

# $(\ell, \delta)$ -Diversity: Linkage-Robustness via a Composition Theorem

Arvind Rameshwar

Dept. of EE, IIT Madras

in collaboration with

Anshoo Tandon

Centre of Data for Public Good, IISc

HyTS 2026



DATA FOR  
PUBLIC GOOD

# Some background: Privacy

- ▶ Privacy of users is a fundamental desideratum in the release of datasets for statistical analysis.

On Taxis and Rainbow Tables: Lessons for researchers and governments from NYC's improperly anonymized taxi logs.

Sema Williams  
July 16th, 2014

1 comment Estimated reading time: 5 minutes

When New York City's Taxi and Limousine Commission made publicly available 20GB worth of trip and fare logs, many welcomed the vast trove of open data. Unfortunately, prior to being widely shared, the personally identifiable information had not been anonymized properly. [Vijay Pandurangan](#) describes the structure of the

## Weaving Technology and Policy Together to Maintain Confidentiality

Latanya Sweeney

Organizations often release and receive medical data with all explicit identifiers, such as name, address, telephone number, and Social Security number (SSN), removed on the assumption that patient confidentiality is maintained because the resulting data look anonymous. However, in most of these cases, the remaining data can be used to reidentify individuals by linking or matching the data to other data bases or by looking at unique characteristics found in the fields and records of the data base itself. When these less apparent aspects are taken into account, each released record can map to many possible people, providing a level of anonymity that the record-holder determines. The greater the number of candidates per record, the more anonymous the data.

I examine three general-purpose computer programs for maintaining patient confidentiality when disclosing electronic medical records: the Scrub System, which locates and suppresses or replaces personally identifying information in letters between doctors and in notes written by clinicians; the Duality System, which generalizes values based on a profile of the data recipient at the time of disclosure;

tion concerning a person's health or treatment that enables someone to identify that person. The expression *personal health information* refers to health information that may or may not identify individuals. As I will show, in many releases of personal health information, individuals can be recognized. *Anonymous personal health information*, by contrast, contains details about a person's medical condition or treatment but the identity of the person cannot be determined.

In general usage, confidentiality of personal information protects the interests of the organization while privacy protects the autonomy of the individual; but, in medical usage, both terms mean privacy. The historical origin and ethical basis of medical confidentiality begins with the Hippocratic Oath, which was written between the sixth century B.C. and the first century A.D. It states:

Whosoever I shall see or hear in the course of my dealings with men, if it be what should not be published abroad, I will never divulge, holding such things to be holy secrets.

- ▶ Differential privacy (DP) [Dwork et al. (2006)] is now widely adopted as the “gold standard” of privacy, for aggregate statistics/learning outcomes release.

# Some background: Privacy research

Journal of Privacy and Confidentiality  
Vol. 9, 2019

Submitted  
April 20, 2019  
Published  
October 2019

## Concentrated Differential Privacy Simplifications, Extensions, and Lov

Mark Bun<sup>1</sup> and Thomas Steinke<sup>1,2</sup>

<sup>1</sup> Microsoft, Internet Explorer  
John A. Paulson School of Engineering and Applied Sciences

## Deep Learning with Differential Privacy

October 25, 2016

Martin Abadi<sup>1</sup>  
H. Brendan McMahan<sup>1</sup>

Andy Chu<sup>1</sup>  
Ilya Mnih<sup>1</sup>  
L. Zhang<sup>1</sup>

Ian C.  
Kirk

### ABSTRACT

Machine learning techniques based on neural networks are achieving remarkable results in a wide variety of domains. However, the training of neural network-based, representative datasets, which may be considered and contain sensitive information. This data should not expose private information to one without the consent of differential privacy. Our implementation and experiments demonstrate that we can train deep neural networks with sensitive information, while a random privacy budget, and at a reasonable cost in terms of accuracy, training efficiency, and model quality.

### 1. INTRODUCTION

Recent progress in neural networks has led to impressive results in a wide range of applications, including image classification, language representation, speech recognition, and so on [16, 20, 26, 30, 32]. These advances are enabled, in part, by the availability of large and representative datasets for training neural networks. These datasets are often considered, and are often contain sensitive

## DIFFERENTIAL PRIVACY IN PRACTICE

CYNTHIA DWORKI, NITEN KOHLE, AND I

340 Harvard University, Harvard University, Cambridge, MA  
e-mail address: dworki@cs.harvard.edu

100 South Hall, UC Berkeley School of Information,  
e-mail address: nitentk@cs.berkeley.edu

Abstract—We present a natural extension of differential privacy to the case of noisy data.

## PATE-GAN: GENERATIVE DIFFERENTIAL PRIVACY

Jason Jordan<sup>1</sup>

Engineering Science Department  
University of Oxford, UK  
jordan.jason@cs.ox.ac.uk

Mehmet Ural<sup>1</sup>  
University of Cambridge, UK  
Mehmet.Ural@eng.ox.ac.uk  
Alex Teyssie<sup>1</sup>  
Alex.Teyssie@eng.ox.ac.uk

## The Discrete Gaussian for Differential Privacy

Ohmri E. Carmi  
IBM Research, Almaden  
ohmiec@us.ibm.com

Gautam Kamath  
University of Waterloo  
gkamath@uwaterloo.ca

Itai  
Dagum

### Abstract

A key tool for building differentially private systems is adding Gaussian noise to the output of a function evaluated on a sensitive dataset. Various continuous distributions present several practical challenges. For finite computers, several exact representations require non-terminating series and have a long tail to handle floating-point precision for a wide range of input values. Moreover, when the underlying data is integer-valued, adding continuous noise makes sense.

With these shortcomings in mind, we introduce and analyze the discrete version of differential privacy. Specifically, we demonstrate that adding discrete Gaussian noise provides essentially the same security guarantees as the addition of continuous Gaussian noise, as an example of efficient algorithms for exact sampling from this distribution is applicable for privately answering counting queries, or more generally, low-sensitivity integer-valued queries.

### 1 Introduction

## Rényi Differential Privacy

Ben Morris

2018-11-15 14:23:13  
DOI: 10.1007/978-1-4939-9836-7\_10

© 2018 Intel Corporation and Applied Mathematics

## UNIVERSALLY UTILITY-MAXIMIZING PRIVACY MECHANISMS\*

ADITHYAN KRISHNAN, TIM MCHUGH, AND MINGHONG SHEN

Abstract. A mechanism for releasing information about a statistical database with sensitive data need not be a trade-off between utility and privacy. Publishing fully accurate information maintains utility while retaining privacy, while providing random noise accomplishes the opposite. Privacy can be globally quantified using the framework of  $\epsilon$ -differential privacy, which requires that a mechanism's output distribution is nearly the same whether a given database entry is included. The goal of this paper is to formalize and provide strong and general utility guarantees, subject to differential privacy. We prove mechanisms that maximize one-optimal utility in every context, independent of the side information included in a prior distribution over query results and performance included in a concrete and reasonable loss function. Our main result is the following: for each fixed entry and differential privacy level, there is a generic mechanism  $M^*$ —a discrete version of the simple and well-studied Laplace mechanism—that side enables some form of Laplace distribution—that is simultaneously optimal loss-maximizing for every possible loss, subject to the differential privacy constraint. This is an extremely strong utility guarantee: every practical loss, no matter what the side information and privacy, derives from such utility loss  $M^*$  as loss-maximizing with  $\epsilon$ -differential privacy mechanism. In fact, it is typically restricted to a Gaussian, for every loss that is an optimal mechanism. We show that there is a non-independent prior that generates mechanisms that are optimal for every possible loss, subject to the same privacy constraint. This is a very general class of algorithms, demonstrating the flexibility of the Laplace mechanism for the distribution from which it is drawn. By adhering with a small amount of knowledge about an individual entry, with high probability, this individual's associated dataset can be seen all

as record contains many attributes (i.e., table schema), which can be viewed as utility means that for the average record, "table" records in the multi-dimensional of the attributes. This capacity is equivalent to a 1-bit and related to the fact an individual transaction and performance include statistically sensitive attributes.

less. Our first contribution is a formal utility function in a statistical database context, which is an extension of the Laplace mechanism for the distribution from which it is drawn. By adhering with a small amount of knowledge about an individual entry, with high probability, this individual's associated dataset can be seen all

as record contains many attributes (i.e., table schema), which can be viewed as utility means that for the average record, "table" records in the multi-dimensional of the attributes. This capacity is equivalent to a 1-bit and related to the fact an individual transaction and performance include statistically sensitive attributes.

Key words. differential privacy, utility maximization, genericity

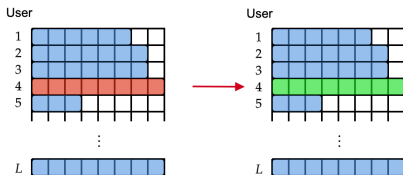
AMS subject classification. 68Q99

DOI: 10.1007/978-1-4939-9836-7\_10

... and continuing to multiply

# DP: A mathematically grounded framework for privacy

- ▶ We say that two datasets  $\mathcal{D}_1, \mathcal{D}_2$  are neighbours if they differ in the **sample values** of a **single user**.

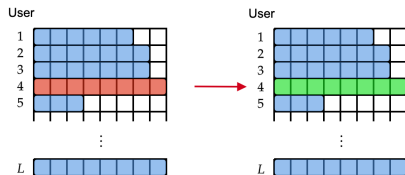


- ▶ A mechanism  $M$  is  $\epsilon$ -DP if for every pair of datasets  $\mathcal{D}_1, \mathcal{D}_2$  that are neighbours, and for every (measurable)  $Y$ ,

$$e^{-\epsilon} \Pr[M(\mathcal{D}_2) \in Y] \leq \Pr[M(\mathcal{D}_1) \in Y] \leq e^{\epsilon} \Pr[M(\mathcal{D}_2) \in Y].$$

# DP: A mathematically grounded framework for privacy

- ▶ We say that two datasets  $\mathcal{D}_1, \mathcal{D}_2$  are neighbours if they differ in the **sample values** of a **single user**.



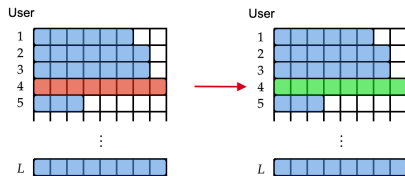
- ▶ A mechanism  $M$  is  $\epsilon$ -DP if for every pair of datasets  $\mathcal{D}_1, \mathcal{D}_2$  that are neighbours, and for every (measurable)  $Y$ ,

$$e^{-\epsilon} \Pr[M(\mathcal{D}_2) \in Y] \leq \Pr[M(\mathcal{D}_1) \in Y] \leq e^{\epsilon} \Pr[M(\mathcal{D}_2) \in Y].$$

Think of  $e^{\epsilon} \approx 1 + \epsilon$ , for  $\epsilon > 0$  small

# DP: A mathematically grounded framework for privacy

- ▶ We say that two datasets  $\mathcal{D}_1, \mathcal{D}_2$  are neighbours if they differ in the **sample values** of a **single user**.



- ▶ A mechanism  $M$  is  $(\epsilon, \delta)$ -DP if for every pair of datasets  $\mathcal{D}_1, \mathcal{D}_2$  that are neighbours, and for every (measurable)  $Y$ ,

$$e^{-\epsilon} (\Pr[M(\mathcal{D}_2) \in Y] - \delta) \leq \Pr[M(\mathcal{D}_1) \in Y] \leq e^{\epsilon} \Pr[M(\mathcal{D}_2) \in Y] + \delta.$$

We must have  $\delta \in (0, 1)$ ; think of  $\delta \approx$  prob. of DP “failure”

## A case for moving away from DP

- ▶ Data owners may often be incentivized to release “microdata,” via the release of **entire datasets**.

One of the following alternatives can then be adopted (in theory):

- ▶ Add noise to **every** entry of the dataset and release  
[**Large noise addition  $\equiv$  Poor accuracy!**]
- ▶ Select a “random” dataset that ensures high utility, via a DP selection mechanism  
[**High computational complexity!**]
- ▶ Generate and release DP **synthetic** datasets  
[**Could run into regulatory issues in practice!**]

## A case for moving away from DP

- ▶ Data owners may often be incentivized to release “microdata,” via the release of **entire datasets**.

One of the following alternatives can then be adopted (in theory):

- ▶ Add noise to **every** entry of the dataset and release  
[**Large noise addition  $\equiv$  Poor accuracy!**]
- ▶ Select a “random” dataset that ensures high utility, via a DP selection mechanism  
[**High computational complexity!**]
- ▶ Generate and release DP **synthetic** datasets  
[**Could run into regulatory issues in practice!**]

In this work, we embrace early, heuristic notions of “**anonymity**”, instead, and seek to formalize them.

## Some background: Anonymization

The following heuristic notions of dataset anonymity are well-known:

- ▶  $k$ -anonymity [Samarati and Sweeney (1998), Samarati (2001)]: Associated with every equivalence class (EC), there are  $\geq k$  sensitive attributes
- ▶  $\ell$ -diversity [Machanavajjhala et al. (2007)]: Associated with every EC, there are  $\geq \ell$  distinct sensitive attributes
- ▶  $t$ -closeness [Li, Li, and Venkatasubramanian (2007)]: Distribution of sensitive attributes in any equivalence class is “close” to that in the overall dataset

We focus on the well-known problem of **dataset linkage**:

How do anonymity parameters **degrade** upon the **linkage/composition** of multiple datasets?

## Dataset linkage: An example

Consider medical records in a hospital.

Gender	Postal Code	Disease
Male	600020	Heart disease
Male	600036	Lung infection
Female	600021	Osteoporosis
Female	600021	Cervical cancer
Female	600021	Osteoporosis

(a) Raw dataset

(Gender, Postal Code)	Disease
(Male, 600020) or (Male, 600036)	Heart disease Lung infection
(Female, 600021)	Osteoporosis Osteoporosis Cervical cancer

(b) Anonymized dataset that respects 2-anonymity and 2-diversity

## Dataset linkage: An example

An adversary now has access to **two, independently anonymized datasets**, which he/she knows that (Female, 600021) participates in.  
public quasi-identifier

(Gender, Postal Code)	Disease
(Male, 600020) or (Male, 600036)	Heart disease Lung infection
(Female, 600021)	Osteoporosis Cervical cancer

Anonymized dataset  $\bar{D}_1$

(Gender, Postal Code)	Disease
(Female, 600025)	Irritable bowel syndrome Lung infection
(Female, 600021)	Cervical cancer Amoebic dysentery

Anonymized dataset  $\bar{D}_2$

## Dataset linkage: An example

An adversary now has access to **two, independently anonymized datasets**, which he/she knows that **(Female, 600021)** participates in.

public quasi-identifier

(Gender, Postal Code)	Disease
(Male, 600020) or (Male, 600036)	Heart disease Lung infection
(Female, 600021)	Osteoporosis Cervical cancer

Anonymized dataset  $\bar{D}_1$

(Gender, Postal Code)	Disease
(Female, 600025)	Irritable bowel syndrome Lung infection
(Female, 600021)	Cervical cancer Amoebic dysentery

Anonymized dataset  $\bar{D}_2$

## Some unavoidable notation

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

## Some unavoidable notation

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

Let  $n^{(\mathbf{q})}(s) = \#$  samples with record  $(\mathbf{q}, s)$ .

## Some unavoidable notation

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

Let  $n^{(\mathbf{q})}(s) = \#$  samples with record  $(\mathbf{q}, s)$ .

### Definition

A dataset  $\mathcal{D}$  satisfies  $\ell$ -diversity, if  $\ell \in [|\mathcal{S}|]$  is the largest integer such that for any vector  $\mathbf{q} \in \mathcal{Q}$  with  $\sum_{s \in \mathcal{S}} n^{(\mathbf{q})}(s) > 0$ , we have  $|\{s \in \mathcal{S} : n^{(\mathbf{q})}(s) > 0\}| \geq \ell$ .

In order to satisfy  $\ell$ -diversity, the data owner bins quasi-identifiers  $\mathbf{q}$  into “equivalence classes”  $\bar{\mathbf{q}}$ , resulting in the anonymized dataset  $\bar{\mathcal{D}}$ .

## Some unavoidable formalism

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

- ▶ Datasets  $\mathcal{D}_1, \dots, \mathcal{D}_t$  (with a common user) are **independently** anonymized to yield  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$ .

## Some unavoidable formalism

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

- ▶ Datasets  $\mathcal{D}_1, \dots, \mathcal{D}_t$  (with a common user) are **independently** anonymized to yield  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$ .
- ▶ An adversary has access to  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$  and knows the quasi-identifier  $\mathbf{q}$  of a particular user, who participates in **all** these datasets.

## Some unavoidable formalism

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

- ▶ Datasets  $\mathcal{D}_1, \dots, \mathcal{D}_t$  (with a common user) are **independently** anonymized to yield  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$ .
- ▶ An adversary has access to  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$  and knows the quasi-identifier  $\mathbf{q}$  of a particular user, who participates in **all** these datasets.
- ▶ The adversary identifies equivalence classes  $\overline{\mathbf{q}}_i, i \in [t]$ , which contain  $\mathbf{q}$ . Let

$$n_{[t]}^{(\mathbf{q})}(s) := \min\{ \underbrace{n^{(\overline{\mathbf{q}}_j)}(s)}_{\#s \text{ samples in } \overline{\mathbf{q}}_j} : j \in [t] \}.$$

## Some unavoidable formalism

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

- ▶ Datasets  $\mathcal{D}_1, \dots, \mathcal{D}_t$  (with a common user) are **independently** anonymized to yield  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$ .
- ▶ An adversary has access to  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$  and knows the quasi-identifier  $\mathbf{q}$  of a particular user, who participates in **all** these datasets.
- ▶ The adversary identifies equivalence classes  $\overline{\mathbf{q}}_i, i \in [t]$ , which contain  $\mathbf{q}$ . Let

$$n_{[t]}^{(\mathbf{q})}(s) := \min\{ \underbrace{n^{(\overline{\mathbf{q}}_j)}(s)}_{\#s \text{ samples in } \overline{\mathbf{q}}_j} : j \in [t] \}.$$

- ▶ The adversary constructs the **linkage**

$$\mathcal{L}^{(\mathbf{q})} := \{(s, n_{[t]}^{(\mathbf{q})}(s)) : n_{[t]}^{(\mathbf{q})}(s) > 0\}.$$

The linkage gives rise to a **post-linkage dataset**  $\overline{\mathcal{D}}_{[t]}^{(\mathbf{q})}$ .

## Some unavoidable formalism

Any dataset  $\mathcal{D} = \{(\mathbf{q}_i, s_i) : i \in [N]\}$ . We focus on  $\ell$ -diversity as our notion of anonymity.

- ▶ Datasets  $\mathcal{D}_1, \dots, \mathcal{D}_t$  (with a common user) are **independently** anonymized to yield  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$ .
- ▶ An adversary has access to  $\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t$  and knows the quasi-identifier  $\mathbf{q}$  of a particular user, who participates in **all** these datasets.
- ▶ The adversary identifies equivalence classes  $\overline{\mathbf{q}}_i$ ,  $i \in [t]$ , which contain  $\mathbf{q}$ . Let

$$n_{[t]}^{(\mathbf{q})}(s) := \min\{ \underbrace{n^{(\overline{\mathbf{q}}_j)}(s)}_{\#s \text{ samples in } \overline{\mathbf{q}}_j} : j \in [t] \}.$$

- ▶ The adversary constructs the **linkage**

$$\mathcal{L}^{(\mathbf{q})} := \{(s, n_{[t]}^{(\mathbf{q})}(s)) : n_{[t]}^{(\mathbf{q})}(s) > 0\}.$$

The linkage gives rise to a **post-linkage dataset**  $\overline{\mathcal{D}}_{[t]}^{(\mathbf{q})}$ . **Is  $\overline{\mathcal{D}}_{[t]}^{(\mathbf{q})}$   $\ell'$ -diverse?**

# Worst-case degradation of diversity

We are interested in computing the worst-case anonymity parameter of  $\overline{\mathcal{D}}_{[t]}^{(\mathbf{q})}$ , i.e.,

$$\ell_{[t]} := \min_{\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t} \{ \widehat{\ell} : \overline{\mathcal{D}}_{[t]}^{(\mathbf{q})} \text{ obeys } \widehat{\ell}\text{-diversity} \}.$$

## Theorem

*Under some very mild regularity conditions, for any  $\ell \leq [|\mathcal{S}|]$ , we have that*

$$\ell_{[t]} = \begin{cases} 1, & \text{if } \ell \leq L \cdot \left(\frac{t-1}{t}\right) + 1, \\ L + 1 - (L - \ell + 1)t, & \text{otherwise.} \end{cases}$$

## Worst-case degradation of diversity

We are interested in computing the worst-case anonymity parameter of  $\overline{\mathcal{D}}_{[t]}^{(\mathbf{q})}$ , i.e.,

$$\ell_{[t]} := \min_{\overline{\mathcal{D}}_1, \dots, \overline{\mathcal{D}}_t} \{ \hat{\ell} : \overline{\mathcal{D}}_{[t]}^{(\mathbf{q})} \text{ obeys } \hat{\ell}\text{-diversity} \}.$$

### Theorem

*Under some very mild regularity conditions, for any  $\ell \leq [|\mathcal{S}|]$ , we have that*

$$\ell_{[t]} = \begin{cases} 1, & \text{if } \ell \leq L \cdot \left(\frac{t-1}{t}\right) + 1, \\ L + 1 - (L - \ell + 1)t, & \text{otherwise.} \end{cases}$$

**Key take-away:** The worst-case anonymity parameter is very poor (equals 1!) for **most**  $\ell$  values.

... We hence need a different notion of anonymity that **degrades gracefully** upon linkage.

## Proof idea

The proof makes use of a simple reduction to the analysis of arrays of binary strings.

- ▶ Given the equivalence classes  $\bar{\mathbf{q}}_j$  in  $\bar{\mathcal{D}}_j$ , one can define

$$\iota_j := \left( \mathbf{1}_{\{n^{(\bar{\mathbf{q}}_j)}(s) > 0\}} : s \in \mathcal{S} \right),$$

for  $j \in [t]$ . Clearly, we have that  $w(\iota_j) \geq \ell$ , for all  $j \in [t]$ .

## Proof idea

The proof makes use of a simple reduction to the analysis of arrays of binary strings.

- ▶ Given the equivalence classes  $\bar{\mathbf{q}}_j$  in  $\bar{\mathcal{D}}_j$ , one can define

$$\iota_j := \left( \mathbb{1}\{n^{(\bar{\mathbf{q}}_j)}(s) > 0\} : s \in \mathcal{S} \right),$$

for  $j \in [t]$ . Clearly, we have that  $w(\iota_j) \geq \ell$ , for all  $j \in [t]$ .

- ▶ One can then define an array  $M \in \{0, 1\}^{t \times L}$  whose  $j^{\text{th}}$  row  $M_j$  is  $\iota_j$ , with

$$\ell_{[t]} = w \left( \odot_{j=1}^t M_j \right).$$

- ▶ Observe that if  $\ell > L \cdot \left( \frac{t-1}{t} \right) + 1$ , then we cannot have at least one 0 in  $L - 1$  columns, resulting in  $\ell_{[t]} > 1$ .

## Proof idea

The proof makes use of a simple reduction to the analysis of arrays of binary strings.

- ▶ Given the equivalence classes  $\bar{\mathbf{q}}_j$  in  $\bar{\mathcal{D}}_j$ , one can define

$$\iota_j := \left( \mathbb{1}_{\{n^{(\bar{\mathbf{q}}_j)}(s) > 0\}} : s \in \mathcal{S} \right),$$

for  $j \in [t]$ . Clearly, we have that  $w(\iota_j) \geq \ell$ , for all  $j \in [t]$ .

- ▶ One can then define an array  $M \in \{0, 1\}^{t \times L}$  whose  $j^{\text{th}}$  row  $M_j$  is  $\iota_j$ , with

$$\ell_{[t]} = w \left( \odot_{j=1}^t M_j \right).$$

- ▶ Observe that if  $\ell > L \cdot \left( \frac{t-1}{t} \right) + 1$ , then we cannot have at least one 0 in  $L - 1$  columns, resulting in  $\ell_{[t]} > 1$ .

The proof then proceeds by arguing a similar “converse” for other  $\ell$  values, and constructing sequences that achieve the stated  $\ell_{[t]}$ .

## $(\ell, \delta)$ -Diversity: A definition

Consider the setting of datasets constructed from samples from **volunteering** users.

- ▶ A good assumption here is that the samples  $(\mathbf{q}_i, s_i)$  are i.i.d. across  $i \in [N]$ , from a joint distribution  $P_{\mathbf{Q}, S} = P_{\mathbf{Q}} \cdot P_{S|\mathbf{Q}}$ .
- ▶ WLOG, we assume that  $P_{\mathbf{Q}} \succ 0$  and  $P_S \succ 0$ .

### Definition

A dataset  $\mathcal{D}$  obeys  $(\ell, \delta)$ -diversity, for some fixed  $\delta \in (0, 1)$ , if it respects  $\ell$ -diversity w.p.  $\geq 1 - \delta$ .

## $(\ell, \delta)$ -Diversity: A definition

Consider the setting of datasets constructed from samples from **volunteering** users.

- ▶ A good assumption here is that the samples  $(\mathbf{q}_i, s_i)$  are i.i.d. across  $i \in [N]$ , from a joint distribution  $P_{\mathbf{Q}, S} = P_{\mathbf{Q}} \cdot P_{S|\mathbf{Q}}$ .
- ▶ WLOG, we assume that  $P_{\mathbf{Q}} \succ 0$  and  $P_S \succ 0$ .

### Definition

A dataset  $\mathcal{D}$  obeys  $(\ell, \delta)$ -diversity, for some fixed  $\delta \in (0, 1)$ , if it respects  $\ell$ -diversity w.p.  $\geq 1 - \delta$ .

1. How must a dataset owner ensure  $(\ell, \delta)$ -diversity?  
[Achievability]
2. How does  $(\ell, \delta)$ -diversity degrade upon linkage?  
[Linkage-resilience]
3. Can we achieve  $(\ell, \delta)$ -diversity while guaranteeing “closeness” to the true dataset? [Utility]

# Achievability

We assume the existence of a **Central Anonymizer** (CA) who knows  $P_{\mathbf{Q},S}$ .

- 
- 1: **procedure** MECH- $(\ell, \delta)$ -DIVERSITY
  - 2: Clients convey  $(\ell, \delta)$  to the CA.
  - 3: CA chooses  $p \in (0, p_\ell]$  and constructs eq. classes  $\bar{\mathbf{q}}$  s.t. for each  $\bar{\mathbf{q}}$ , we have  $|\mathcal{S}_{\bar{\mathbf{q}}}| := |\{s : P(\bar{\mathbf{q}}, s) \geq p\}| \geq \ell$ .
  - 4: CA broadcasts  $N = \frac{\ln(\frac{m\ell}{\delta})}{\ln(\frac{1}{1-p})}$  and  $\{\bar{\mathbf{q}}\}$  to the data owners (DOs).
  - 5: DOs obtain  $N$  samples and construct  $\bar{\mathcal{D}}_1, \dots, \bar{\mathcal{D}}_t$  using  $\{\bar{\mathbf{q}}\}$ .
- 

Here,  $p_i = i^{\text{th}}$ -largest  $P_S(s)$ ; and  $m := \min \left\{ |\mathcal{Q}|, \frac{1}{\ell p}, \frac{\sum_{i=\ell}^{|S|} P_i}{p} \right\}$ .

# Achievability

We assume the existence of a **Central Anonymizer** (CA) who knows  $P_{\mathbf{Q},S}$ .

- 
- 1: **procedure** MECH- $(\ell, \delta)$ -DIVERSITY
  - 2: Clients convey  $(\ell, \delta)$  to the CA.
  - 3: CA chooses  $p \in (0, p_\ell]$  and constructs eq. classes  $\bar{\mathbf{q}}$  s.t. for each  $\bar{\mathbf{q}}$ , we have  $|\mathcal{S}_{\bar{\mathbf{q}}}| := |\{s : P(\bar{\mathbf{q}}, s) \geq p\}| \geq \ell$ .
  - 4: CA broadcasts  $N = \frac{\ln(\frac{m\ell}{\delta})}{\ln(\frac{1}{1-p})}$  and  $\{\bar{\mathbf{q}}\}$  to the data owners (DOs).
  - 5: DOs obtain  $N$  samples and construct  $\bar{\mathcal{D}}_1, \dots, \bar{\mathcal{D}}_t$  using  $\{\bar{\mathbf{q}}\}$ .
- 

## Theorem

Any dataset  $\bar{\mathcal{D}}$  obtained via Mech- $(\ell, \delta)$ -diversity satisfies  $(\ell, \delta)$ -diversity.

Here,  $p_i = i^{\text{th}}$ -largest  $P_S(s)$ ; and  $m := \min \left\{ |\mathcal{Q}|, \frac{1}{\ell p}, \frac{\sum_{i=\ell}^{|\mathcal{S}|} P_i}{p} \right\}$ .

## Proof idea

- ▶ The proof first argues that  $|\overline{Q}| \leq m$ .
- ▶ Next, note that for any equivalence class  $\overline{\mathbf{q}}$ ,

$$\Pr[N(\overline{\mathbf{q}}, s) = 0] \leq (1 - p)^N.$$

## Proof idea

- ▶ The proof first argues that  $|\overline{Q}| \leq m$ .

- ▶ Next, note that for any equivalence class  $\overline{\mathbf{q}}$ ,

$$\Pr[N(\overline{\mathbf{q}}, s) = 0] \leq (1 - p)^N.$$

- ▶ Let  $M_{\overline{\mathbf{q}}}$  denote # distinct sensitive attributes in  $\overline{\mathbf{q}}$ . Then,

$$\Pr[M_{\overline{\mathbf{q}}} \geq \ell] \geq 1 - \ell \cdot (1 - p)^N.$$

## Proof idea

- ▶ The proof first argues that  $|\overline{\mathcal{Q}}| \leq m$ .

- ▶ Next, note that for any equivalence class  $\overline{\mathbf{q}}$ ,

$$\Pr [N(\overline{\mathbf{q}}, s) = 0] \leq (1 - \rho)^N.$$

- ▶ Let  $M_{\overline{\mathbf{q}}}$  denote # distinct sensitive attributes in  $\overline{\mathbf{q}}$ . Then,

$$\Pr [M_{\overline{\mathbf{q}}} \geq \ell] \geq 1 - \ell \cdot (1 - \rho)^N.$$

- ▶ Finally, via a union bound again, we have

$$\begin{aligned} \Pr [M_{\overline{\mathbf{q}}} \geq \ell, \text{ for all } \overline{\mathbf{q}} \in \overline{\mathcal{Q}}] &\geq 1 - \ell |\overline{\mathcal{Q}}| \cdot (1 - \rho)^N \\ &\geq 1 - m\ell \cdot (1 - \rho)^N = 1 - \delta, \end{aligned}$$

by the choice of  $N$ .

# Linkage resilience

The following [composition theorem](#) (à la the Basic Composition Theorem of DP [[Dwork-Roth \(2014\)](#)]) provides the degradation of anonymity on linkage:

## Theorem

*The overall post-linkage dataset  $\overline{\mathcal{D}}_{[t]}$  obtained via anonymization using MECH- $(\ell, \delta)$ -DIVERSITY, satisfies  $(\ell, t\delta)$ -diversity.*

The proof is via a simple application of the [union bound](#).

# Linkage resilience

The following [composition theorem](#) (à la the Basic Composition Theorem of DP [[Dwork-Roth \(2014\)](#)])) provides the degradation of anonymity on linkage:

## Theorem

*The overall post-linkage dataset  $\overline{\mathcal{D}}_{[t]}$  obtained via anonymization using MECH- $(\ell, \delta)$ -DIVERSITY, satisfies  $(\ell, t\delta)$ -diversity.*

The proof is via a simple application of the [union bound](#).

Observe that the Theorem asserts that the [entire post-linkage dataset](#), in **all** equivalence classes, satisfies  $(\ell, t\delta)$ -diversity.

# Linkage resilience

The following **composition theorem** (à la the Basic Composition Theorem of DP [Dwork-Roth (2014)]) provides the degradation of anonymity on linkage:

## Theorem

*The overall post-linkage dataset  $\overline{\mathcal{D}}_{[t]}$  obtained via anonymization using MECH- $(\ell, \delta)$ -DIVERSITY, satisfies  $(\ell, t\delta)$ -diversity.*

The proof is via a simple application of the **union bound**.

We shall next make explicit the “**generalization**” strategy for obtaining eq. classes to ensure “**high utility**,” in a special setting of practical relevance.

# Utility - I

Consider the setting where there exists a total order on quasi-identifiers:

$$q_1 < q_2 < \dots < q_{|\mathcal{Q}|}.$$

- ▶ We consider the natural class of “generalization” (or binning) strategies that give rise to eq. classes with contiguous quasi-identifiers.

[e.g., contiguous PIN codes, contiguous ages, etc.]

- ▶ A natural notion of “utility” hence is the number of eq. classes: large  $|\overline{\mathcal{Q}}| \equiv$  high utility.

**Q:** What is the structure of the optimal (highest utility) contiguous generalization strategy that guarantees  $\ell$ -diversity?

# Utility - I

Consider the setting where there exists a total order on quasi-identifiers:

$$\mathbf{q}_1 < \mathbf{q}_2 < \dots < \mathbf{q}_{|\mathcal{Q}|}.$$

- ▶ We consider the natural class of “**generalization**” (or binning) strategies that give rise to eq. classes with **contiguous** quasi-identifiers.

[e.g., **contiguous PIN codes**, **contiguous ages**, etc.]

- ▶ A natural notion of “**utility**” hence is the number of eq. classes: large  $|\overline{\mathcal{Q}}| \equiv$  high utility.

**Q:** What is the structure of the optimal (highest utility) **contiguous** generalization strategy that guarantees  $\ell$ -diversity?

**A:** The greedy procedure that **iteratively adds quasi-identifiers to eq. classes** until  $|\overline{S_{\overline{q}}}| \geq \ell$  is optimal (proof via induction).

# Utility - I

Consider the setting where there exists a total order on quasi-identifiers:

$$\mathbf{q}_1 < \mathbf{q}_2 < \dots < \mathbf{q}_{|\mathcal{Q}|}.$$

- ▶ We consider the natural class of “**generalization**” (or binning) strategies that give rise to eq. classes with **contiguous** quasi-identifiers.  
[e.g., **contiguous PIN codes**, **contiguous ages**, etc.]
- ▶ A natural notion of “**utility**” hence is the number of eq. classes: large  $|\overline{\mathcal{Q}}| \equiv$  high utility.

**Q:** What is the structure of the optimal (highest utility) **contiguous** generalization strategy that guarantees  $\ell$ -diversity?

**A:** The greedy procedure that **iteratively adds quasi-identifiers to eq. classes** until  $|\overline{S_{\overline{\mathcal{Q}}}}| \geq \ell$  is optimal (proof via induction).

An explicit variant of this algorithm can be derived when  $P_{\mathbf{Q},S} = P_{\mathbf{Q}} \cdot P_S$ .

## Numerical results

We present numerical results on the # samples required by Mech- $(\ell, \delta)$ -Diversity.

We consider the setting of hospital records, where each record is

$$\mathbf{r} \in \underbrace{\text{POSTAL-CODE} \times \text{GENDER} \times \text{AGE}}_{\text{quasi-identifiers}} \times \underbrace{\text{DISEASE}}_{\text{sens. attr.}},$$

with some fixed cardinalities of the underlying sets.

## Numerical results

We present numerical results on the # samples required by Mech- $(\ell, \delta)$ -Diversity.

We consider the setting of hospital records, where each record is

$$\mathbf{r} \in \underbrace{\text{POSTAL-CODE} \times \text{GENDER} \times \text{AGE}}_{\text{quasi-identifiers}} \times \underbrace{\text{DISEASE}}_{\text{sens. attr.}},$$

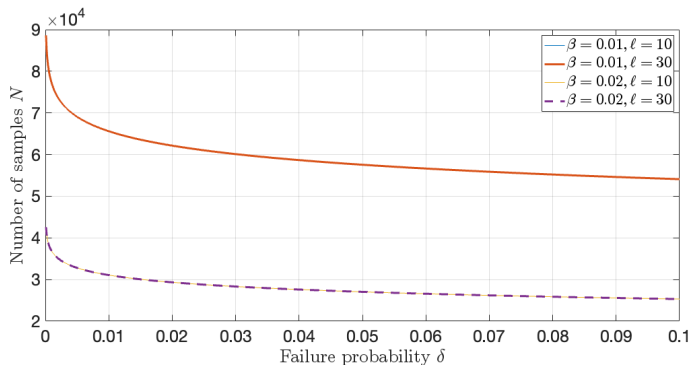
with some fixed cardinalities of the underlying sets.

We consider the following marginal distributions  $P_S$ :

1. **Uniform marginals:** Here, we let  $P_S(s) = 1/|S|$ , for all  $s \in [|S|]$ ; observe that  $p_i = 1/|S|$ , for all  $i \in [|S|]$ .
2. **Geometric marginals:** Here, we let  $P_S(s) = p_1 \rho^{s-1}$ , for suitable  $p_1, \rho \in (0, 1)$ ; observe that  $p_i = P_S(i)$ , for all  $i \in [|S|]$ .

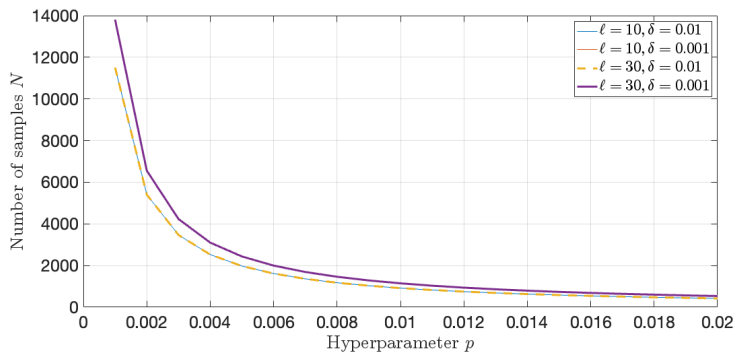
The exact parameters can be found in <https://arxiv.org/abs/2506.18405>.

# Plots - I



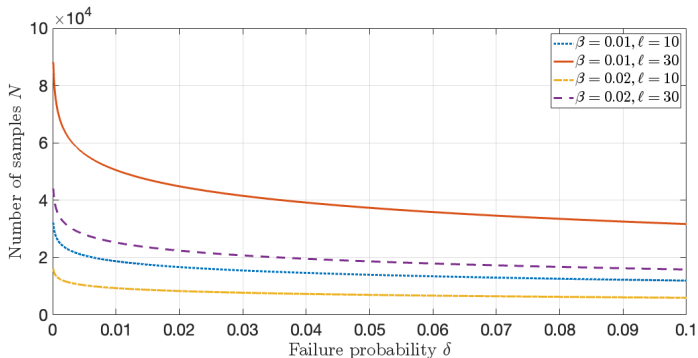
Plots showing the variation of the number of samples  $N$  with the parameter  $\delta$ , under **uniform marginals** of sensitive attributes. Here,  $\beta := p/p_\ell$ .

## Plots - II



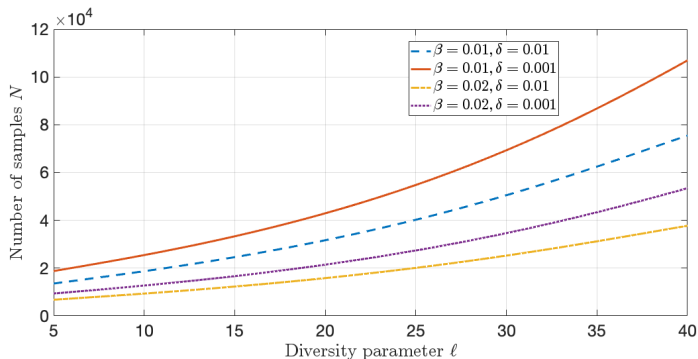
Plots showing the variation of the number of samples  $N$  with the parameter  $p$ , under **uniform marginals** of sensitive attributes.

## Plots - III



Plots showing the variation of the number of samples  $N$  with the parameter  $\delta$ , under **geometric marginals** of sensitive attributes. Here,  $\beta := p/p_\ell$ .

## Plots - IV



Plots showing the variation of the number of samples  $N$  with the parameter  $\ell$ , under **geometric marginals** of sensitive attributes. Here,  $\beta := p/p_\ell$ .

# Ongoing and future research directions

- ▶ Extend our analysis for the setting of datasets with **correlated** (possibly Markovian) **samples**
- ▶ Design **achievability schemes** (generalization strategies) without the need for a **Central Anonymizer (CA)**
- ▶ Extend our analysis techniques to other notions of anonymity such as **entropy  $\ell$ -diversity**

Thank You!