#### **Data privacy: Introduction**

Vicenç Torra

January 2025

Umeå University, Sweden

V. Torra (2022) A guide to data privacy, Springer (Chapter 1)

## Outline

#### 1. Introduction

- A context
- Privacy for machine learning and statistics
- Two motivating examples
- Data privacy: the core
- Data privacy: difficulties
- Privacy and society

#### 2. Concepts

- Terminology
- Transparency
- Interim discussion
- Privacy by design
- 3. Summary
- 4. References

# Introduction

 ${\sf Introduction} >$ 

### Introduction to data privacy

#### • Chapter 1. Introduction

- 1 Introduction
- 1.1 Motivations for Data Privacy
- 1.1.1 Privacy, Security and Inference
- 1.2 Two Motivating Examples
- 1.2.1 Sharing a Database
- 1.2.2 Sharing a Computation
- 1.2.3 Privacy Leakages and Risk
- 1.3 Privacy and Society
- 1.4 Terminology
- 1.4.1 The Framework
- 1.4.2 Anonymity and Unlinkability
- 1.4.3 Disclosure
- 1.4.4 Dalenius' Definitions for Attribute and Identity Disclosure
- 1.4.5 Plausible Deniability
- 1.4.6 Undetectability and Unobservability
- 1.4.7 Pseudonyms and Identity
- 1.4.8 Transparency
- 1.5 Privacy and Disclosure
- 1.6 Privacy by Design
- 1.7 Bibliographical Notes

## A context:

#### **Data-driven machine learning/statistical models**<sup>1</sup>

<sup>1</sup>This is not the only type of scenario/application context we will discuss

Vicenç Torra; Data privacy: Introduction

## Prediction using (machine learning/statistical) models

 Data is collected to be used (otherwise, better not to collect them<sup>2</sup>)

<sup>2</sup>Concept: Data minimization

## Prediction using (machine learning/statistical) models

 Application of a model for decision making data ⇒ prediction/decision



• Example: predict the length-of-stay at admission

## Data-driven machine learning/statistical models

- From (huge) databases, build the "decision maker"
  - Use (logistic) regression, deep lerning, neural networks, . . . classification algorithms, decision trees, . . .



• Example: build a predictor from hospital historical data about lengthof-stay at admission

## **Privacy for machine learning and statistics:**

Data-driven machine learning/statistical models

### Data is sensitive

- Who/how is going to create this model (this "decision maker")?
- Case #1. Sharing (part of the data)



### Data is sensitive

- Who/how is going to create this model (this "decision maker")?
- Case #2. Not sharing data, only querying data



## **Two motivating examples**

- Data privacy: core
  - Someone needs to access to data to perform authorized analysis, but access to the data and the result of the analysis should avoid disclosure.



E.g., you are authorized to compute the average stay in a hospital, but maybe you are not authorized to see the length of stay of your neighbor.

- Case #1. Sharing (part of the data)
- Q: How different children ages and diagnoses affect this length of stay? Average length of stay is decreasing in the last years due to new hospital policies?
- Data: Existing database with previous admissions (2010–2019). To avoid disclosure a view of the DB restricting records to children born before 2019 and only providing for these records year of birth, town, year of admission, illness, and length of stay.

### Data is sensitive

### • Case #1. Sharing (part of the data)

Year birth	Year Admission	Town	Illness	Length stay (days)
2017	2019	Umeå	а	3
2015	2020	Umeå	b	2
2011	2020	Luleå	С	5
2017	2019	Luleå	а	2
2016	2020	Dorotea	b	4
2016	2020	Holmöns	d	2
2015	2019	Täfteå	е	4
2015	2019	Täfteå	е	4
2015	2018	Täfteå	е	4
2015	2018	Täfteå	е	4

#### • Is this data safe?

## Data is sensitive

### • Case #1. Sharing (part of the data)

Year birth	Year Admission	Town	Illness	Length stay (days)
2017	2019	Umeå	а	3
2015	2020	Umeå	b	2
2011	2020	Luleå	С	5
2017	2019	Luleå	а	2
2016	2020	Dorotea	b	4
2016	2020	Holmöns	d	2
2015	2019	Täfteå	е	4
2015	2019	Täfteå	е	4
2015	2018	Täfteå	е	4
2015	2018	Täfteå	е	4

#### • Is this data safe?

Holmöns 63, Täfteå 1383, Luleå 49123, Umeå 83249, Dorotea 2366

- Case #1. Sharing (part of the data). Example 2
  - $\,\circ\,$  Q: stress influenced by studies and commuting distance ?

```
    Data: DB/view = (where students live, what study, got sick?)
    ( Umeå, CS, no )
```

Ú Umeå, CS, yes Umeå, . . . , . . .

```
(Holmsund, CS, no
(Holmsund, CS, no
(Holmsund, CS, yes
(Holmsund, ..., ...
```

( Norrbyn, XXXX, yes )
O No "personal data", is this ok ? NO!!

 $\Rightarrow$  We learn that our friend<sup>3</sup> is sick !!

<sup>&</sup>lt;sup>3</sup>Norrbyn 168 (2023)

- Case #2. Sharing a computation.
  - Q: Mean income of admitted to hospital unit (e.g., psychiatric unit) for a given Town (Bunyola)?

- Case #2. Sharing a computation.
  - Q: Mean income of admitted to hospital unit (e.g., psychiatric unit) for a given Town (Bunyola)?
  - Mean income is not "personal data", is this ok ? NO!!:

- Case #2. Sharing a computation.
  - Q: Mean income of admitted to hospital unit (e.g., psychiatric unit) for a given Town (Bunyola)?
  - Mean income is not "personal data", is this ok ? NO!!:
  - $\circ$  Example 1000 2000 3000 2000 1000 6000 2000 10000 2000 4000  $\Rightarrow$  mean = 3300

- Case #2. Sharing a computation.
  - Q: Mean income of admitted to hospital unit (e.g., psychiatric unit) for a given Town (Bunyola)?
  - Mean income is not "personal data", is this ok ? NO!!:
  - $\circ$  Example 1000 2000 3000 2000 1000 6000 2000 10000 2000 4000  $\Rightarrow$  mean = 3300
  - Adding Ms. Rich's salary 100,000 Eur/month: mean = 12090,90 !

     (a extremely high salary changes the mean significantly)
     ⇒ We infer Ms. Rich from Town was attending the unit

• Case #2. Sharing a computation. Example 2



 Regression of income with respect to age with (right) and without (left) the record of Dona Obdúlia

• income = -4524.2 + 207.5 age (without Ms. Rich = Dona Obdúlia)

 $\circ$  income = -54307 + 1652 age (with Ms. Rich = Dona Obdúlia)

## Data privacy: the core

### Data is sensitive: Data privacy

- Privacy leaks
  - Transmissions of unencrypted e-mails: senders, receivers, content
  - Unauthorized access to an account or service: sensitive information
  - AOL<sup>4</sup>: linkage of search log queries
    - $\Rightarrow$  User No. 4417749, hundreds of searches over a three-month

period including queries 'landscapers in Lilburn, Ga'

- → Thelma Arnold identified!
- Netflix case<sup>5</sup>: linkage of movie ratings

 $\Rightarrow$  individual users matched with film ratings on the Internet Movie Database.

• Similar with credit card payments, shopping carts, ...

<sup>&</sup>lt;sup>4</sup>http://www.nytimes.com/2006/08/09/technology/09aol.html <sup>5</sup>Narayanan and Shmatikov, 2008

## Data is sensitive: Data privacy

- A personal view of core and boundaries of data privacy: boundaries
  - $\circ\,$  Database in a computer or in a removable device
    - $\Rightarrow$  access control to avoid unauthorized access
    - $\implies$  Access to address (admissions), Access to blood test (admissions?)
  - $\circ$  Data is transmitted
    - $\Rightarrow$  security technology to avoid unauthorized access
    - $\implies$  Data from blood glucose meter sent to hospital. Network sniffers
      - Transmission is sensitive: Near miss/hit report to car manufacturers
- Data privacy: disclosure because of inference



## Data is sensitive: Data privacy

- A personal view of core and boundaries of data privacy: core
- Someone needs to access to data to perform authorized analysis, but access to the data and the result of the analysis should avoid disclosure.
  - $\circ$  data uses / rellevant techniques
    - ▷ Data to be used for data analysis
      - $\Rightarrow$  compute indices, find patterns, build models
      - $\Rightarrow$  statistics, machine learning, data mining
    - ▷ Data is transmitted
      - ⇒ communication and inform. retrieval (identities/data)
      - $\Rightarrow$  communications



## **Data privacy: difficulties**

• Naive anonymization does not work

Passenger manifest for the Missouri, arriving February 15, 1882; Port of Boston<sup>6</sup> Names, Age, Sex, Occupation, Place of birth, Last place of residence, Yes/No, condition (healthy?)

AT TEAM & CO. LIN SCHEMERS, LEW SCHEMERS, LEW SCHEMERS, LEW SCHEMERS, Standard State of Marrielle & Response of Mar	ST C	DF I BES taken In Manna Or of Manna exit	ASSEN tordine Jordine Martine	GERS. J. J. J	ting to	nden ut	HE GI and the count of the Boston the Boston the count of the the count of the file the file the count of the count of the file the count of the file t
NAMER States	4 Chul	ARL.	OCCUPATION.	PLACE OF BURTH.	Last Part of Bestleme	Hits Assettan	CONDITION
	PARTE STATIS	Male	Gullenan Saite Gangeran kalilenan	Hochers bud Radrat chile Souten Boton Hi. Soton Hi. Cours bug Marter His Marter His	Childrent big Settland big Constra Boston H.L. Boston H.L. Const King Brechn H.L. Brechn H.L.	que se an	Arethy.

<sup>6</sup>https://www.sec.state.ma.us/arc/arcgen/genidx.htm

- Difficulties: highly identifiable data
  - (Sweeney, 1997) on USA population
    - ▷ 87.1% (216 million/248 million) were likely made them unique based on

5-digit ZIP, gender, date of birth,

- ▷ 3.7% (9.1 million) had characteristics that were likely made them unique based on
  - 5-digit ZIP, gender, Month and year of birth.

- Difficulties: highly identifiable data and high dimensional data
  - Data from mobile devices:
    - ⇒ two positions can make you unique (home and working place)
  - $\circ$  AOL<sup>7</sup> and Netflix cases (search logs and movie ratings)
    - $\Rightarrow$  User No. 4417749, hundreds of searches over a three-month period including queries 'landscapers in Lilburn, Ga'
      - $\longrightarrow$  Thelma Arnold identified!
    - $\Rightarrow$  individual users matched with film ratings on the Internet Movie Database.
  - Similar with credit card payments, shopping carts, ...

<sup>&</sup>lt;sup>7</sup>http://www.nytimes.com/2006/08/09/technology/09aol.html

- Difficulties: highly identifiable data and high dimensional data
  - Ex1: Sickness influenced by studies and commuting distance ?
  - Ex2: Mean income of admitted to hospital unit (e.g., psychiatric unit) for a given Town?
  - Ex3: Driving behavior in the morning
    - ▷ Automobile manufacturer uses (data from vehicles)
    - ▷ Data: First drive after 6:00am
      - (GPS origin + destination, time)  $\times$  30 days
    - ▷ No "personal data", is this ok?: NO!!!:
    - How many cars from your home to your work? Are you exceeding the speed limit? Are you visiting a psychiatric clinic every tuesday?

- Difficulties: highly identifiable data and high dimensional data
  - Data from people: In some cases, it is the combination of characteristics that make you unique. (Search logs, Market Basket Analysis)
    - ⊳ Manga
    - ▷ K-pop music
    - ▷ Opera
    - ▷ IF Björklöven
    - ▷ ...

 $\rightarrow |\{ \text{ Manga} \cap k\text{-pop} \cap \text{Opera} \cap \text{IF Björklöven } \}| = 1 = 1$ 

- Data privacy is "impossible", or not? challenging
  - Privacy vs. utility
  - Privacy vs. security
  - Computationally feasible

## **Privacy and society**

- Privacy a fundamental right. (Ch. 1.1)
  - Universal Declaration of Human Rights (UN). European Convention on Human Rights (Council of Europe). General Data Protection Regulation - GDPR (EU). National regulations.
- Regulations
  - EU: General Data Protection Regulation (GDPR)
  - USA: Health Insurance Protability and Accountability Act (HIPAA, 1996), Patriot Act (2001), Homeland Security Act (2002), Children's Online Privacy Protection Act (COPPA, 2000).
  - California: California Consumer Privacy Act (CCPA)

## Legislation and motivation

### • Enforcement

- Obligations with respect to data processing
- Requirement to report personal data breaches
- Grant individual rights (to be informed, to access, to rectification, to erasure, ...)
- Data protection officer<sup>8</sup>
- $\circ$  Data protection impact assessment (DPIA) / privacy impact assessment (PIA)  $^9$
- Companies own interest.
  - Competitors can take advantage of information.
- Avoiding privacy breach. Several well known cases.

<sup>&</sup>lt;sup>8</sup>https://edps.europa.eu/data-protection/data-protection/reference-library/data-protection/reference-library/data-protection/reform/rules-business-au

- Privacy and society
  - Not only a computer science/technical problem
    - ▷ Social roots of privacy
    - Multidisciplinary problem
  - Social, legal, philosophical questions
  - Culturally relative?
    - I.e., the importance of privacy is the same among all people ?
  - Are there aspects of life which are inherently private or just conventionally so?

- Privacy and society. Is this a new problem? Yes and not
  - $\circ$  No side. See the following:

Instantaneous photographs and newspaper enterprise have invaded the sacred precincts of private and domestic life; and numerous mechanical devices threaten to make good the prediction that "what is whispered in the closet shall be proclaimed from the house-tops." (...)

Gossip is no longer the resource of the idle and of the vicious, but has become a trade, which is pursued with industry as well as effrontery (...) To occupy the indolent, column upon column is filled with idle gossip, which can only be procured by intrusion upon the domestic circle.

(S. D. Warren and L. D. Brandeis, 1890)<sup>10</sup>

• Yes side: big data, storage, surveillance/CCTV, RFID, IoT

<sup>&</sup>lt;sup>10</sup>https://www.jstor.org/stable/1321160

# Concepts

#### Technical solutions

- Statistical disclosure control (SDC)
- Privacy preserving data mining (PPDM)
- Privacy enhancing technologies (PET)
- Socio-technical aspects
  - Technical solutions are not enough
  - Implementation/management of solutions for achieving data privacy need to have a holistic perspective of information systems
  - E.g., employees and customers: how technology is applied

## Terminology

• Terminology using as framework a communication network with senders (actors) and receivers (actees)



- Attacker, adversary, intruder
  - $\circ$  the set of entities working against some protection goal
  - increase their knowledge (e.g., facts, probabilities, ...)
     on the items of interest (lol) (senders, receivers, messages, actions)

- Anonymity set. Anonymity of a subject means that the subject is not identifiable within a set of subjects, the anonymity set. Not distinguishable!
- Unlinkability. Unlinkability of two or more lol, the attacker cannot sufficiently distinguish whether these lols are related or not.
   ⇒ Unlinkability with the sender implies anonymity of the sender.
  - Linkability but anonymity. E.g., an attacker links all messages of a transaction, due to timing, but all are encrypted and no information can be obtained about the subjects in the transactions: anonymity not compromised.

(region of the anonymity box outside unlinkability box)

- Examples of anonymity in communications (definition of IoI):
  - Sender anonymity. No link between a message and the sender.
  - **Recipient anonymity.** No link between a message and the receiver.
  - Relationship anonymity. No link between a message and both sender and receiver.



- Disclosure. Attackers take advantage of observations to improve their knowledge on some confidential information about an IoI.
   ⇒ SDC/PPDM: Observe DB, ∆ knowledge of a particular subject (the respondent in a database)
  - Identity disclosure (entity disclosure). Linkability. Finding Mary in the database.
  - Attribute disclosure. Increase knowledge on Mary's salary.
     also: learning that someone is in the database, although not found.

- Disclosure. Discussion.
  - Identity disclosure. Avoid.
  - Attribute disclosure. A more complex case. Some attribute disclosure is expected in data mining.

At the other extreme, any improvement in our knowledge about an individual could be considered an intrusion. The latter is particularly likely to cause a problem for data mining, as the goal is to improve our knowledge. (J. Vaidya et al., 2006, p. 7.)

- Identity disclosure vs. attribute disclosure
  - Usually, identity disclosure implies attribute disclosure

Find record (HYU, Tarragona, 58), learn variable (Heart Attack)

Respondent	City	Age	Illness
ABD	Barcelona	30	Cancer
COL	Barcelona	30	Cancer
GHE	E Tarragona		AIDS
CIO	CIO Tarragona		AIDS
HYU	Tarragona	58	Heart attack

- Identity disclosure vs. attribute disclosure
  - Usually, identity disclosure implies attribute disclosure

Respondent	City	Age	Illness
ABD	Barcelona	30	Cancer
COL	Barcelona	30	Cancer
GHE	Tarragona	60	AIDS
CIO	Tarragona	60	AIDS
HYU	Tarragona	60	Heart attack

• Identity disclosure without attribute disclosure. Use all attributes (*Tarragona*, 60, *Heart attack*).

- Identity disclosure vs. attribute disclosure
  - Usually, identity disclosure implies attribute disclosure

Respondent	City	Age	Illness
ABD	Barcelona	30	Cancer
COL	Barcelona	30	Cancer
GHE	Tarragona	60	AIDS
CIO	Tarragona	60	AIDS
HYU	Tarragona	60	Heart attack

- Identity disclosure without attribute disclosure. Use all attributes (*Tarragona*, 60, *Heart attack*).
- $\circ$  Attribute disclosure without identity disclosure. k-anonymity (ABD, Barcelona, 30) not reidentified but learn Cancer

- Identity disclosure vs. attribute disclosure
  - Usually, identity disclosure implies attribute disclosure

Respondent	City	Age	Illness
ABD	Barcelona	30	Cancer
COL	Barcelona	30	Cancer
GHE	Tarragona	60	AIDS
CIO	Tarragona	60	AIDS
HYU	Tarragona	60	Heart attack

- Identity disclosure without attribute disclosure. Use all attributes (*Tarragona*, 60, *Heart attack*).
- Attribute disclosure without identity disclosure. k-anonymity (*ABD*, *Barcelona*, 30) not reidentified but learn *Cancer*
- Neither identity nor attribute disclosure

- Identity disclosure and anonymity are exclusive.
  - Identity disclosure implies non-anonymity
  - Anonymity implies no identity disclosure.

Identity	Disclosure	Attribute Disclosure
	Anonymity	
	Unlin	kability

### • So,

- Identity and attribute disclosure
- Anonymity, unlinkability
- Privacy try to provide guarantees / algorithms for avoiding disclosure

- Undetectability and unobservability
  - Undetectability of an IoI. The attacker cannot sufficiently distinguish whether IoI exists or not.
    - E.g. Intruders cannot distinguish messages from random noise
    - $\Rightarrow$  Steganography
  - Unobservability of an IoI means
    - ▷ undetectability of the lol against all subjects uninvolved in it and
    - > anonymity of the subject(s) involved in the lol even against the other subject(s) involved in that lol.

Unobservability pressumes undetectability but at the same time it also pressumes anonymity in case the items are detected by the subjects involved in the system. From this definition, it is clear that unobservability implies anonymity and undetectability.

## Terminology

### • Plausible deniability

- I have nothing to do with this database, model, etc
- $\circ\,$  Is this statement credible?
- For a database
  - at record level: This record is not mine!
  - at database level: I am not in this database!

### • Plausible deniability

- I have nothing to do with this database, model, etc
- Is this statement credible?
- For a database
  - at record level: This record is not mine!
    at database level: I am not in this database!
- We will see that some privacy models provide guarantees for plausible deniability

- Connections between plausible deniability and anonymity set
  - Plausible deniabilty: perspective of the individual
  - $\circ$  Anonymity set: perspective of the intruder

## Transparency

### • Transparency

- DB is published: give details on how data has been produced.
   Description of any data protection process and parameters
- Positive effect on data utility. Use information in data analysis.
- Negative effect on risk. Intruders use the information to attack.

### • The transparency principle in data privacy<sup>11</sup>

Given a privacy model, a masking method should be compliant with this privacy model even if everything about the method is public knowledge. (Torra, 2017, p17)

<sup>&</sup>lt;sup>11</sup>Similar to the Kerckhoffs's principle (Kerckhoffs, 1883) in cryptography: a cryptosystem should be secure even if everything about the system is public knowledge, except the key

## **Interim discussion**

- Privacy and disclosure
  - Identity disclosure appears when information about individuals make them unique: avoid identity disclosure
    - ▷ A few attributes, very distinctive attributes
      - $\Rightarrow$  highly identifiable data
    - ▷ A large number of attributes, not distinctive
      - $\Rightarrow$  high dimensional data
  - Same about "interests", multifaceted individuals
    - ▷ FC Barcelona, Umeå IK
    - ▷ Vegetarian, beef lover
    - ⊳ etc
  - Attribute disclosure: some attribute disclosure is usually expected, difficult to define/protect

## Privacy by design

- Privacy by design (Cavoukian, 2011)
  - Privacy "must ideally become an organization's default mode of operation" (Cavoukian, 2011) and thus, not something to be considered a posteriori. In this way, privacy requirements need to be specified, and then software and systems need to be engineered from the beginning taking these requirements into account.
  - In the context of developing IT systems, this implies that privacy protection is a system requirement that must be treated like any other functional requirement. In particular, privacy protection (together with all other requirements) will determine the design and implementation of the system (Hoepman, 2014)

- Privacy by design principles (Cavoukian, 2011)
  - 1. Proactive not reactive; Preventative not remedial.
  - 2. Privacy as the default setting.
  - 3. Privacy embedded into design.
  - 4. Full functionality positive-sum, not zero-sum.
  - 5. End-to-end security full lifecycle protection.
  - 6. Visibility and transparency keep it open.
  - 7. Respect for user privacy keep it user-centric.

- Data minimization: a cornerstone for privacy
  - The principle of "data minimisation" means that a data controller should limit the collection of personal information to what is directly relevant and necessary to accomplish a specified purpose. They should also retain the data only for as long as is necessary to fulfil that purpose. In other words, data controllers should collect only the personal data they really need, and should keep it only for as long as they need it.

The data minimisation principle is expressed in Article 5(1)(c) of the GDPR and Article 4(1)(c) of Regulation (EU) 2018/1725, which provide that personal data must be "adequate, relevant and limited to what is necessary in relation to the purposes for which they are processed".<sup>12</sup>

<sup>&</sup>lt;sup>12</sup>https://edps.europa.eu/data-protection/data-protection/glossary/d\_en

# Summary

#### • Concepts

- What is data privacy?
- Multidisciplinary problem and socio-technical aspects to be considered
- Difficulties of data privacy: naive annonymization does not work
- Linkability and anonymity set
- Identity and attribute disclosure
- Plausible deniability
- Transparency
- Privacy by design
- Data minimization

## References

- Torra, V. (2022) Guide to data privacy, Springer.
- Cavoukian, A. (2011) Privacy by design. The 7 foundational principles in Privacy by Design. Strong privacy protection now, and well into the future.
- Narayanan, A., Shmatikov, V. (2008) Robust De-anonymization of Large Sparse Datasets, Proc. of the 2008 IEEE Symposium on Security and Privacy (SP '08), 111-125.
- Sweeney, L. (1997) Simple Demographics Often Identify People Uniquely, Carnegie Mellon University, Data Privacy Working Paper 3. Pittsburgh 2000.
- Dalenius, T. (1977) Towards a methodology for statistical disclosure control, Statistisk Tidskrift 5 429-444.
- D'Acquisto, G., Domingo-Ferrer, J., Kikiras, P., Torra, V., de Montjoye, Y.-A., Bourka, A. (2015) Privacy by design in big data: An overview of privacy enhancing technologies in the era of big data analytics, ENISA Report.
- V. Torra, G. Navarro-Arribas (2016) Big Data Privacy and Anonymization, Privacy and Identity Management 15-26. https://doi.org/10.1007/978-3-319-55783-0\_2 (open access)