

# Data privacy: Motivation

Vicenç Torra

February, 2018

SAIL + PICS, School of Informatics, University of Skövde, Sweden

# Outline

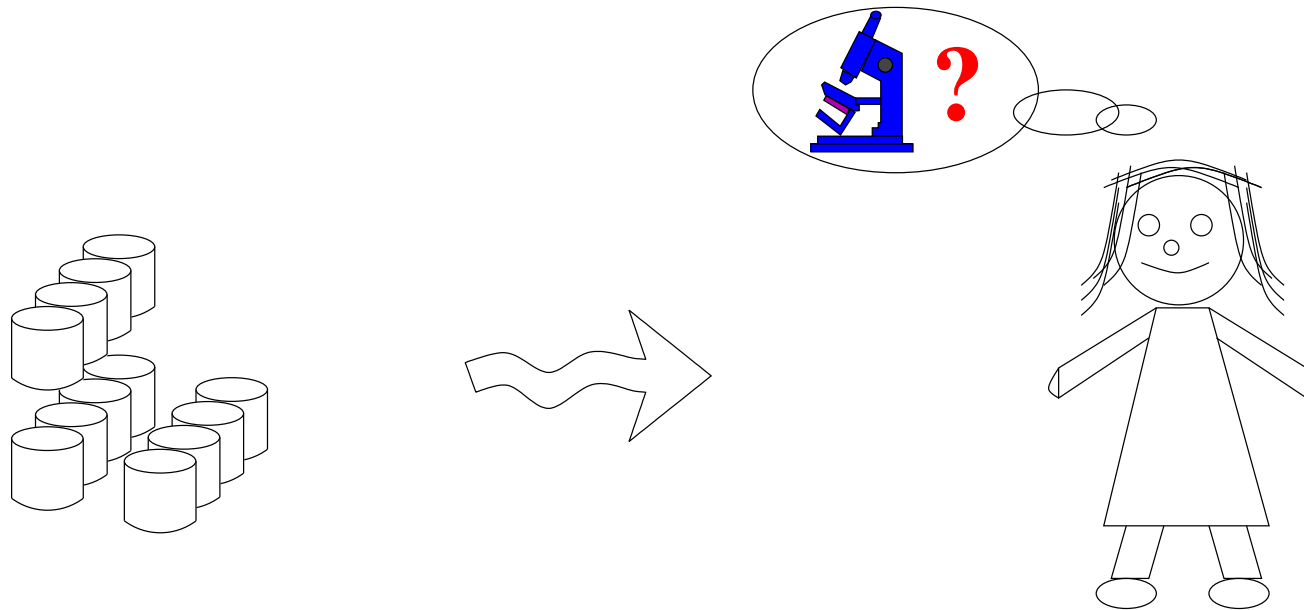
---

## 1. Motivation

# Motivation

# Introduction

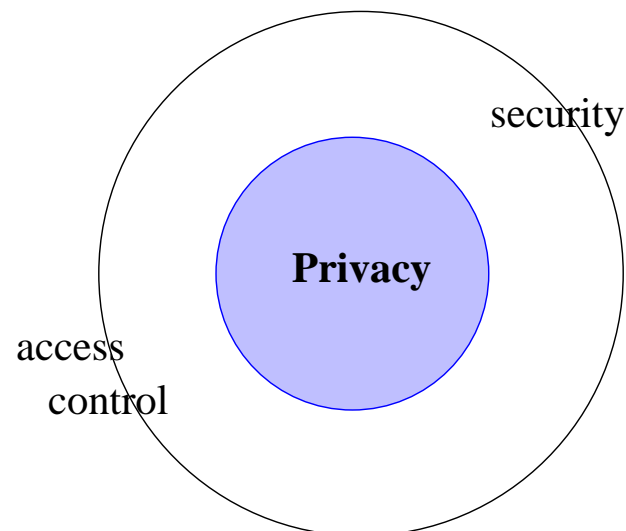
- Data privacy: core
  - Someone needs to access to data to perform **authorized analysis**, but **access to the data** and the **result of the analysis** should avoid **disclosure**.



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.

# Introduction

- Data privacy: boundaries
  - Database in a computer or in a removable device
    - ⇒ access control to avoid unauthorized access
      - ⇒⇒ Access to address (admissions), Access to blood test (admissions?)
  - Data is transmitted
    - ⇒ security technology to avoid unauthorized access
      - ⇒⇒ Data from blood glucose meter sent to hospital. Network sniffers
      - Transmission is sensitive: Near miss/hit report to car manufacturers



# Motivation

---

- Legislation. 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 (EU). National regulations.
- Companies own interest.
  - Competitors can take advantage of information.
- Avoiding privacy breach. Several well known cases (see later).

# Difficulties

- Difficulties: Naive anonymization **does not work**

Passenger manifest for the Missouri, arriving February 15, 1882; Port of Boston<sup>1</sup>

Names, Age, Sex, Occupation, Place of birth, Last place of residence, Yes/No, condition (healthy?)

JOHN WARD & CO.,  
GROG MERCHANTS,  
47E ST. BOSTON.

#61

**LIST OF PASSENGERS.**

REPORT AND LIST OF PASSENGERS taken on board the *S. S. Hooper* of *London*,  
whom *Fredrick Murrell* is Master, berthing from the Port of *London* to *Boston*.

I, *Fredrick Murrell* Master of the *S. S. Hooper* from *London* do solemnly swear that the Report herewith made, in conformity with the Laws of the Commonwealth of Massachusetts, relating to these Passengers, is true and correct, to the best of my knowledge and belief. As sworn to at *Boston*, this *16* day of *April* 1882.

Before me, *Amos Hildreth* Justice of the Peace.

*F. Murrell* Master.

	NAME	AGE	SEX	OCCUPATION	PLACE OF BIRTH	Last Place of Residence	If in American Service		CONDITION
							Yes	No	
1	<i>George W. Smith</i>	18	Male	<i>Bookbinder</i>	<i>London, Eng.</i>	<i>London, Eng.</i>			<i>Healthy</i>
2	<i>George W. Smith</i>	24	-	<i>Seaman</i>	<i>Madras India</i>	<i>London, Eng.</i>			
3	<i>George W. Smith</i>	21	-	<i>Tailor</i>	<i>London</i>	<i>London</i>			
4	<i>George W. Smith</i>	21	-	<i>Seaman</i>	<i>London</i>	<i>London</i>			
5	<i>George W. Smith</i>	25	-	<i>Bookbinder</i>	<i>Boston, U.S.</i>	<i>Boston, U.S.</i>			
6	<i>George W. Smith</i>	18	-	-	<i>Island</i>	<i>Boston, U.S.</i>			
7	<i>George W. Smith</i>	16	-	-	<i>London (I)</i>	<i>London</i>			
8	<i>George W. Smith</i>	42	-	-	<i>Cross, Eng.</i>	<i>Cross, Eng.</i>			
9	<i>George W. Smith</i>	28	-	-	<i>Albany Mass</i>	<i>Boston, U.S.</i>			
10	<i>George W. Smith</i>	25	-	<i>Bookbinder</i>	<i>Boston, U.S.</i>	<i>Boston, U.S.</i>			

73

<sup>1</sup><https://www.sec.state.ma.us/arc/gen/genidx.htm>

# Difficulties

---

- 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

---

- Difficulties: **highly identifiable data**
  - Data from mobile devices:
    - ★ two positions can **make you unique** (home and working place)
  - AOL<sup>2</sup> and Netflix cases (search logs and movie ratings)
    - ⇒ User No. 4417749, hundreds of searches over a three-month period including queries 'landscapers in Lilburn, Ga' ⇒ Thelma Arnold identified!
    - ⇒ individual users matched with film ratings on the Internet Movie Database.
  - Similar with credit card payments, shopping carts, ...  
(i.e., **high dimensional data**)

---

<sup>2</sup><http://www.nytimes.com/2006/08/09/technology/09aol.html>

# Difficulties

---

- Difficulties: highly identifiable data
  - Example #1:
    - ★ University goal: know how sickness is influenced by studies and by commuting distance
    - ★ Data: where students live, what they study, if they got sick
    - ★ No “personal data”, is this ok ?

# Difficulties

---

- Difficulties: highly identifiable data
  - Example #1:
    - ★ University goal: know how sickness is influenced by studies and by commuting distance
    - ★ Data: where students live, what they study, if they got sick
    - ★ No “personal data”, is this ok ?
    - ★ **NO!!!**: How many in your degree live in your town ?

# Difficulties

---

- Difficulties: highly identifiable data
  - Example #1:
    - ★ University goal: know how sickness is influenced by studies and by commuting distance
    - ★ Data: where students live, what they study, if they got sick
    - ★ No “personal data”, is this ok ?
    - ★ **NO!!!**: How many in your degree live in your town ?
  - Example #2:
    - ★ Car company goal: Study driving behaviour in the morning
    - ★ Data: First drive (GPS origin + destination, time) × 30 days
    - ★ No “personal data”, is this ok?

# Difficulties

---

- Difficulties: highly identifiable data
  - Example #1:
    - ★ University goal: know how sickness is influenced by studies and by commuting distance
    - ★ Data: where students live, what they study, if they got sick
    - ★ No “personal data”, is this ok ?
    - ★ **NO!!!**: How many in your degree live in your town ?
  - Example #2:
    - ★ Car company goal: Study driving behaviour in the morning
    - ★ Data: First drive (GPS origin + destination, time) × 30 days
    - ★ No “personal data”, is this ok?
    - ★ **NO!!!!**: How many (cars) go from your parking to your university every morning ? Are you exceeding the speed limit ? Are you visiting a psychiatrist every tuesday ?

# Difficulties

---

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