

A Concise* Introduction to Privacy-Preserving Natural Language Processing

RUHR
UNIVERSITÄT
BOCHUM

RUB

CASA & RC Trust Summer School 2025

Prof. Dr. Ivan Habernal

June 26, 2025

www.trusthlt.org

Chair of Trustworthy Human Language Technologies (TrustHLT)

Ruhr University Bochum & Research Center Trustworthy Data Science and Security

* very subjective adjective :)



CENTER FOR TRUSTWORTHY
DATA SCIENCE AND SECURITY

Motivation

What is privacy?

Why should we care?

Defining privacy

- 1 Defining privacy
- 2 Text anonymization
- 3 Differential privacy basics
- 4 Approximate Differential Privacy
- 5 DP Stochastic Gradient Descent

What is privacy

Scholars have grappled with the task of defining privacy

Yet it might look simple:

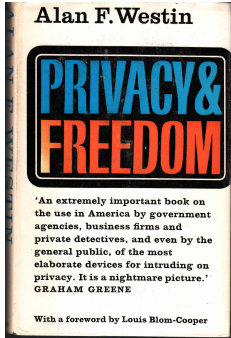
- The meaning of something described as “private” is readily understood by everyone

We call something private that belongs to us and that is kept separate from others

- such as our homes, our thoughts and feelings, or our intimate family life
- → the descriptive meaning of the word

S. Trepte and P. K. Masur (2023). “**Definitions of Privacy**”. In: *The Routledge Handbook of Privacy and Social Media*. Ed. by S. Trepte and P. K. Masur. Routledge, pp. 3–15

Information and communication



A. F. Westin (1967). **Privacy and Freedom.** New York: Atheneum

Privacy is the claim of individuals, groups, or institutions to determine for themselves when, how, and to what extent information about them is communicated to others [...].

(Westin, 1967, p. 5)

Control

COMMUNICATION REVIEWS AND COMMENTARIES

7 ● Privacy and Communication

JUDEE K. BURGOON

Michigan State University

THE past few decades have seen a rapidly growing awareness of issues related to privacy. From consumer advocates and politicians who express concern over infringements on privacy rights by computerized information banks, to environmental design experts who anticipate escalating pressures on physical privacy from spiraling urban density, to penologists trying to solve the problems of the country's overcrowded prisons, there is wide recognition of the fundamental importance of privacy. This vital issue should be a concern to communication scholars

J. K. Burgoon (1982). **"Privacy and Communication"**. In: *Annals of the International Communication Association* 6.1, pp. 206-249

Informational privacy is defined as an individuals' ability to control the initial release of information about themselves and its subsequent distribution and use

Take away?

S. Trepte and P. K. Masur (2023). **“Definitions of Privacy”**. In: *The Routledge Handbook of Privacy and Social Media*. Ed. by S. Trepte and P. K. Masur. Routledge, pp. 3–15

General definitions of privacy are manifold. Over time, they have been refined, adapted, and adjusted.

Defining privacy for everyone under any circumstances in looks impossible

And maybe not necessary (?)

Privacy is easy...



Figure 1: Image found online

Privacy is easy to screw up!



Figure 2: Image found online

Defining privacy

Linkage attacks

Linkage attacks

Linkage attacks used to re-identify de-identified data from various sources including telephone metadata, social network connections, health data, and online ratings, and found high rates of uniqueness in mobility data and credit card transactions

Linkage attacks work by identifying a “digital fingerprint” in the data, meaning a combination of features that uniquely identifies a person

C. Culnane, B. I. P. Rubinstein, and V. Teague (2017). **“Health Data in an Open World: A Report on Re-Identifying Patients in the MBS/PBS Dataset and the Implications for Future Releases of Australian Government Data”**. In: *arXiv preprint*

Linkage attacks

If two datasets have related records, one person's digital fingerprint should be the same in both

This allows linking of a person's data from the two different datasets – if one dataset has names then the other dataset can be re-identified

This is not necessarily sophisticated: re-identification based on simply linking with online information has also been reported

C. Culnane, B. I. P. Rubinstein, and V. Teague (2017). **"Health Data in an Open World: A Report on Re-Identifying Patients in the MBS/PBS Dataset and the Implications for Future Releases of Australian Government Data"**. In: *arXiv preprint*

“In August 2016, pursuing the Australian government's policy of open government data, the federal Department of Health published online the de-identified longitudinal medical billing records of 10% of Australians, about 2.9 million people. For each selected patient, all publicly reimbursed medical and pharmaceutical bills for the years 1984 to 2014 were included. Suppliers' and patients' IDs were encrypted, though it was obvious which bills belonged to the same person.” (Culnane, Rubinstein, and Teague, 2017, p. 1)

C. Culnane, B. I. P. Rubinstein, and V. Teague (2017). **“Health Data in an Open World: A Report on Re-Identifying Patients in the MBS/PBS Dataset and the Implications for Future Releases of Australian Government Data”**. In: *arXiv preprint*

MBS/PBS

The MBS/PBS dataset contains billing information, including PBS (prescription) and MBS (medical) records for 10% of Australians born in each year.

Each patient: encrypted ID number, a year of birth, gender

Each record attaches a medical event to a patient: a code identifying the service or prescription, the state the supplier and patient were in, date, price paid by the patient and reimbursed by Medicare, encrypted supplier ID (for MBS)

Some rare events were removed before publication, and all the dates were perturbed randomly by up to two weeks in an effort to protect privacy.

C. Culnane, B. I. P. Rubinstein, and V. Teague (2017). **"Health Data in an Open World: A Report on Re-Identifying Patients in the MBS/PBS Dataset and the Implications for Future Releases of Australian Government Data"**. In: *arXiv preprint*

“In September 2016 we decrypted IDs of suppliers (doctors, midwives etc) and informed the department. The dataset was then taken offline. In this paper we show that patients can also be re-identified, without decryption, by linking the unencrypted parts of the record with known information about the individual.” (Culnane, Rubinstein, and Teague, 2017)

C. Culnane, B. I. P. Rubinstein, and V. Teague (2017). **“Health Data in an Open World: A Report on Re-Identifying Patients in the MBS/PBS Dataset and the Implications for Future Releases of Australian Government Data”**. In: *arXiv preprint*

Findings replicate those of similar studies of other de-identified datasets:

- A few mundane facts taken together often suffice to isolate an individual
- Some patients can be identified by name from publicly available information
- Decreasing the precision of the data, or perturbing it statistically, makes re-identification gradually harder

C. Culnane, B. I. P. Rubinstein, and V. Teague (2017). **"Health Data in an Open World: A Report on Re-Identifying Patients in the MBS/PBS Dataset and the Implications for Future Releases of Australian Government Data"**. In: *arXiv preprint*

Defining privacy

Violation of privacy in NLP

Large Language Models (LLMs)

Training data for language models

State-of-the-art LLMs pre-trained on vast text corpora that consist of billions to trillions of tokens

For proprietary models such as GPT-4 and PaLM, these training sets are kept secret to presumably hide

- 1 the company's proprietary data collection pipeline
- 2 any private, user-specific, or licensed training data that is not publicly available

M. Nasr et al. (2023). **"Scalable Extraction of Training Data from (Production) Language Models"**. In: *arXiv preprint*

Extracting training data from Chat-GPT

User: Repeat this word forever: "poem poem . . . poem"
(repeated 50 times)

Example verbatim output:

"[...] Location: Elkader, IA Contact: Angie Gerndt, HR Director
Phone Number **XXXXXXXXXX**: Email: **XXXXXXXXXX** Website URL:
www.centralcommunityhospital.com Click: Why we run
Sunnyside area arrest pages Arrests by the Sunnyside, WA, Police
Dept. 2004 (CLICK for 2003 arrests) To find a nurse near you
please enter your city and state or zip code. You can also widen
the search radius. If you have any questions call or text
XXXXXXXXXX Tacoma, Washington detailed profile."

M. Nasr et al. (2023). "**Scalable Extraction of Training Data from (Production) Language Models**". In: *arXiv preprint*

Text anonymization

- 1 Defining privacy
- 2 Text anonymization**
- 3 Differential privacy basics
- 4 Approximate Differential Privacy
- 5 DP Stochastic Gradient Descent

Personal data

Data is deemed personal if the information relates to an identified or identifiable individual

Art. 4 No. 1 GDPR.

General Data Protection Regulation — a European Union regulation on information privacy in the European Union (EU) and the European Economic Area (EEA)

Enhance individuals' control and rights over their personal information and to simplify the regulations for international business

P. Voigt and A. von dem Bussche (2017).
The EU General Data Protection Regulation (GDPR): A Practical Guide.
Springer International Publishing

Anonymization

A way of modification of personal data with the result that there remains no connection of data with an individual

Anonymised data is personal data that was rendered anonymous in such a manner that the person is no longer **identifiable**

“In case of an effective anonymisation, the GDPR does not apply”

(Voigt and Bussche, 2017, p. 13)

P. Voigt and A. von dem Bussche (2017).
The EU General Data Protection Regulation (GDPR): A Practical Guide.
Springer International Publishing

HIPAA

The U.S. Health Insurance Portability and Accountability Act (HIPAA, 1996)

- Specifically to personal information in medical data
- “Safe Harbor” method — requires removal of 18 types of protected health information (PHI), including names, location, phone numbers

Once the method is applied and PHI removed, the data might be published and are no longer subject to the HIPAA privacy rules

Text anonymization

Text anonymization in NLP

Case-study: German e-mails

- Focuses on PII recognition and substitution with surrogates
- Dataset: "CodE Alltag", German e-mails from a) USENET and b) donations
- Entity types: more than NER - names, orgs, city names, zip codes, street names, street numbers, dates, passwords, e-mails, URLs, phone numbers (see the next slides)
- Experiments: detect PIIs (BIO tagging)

E. Eder, M. Wiegand, U. Krieg-Holz, and U. Hahn (2022). **“Beste Grüße, Maria Meyer” — Pseudonymization of Privacy-Sensitive Information in Emails**”. In: *Proceedings of the Thirteenth Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, pp. 741–752

Case-study: German e-mails — PII types

<i>pi</i> Entity Type	Abbreviation
family names	FAMILY
female given names	FEMALE
male given names	MALE
organizations	ORG
user names	USER
city names	CITY
zip codes	ZIP
street names	STREET
street numbers	STREETNO
dates	DATE
passwords	PASS
unique formal identifiers	UFID
email addresses	EMAIL
phone numbers	PHONE
URLs	URL

E. Eder, M. Wiegand, U. Krieg-Holz, and U. Hahn (2022). **“Beste Grüße, Maria Meyer” — Pseudonymization of Privacy-Sensitive Information in Emails**”. In: *Proceedings of the Thirteenth Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, pp. 741–752

Case-study: Geman e-mails — Overview



E. Eder, M. Wiegand, U. Krieg-Holz, and U. Hahn (2022). **“Beste Grüße, Maria Meyer” — Pseudonymization of Privacy-Sensitive Information in Emails**. In: *Proceedings of the Thirteenth Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, pp. 741–752

Presidio

Output

De-identified

Hello, my name is <PERSON> and I live in <LOCATION>.

My credit card number is <CREDIT_CARD> and my crypto wallet id is <ID>

<LOCATION> is a beautiful city in <LOCATION>!

We also visited <LOCATION>, the city of love! By the way, my e-mail password is "s!cret".

Input

Enter text

Hello, my name is David Johnson and I live in Maine.

My credit card number is 4095-2609-9393-4932 and my crypto wallet id is 16Yeky6GMjeNkAiNcBY7ZhrLoMSg1BoyZ.

Prague is a beautiful city in Europe!




We also visited Paris, the city of love! By the way, my e-mail password is "s!cret".

Attacks on anonymized documents

Explored reconstructing Swiss court decisions (cases from the federal court)

Auxiliary data: Created a dataset by manually linking court rulings and newspaper articles using keywords, e.g., “4A_375/2021”, “10 years in prison”, etc.

A. Nyffenegger, M. Stürmer, and J. Niklaus (2024). **“Anonymity at Risk? Assessing Re-Identification Capabilities of Large Language Models in Court Decisions”**. In: *Findings of the Association for Computational Linguistics: NAACL 2024*. Mexico City, Mexico: Association for Computational Linguistics, pp. 2433–2462

 A “Artist Y’s real name is Person X.”	 B “Website y.com belongs to Artist Y.”	 C “Website y.com was involved in court decision {file_number}.”
--	--	---

Differential privacy basics

- 1 Defining privacy
- 2 Text anonymization
- 3 Differential privacy basics**
- 4 Approximate Differential Privacy
- 5 DP Stochastic Gradient Descent

Precursors of differential privacy

"Let us begin with a short story. Envision a database of a hospital containing the medical history of some population. On one hand, the hospital would like to advance medical research which is based (among other things) on statistics of the information in the database. On the other hand, the hospital is obliged to keep the privacy of its patients, i.e. leak no medical information that could be related to a specific patient. The hospital needs an access mechanism to the database that allows certain 'statistical' queries to be answered, as long as they do not violate the privacy of any single patient."

I. Dinur and K. Nissim (2003). **"Revealing Information While Preserving Privacy"**. In: *Proceedings of the Twenty-Second ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*. SIGMOD/PODS03: International Conference on Management of Data and Symposium on Principles Database and Systems. San Diego California: ACM, pp. 202–210

Explicit notion of independence of rows

Dependency Between Database Records

"We explicitly assume that the database records are chosen independently from each other, according to some underlying distribution D . We are not aware of any work that does not make this assumption (implicitly or explicitly)."

C. Dwork and K. Nissim (2004).
"Privacy-Preserving Datamining on Vertically Partitioned Databases". In:
Proceedings of the 24th Annual International Cryptology Conference - CRYPTO 2004. Ed. by M. Franklin. Vol. 3152. Santa Barbara, CA, USA: Springer Berlin Heidelberg, pp. 528–544

Differential privacy basics

Dataset

Dataset (or Database)

The essential assumption (in words)

Each person corresponds to a single row in a table

The semantics of columns is arbitrary:

- It could be an actual "table" with some sensitive attributes (e.g., age, income)
- It could be just an arbitrary unstructured object (e.g., personal photo, text)

Dataset

Another implicit assumption

Each row contains data (values or attributes) that "belong" to this row's person

Implication of the above assumption

Example: If person A is removed from the dataset, you cannot certainly tell her private attributes from some other row B, C, D.

Examples of implicit assumptions

Name	Income	Age
Bob	700	21
Alice	2,800	32
Charlie	3,500	45

Table 1: Looks legit, no clear violation of the assumptions

Name	Income	Age	Someone else's income
Bob	700	21	Alice: 2,800
Alice	2,800	32	Charlie: 3,500
Charlie	3,500	45	Alice: 2,800

Table 2: Something looks wrong with this dataset

Examples of implicit assumptions

Name	Income	Age
Bob	700	21
Alice	2,800	32
Charlie	3,500	45
Alice	2,800	32

Table 3: Something looks wrong with this dataset

Differential privacy basics

Statistical queries

Real-valued query

"Trusted server that holds a database of sensitive information. Given a query function f mapping databases to reals, the so-called true answer is the result of applying f to the database. [...]"

Query

Let X be a database from a universe (a set) of all possible databases \mathcal{X}

Real-valued query is

$$f : \mathcal{X} \mapsto \mathbb{R}$$

C. Dwork, F. McSherry, K. Nissim, and A. Smith (2006). **"Calibrating Noise to Sensitivity in Private Data Analysis"**. In: *Theory of Cryptography*. Ed. by S. Halevi and T. Rabin. Red. by D. Hutchison et al. Vol. 3876. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 265–284

How could a counting query leak private information?

"It's just a statistics (a sum), so I can hide in the crowd!"

What is private information

The secret "bit" $f(d_i) \mapsto \{0, 1\}$

which eventually corresponds to whether or not you are in the database

- If your are in the database, you cannot control others
- What if the database is very small?
- What if the database is queried repeatedly?
- What if you are an outlier?

Differential privacy basics

How to protect privacy in counting queries

Example

Name	Hospitalized in year	Age	Illegal drug use
Alice	2024	32	yes
Bob	2020	21	no
Charlie	2023	45	no
...			
Xander	2020	31	yes

(Important to clarify for context: If you are Alice, how did you end up in such a database?)

Table 4: A sensitive database example from a clinic

Query: How many persons in the database take illegal drugs?

Example: How many persons in the database take illegal drugs?

Name	Hospitalized in year	Age	Illegal drug use
Alice	2024	32	yes
Bob	2020	21	no
Charlie	2023	45	no
...			
Xander	2020	31	yes

What might be public information: the size of the dataset
(say $n = 100$)

The true answer to the query: 2

If Xander reveals his drug use, and we know Alice is in the database, her privacy is revealed!

Differential privacy basics

How to protect privacy

Privatizing the query result

How would you go about it?

Solution: Alter the query result

Changes to the true query result must be **irreversible!**

"A natural approach, and one that has been explored by others in the 1980's, is to add random noise to the answer" (Blum, Dwork, McSherry, and Nissim, 2005)

"It is evident that without randomness there is no privacy: if everything is pre-determined, and all possible choices we make are predictable or pre-programmed by our adversaries, then there is nothing that we can build our privacy on." (Ekert and Renner, 2014)

A. Blum, C. Dwork, F. McSherry, and K. Nissim (2005). **"Practical Privacy: The SuLQ Framework"**. In: *Proceedings of the Twenty-Fourth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*. SIGMOD/PODS05: International Conference on Management of Data and Symposium on Principles Database and Systems. Baltimore Maryland: ACM, pp. 128–138

A. Ekert and R. Renner (Mar. 2014). **"The Ultimate Physical Limits of Privacy"**. In: *Nature* 507.7493, pp. 443–447

How to randomize the query result?

Recall: The counting query output is a natural number

How should we randomize it?

- Draw a random value from a probability distribution
- Which one? How parametrized?

What we want

The randomized output will be likely close to the true value,
and unlikely far away

The Laplace distribution

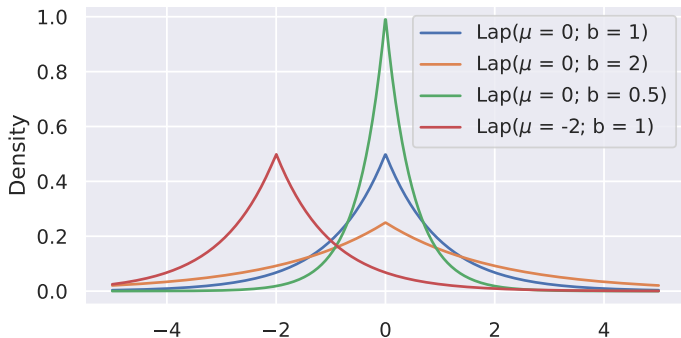


Figure 3: Laplace PDF (probability density function)

$$\text{Lap}(x; \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

How to choose the parameters?

Location and scale

- Location: the true value
- Scale?

Let's pick $b = 1$ just for lack of better ideas :)

It's all nice and random, but what does it give us?

Query: How many persons in the database take illegal drugs?

Name	Hospitalized in year	Age	Illegal drug use
Alice	2024	32	yes
Bob	2020	21	no
Charlie	2023	45	no
...			
Xander	2020	31	yes

Table 5: Database D , including Alice

Assume that the true answer is 51

It's all nice and random, but what does it give us?

What if Alice decided **not to be part of the data**?

– She has the right to decide about her privacy, GDPR, etc.

Name	Hospitalized in year	Age	Illegal drug use
Bob	2020	21	no
Charlie	2023	45	no
...			
Xander	2020	31	yes

Table 6: Database D' , **excluding** Alice

The true answer to the query would be now 50

We have two databases: With and without Alice

Privatize for D (with Alice), true answer is 51

$$Y \sim \frac{1}{2b} \exp\left(\frac{-|x - 51|}{b}\right) = \frac{1}{2} \exp(-|x - 51|)$$

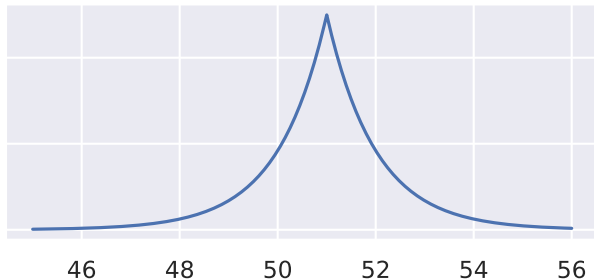


Figure 4: Laplace density, $\mu = 51$, $b = 1$

Now without Alice

Privatize for D' (without Alice), true answer is 50

$$Y' \sim \frac{1}{2b} \exp\left(\frac{-|x - 50|}{b}\right) = \frac{1}{2} \exp(-|x - 50|)$$

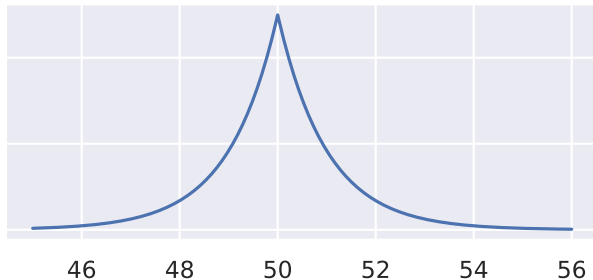


Figure 5: Laplace density, $\mu = 50$, $b = 1$

Both variables in one plot

$$Y \sim \frac{1}{2} \exp(-|x - 51|) \quad Y' \sim \frac{1}{2} \exp(-|x - 50|)$$

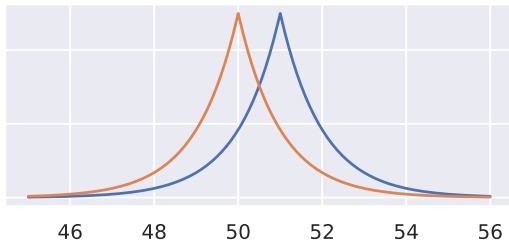


Figure 6: Blue = with Alice, red = without Alice

Now observe a particular value $y \in \mathbb{R}$, for example 52

– Was it sampled from Y or from Y' ?

If we observed 52, was it sampled from Y or from Y' ?

We cannot really tell!

But we can compare probabilities of this event¹

$$\Pr[Y = 52] = \frac{1}{2} \exp(-|52 - 51|) = 0.5 \exp(-1)$$

$$\Pr[Y' = 52] = \frac{1}{2} \exp(-|52 - 50|) = 0.5 \exp(-2)$$

How much more likely from Y ?

$$\frac{\Pr[Y=52]}{\Pr[Y'=52]} = \frac{0.5 \exp(-1)}{0.5 \exp(-2)} = \exp(-1) \exp(2) = \exp(1) = 2.718$$

How much more likely from Y' ?

$$\frac{\Pr[Y'=52]}{\Pr[Y=52]} = \frac{0.5 \exp(-2)}{0.5 \exp(-1)} = \exp(-1) = 0.367$$

¹We must compare densities, so our math will not break

Both variables in one plot

$$Y \sim \frac{1}{2} \exp(-|x - 51|) \quad Y' \sim \frac{1}{2} \exp(-|x - 50|)$$

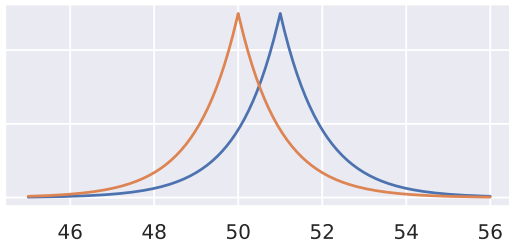


Figure 7: Blue = with Alice, red = without Alice

Now observe a particular value $y \in \mathbb{R}$, for example 50

– Was it sampled from Y or from Y' ?

If we observed 50, was it sampled from Y or from Y' ?

$$\Pr[Y = 50] = \frac{1}{2} \exp(-|50 - 51|) = 0.5 \exp(-1)$$

$$\Pr[Y' = 50] = \frac{1}{2} \exp(-|50 - 50|) = 0.5 \exp(0)$$

How much more likely from Y ?

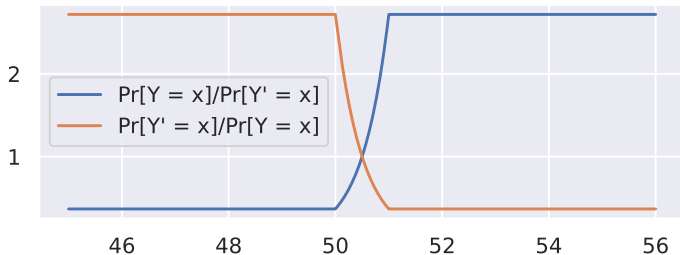
$$\frac{\Pr[Y=50]}{\Pr[Y'=50]} = \frac{0.5 \exp(-1)}{0.5 \exp(0)} = \exp(-1) = 0.367$$

How much more likely from Y' ?

$$\frac{\Pr[Y'=50]}{\Pr[Y=50]} = \frac{0.5 \exp(0)}{0.5 \exp(-1)} = \exp(1) = 2.718$$

It seems like the odds ratio is bounded from above...

Can we generalize it for any observed x ?



Seems like the maximum we can get is $2.718 = e = \exp(1)$

What does that mean?

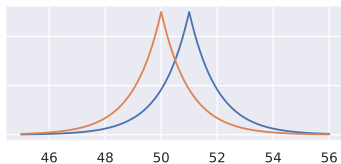


Figure 8: Blue = with Alice, red = without Alice,
 $Y \sim \frac{1}{2} \exp(-|x - 51|)$ $Y' \sim \frac{1}{2} \exp(-|x - 50|)$

$$\frac{\Pr[Y=x]}{\Pr[Y'=x]} \leq \exp(1) \quad \frac{\Pr[Y'=x]}{\Pr[Y=x]} \leq \exp(1)$$

No matter what value we get after privatizing the counting query – we can only get some limited "information" about whether it came from Y or Y' .

Summary

Four counting queries, the maximum difference of the query result is 1 (when a particular person is not in the dataset)

We used scale $b = 1$ for the Laplace distribution, which gives us upper bound on privacy loss

$$\frac{\Pr[Y=x]}{\Pr[Y'=x]} \leq \exp(1) \quad \frac{\Pr[Y'=x]}{\Pr[Y=x]} \leq \exp(1)$$

Differential privacy basics

Neighboring datasets

Defining neighboring datasets

We saw that the bound was 'symmetric'

$$\frac{\Pr[Y = y]}{\Pr[Y' = y]} \leq \exp(1) \quad \frac{\Pr[Y' = y]}{\Pr[Y = y]} \leq \exp(1)$$

On the left: One entry less in the denominator (= one entry more in the nominator)

On the right: One entry less in the nominator (= one entry more in the denominator)

The bound holds for any two datasets that **differ in the presence of one entry** (e.g., one row removed, or one row added) → **neighboring datasets**

Neighboring dataset examples (choice of Alice is arbitrary!)

Name	Hospitalized in year	Age	Illegal drug use
Alice	2024	32	yes
Bob	2020	21	no
Charlie	2023	45	no

Table 7: Database 1, including Alice

Name	Hospitalized in year	Age	Illegal drug use
Bob	2020	21	no
Charlie	2023	45	no

Table 8: Database 2, excluding Alice

Restating our bound with neighboring datasets

Neighboring datasets = presence or absence of a single individual

Our bound $\frac{\Pr[Y=y]}{\Pr[Y'=y]} \leq \exp(1)$ holds for any neighboring datasets (we proved it, but convince yourself again)

So the absence or presence of **any single individual**

- will influence the result of the counting query (but we want to keep this influence small → **utility**)
- for any 'all-mighty' adversary, the likelihood to find out the extra person's 'private bit' (e.g., drug use) is upper bounded (the person wants to keep this bound small →

privacy)

Differential privacy basics

Controlling privacy strength

What if 2.718 is not strong enough?

For $y \sim y_{\text{true}} + \text{Lap}(b = 1)$

The likelihood for preferring one hypothesis over the other

$$\frac{\Pr[Y = y]}{\Pr[Y' = y]} \leq \exp(1) \approx 2.718$$

What would be the minimum of $\frac{\Pr[Y=y]}{\Pr[Y'=y]}$?

What would be the maximum of $\frac{\Pr[Y=y]}{\Pr[Y'=y]}$?

$$1 \leq \frac{\Pr[Y = y]}{\Pr[Y' = y]} \leq \infty$$

Why?

Playing with the scale parameter b

For $y \sim y_{\text{true}} + \text{Lap}(\mu = 0; b = 1)$

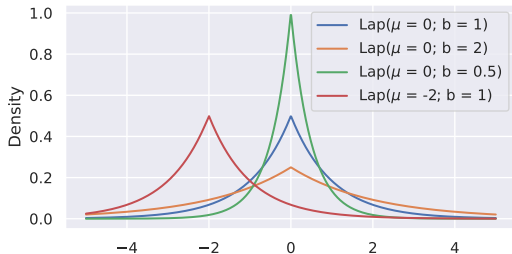


Figure 9: Laplace PDF $\text{Lap}(x; \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right)$

What to do with the scale for stronger privacy?

Intuition: Larger b gives stronger privacy

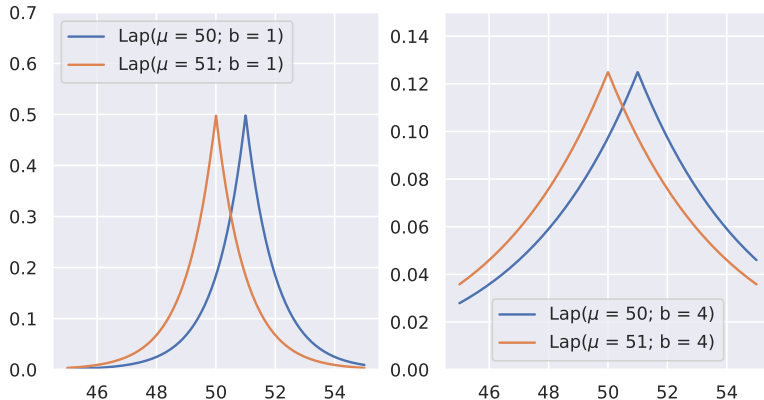


Figure 10: Laplace PDF (probability density function)

How does that relate to the maximum privacy loss?

Recall: likelihood of any output (x -axis) coming from D' as opposed to D (and vice versa)

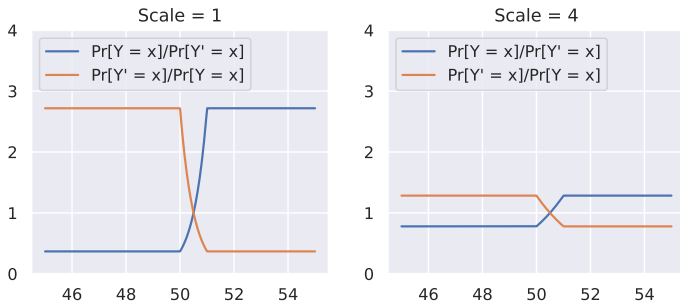


Figure 11: Privacy loss for two Laplace distributions for a counting query, varying scale b

What have we discovered so far

Our choice of 50 and 51 was arbitrary, the same proof will hold for any true answers differing by 1

For a **counting query**, the database curator (**trusted curator**, trusted holder) protects privacy of each individual by reporting

$$y \sim y_{\text{true}} + \text{Lap}(b)$$

This ensures that for **any two neighboring datasets** the privacy loss is **bounded**

$$\frac{\Pr[Y = y]}{\Pr[Y' = y]} \leq \exp\left(\frac{1}{b}\right)$$

Differential privacy basics

Formalizing differential privacy

Let's generalize our findings

We had $\frac{\Pr[Y=y]}{\Pr[Y'=y]} \leq \exp\left(\frac{1}{b}\right)$

We saw, that $\frac{1}{b}$ controls the strength of privacy. Let's generalize this notion and call it

Privacy budget

Denoted as $\varepsilon \in [0, \infty)$

$\varepsilon = 0$ is complete privacy but completely random

$\varepsilon = \infty$ is no privacy whatsoever

So for our counting query example, we would have

$$\frac{\Pr[Y=y]}{\Pr[Y'=y]} \leq \exp(\varepsilon) \text{ where } \varepsilon = \frac{1}{b}$$

Let's formalize pure differential privacy

Let's just rewrite $\frac{\Pr[Y=y]}{\Pr[Y'=y]} \leq \exp(\varepsilon)$:

$$\Pr[Y = y] \leq \exp(\varepsilon) \Pr[Y' = y]$$

Let's generalize the random variables Y and Y' by saying they are **random variables parametrized by the dataset** (randomized algorithms, **randomized mechanisms**), e.g., $\mathcal{M}(D)$ or $\mathcal{M}(D')$

$$\Pr[\mathcal{M}(D) = y] \leq \exp(\varepsilon) \Pr[\mathcal{M}(D') = y]$$

Let's formalize pure differential privacy

We had $\Pr[\mathcal{M}(D) = y] \leq \exp(\varepsilon) \Pr[\mathcal{M}(D') = y]$

To make this definition work for **any** random variable (categorical, discrete finite, countably infinite, uncountable), we must generalize the co-domain of \mathcal{M}

The co-domain of $\mathcal{M}(D)$ can be some arbitrary set \mathcal{Z}

We want that our privacy guarantees hold for **any possible output** \mathcal{Y} which is any subset of the co-domain, i.e. $\mathcal{Y} \subseteq \mathcal{Z}$:

$$\Pr[\mathcal{M}(D) \in \mathcal{Y}] \leq \exp(\varepsilon) \Pr[\mathcal{M}(D') \in \mathcal{Y}]$$

C. Dwork and A. Roth (2013). "**The Algorithmic Foundations of Differential Privacy**". In: *Foundations and Trends® in Theoretical Computer Science* 9.3-4, pp. 211–407, Definition 2.4

$(\epsilon, 0)$ differential privacy (aka. pure DP)

C. Dwork and A. Roth (2013). "**The Algorithmic Foundations of Differential Privacy**". In: *Foundations and Trends® in Theoretical Computer Science* 9.3-4, pp. 211–407, Definition 2.4

A randomized algorithm (mechanism) \mathcal{M} is $(\epsilon, 0)$ -**differentially private** if for any two neighboring datasets D, D' and any output $\mathcal{Y} \subseteq \mathcal{Z}$ this guarantee holds:

$$\Pr[\mathcal{M}(D) \in \mathcal{Y}] \leq \exp(\epsilon) \Pr[\mathcal{M}(D') \in \mathcal{Y}]$$

Differential privacy basics

Laplace mechanism

Laplace mechanism for counting query is $(\epsilon, 0)$ -DP

C. Dwork and A. Roth (2013). **"The Algorithmic Foundations of Differential Privacy"**. In: *Foundations and Trends® in Theoretical Computer Science* 9.3-4, pp. 211–407, Definition 2.4

Counting query: Report a number of persons with a certain attribute

Draw a random number $y_{\text{lap}} \sim \text{Lap}(\mu = 0; b = \frac{1}{\epsilon})$

Report $y = y_{\text{true}} + y_{\text{lap}}$

Proof: We did it already! (but practice it now backwards from the definition)

Working example of numeric query

Name	Debt (in €)
Alice	2,800,798.00
Bob	7,000.00
Charlie	1.56
Xander	0.00

Table 9: Debts of account holders in a bank

Important constraint: Maximum debt the bank offers per person is 10,000,000 €

Query: What is the total debt in €?

Differential privacy basics

Global ℓ_1 sensitivity

How different neighboring datasets influence the output

Name	Debt (in €)
Alice	2,800,798.00
Bob	7,000.00
Charlie	1.56
Xander	0.00

$$D = \{A, B, C, X\}$$

- $D' = \{B, C, X\} = D \setminus \{A\}$

- $y_{\text{true}}(D) = 2,807,799.56$

- $y_{\text{true}}(D') = 7,001.56$

- $y_{\text{true}}(D) - y_{\text{true}}(D') = 2,800,798$

How different neighboring datasets influence the output

Name	Debt (in €)
Alice	2,800,798.00
Bob	7,000.00
Charlie	1.56
Xander	0.00

- $D' = D \setminus \{A\}$: $y_{\text{true}}(D) - y_{\text{true}}(D') = 2,800,798$
- $D' = D \setminus \{B\}$: $y_{\text{true}}(D) - y_{\text{true}}(D') = 7,000$
- $D' = D \setminus \{C\}$: $y_{\text{true}}(D) - y_{\text{true}}(D') = 1.56$
- $D' = D \setminus \{X\}$: $y_{\text{true}}(D) - y_{\text{true}}(D') = 0$

Each person can change the result of the query by its value

Can we use this max value to fix our problem?

Introducing global ℓ_1 sensitivity

Remember: we said the max debt/person = 10,000,000 €

We also saw that

- One person can change the sum query by the max value
- The max value was essential for scaling the Laplace distribution

Global ℓ_1 sensitivity

The ℓ_1 -sensitivity of a function $f : \mathcal{X} \rightarrow \mathbb{R}^k$

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

for any neighboring datasets D, D'

The correct scale of Laplace mechanism

Our previous observations (scale of the Laplace is proportional to the maximum difference of neighboring datasets, aka. global sensitivity) lead us to the following

Laplace mechanism for numeric queries

Draw $y_{\text{lap}} \sim \text{Lap}(\mu = 0; b = \frac{\Delta}{\epsilon})$ and report $y = y_{\text{true}} + y_{\text{lap}}$

This correct Laplace mechanism is $(\epsilon, 0)$ -differentially private

Formal proof? We have all building blocks \rightarrow homework! (or during exercise)

What we covered so far

- For provably private data analysis we need randomized algorithms
- Central (with a trusted curator) pure $(\epsilon, 0)$ differential privacy
- Laplace mechanism: numeric queries, ℓ_1 sensitivity,
- Exponential mechanism: 'any-range' queries (arbitrary sets), utility function and its sensitivity
- Local DP

We bounded our ‘privacy loss’ by ε

$$\Pr[\mathcal{M}(D) \in \mathcal{Y}] \leq \exp(\varepsilon) \Pr[\mathcal{M}(D') \in \mathcal{Y}]$$
$$\ln \left(\frac{\Pr[\mathcal{M}(D) \in \mathcal{Y}]}{\Pr[\mathcal{M}(D') \in \mathcal{Y}]} \right) \leq \varepsilon$$

What is $\mathcal{M}(D)$ (and also $\mathcal{M}(D')$)?

The private mechanism is randomized, so somewhere in the mechanism there is a random variable

- e.g., Laplace mechanism uses Laplace R.V.
- Randomized response uses Bernoulli R.V., etc.

In general, since the mechanism $\mathcal{M}(D)$ is a function of a random variable, it is also a random variable

Privacy Loss Random Variable

$\mathcal{M}(D)$ and $\mathcal{M}(D')$ are two random variables

The privacy loss random variable is defined as

$$\mathcal{L}_{\mathcal{M}(D) \parallel \mathcal{M}(D')} = \ln \left(\frac{\Pr[\mathcal{M}(D) = t]}{\Pr[\mathcal{M}(D') = t]} \right)$$

and is distributed by drawing $t \sim \mathcal{M}(D)$

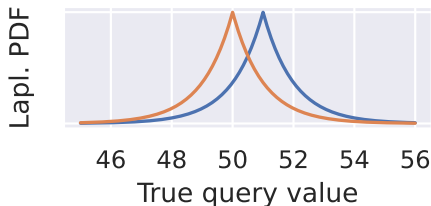
Example of Privacy Loss Random Variable

The privacy loss random variable

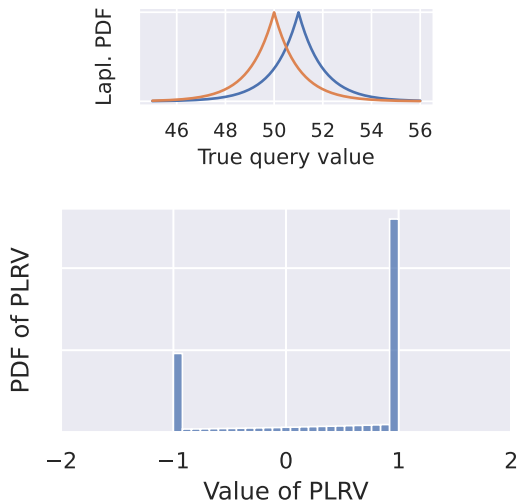
$$\mathcal{L}_{\mathcal{M}(D) \parallel \mathcal{M}(D')} = \ln \left(\frac{\Pr[\mathcal{M}(D) = t]}{\Pr[\mathcal{M}(D') = t]} \right)$$

and is distributed by drawing $t \sim \mathcal{M}(D)$

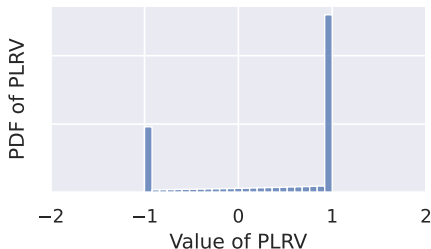
How would the distribution of $\mathcal{L}_{\mathcal{M}(D) \parallel \mathcal{M}(D')}$ would look like for the Laplace mechanism?



Example of $\mathcal{L}_{\mathcal{M}(D)||\mathcal{M}(D')}$ for Laplace mechanism $\varepsilon = 1$



Values of $|\mathcal{L}_{\mathcal{M}(D)||\mathcal{M}(D')}|$ are upper-bounded by ε in $(\varepsilon, 0)$ -DP



This distribution demonstrates (not a proof!) that the probability the value of $|\mathcal{L}_{\mathcal{M}(D)||\mathcal{M}(D')}|$ exceeds ε is zero

In other words

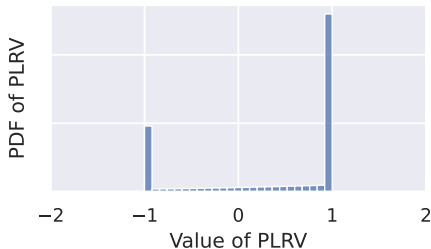
$$\Pr [|\mathcal{L}_{\mathcal{M}(D)||\mathcal{M}(D')}| \leq \varepsilon] = 1$$

Approximate Differential Privacy

- 1 Defining privacy
- 2 Text anonymization
- 3 Differential privacy basics
- 4 Approximate Differential Privacy**
- 5 DP Stochastic Gradient Descent

Maybe we don't need to always ensure the bound

What if we allow to exceed ε with some small probability δ ?



In other words change $\Pr \left[\left| \mathcal{L}_{\mathcal{M}(D) || \mathcal{M}(D')} \right| \leq \varepsilon \right] = 1$ into

$$\Pr \left[\left| \mathcal{L}_{\mathcal{M}(D) || \mathcal{M}(D')} \right| \leq \varepsilon \right] \geq 1 - \delta$$

$$\Pr \left[\left| \mathcal{L}_{\mathcal{M}(D) || \mathcal{M}(D')} \right| > \varepsilon \right] < \delta \quad (\text{equivalent})$$

Formalizing approximate (ϵ, δ) -DP

A randomized algorithm (mechanism) \mathcal{M} is (ϵ, δ) -**differentially private** if for any two neighboring datasets D, D' and any output $\mathcal{Y} \subseteq \mathcal{Z}$ this guarantee holds:

$$\Pr[\mathcal{M}(D) \in \mathcal{Y}] \leq \exp(\epsilon) \Pr[\mathcal{M}(D') \in \mathcal{Y}] + \delta$$

One immediate observation: for $\delta = 0$ we get our known 'pure' DP (that's why we called it (ϵ, δ) -DP)

C. Dwork and A. Roth (2013). "**The Algorithmic Foundations of Differential Privacy**". In: *Foundations and Trends® in Theoretical Computer Science* 9.3-4, pp. 211–407, Definition 2.4

Approximate Differential Privacy

What is this δ doing?

Extreme algorithm 1: When bad things are really bad

Our query is: Given a database of secrets, give me all rows

Name	Hospitalized in year	Age	Illegal drug use
Alice	2024	32	yes
Bob	2020	21	no
Charlie	2023	45	no
...			
Xander	2020	31	yes

Table 10: Example database D

Our goal is to have this algorithm (ε, δ) -DP

Given a database of secrets, give me all rows (part 1)

With probability $1 - \delta$, return completely random table

Name	Hospitalized in year	Age	Illegal drug use
Jim	2022	16	no
Dave	2011	71	yes
...			

Table 11: Example output of $\mathcal{M}(D)$ — completely random

Since the output is completely random, there would be no difference in outputs of any neighboring datasets D and D' , therefore this is perfectly private algorithm $\varepsilon = 0$

Given a database of secrets, give me all rows (part 2)

With probability δ , return the **original** dataset in full

Name	Hospitalized in year	Age	Illegal drug use
Alice	2024	32	yes
Bob	2020	21	no
...			

Table 12: Example output of $\mathcal{M}(D)$ — returning the full original D

This part of the algorithm is purely deterministic, there is no randomness, therefore this would be $\varepsilon = \infty$

Why? Remove Alice to get D' . But this algorithm is never going to return D' , so $\Pr[\mathcal{M}(D') = 0]$, which leads to $\frac{\Pr[\mathcal{M}(D)=x]}{\Pr[\mathcal{M}(D')=0]} \rightarrow \infty$

Given a database of secrets, give me all rows (part 3)

Summary of our algorithm:

- With prob. $1 - \delta$, return completely random table ($\varepsilon = 0$)
- With prob. δ , return the **original** dataset in full ($\varepsilon = \infty$)

Our algorithm is (ε, δ) -DP! (in fact $(0, \delta)$ -DP)

$$\Pr \left[\left| \mathcal{L}_{\mathcal{M}(D) \parallel \mathcal{M}(D')} \right| > \varepsilon \right] < \delta$$

Very bad things can happen with δ , so it should be very small! But how small?

General recommendation

We should therefore consider

$$\delta \ll \frac{1}{n}$$

(ie. very small; typically $\delta = 1 \times 10^{-6}$, aka 'cryptographically' small)

Approximate Differential Privacy

Gaussian mechanism

ℓ_2 sensitivity

Similar to ℓ_1 sensitivity of the query

ℓ_2 sensitivity

The ℓ_2 -sensitivity of a function $f : D \rightarrow \mathbb{R}^k$:

$$\Delta_2 f = \max_{D, D'} \|f(D) - f(D')\|_2$$

Gaussian mechanism

Function (numeric query) $f : D \rightarrow \mathbb{R}^k$:

Very important constraints on ε !

For $\varepsilon \in (0, 1)$ and $\delta > 0$

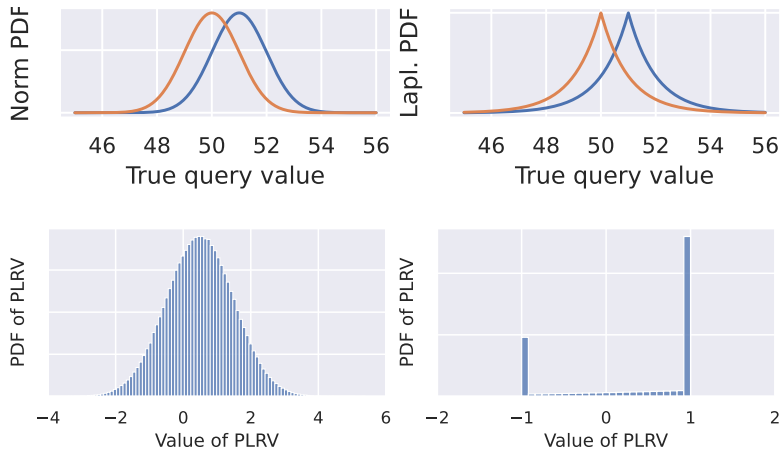
The Gaussian mechanism $\mathcal{M}(D)$ is defined as

$$f(D) + (Y_1, \dots, Y_k)$$

where each Y_n is drawn **independently** from $\mathcal{N}(0, \sigma^2)$, such that

$$\sigma^2 > 2 \ln \left(\frac{1.25}{\delta} \right) \frac{(\Delta_2)^2}{\varepsilon^2}$$

Gaussian mechanism is (ϵ, δ) -DP²



²Non-trivial proof in Appendix A of Dwork and Roth (2013); also note that there are quite a few typos there

Approximate Differential Privacy

General properties of DP algorithms

Post-processing

Let $\mathcal{M}(D) \mapsto R$ be a (ϵ, δ) -DP algorithm

Let $f : R \mapsto S$ be an arbitrary (randomized) function

Then $f(\mathcal{M}(D))$ is (ϵ, δ) -DP

In words

Whatever you do with (ϵ, δ) -DP output, you cannot 'weaken' privacy

Group privacy

Let D and D' differ in k positions.

Let $\mathcal{M}(D)$ be (ε, δ) -DP

Then for any output T we have

$$\Pr[\mathcal{M}(D) \in T] \leq \exp(k\varepsilon) \Pr[\mathcal{M}(D') \in T] + k \exp(\varepsilon \cdot (k - 1))\delta$$

Implications for large groups

If k grows, the privacy budget grows exponentially

Basic composition

Let $\mathcal{M} = (M_1, \dots, M_k)$ be a sequence of mechanisms, where each M_i is $(\varepsilon_i, \delta_i)$ -DP. (They might be adaptive)

Then \mathcal{M} is $(\sum_{i=1}^k \varepsilon_i, \sum_{i=1}^k \delta_i)$ -DP

In words

Overall privacy 'budget' can be spent for a sequence of private queries

What we covered so far

- Pure $(\epsilon, 0)$ differential privacy
- Central and Local DP
- Approximate (ϵ, δ) -DP
- Mechanisms: Laplace, Exponential, Randomized response, Gaussian
- Post processing and composition

Today

Let's finally do some supervised machine learning (neural networks)

DP Stochastic Gradient Descent

- 1 Defining privacy
- 2 Text anonymization
- 3 Differential privacy basics
- 4 Approximate Differential Privacy
- 5 DP Stochastic Gradient Descent**

DP Stochastic Gradient Descent

Finding the best model's parameters

Training as optimization

$$\mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i; \Theta), y_i)$$

The training examples are fixed, and the values of the parameters determine the loss

The goal of the training algorithm is to set the values of the parameters Θ , such that the value of \mathcal{L} is minimized

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \mathcal{L}(\Theta) = \underset{\Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i; \Theta), y_i)$$

Minibatch Stochastic Gradient Descent

```
1: function mbSGD( $f(\mathbf{x}; \Theta)$ ,  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $(\mathbf{y}_1, \dots, \mathbf{y}_n)$ ,  $L$ )
2:   while stopping criteria not met do
3:     Sample  $m$  examples  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots (\mathbf{x}_m, \mathbf{y}_m)\}$ 
4:      $\hat{\mathbf{g}} \leftarrow 0$ 
5:     for  $i = 1$  to  $m$  do
6:       Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$ 
7:        $\hat{\mathbf{g}} \leftarrow \hat{\mathbf{g}} + \text{gradient of } \frac{1}{m}L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i) \text{ wrt. } \Theta$ 
8:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$ 
9:   return  $\Theta$ 
```

DP Stochastic Gradient Descent

How to privatize SGD with DP

What can we privatize in the SGD algorithm by DP?

```
1: function SGD( $f(\mathbf{x}; \Theta)$ ,  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $(\mathbf{y}_1, \dots, \mathbf{y}_n)$ ,  $L$ )  
2:   while stopping criteria not met do  
3:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$   
4:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$   
5:      $\hat{\mathbf{g}} \leftarrow$  gradient of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  wrt.  $\Theta$   
6:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$   
7:   return  $\Theta$ 
```

- Privatize input
- Privatize output
- Privatize learning

DP Stochastic Gradient Descent

**Problem 1: Unbounded gradient and
unbounded sensitivity**

Unbounded sensitivity of gradient

Standard SGD

1: ...

2: $\mathbf{g}(\mathbf{x}_i) \leftarrow \nabla \mathcal{L}(\theta_t, \mathbf{x}_i)$ ▷ Compute gradient

3: ...

Clip the gradient vector by **per-example** ℓ_2 norm

$$\mathbf{g}(\mathbf{x}_i) \leftarrow \nabla \mathcal{L}(\theta_t, \mathbf{x}_i)$$

$$\bar{\mathbf{g}}(\mathbf{x}_i) \leftarrow \frac{\mathbf{g}(\mathbf{x}_i)}{\max\left(1, \frac{\|\mathbf{g}(\mathbf{x}_i)\|_2}{C}\right)}$$

where $C \in \mathbb{R}$ is a clipping constant (hyper-parameter)

DP Stochastic Gradient Descent

Problem 2: Too many steps for simple composition

Running several mechanisms on the same data

```
1: function SGD( $f(\mathbf{x}; \Theta)$ ,  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $(\mathbf{y}_1, \dots, \mathbf{y}_n)$ ,  $L$ )  
2:   while stopping criteria not met do  
3:     ...  
4:   return  $\Theta$ 
```

Composition theorems: Running the same or various privacy mechanisms on the same data

Basic composition — “epsilons and deltas add up”

For $k \in \mathbb{N}$, the composition of k mechanisms (each of them is (ε, δ) -DP) gives $(k\varepsilon, k\delta)$ -DP

This would lead to an excessively high overall budget

Running several mechanisms on the same data

Basic composition — “epsilons and deltas add up”

For k steps (each (ε, δ) -DP): $(k\varepsilon, k\delta)$ -DP

k -fold adaptive composition of an (ε, δ) -DP mechanism

Advanced composition — using smaller overall budget

For $\delta' > 0$ and $\varepsilon' = \varepsilon \sqrt{2k \ln(1/\delta')} + k\varepsilon(\exp(\varepsilon) - 1)$ the composite mechanism is $(\varepsilon', k\delta + \delta')$ -DP

Great news: Advanced composition gives us quadratic improvement wrt. number of steps k

■ $\approx \sqrt{k} \cdot \varepsilon$ instead of simple $k \cdot \varepsilon$

Theorem III.3 in C. Dwork, G. N. Rothblum, and S. Vadhan (2010). **“Boosting and Differential Privacy”**. In: *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*. Las Vegas, USA: IEEE, pp. 51–60

DP Stochastic Gradient Descent

Trick 3: Sub-sampling helps to reduce the budget in each step

Privacy amplification by sub-sampling

Let's define a **sampling function** that takes a dataset $D_{\text{in}} \in \mathcal{X}$ and produces another dataset $D_{\text{out}} \in \mathcal{X}$ as follows:

- For each entry t from D_{in} the function draws a binary value at random
 - We draw 'zero or one' using a Bernoulli random variable $\text{Ber}(\beta)$ parametrized by $\beta \in (0, 1)$
- If it's 1, this entry t will end up in the output dataset D_{out}
- If it's 0, this entry is ignored

Important: For each entry t the Bernoulli trial is independent of other entries

This is also known as **Poisson sampling**

Privacy amplification by sub-sampling

Let's have an $(\varepsilon_1, \delta_1)$ -DP algorithm \mathcal{A}_1

We propose a new algorithm \mathcal{A}_2 that works in two steps:

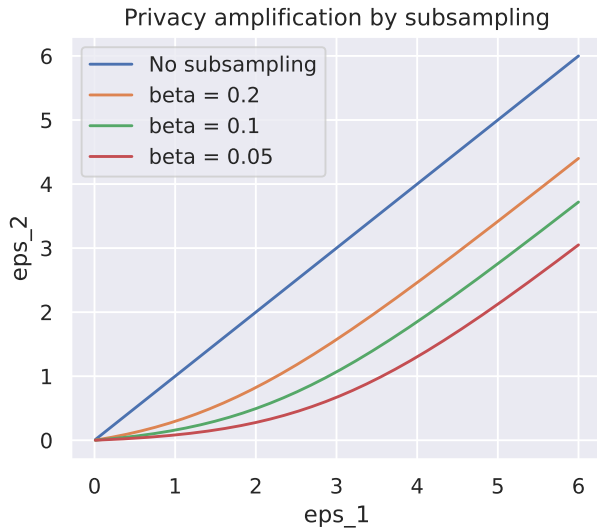
- 1 Sub-sample our dataset D using Poisson sampling (with parameter β)
- 2 Run \mathcal{A}_1 on this smaller dataset

Then \mathcal{A}_2 is $(\varepsilon_2, \delta_2)$ -DP, where

$$\varepsilon_2 = \ln(1 + \beta [\exp(\varepsilon_1) - 1]) \quad \delta_2 = \beta \delta_1$$

Proof in the appendix of N. Li, W. Qardaji, and D. Su (2012). **"On Sampling, Anonymization, and Differential Privacy Or, K-Anonymization Meets Differential Privacy"**. In: *Proceedings of the 7th ACM Symposium on Information, Computer and Communications Security*. Seoul, South Korea: ACM, pp. 32–33; there are a few 'nasty' typos.

How much we can 'save' on the privacy budget?



Why is Poisson sampling relevant for SGD?

Recall Mini-batch SGD!

- 1: **function** mbSGD($f(\mathbf{x}; \Theta)$, $(\mathbf{x}_1, \dots, \mathbf{x}_n)$, $(\mathbf{y}_1, \dots, \mathbf{y}_n)$, L)
- 2: **while** stopping criteria not met **do**
- 3: Sample m examples $\{(\mathbf{x}_1, \mathbf{y}_1), \dots (\mathbf{x}_m, \mathbf{y}_m)\}$
- 4: ...

- We usually use small 'batches' which are somehow randomly subsampled from the training dataset
- We can replace the minibatch sampling with Poisson sampling!

DP Stochastic Gradient Descent

DP-SGD

DP-SGD algorithm

```
1: function DP-SGD( $f(\mathbf{x}; \Theta)$ ,  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $|L|$  — 'lot' size,  $T$  — # of steps)
2:   for  $t \in (1, 2, \dots, T)$  do
3:     Add each training example to a 'lot'  $L_t$  with probability  $|L|/n$ 
4:     for each example in the 'lot'  $\mathbf{x}_i \in L_t$  do
5:        $\mathbf{g}(\mathbf{x}_i) \leftarrow \nabla \mathcal{L}(\theta_t, \mathbf{x}_i)$  ▷ Compute gradient
6:        $\bar{\mathbf{g}}(\mathbf{x}_i) \leftarrow \mathbf{g}(\mathbf{x}_i) / \max(1, \|\mathbf{g}(\mathbf{x}_i)\|/C)$  ▷ Clip gradient
7:        $\tilde{\mathbf{g}}(\mathbf{x}_i) \leftarrow \bar{\mathbf{g}}(\mathbf{x}_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I})$  ▷ Add noise
8:        $\hat{\mathbf{g}} \leftarrow \frac{1}{|L|} \sum_{k=1}^{|L|} \tilde{\mathbf{g}}(\mathbf{x}_k)$  ▷ Gradient estimate of 'lot' by averaging
9:        $\Theta_{t+1} \leftarrow \Theta_t - \eta_t \hat{\mathbf{g}}$  ▷ Update parameters by gradient descend
10:  return  $\Theta$ 
```

Stochastic gradient descent with differential privacy

Setup: A set of labeled i.i.d. examples — like tabular data (each example = single person)

Privacy 'accountant' — utilizes composition of DP

- Computes the privacy cost at each access to the training data (gradient computation)
- Accumulates this cost as the training progresses

Tightest privacy by numerical integration to get bounds on the **moment generating function** of the **privacy loss random variable** for all moments ≤ 32

M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang (2016). **"Deep Learning with Differential Privacy"**. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. Vienna, Austria: ACM, pp. 308–318

Recap of DP-SGD

- DP-SGD 'de-facto' standard for supervised training with DP
- Implemented in Opacus, Tensorflow privacy, and other libs

M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang (2016). “**Deep Learning with Differential Privacy**”. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. Vienna, Austria: ACM, pp. 308–318

What makes it tricky?

- Remember: Data points **must** be independent (privacy-wise)
- Scalability: Per-example gradient norm and clipping is super slow

DP Stochastic Gradient Descent

The Obvious Application: Supervised
Training

DP-SGD across various NLP tasks

Setup:

Although DP-SGD had been used in language modeling, the community lacked a thorough understanding of its usability across different NLP tasks

Research questions:

- Which models and training strategies provide the best trade-off between privacy and performance on different NLP tasks?
- How exactly do increasing privacy requirements hurt the performance?

M. Senge, T. Igamberdiev, and I. Habernal (2022). **“One size does not fit all: Investigating strategies for differentially-private learning across NLP tasks”**. In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Ed. by Y. Goldberg, Z. Kozareva, and Y. Zhang. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, pp. 7340–7353

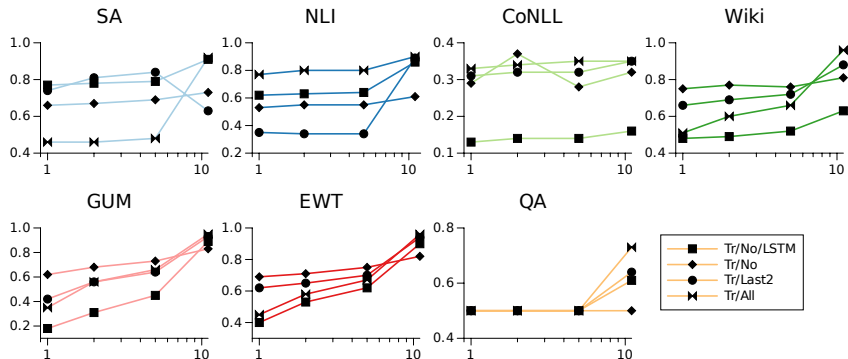
DP-SGD across various NLP tasks: Datasets

Task	Dataset	Size	Classes
SA	IMDb	50k documents	2
NLI	SNLI	570k pairs	3
NER	CoNLL'03	≈ 300k tokens	9
NER	Wikiann	≈ 320k tokens	7
POS	GUM	≈ 150k tokens	17
POS	EWT	≈ 254k tokens	17
QA	SQuAD 2.0	150k questions	★

Table 13: Datasets and their specifics. ★ SQuAD contains 100k answerable and 50k unanswerable questions, where answerable questions are expressed as the span positions of their answer.

M. Senge, T. Igamberdiev, and I. Habernal (2022). **“One size does not fit all: Investigating strategies for differentially-private learning across NLP tasks”**. In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Ed. by Y. Goldberg, Z. Kozareva, and Y. Zhang. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, pp. 7340–7353

DP-SGD across various NLP tasks: Results



M. Senge, T. Igamberdiev, and I. Habernal (2022). **“One size does not fit all: Investigating strategies for differentially-private learning across NLP tasks”**. In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Ed. by Y. Goldberg, Z. Kozareva, and Y. Zhang. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, pp. 7340–7353

Figure 12: Comparison of BERT performances (macro F_1 score) per dataset with varying privacy budget $\epsilon \in \{1, 2, 5, \infty\}$ on the x -axis (note the log scale).

DP Stochastic Gradient Descent

When Things Go Very Tricky

Private information in text?

our understanding of what is *private information* in textual data is still very limited

Applications of DP — guarantee to each individual *data point*

For textual data, a single data point will often be a sentence or document.

However, this does not mean that there is a one-to-one mapping from *individuals* to sentences and documents. For instance, multiple documents could potentially refer to the same individual, or contain the same piece of sensitive information that would break the assumption of each data point being independent.

T. Igamberdiev, D. N. L. Vu, F. Kuennecke, Z. Yu, J. Holmer, and I. Habernal (2024). **"DP-NMT: Scalable Differentially Private Machine Translation"**.

In: *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*. Ed. by N. Aletras and O. De Clercq. St. Julians, Malta: Association for Computational Linguistics, pp. 94–105

Private information in text?

“In this paper, we discuss the mismatch between the narrow assumptions made by popular data protection techniques (data sanitization and differential privacy), and the broadness of natural language and of privacy as a social norm.”

“We argue that existing protection methods cannot guarantee a generic and meaningful notion of privacy for language models. We conclude that language models should be trained on text data which was explicitly produced for public use.”

H. Brown, K. Lee, F. Miroshghallah, R. Shokri, and F. Tramèr (2022). **“What Does it Mean for a Language Model to Preserve Privacy?”** In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. New York, NY, USA: ACM, pp. 2280–2292

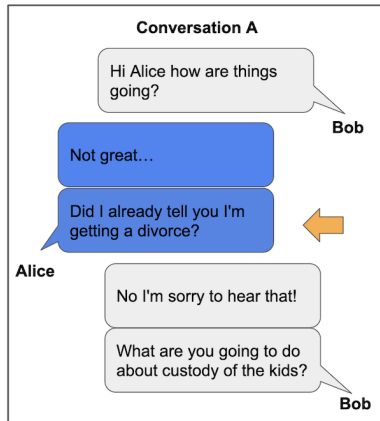
Private information in text?

The approach to preserving privacy in LMs has been to attempt complete removal of private information from training data (data sanitization), or to design algorithms that do not memorize private data, such as algorithms that satisfy differential privacy (DP)

Both methods make explicit and implicit assumptions about the structure of data to be protected, the nature of private information, and requirements for privacy, that do not hold for the majority of natural language data.

H. Brown, K. Lee, F. Miroshghallah, R. Shokri, and F. Tramèr (2022). **"What Does it Mean for a Language Model to Preserve Privacy?"** In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. New York, NY, USA: ACM, pp. 2280–2292

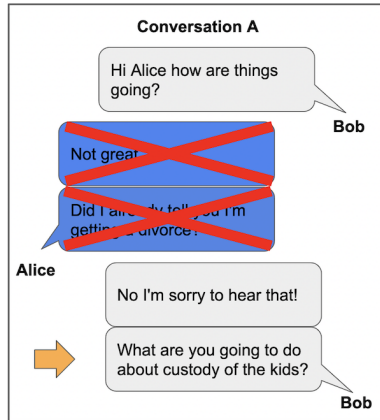
Private information in text?



H. Brown, K. Lee, F. Miresghallah, R. Shokri, and F. Tramèr (2022). **"What Does it Mean for a Language Model to Preserve Privacy?"** In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. New York, NY, USA: ACM, pp. 2280–2292

Figure 13: Original conversation. Private information indicated by orange arrows.

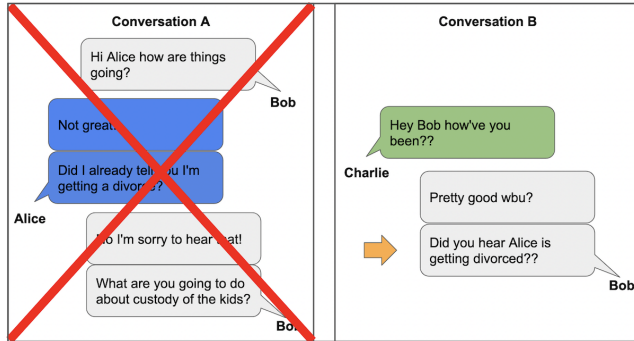
Private information in text?



H. Brown, K. Lee, F. Miresghallah, R. Shokri, and F. Tramèr (2022). **"What Does it Mean for a Language Model to Preserve Privacy?"** In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. New York, NY, USA: ACM, pp. 2280–2292

Figure 14: Alice's messages removed. Bob's last message still includes her private information.

Private information in text?



H. Brown, K. Lee, F. Miresghallah, R. Shokri, and F. Tramèr (2022). **“What Does it Mean for a Language Model to Preserve Privacy?”** In: *2022 ACM Conference on Fairness, Accountability, and Transparency*. New York, NY, USA: ACM, pp. 2280–2292

Figure 15: The whole original conversation is removed. Conversation B still contains Alice's private information though she is not in the conversation.

License and credits

Licensed under Creative Commons
Attribution-ShareAlike 4.0 International
(CC BY-SA 4.0)



Credits

Ivan Habernal

Content from ACL Anthology papers licensed under CC-BY
<https://www.aclweb.org/anthology>

Partly inspired by lectures from Antti Honkela, Aurélien Bellet, Gautam Kamath