
dedication (optional)

Summary

Write your summary here...

Preface

Write your preface here...

Table of Contents

Summary	i
Preface	ii
Table of Contents	iv
List of Tables	v
List of Figures	vii
Abbreviations	viii
1 Introduction	1
1.1 Equations	1
1.2 Figures	1
1.3 References	2
1.4 Tables	2
2 Background and Motivation	5
2.1 Concepts and Expressions	5
2.2 Privacy Breaches	5
2.2.1 Netflix Prize Competition	5
2.2.2 Group Insurance Commission	6
2.2.3 New York Taxi dataset	6
3 Basic Theory	7
3.1 Differential Privacy	7
3.1.1 Definition of Differential Privacy	8
3.1.2 Noise Mechanisms	8
4 Experiment	11

5	Analysis	13
6	Conclusion	15
	Bibliography	17
	Appendix	19

List of Tables

1.1	Table 1.	3
-----	------------------	---

List of Figures

1.1	Pikachu.	2
-----	------------------	---

Abbreviations

Symbol = definition

Chapter 1

Introduction

1.1 Equations

To write an equation

```
\begin{eqnarray}\label{eq1}  
F = m \times a  
\end{eqnarray}
```

This will produceasdasd asdf asdf asdf

$$F = m \times a \tag{1.1}$$

To refer to the equation

```
\eqref{eq1}
```

This will produce (1.1).

1.2 Figures

To create a figure

```
\begin{figure}[h!]  
  \centering  
    \includegraphics[width=0.5\textwidth]{fig/pikachu}  
    \caption{Pikachu.}  
  \label{fig1}  
\end{figure}
```

To refer to the figure

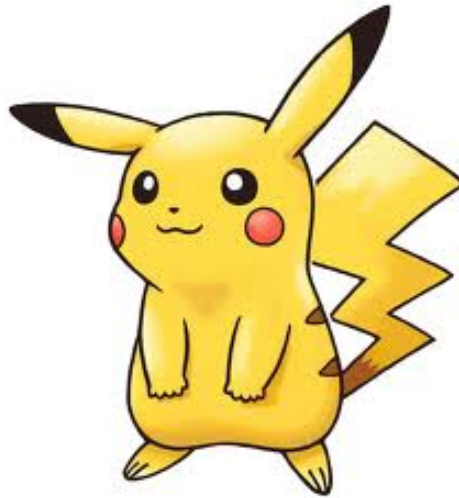


Figure 1.1: Pikachu.

```
\textbf{Fig. \ref{fig1}}
```

This will produce **Fig. 1.1**

1.3 References

To cite references

```
\cite{1,2,3}
```

or

```
\citep{1,2,3}
```

This will produce: [? ? ?] or [? ? ?], respectively.

1.4 Tables

To creat a table

```
\begin{table}[!h]
\begin{center}
\begin{tabular}{| l | l | l | l |}
\hline
\textbf{No.} & \textbf{Data 1} & \textbf{Data 2} & \\ \hline
1 & a1 & b1 & \\ \hline
\end{tabular}
\end{center}
\end{table}
```

```
2 & a2 & b2 \\ \hline
\end{tabular}
\end{center}
\caption{Table 1.}
\label{Tab1}
\end{table}
```

This will produce

No.	Data 1	Data 2
1	a1	b1
2	a2	b2

Table 1.1: Table 1.

To refer to the table

```
\textbf{Table. \ref{Tab1}}
```

This will produce **Table. 1.1.**

Background and Motivation

In this section we will first explain some basic concepts and expressions that are used in the privacy context such as anonymization operations and Personally Identifiable Information (PII). Then we will have a look at some classic examples of failure to preserve privacy when data publishing and how these attacks motivated us to choose our topic for this thesis.

2.1 Concepts and Expressions

2.2 Privacy Breaches

In the recent years there have been many failures in privacy preserving data publishing. Many companies have been faced with a PR disaster after releasing data about their customers thinking them being anonymized, only to have people de-anonymize their data and breaching the privacy of the datasets' participants. In this section we will have a look at some of these privacy failures.

2.2.1 Netflix Prize Competition

Netflix, the world's largest online movie streaming website, decided in 2006 to crowd-source a new movie suggestion algorithm and offered a cash prize of 1 million dollar for the most efficient algorithm. To help the research, they released 100 million supposedly anonymized movie ratings from their own database. Each record included an anonymized user ID, some information about the movie and the user's rating information. Two researchers from the University of Texas demonstrated that an adversary who knows only a little bit about an individual subscriber can easily identify this subscriber's record in the dataset [6]. Using the publicly available dataset from the Internet Movie Database (IMDB) as the source of background knowledge, they matched certain subscribers with their Netflix records, and uncovered their apparent political preferences and other potentially sensitive information[7].

2.2.2 Group Insurance Commission

In 1997, Latanya Sweeney wrote a paper on how she had identified the medical records of Massachusetts governor William Weld based on publicly available information from the database of Group Insurance Commission. She achieved this analyzing data from a public voter list, and linked it with patient-specific medical data through a combination of birth date, zip code, and gender[9]. As these columns were similar in both databases, their combination could be used to identify medical records that belong to either one person, or a small group of people. Sweeney hypothesized that 87% of the US population could be identified by having the combination of the three aforementioned records. It's worth noting here that this theory is not conclusive. A paper by Daniel Barth-Jones suggests that the re-identification of Weld may have been a fluke due to his public figure, and that ordinary people risk of identification is much lower[1].

2.2.3 New York Taxi dataset

The New York City Taxi and Limousine Commission released a dataset in 2013 containing details about every taxi ride that year, including pickup and dropoff times, location, fare, as well as anonymized (hashed) versions of the taxi's license and medallion numbers. Vijay Punduranga, a researcher for Google, wrote a blog-post where he showed how he exploited a vulnerability in the hashing-function to re-identify the drivers. He then showed how this could be potentially used to calculate any driver's personal income[8].

Another researcher, called Anthony Tockar, wrote an article during his internship at Neustar Research where he proved that the dataset also contained an inherent privacy risk to the passengers which had been riding New York Taxis. Even though there was no information in the dataset on who had been riding the taxis, Tockar showed that by using auxiliary information such as timestamped pictures, he could stalk celebrities and figure out to where they were driving, and how much they tipped the driver. He also used map data from Google Maps to create a map of dropoff locations for people that had exited a late night visit from gentleman's club and taken a cab home. He then used websites like Spokeo and Facebook to find the cab customer's ethnicity, relationship status, court records, and even a profile picture[10].

Basic Theory

In a world where massive amounts of sensitive personal data are being collected, attacks on the individual's privacy are becoming more and more of a threat. One type of attack is the identification of an individual's personal information from massive data sets, such as people's movie ratings from the Netflix data set[7], and the medical records of a former governor of Massachusetts[1]. These types of privacy breaches may lead to the unwanted discovery of a person's embarrassing information, and could also lead to the theft of an individual's private data or identity.

Many different approaches have been tried by data custodians to privatize the data they hold, such as removing any columns containing Personally Identifiable Information (PII), anonymizing the data by providing k-anonymity protection[9], or perform group based anonymization through l-diversity[4]. All of these methods mentioned have been proved to be susceptible in some way or form to attacks [3]. Motivated by these shortcomings, a researcher at Microsoft came up with a data theoretical framework called differential privacy, which operates off a solid mathematical foundation and have strong theoretical guarantees on the privacy and utility of the released data.

3.1 Differential Privacy

The term "differential privacy" was defined by Dwork as a description of a promise, made by a data holder to a data subject: "You will not be affected, adversely or otherwise, by allowing your data to be used in any study or analysis, no matter what other studies, data sets, or information sources, are available." [2] In an ideal situation, databases which implement differential privacy mechanisms can make confidential data widely available for accurate data analysis, without resorting to data usage agreements, data protection plans, or restricted views. Nevertheless, the Fundamental Law of Information Recovery states that overly accurate answers to too many questions will destroy privacy in a spectacular way [2], meaning that data utility will eventually be consumed.

3.1.1 Definition of Differential Privacy

The classic example for explaining a security breach is the case of Mr White: Suppose you have access to a database that allows you to compute the income of all residents in a specified area. If you knew that Mr White was going to move, simply querying the database before and after his relocation would allow you to deduce his income.

Definition 1: a mechanism \tilde{f} is a random function that takes a dataset D as input, and outputs a random variable $\tilde{f}(D)$.

Definition 2: the distance of two datasets, $d(D, D')$, denotes the minimum number of sample changes that are required to change D into D' .

Formally, differential privacy is defined as follows: A randomized function f gives ϵ -differential privacy if for all data sets D and D' differing on at most one row, and all $S \subseteq \text{Range}(f)$,

$$\Pr[\tilde{f}(D) \in S] \leq e^{\epsilon} \times \Pr[\tilde{f}(D') \in S]$$

What this means is that the risk to an individual's privacy should not be substantially increase as a result of participating in a statistical database, as the risk is bounded by the parameter ϵ . Therefore an attacker should not be able to learn anything about any participant that they would not have learned if the participant had opted out of participating.

3.1.2 Noise Mechanisms

This is guaranteed by applying noise to the result of an query to the dataset, by using the function $f()$.

Definition 3: A query f is a function that takes a dataset as input. The answer to the query f is denoted $f(D)$.

There are many different mechanisms for applying this noise, but the two most common are the Laplace mechanism and the Exponential mechanism.

Laplace Mechanism

The Laplace mechanism involves adding random noise which follows the Laplace statistical distribution. The most common question that needs to be answered before doing research with differentially private data, is how we should define our Laplace random variable, i.e how much noise needs to be added. The Laplace distribution centered around zero has only one parameter, its scale, and this is proportional to its standard deviation. The scale is naturally dependent on the privacy parameter ϵ , and also on the risk of the most different individual having their private information leaked from the data. This risk is called the sensitivity of the query, and is defined mathematically as:

$$\Delta f = \max_{D, D'} ||f(D) - f(D')||_1$$

This equation states that the maximum difference in the values that the query f may take is on a pair of databases that differ on only one row. Dwork proved that adding a random Laplace variable, $(\Delta f / \epsilon)$, to a query you could guarantee ϵ -differential privacy[2].

Exponential Mechanism

The Exponential mechanism proposed by McSherry and Talwar [5] is a method for selecting one element from a set, and is commonly used if a non-numeric value query is used. An example would be: "What is the most common eye color in this room?". Here it would not make sense to perturb the answer by adding noise drawn from the Laplace distribution. The idea of the mechanism is to select the output from all the possible answers at random, with the probability of selecting a particular output being higher for those outputs that are "closer" to the true output.

More formally, let A be the range of possible outputs for the query function f . Also, let $u_f(D, a)$ be a utility function that measures how good an output $a \in A$ is as an answer to the query function f given that the input dataset is D (Note that higher values of u_f represents better outputs). The sensitivity function will then be defined as the maximum possible change in the utility function's value u_f due to the addition or removal of one person's data from the input, i.e:

Definition 4: the sensitivity of score function H is defined as

$$S(u_f) = \max_{(D, D')=1, a \in A} ||u_f(D, a) - u_f(D', a)||$$

Chapter 4

Experiment

Chapter 5

Analysis

Chapter 6

Conclusion

Bibliography

- [1] Barth-Jones, D. C., 2012. The 're-identification' of governor william weld's medical information: a critical re-examination of health data identification risks and privacy protections, then and now. Then and Now (June 4, 2012).
- [2] Dwork, C., Roth, A., 2013. The algorithmic foundations of differential privacy. *Theoretical Computer Science* 9 (3-4), 211–407.
- [3] Ganta, S. R., Kasiviswanathan, S. P., Smith, A., 2008. Composition attacks and auxiliary information in data privacy. In: *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, pp. 265–273.
- [4] Machanavajjhala, A., Kifer, D., Gehrke, J., Venkitasubramaniam, M., 2007. l-diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 1 (1), 3.
- [5] McSherry, F., Talwar, K., Oct 2007. Mechanism design via differential privacy. In: *Foundations of Computer Science, 2007. FOCS '07. 48th Annual IEEE Symposium on*. pp. 94–103.
- [6] Narayanan, A., Shmatikov, V., 2006. How to break anonymity of the netflix prize dataset. *CoRR abs/cs/0610105*.
URL <http://arxiv.org/abs/cs/0610105>
- [7] Narayanan, A., Shmatikov, V., 2008. Robust de-anonymization of large sparse datasets. In: *Security and Privacy, 2008. SP 2008. IEEE Symposium on*. IEEE, pp. 111–125.
- [8] Pandurangan, V., Jun. 2014. On taxis and rainbows: Lessons from nycs improperly anonymized taxi logs.
URL <https://medium.com/@vijayp/of-taxis-and-rainbows-f6bc289679a1>
- [9] Sweeney, L., 2002. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10 (05), 557–570.

[10] Tockar, A., Sep. 2014. Riding with the stars: Passenger privacy in the nyc taxicab dataset.

URL <http://research.neustar.biz/2014/09/15/>

[riding-with-the-stars-passenger-privacy-in-the-nyc-taxicab-dataset/](http://research.neustar.biz/2014/09/15/riding-with-the-stars-passenger-privacy-in-the-nyc-taxicab-dataset/)

Appendix

Write your appendix here...