

Privacy-Preserving Natural Language Processing

RUHR
UNIVERSITÄT
BOCHUM

RUB

Lecture 1 — Introduction

Prof. Dr. Ivan Habernal

April 10, 2025

www.trusthlt.org

Chair of Trustworthy Human Language Technologies (TrustHLT)

Ruhr University Bochum & Research Center Trustworthy Data Science and Security



CENTER FOR TRUSTWORTHY
DATA SCIENCE AND SECURITY

Motivation

What is privacy?

Why should we care?

Defining privacy

- 1 Defining privacy
- 2 When things go wrong
- 3 Linkage attacks
- 4 And now something completely different...
- 5 Violation of privacy in NLP
- 6 Outlook and course logistics

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

What is privacy (contd.)

We also seek to emphasize that the private thing **should not be accessed or known** by others, certainly not by the general public

- If we consider something private, it belongs to us and we get to decide what happens with it
- → normative¹ meaning of the word

¹“Relating to an ideal standard or model; how things should be, how to value them

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

Privacy as the right to be let alone

S. D. Warren and L. D. Brandeis (1890).
"The Right to Privacy". In: *Harvard Law Review* 4.5, pp. 193–220

HARVARD LAW REVIEW.

VOL. IV.

DECEMBER 15, 1890.

NO. 5.

THE RIGHT TO PRIVACY.

"It could be done only on principles of private justice, moral fitness, and public convenience, which, when applied to a new subject, make common law without a precedent; much more when received and approved by usage."

WILLES, J., in *Millar v. Taylor*, 4 Burr. 2303, 2312.

THAT the individual shall have full protection in person and in property is a principle as old as the common law; but

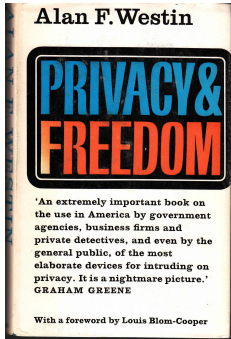
Privacy as the right to be let alone

S. D. Warren and L. D. Brandeis (1890).
"The Right to Privacy". In: *Harvard Law Review* 4.5, pp. 193–220

Warren and Brandeis (1890) paved the way to defining the actual term privacy

- Demanded 'the right to be let alone', meaning that privacy translates to freedom from intrusion by the press
- Their essay inspired an international conversation on how to define privacy in law and beyond

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)

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 over-crowded 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

When things go wrong

- 1 Defining privacy
- 2 When things go wrong**
- 3 Linkage attacks
- 4 And now something completely different...
- 5 Violation of privacy in NLP
- 6 Outlook and course logistics

The Netflix competition

On October 2, 2006, Netflix announced the 1-million Netflix Prize for improving their movie recommendation service

A. Narayanan and V. Shmatikov (2008).
“How To Break Anonymity of the Netflix Prize Dataset”. In: *arXiv preprint*

- Netflix publicly released a dataset containing 100,480,507 movie ratings, created by 480,189 Netflix subscribers between December 1999 and December 2005
- Almost 15% of all their customers
- The ratings data appear to not have been perturbed to any significant extent

Microdata

Datasets containing “micro-data,” that is, information about specific individuals, are increasingly becoming public—both in response to “open government” laws, and to support data mining research

Privacy risks of publishing micro-data are well-known

Even if identifying information such as names, addresses, and Social Security numbers has been removed, the adversary can use contextual and background knowledge, as well as cross-correlation with publicly available databases, to re-identify individual data records.

A. Narayanan and V. Shmatikov (2008).
“How To Break Anonymity of the Netflix Prize Dataset”. In: *arXiv preprint*

Microdata

Micro-data are characterized by high dimensionality and sparsity

Many attributes, each of which can be viewed as a dimension (an attribute can be thought of as a column in a database schema)

Sparsity means that a pair of random records are located far apart in the multi-dimensional space defined by the attributes (individual transaction and preference records tend to include statistically rare attributes)

A. Narayanan and V. Shmatikov (2008).
"How To Break Anonymity of the Netflix Prize Dataset". In: *arXiv preprint*

What is the database

Database D is an $N \times M$ matrix

- Each row is a record associated with some individual
- Columns are attributes
- Values might be boolean, int, date, etc.

The number of columns reflects the total number of items (e.g., a few thousands for movies)

Each column = a dimension, each individual record as a point in the multidimensional attribute space

“Collaborative filtering” — predicting a customer’s future choices using the knowledge of what similar customers did

A. Narayanan and V. Shmatikov (2008).
“How To Break Anonymity of the Netflix Prize Dataset”. In: *arXiv preprint*

The Netflix competition

Removing the identifying information from the records is not sufficient for anonymity.

The Netflix competition

Removing the identifying information from the records is not sufficient for anonymity.

Auxiliary information about some subscriber's movie preferences: the titles of a few of the movies that this subscriber watched, whether she liked them or not, maybe even approximate dates when she watched them.

The Netflix competition

Removing the identifying information from the records is not sufficient for anonymity.

Auxiliary information about some subscriber's movie preferences: the titles of a few of the movies that this subscriber watched, whether she liked them or not, maybe even approximate dates when she watched them.

How much does the adversary need to know about a Netflix subscriber in order to identify her record in the dataset, and thus learn her complete movie viewing history?

A. Narayanan and V. Shmatikov (2008).
"How To Break Anonymity of the Netflix Prize Dataset". In: *arXiv preprint*

Adversary and results of the attack

Adversary model

The adversary's goal is to de-anonymize an anonymous record r from the public database.

A. Narayanan and V. Shmatikov (2008).
"How To Break Anonymity of the Netflix Prize Dataset". In: *arXiv preprint*

Results of de-anonymization

- Very little auxiliary information is needed for de-anonymize an average subscriber record
- With 8 movie ratings (of which 2 may be completely wrong) and dates that may have a 14-day error, 99% of records be uniquely identified in the dataset

Does privacy of Netflix ratings matter?

The privacy question is not “Does the average Netflix subscriber care about the privacy of his movie viewing history?,” but “Are there any Netflix subscribers whose privacy can be compromised by analyzing the Netflix Prize dataset?”

The answer to the latter question is, undoubtedly, yes.

Experiments with cross-correlating non-anonymous records from the Internet Movie Database with anonymized Netflix records, it is possible to learn sensitive non-public information about a person's political or even sexual preferences.

A. Narayanan and V. Shmatikov (2008).
“How To Break Anonymity of the Netflix Prize Dataset”. In: *arXiv preprint*

Does privacy of Netflix ratings matter?

A. Narayanan and V. Shmatikov (2008).
“How To Break Anonymity of the Netflix Prize Dataset”. In: *arXiv preprint*

Consider the information that we have been able to deduce by locating **one of these users' entire movie viewing history** in the Netflix dataset and that cannot be deduced from his public IMDb ratings.

De-anonymization example

“First, we can immediately find his political orientation based on his strong opinions about “Power and Terror: Noam Chomsky in Our Times” and “Fahrenheit 9/11.” Strong guesses about his religious views can be made based on his ratings on “Jesus of Nazareth” and “The Gospel of John.” He did not like “Super Size Me” at all; perhaps this implies something about his physical size? Both items that we found with predominantly gay themes, “Bent” and “Queer as folk” were rated one star out of five. He is a cultish follower of “Mystery Science Theater 3000.” This is far from all we found about this one person, but having made our point, we will spare the reader further lurid details.”

A. Narayanan and V. Shmatikov (2008).
“How To Break Anonymity of the Net-
flix Prize Dataset”. In: *arXiv preprint*

Linkage attacks

- 1 Defining privacy
- 2 When things go wrong
- 3 Linkage attacks**
- 4 And now something completely different...
- 5 Violation of privacy in NLP
- 6 Outlook and course logistics

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*

And now something completely different...

- 1 Defining privacy
- 2 When things go wrong
- 3 Linkage attacks
- 4 And now something completely different...**
- 5 Violation of privacy in NLP
- 6 Outlook and course logistics

"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*

Violation of privacy in NLP

- 1 Defining privacy
- 2 When things go wrong
- 3 Linkage attacks
- 4 And now something completely different...
- 5 Violation of privacy in NLP**
- 6 Outlook and course logistics

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*

Training data memorization

Neural networks, especially ones with many parameters, can memorize their training data

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

This can be exploited by adversaries via **membership inference attacks** that infer whether an example was in the training set, or more powerful data extraction attacks that **recover full training examples**

Extractable memorization

Given a model with a generation routine Gen , an example \mathbf{x} from the training set \mathbb{X} is extractably memorized if an adversary (without access to \mathbb{X}) can construct a prompt \mathbf{p} that makes the model produce \mathbf{x} (i.e., $\text{Gen}(\mathbf{p}) = \mathbf{x}$).

Extracting attack by Carlini et al. (2020)

Recovered training examples from GPT-2

- by prompting it with short strings sampled from the public Internet
- then manually checking whether these strings can also be found with a Google search

Simply prompting the model with data sampled from the model's training distribution (GPT-2 was trained on some unknown text sampled from the Internet), and (reasonably) assuming that any string memorized by the model is also contained in Google's search index

N. Carlini et al. (2020). **"Extracting Training Data from Large Language Models"**. In: *arXiv preprint*

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

Attack on open-source LLMs by Nasr et al. (2023)

- 1 Download 108 bytes of data from Wikipedia, and generate prompts by randomly sampling (with replacement) hundreds of millions of continuous 5-token blocks from this dataset
- 2 Perform an independent generation for each prompt store the results

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

Open-source models: able to directly query the model's training data to evaluate whether or not any generated sample is memorized

Extraction is successful if the model outputs text that contains a substring of length at least 50 tokens that is contained verbatim in the training set

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*

Outlook and course logistics

- 1 Defining privacy
- 2 When things go wrong
- 3 Linkage attacks
- 4 And now something completely different...
- 5 Violation of privacy in NLP
- 6 Outlook and course logistics**

Course content

We should cover the following broad topics

- PII, text redaction, text anonymization
- Differential privacy
- Unlearning

Course logistics

Lectures (Thursdays)

- Mostly theoretical concepts

Exercises (Thursdays)

- Mix of theoretical (e.g., proofs) and practical (e.g., programming)
- Format to be discussed

Exam and grading

Written exam only

No homework assignment (scalability issue without student assistants), no bonus points

FAQ: “What is exam-relevant?”

Everything we discuss here and in the exercises.

Mostly I will **highlight** very important things.

Some things will be optional.

Communication

Important stuff: e-mail!

`ivan.habernal@ruhr-uni-bochum.de`

Announcements: Moodle

Informally: Discord?

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