plum	Hobart	25	M	D	2.9	no
Mustard	Darwin	54	M	W	4.2	no
White	Dubbo	51	F	D	3.8	no
Peacock	Sydney	40	NB	M	4.1	no
Black	Ballina	68	F	W	3.6	no
Violet	Dubbo	33	F	M	2.7	no
Aureate	Sydney	28	F	S	3.5	no

	ID - unique identifier ID={Name}
	QID - Quasi identifier QID={City, Age, Sex }
	Non-sensitive attributes
	Sensitive Attribute

* Strictly speaking, name itself can rarely be considered to be a unique identifier.

Attack Models

We can classify the main attack types into 2 broad categories:

1. Linkage Attack Models:

- 1) **Record linkage**, where an intruder is able to link an individual to a record in the published data table.
- 2) **Attribute linkage**, where an intruder is able to link an individual to a sensitive value in the published data table.
- 3) **Table linkage**, where an intruder is able to link an individual to the published data table itself.

2. **Probabilistic attack**. Ideally, the published data should reveal to an intruder as little additional knowledge about individuals as possible, beyond what he/she already knew before seeing the data (background knowledge, or supplementary knowledge). *Probabilistic attack* occurs when the difference between the prior and the posterior knowledge regarding an individual is “significant”.

Record Linkage

The intruder is able to link an individual to a record in the published data table.

Recall that in published data tables Unique Identifiers (UIs) are typically removed, so record linkage typically relies on QIDs.

Suppose that an individual A , which the intruder is after, has a value qid of the QID, and that the value qid is known to the intruder.

In general, qid identifies a group of records in the table. If the size of the group is 1, we have record linkage.

If the size of the group is more than 1, an intruder may still be able to uniquely identify A with the help of additional knowledge.

Record Linkage Example 1

<i>Name *</i>	<i>City</i>	<i>Age</i>	<i>Gender</i>	<i>Status</i>	<i>Post-traumatic stress disorder</i>	<i>Attempted suicide</i>
<i>White</i>	<i>Sydney</i>	<i>34</i>	<i>F</i>	<i>W</i>	<i>4.1</i>	<i>no</i>
<i>Scarlet</i>	<i>Dubbo</i>	<i>27</i>	<i>F</i>	<i>D</i>	<i>3.9</i>	<i>no</i>
<i>Brown</i>	<i>Sydney</i>	<i>45</i>	<i>M</i>	<i>M</i>	<i>4.3</i>	<i>no</i>
<i>Mustard</i>	<i>Perth</i>	<i>32</i>	<i>M</i>	<i>S</i>	<i>2.1</i>	<i>yes</i>
<i>Green</i>	<i>Ballina</i>	<i>76</i>	<i>NB</i>	<i>M</i>	<i>4.8</i>	<i>no</i>
<i>Green</i>	<i>Darwin</i>	<i>32</i>	<i>F</i>	<i>M</i>	<i>4.6</i>	<i>no</i>
<i>Plum</i>	<i>Hobart</i>	<i>25</i>	<i>M</i>	<i>D</i>	<i>2.9</i>	<i>no</i>
<i>Mustard</i>	<i>Darwin</i>	<i>24</i>	<i>M</i>	<i>W</i>	<i>4.2</i>	<i>no</i>
<i>White</i>	<i>Dubbo</i>	<i>51</i>	<i>F</i>	<i>D</i>	<i>3.8</i>	<i>no</i>
<i>Peacock</i>	<i>Sydney</i>	<i>40</i>	<i>NB</i>	<i>M</i>	<i>4.1</i>	<i>no</i>
<i>Black</i>	<i>Ballina</i>	<i>68</i>	<i>F</i>	<i>W</i>	<i>3.6</i>	<i>no</i>
<i>Violet</i>	<i>Dubbo</i>	<i>33</i>	<i>F</i>	<i>M</i>	<i>2.7</i>	<i>no</i>
<i>Aureate</i>	<i>Sydney</i>	<i>28</i>	<i>F</i>	<i>S</i>	<i>3.5</i>	<i>no</i>

$QID = \{City, Age, Sex\}$

$ID(A) = \text{Scarlet}$

$qid = QID(A) = \{Dubbo, 27, F\}$

Record Linkage Example 2

Table 2.2: Original patient data

Job	Sex	Age	Disease
Engineer	Male	35	Hepatitis
Engineer	Male	38	Hepatitis
Lawyer	Male	38	HIV
Writer	Female	30	Flu
Writer	Female	30	HIV
Dancer	Female	30	HIV
Dancer	Female	30	HIV

A hospital intends to release the Patient Data table (Table 2.2.) to a researcher.

The researcher already has access to an external data (Table 2.3.). Additionally, they know that every patient in Table 2.2. also has a record in Table 2.3.

Table 2.3: External data

Name	Job	Sex	Age
Alice	Writer	Female	30
Bob	Engineer	Male	35
Cathy	Writer	Female	30
Doug	Lawyer	Male	38
Emily	Dancer	Female	30
Fred	Engineer	Male	38
Gladys	Dancer	Female	30
Henry	Lawyer	Male	39
Irene	Dancer	Female	32

What are ID, QID, non-sensitive and sensitive attributes in each table?

What can the researcher learn by linking these two tables?

k-anonymity [Sweeney, 2002]

In her famous 2002 paper [Sweeney, 2002], Sweeny showed that 87% of respondents in 1990 US census (216,000,000) can be uniquely identified using only 3 attributes:

- ZIP code
- Date of Birth
- Gender

In the same paper she famously demonstrated how linking different data sets can be used to compromise sensitive information about individuals.

The National Association of Health Data Organizations (NAHDO) reported that 37 states in the USA have legislative mandates to collect hospital level data and that 17 states have started collecting ambulatory care data from hospitals, physicians offices, clinics, and so forth [2]. The leftmost circle in Figure 1 contains a subset of the fields of information, or *attributes*, that NAHDO recommends these states collect; these attributes include the patient's ZIP code, birth date, gender, and ethnicity.

k-anonymity [Sweeney, 2002]

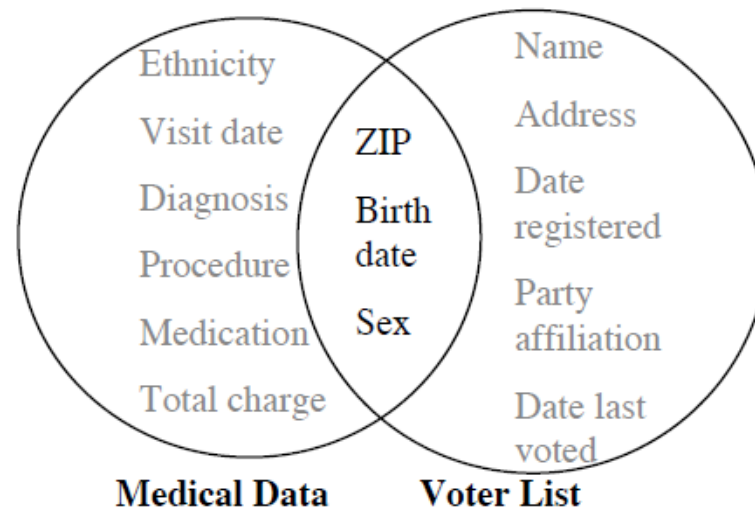


Figure 1 Linking to re-identify data

In Massachusetts, the Group Insurance Commission (GIC) is responsible for purchasing health insurance for state employees. GIC collected patient-specific data with nearly one hundred attributes per encounter along the lines of the those shown in the leftmost circle of Figure 1 for approximately 135,000 state employees and their families. Because the data were believed to be anonymous, GIC gave a copy of the data to researchers and sold a copy to industry [3].

k-anonymity

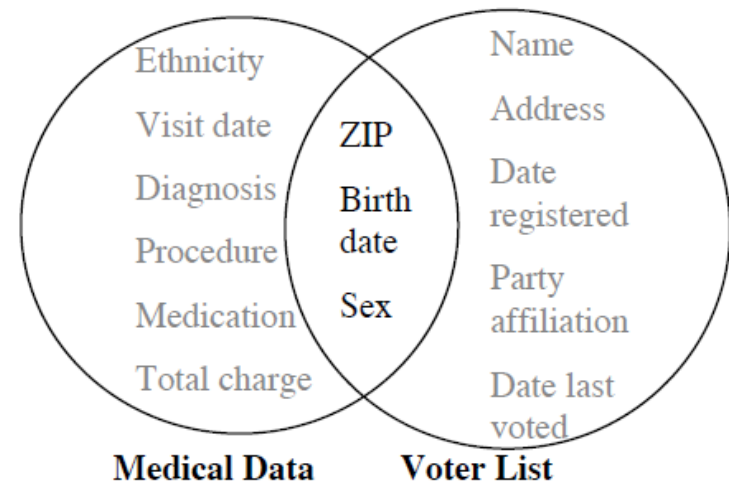


Figure 1 Linking to re-identify data

For twenty dollars I purchased the voter registration list for Cambridge Massachusetts and received the information on two diskettes [4]. The rightmost circle in Figure 1 shows that these data included the name, address, ZIP code, birth date, and gender of each voter. This information can be linked using ZIP code, birth date and gender to the medical information, thereby linking diagnosis, procedures, and medications to particularly named individuals.

For example, William Weld was governor of Massachusetts at that time and his medical records were in the GIC data. Governor Weld lived in Cambridge Massachusetts. According to the Cambridge Voter list, six people had his particular birth date; only three of them were men; and, he was the only one in his 5-digit ZIP code.

k-anonymity

In a series of papers ([Sweeney, 2002], [Samarati et al., 1998]) Samarati and Sweeney proposed *k-anonymity* in order to prevent record linkage.

For each value *qid* of QID that exist in the data table, there are at least *k* record having value *qid* in QID.

If a table satisfies this requirement, we say it is *k-anonymous*.

In a *k-anonymous* table, a probability of successfully linking a record to another table on QID is at most $\frac{1}{k}$.

k-anonymity: Example [Fung et al., 2011]

Table 2.2: Original patient data

Job	Sex	Age	Disease
Engineer	Male	35	Hepatitis
Engineer	Male	38	Hepatitis
Lawyer	Male	38	HIV
Writer	Female	30	Flu
Writer	Female	30	HIV
Dancer	Female	30	HIV
Dancer	Female	30	HIV

Table 2.3: External data

Name	Job	Sex	Age
Alice	Writer	Female	30
Bob	Engineer	Male	35
Cathy	Writer	Female	30
Doug	Lawyer	Male	38
Emily	Dancer	Female	30
Fred	Engineer	Male	38
Gladys	Dancer	Female	30
Henry	Lawyer	Male	39
Irene	Dancer	Female	32

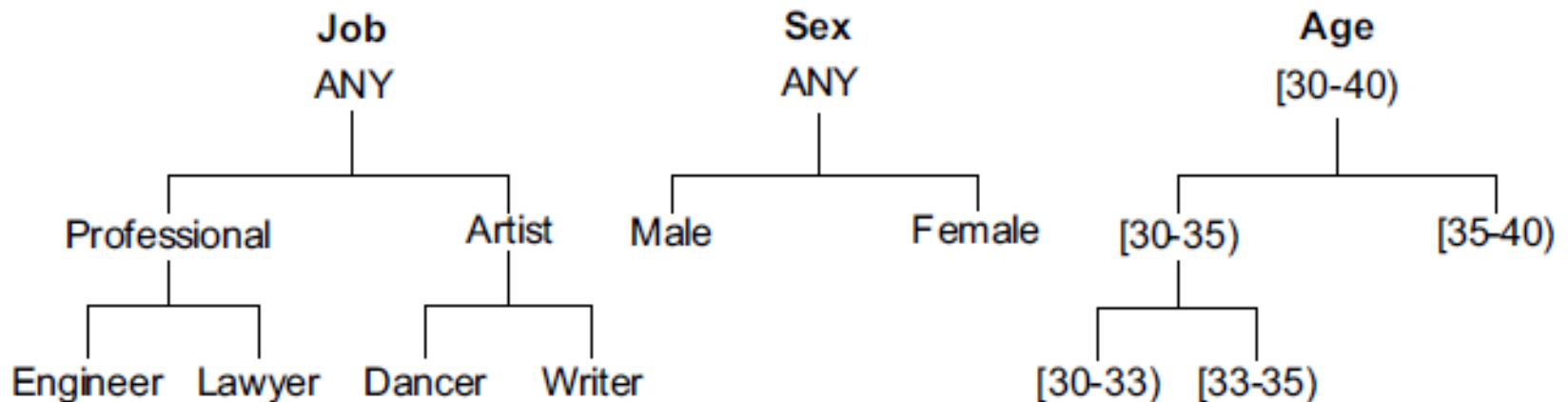


FIGURE 2.1: Taxonomy trees for *Job*, *Sex*, *Age*

k-anonymity: Example [Fung et al., 2011]

Table 2.2: Original patient data

Job	Sex	Age	Disease
Engineer	Male	35	Hepatitis
Engineer	Male	38	Hepatitis
Lawyer	Male	38	HIV
Writer	Female	30	Flu
Writer	Female	30	HIV
Dancer	Female	30	HIV
Dancer	Female	30	HIV

Table 2.4: 3-anonymous patient data

Job	Sex	Age	Disease
Professional	Male	[35-40)	Hepatitis
Professional	Male	[35-40)	Hepatitis
Professional	Male	[35-40)	HIV
Artist	Female	[30-35)	Flu
Artist	Female	[30-35)	HIV
Artist	Female	[30-35)	HIV
Artist	Female	[30-35)	HIV

Table 2.3: External data

Name	Job	Sex	Age
Alice	Writer	Female	30
Bob	Engineer	Male	35
Cathy	Writer	Female	30
Doug	Lawyer	Male	38
Emily	Dancer	Female	30
Fred	Engineer	Male	38
Gladys	Dancer	Female	30
Henry	Lawyer	Male	39
Irene	Dancer	Female	32

Table 2.5: 4-anonymous external data

Name	Job	Sex	Age
Alice	Artist	Female	[30-35)
Bob	Professional	Male	[35-40)
Cathy	Artist	Female	[30-35)
Doug	Professional	Male	[35-40)
Emily	Artist	Female	[30-35)
Fred	Professional	Male	[35-40)
Gladys	Artist	Female	[30-35)
Henry	Professional	Male	[35-40)
Irene	Artist	Female	[30-35)

ϵ -differential privacy

In 2006 Dwork proposed a new model that she termed “different privacy”, based on the following principle: “the risk to the record owner’s privacy should not substantially increase as a result of participating in a statistical database”.

Importantly, unlike other models, differential privacy does not consider prior and posterior knowledge.

Instead, Dwork argues if the responses from the data table are the same or very similar with and without one particular record, there is no (are at least not significant) privacy risk for that individual.

In other words, if a data set is protected using differential privacy model, then removing or adding a single record will not affect the results of any analysis very much. Therefore, data linkage will not pose a risk to privacy.

Therefore, if an individual chooses not to provide their record, the responses from the data set will remain largely unchanged.

ϵ -differential privacy

It is important to note that ϵ -differential privacy does not prevent record and attribute linkages which are prevented by k -anonymity and similar methods.

On the other hand, ϵ -differential privacy signals to an individual that if they allow their record to be added to the dataset, then nothing, or almost nothing, can be discovered from the dataset that could not have been discovered without their record.

From this we see that ϵ -differential privacy prevents table linkage.

References

[Acquisti et al., 2015] A. Acquisti, L. Brandimarte and G. Loewenstein. “Privacy and human behavior in the age of information”, *Science*, 347(6621), 2015.

[Dwork, 2006] C. Dwork. “Differential privacy”, *In Proc. of the 33rd International Colloquium on Automata, Languages and Programming (ICALP)*, pages 1–12, Venice, Italy, 2006.

[Fung et al., 2011] B. C. M. Fung, K. Wang, A. W.-C. Fu and P. S. Yu, “Introduction to Privacy-Preserving Data Publishing - Concepts and Techniques”, *CRC Press*, Tylor & Francis Group, 2011.

[Jourard, 1966] S. M. Jourard. “Some Psychological Aspects of Privacy”, *Law and Contemporary Problems*, vol. 31, no. 2, 1966, pp. 307–18. JSTOR, <https://doi.org/10.2307/1190673>. Accessed 12 Sept. 2023.

[Samarati et al., 1998] P. Samarati and L. Sweeney. “Generalizing data to provide anonymity when disclosing information”, *In Proc. of the 17th ACM SIGACTSIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, page 188, Seattle, WA, June 1998.

[Sweeney, 2002] L. Sweeney, “k-anonymity: a model for protecting privacy”, *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 10(5), 557-570, 2002.