



HIROSHIMA UNIVERSITY

Fundamental Data Science (30104001)

Lecture 2 — Data acquisition and open data.
Data science ethics

Jorge N. Tendeiro

Hiroshima University

Today

- Data acquisition and open data.
- Data science ethics.

Data acquisition and open data

Data

Data is the foundation of data science:

Data is the material on which all the analyses are based.

(*Source*)

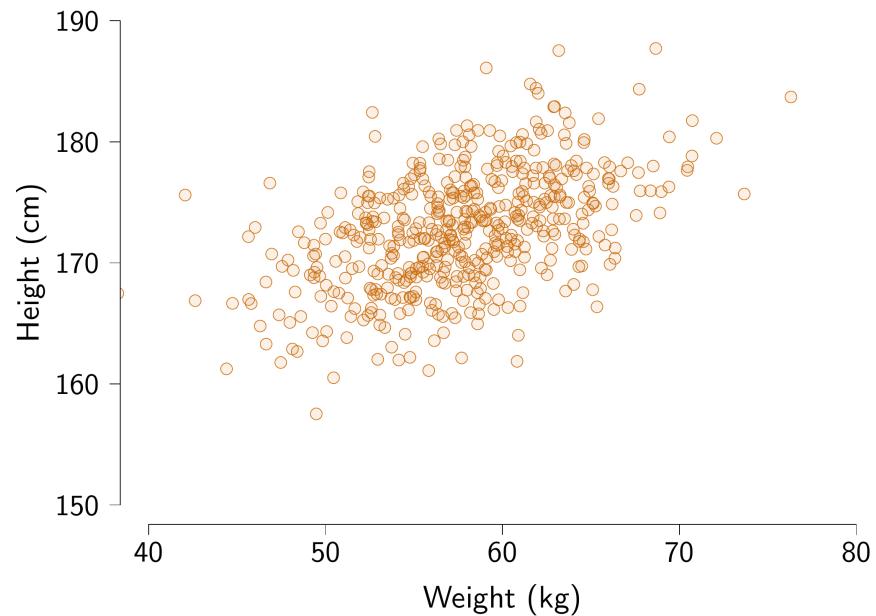
Data are supposed to capture relevant **features** of what we are studying.

Some features are easier to observe and measure (e.g., age), whereas others are quite difficult (e.g., emotions).

Data typically must be **digitalized** so that they can be analyzed with the help of computers.

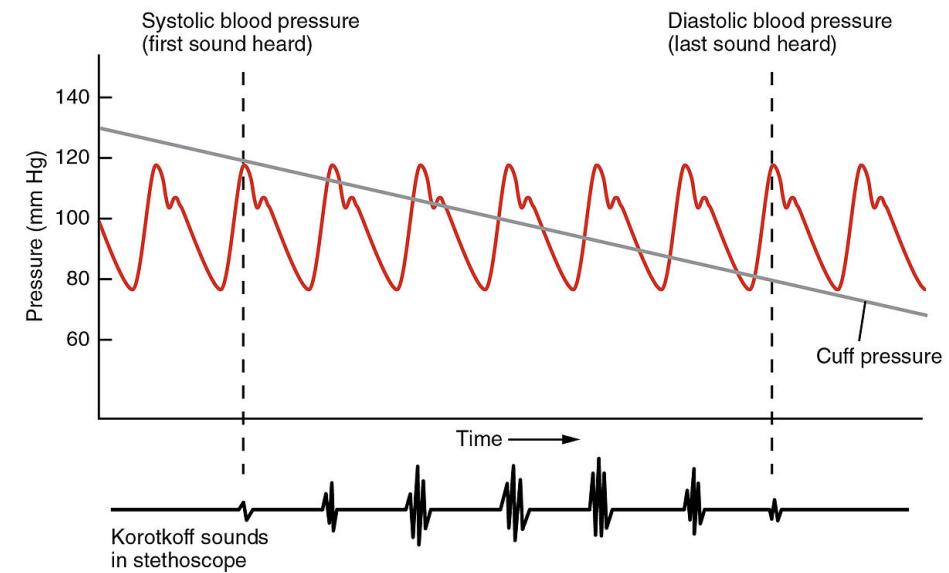
Examples of data

Height, weight:



Source: Kaggle.

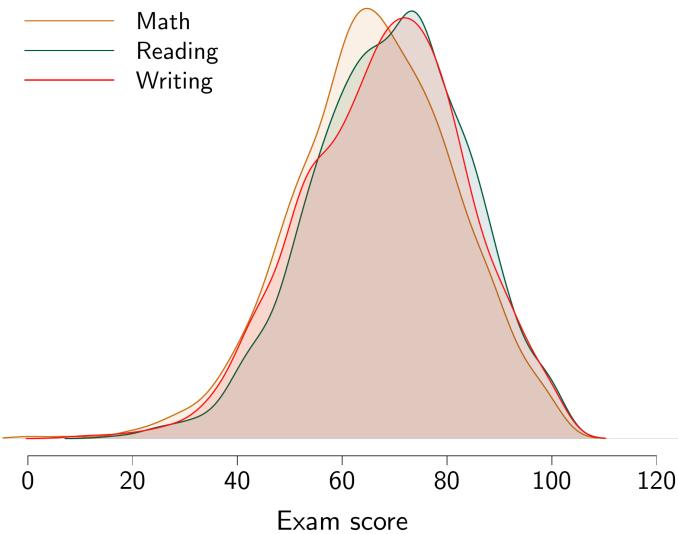
Blood pressure:



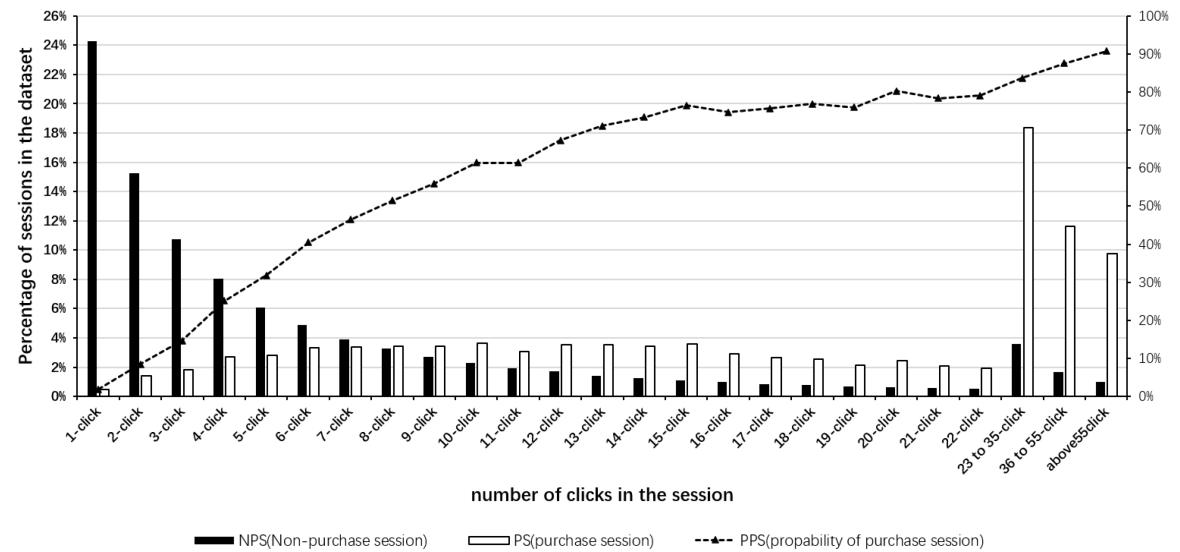
Source: Wikipedia.

Examples of data

Exam scores:



Marketing:

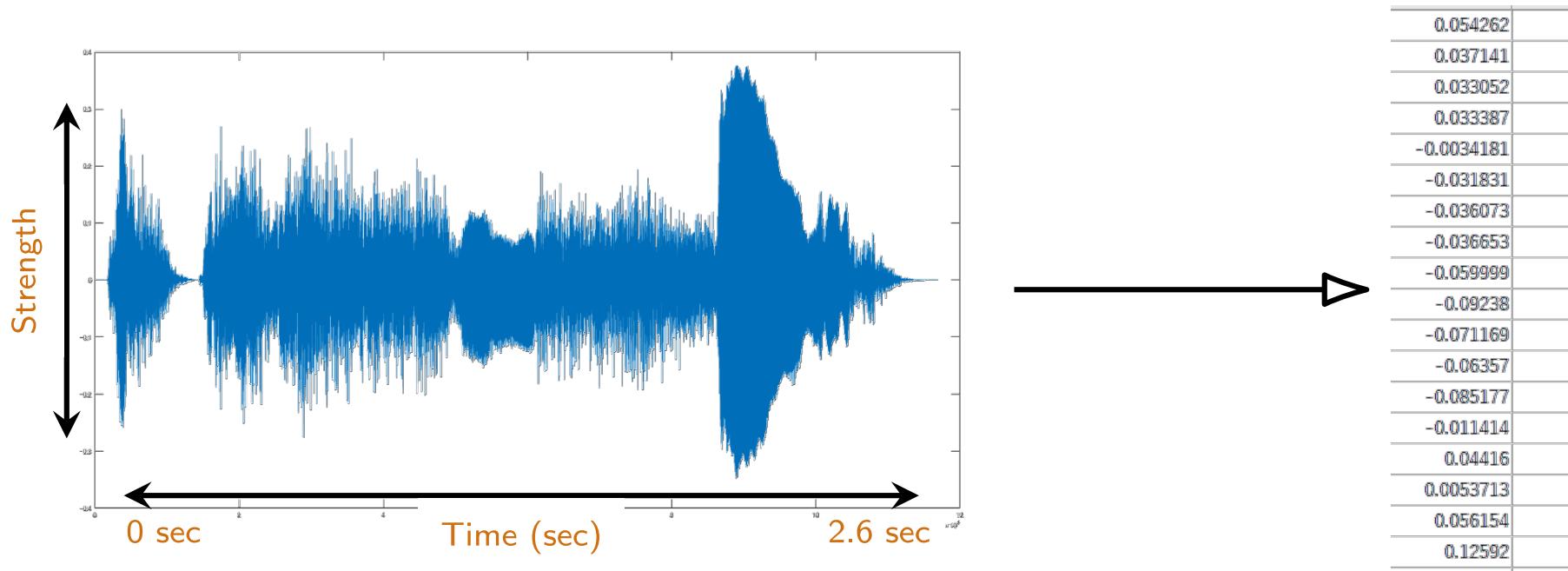


Source: Kaggle.

Source: Wen et al. (2023).

Examples of data (audio)

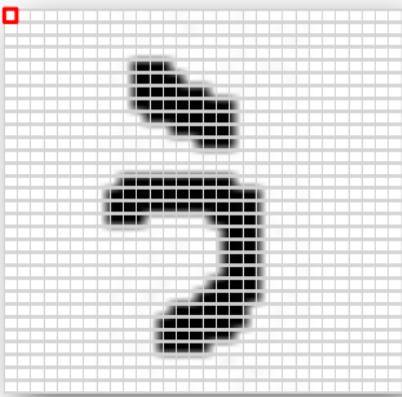
Audio data expressed as a wave:



Numeric data of strength (about 120,000 numeric values in 2.6 sec).

Examples of data (image)

Hand-written letter:



Split an image into pixels (small squares), and quantify the shading of the color in each pixel.

Numerical data (0: black ~ 255: white):

Obtaining data

There are various methods available to collect data.

- Collect the data by **yourself**.

E.g.: Distribute questionnaires, perform experiment/observation.

- ✓ You collect the data that you want.
- ✗ It can take a lot of time and effort.

- Use data collected by **someone else**.

E.g.: Request materials, download publicly available data online.

- ✓ Easier to obtain.
- ✗ The data is not always available or shared.
- ✗ The data need not exactly match what you want.

Open data

Open data is data that is free to access, use, modify, and share by anyone – subject, at most, to measures that preserve provenance and openness.

Open Definition

Requirements that open data must meet ([source](#)):

- ***Open license***
Data must be available in the public domain or provided under an open license allowing secondary use.
- ***Access free of charge***
Data must be provided online and free of charge
(or at most at no more than a reasonable one-time reproduction cost).
- ***Machine readability***
Data must be provided in a form readily processable by a computer (see next page for details).
- ***Open format***
Data must be provided in an open format, placing no restrictions, monetary or otherwise, upon its use.
Data can be fully processed with at least one open-source software tool.

Machine readability

Machine readability: Ease with which files can be imported or edited by computer programs.



In data analysis, tabular files (e.g., XLS, TXT, CSV) are easy to work with.
CSV files are often used in R.

Important

Even with CSV files, you need to first clean data before analyzing them!

It is a part of data science task to convert data into a format which is machine readable.

Benefits to open data

1. Solving various issues and revitalizing the economy through the promotion of public participation and public-private collaboration.

→ *Creation of new businesses and services.*

2. Sophistication and efficiency of administration.

→ *Enhancement of public services.*

3. Improve transparency and reliability.

Open data — Also good to keep in mind

- **Reproducibility**

Results can be reinspected after publication. Mistakes are easier to find.

- **Accountability**

We do not need to "take someone's word for it" (i.e., to believe without proof).

- **Trust**

Open data is a means to show that our work is trustworthy and honest.

- **Contribution**

Open data is a way to give back to the society.

- **Preserve the scientific record**

Open data is a means to not lose valuable information.

- **Save money**

Open data facilitate future data collections, minimizing waste.

Sources: [here](#), [here](#).

Open data — Examples

There is a boom of projects which rely on access to open data.

Below are a couple of interesting sources to read about this:

- <https://citizens-guide-open-data.github.io/guide/3-open-datas-use>.
- <https://www.luzmo.com/blog/examples-of-open-data-in-government>.

I will highlight three examples.

Open data — Example 1: Public health and environment

CurieuzeNeuzen is a citizen participation project in Belgium.

Goal

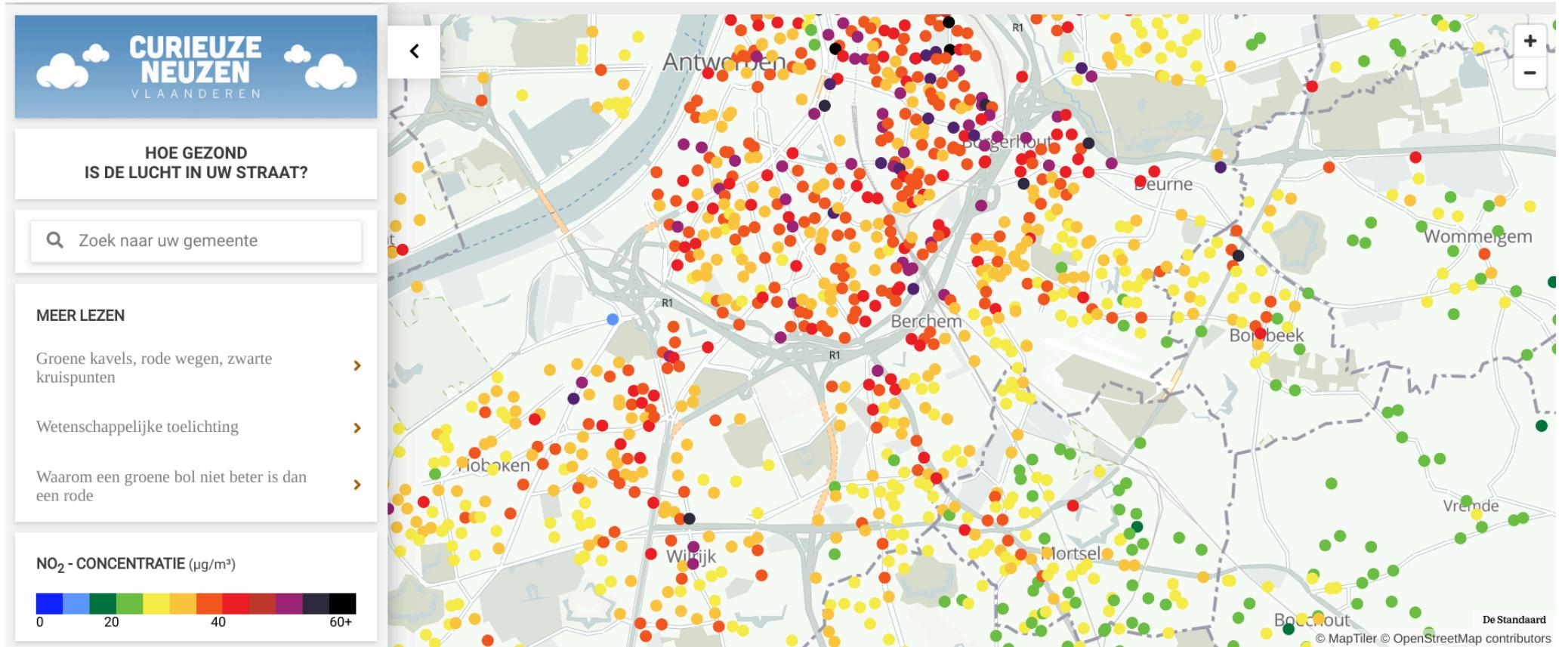
Participate in a large-scale research of air quality in Flanders.

Outcome

They developed an interactive app that monitors the air quality in Flanders.

Data are contributed by each citizen!

Open data – Example 1: Public health and environment



Open data — Example 2: Helping the police

See the full blog post here:

[Parking Data, Ben Wellington.](#)

Ben Wellington used the parking ticket data set from the New York Police Department's [open data portal](#) to identify a problem:

*| The police had been giving parking tickets for cars that were **legally** parked.*

Mr. Wellington used the [open parking ticket data](#) to prove that the police officers were not following a new traffic rule.

This saved the community from paying millions of dollars of unfair fines!

Open data – Example 2: Helping the police

Follow iquantny

I Quant NY

[MAILING LIST](#) [RSS](#) [ARCHIVE](#)

Quantitative Analysis of NYC Open Data: Every data set that the city releases tells a story. This blog is all about telling those stories, one data set at a time.

[About Me](#) [About You](#) [Interviews](#) [Press](#) [Topics](#) [Subscribe](#)

MAY 11, 2016

The NYPD Was Systematically Ticketing Legally Parked Cars for Millions of Dollars a Year- Open Data Just Put an End to It

New York City is a complex place to drive. And when it comes to parking, there are plenty of rules and regulations to follow. It's no wonder that sometimes people get confused and end up getting their cars ticketed or towed.

But in all of these rules, there is one thing that very few drivers seem to know. As of late 2008, **in NYC you can park in front of a sidewalk pedestrian ramp, as long as it's not connected to a crosswalk**. It's all written up in the NYC Traffic [Rules](#), and for more detail, take a look at [this](#) article. The local legislation making these parking spots legal was proposed by Council Member Gentile, and adopted by the Department of Transportation before it ever made it for a vote. Though few people seem to know about the change.

Open data – Example 3: Spotify and MusicBrainz



Spotify is a popular digital music, podcast, and video streaming service.



Spotify relies on **MusicBrainz** to collect music metadata (like the name of the artists, album, year, music titles, etc.).

MusicBrainz is an **open music encyclopedia** that makes all of its data available to the public.

Also **you** can use MusicBrainz, for example, when setting up your own home media center (using for instance **Jellyfin**, **Kodi**, or **Plex**)!

Where to get open data?

There are many repositories.

In Japan:

- Websites to obtain open data:



<https://www.e-stat.go.jp>



<https://www.data.go.jp>



<https://data.e-gov.go.jp/info/ja>

- Example of analyzing open data:

◦ Open data 100

(Government CIO's portal Japan).

<https://cio.go.jp/opendata100>

In the World:

- Websites to obtain open data:



<https://abcnews.go.com/538>



<https://www.kaggle.com/datasets>



Open Science Framework

<https://osf.io/>

◦ ...

Data science ethics

Collecting data



I want to broadcast about the approval rate of the governor of Hiroshima among Hiroshima citizens during tomorrow's news.

Let's ask people at Hiroshima station. We tell upfront that the survey is to broadcast the approval rating of the governor at my news.

(After broadcast) To prove that I didn't fake it, let's publish the collected data on our website!

ID	Name	Sex	Approve (0=no, 1=yes)
1	Mike Brown	male	1
2	Sue Allen	female	0
3	John Watson	male	1
4	Tom Hanks	male	1
5	Norah Jones	female	0
6	Scarlett Johansson	female	0
7	Anne Hathaway	female	0
8	Meryl Streep	female	1

Looking at the published data, do you see any problems?

There are **many problems!!!**

Ethics in data collection and utilization

Ethical considerations in data collection and utilization:

- It is important to obtain **informed consent** after explaining the purpose of data collection and how the data will be utilized.

Informed consent (medical):

A process in which patients are given important information, including possible risks and benefits, about a medical procedure, so that patients can decide if they want to take part in the trial.

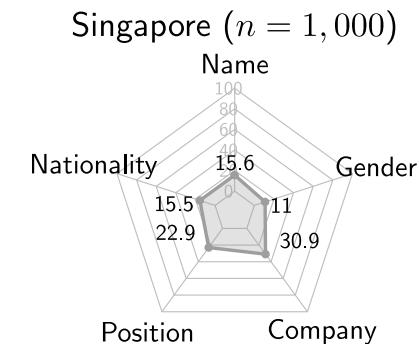
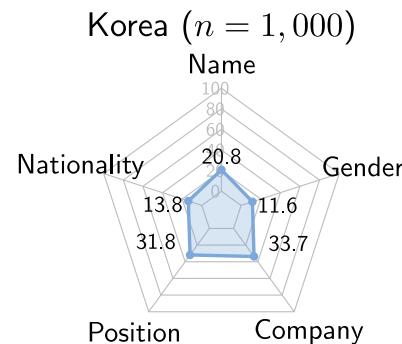
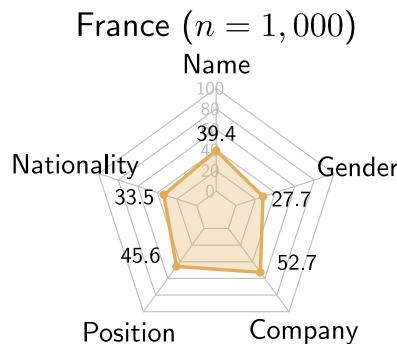
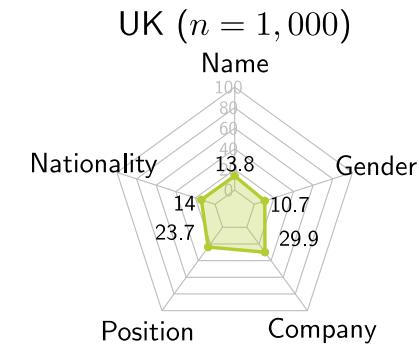
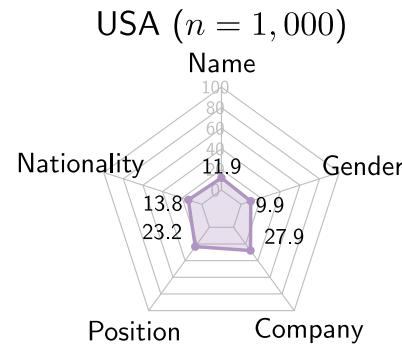
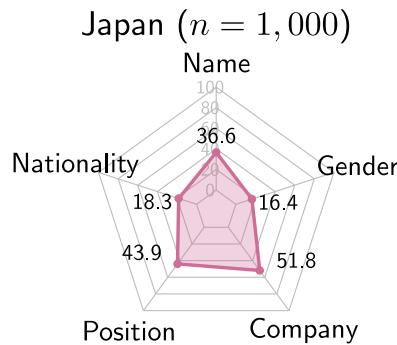
(*National Cancer Institute, NIH*)

- Build a system that allows participants to withdraw their data easily, if they are willing to.
- Note that each participant has a **different background** and views on which data to be kept personal.
→ When collecting data, we need to consider the *social and cultural background* of the data subjects.

Also, keep in mind that the information that people don't want to disclose to others varies across regions or countries.

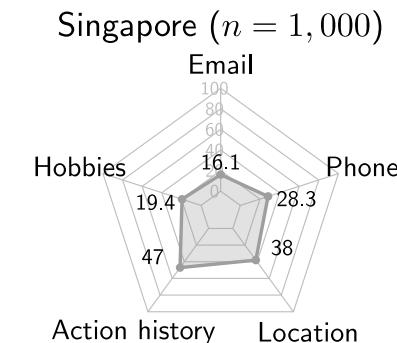
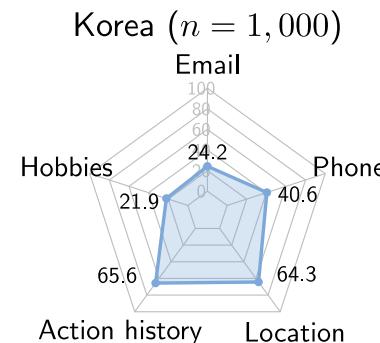
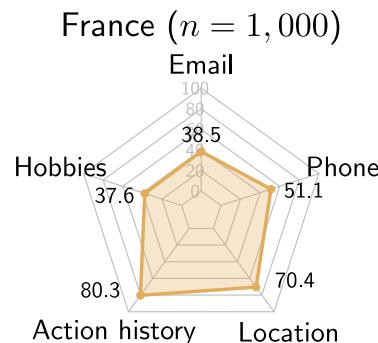
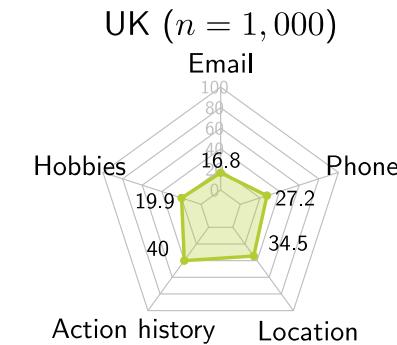
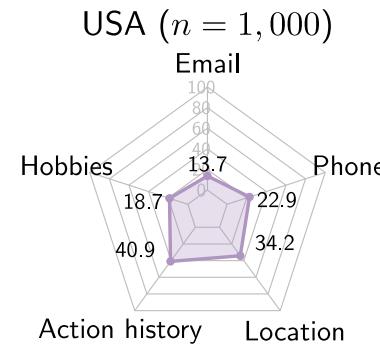
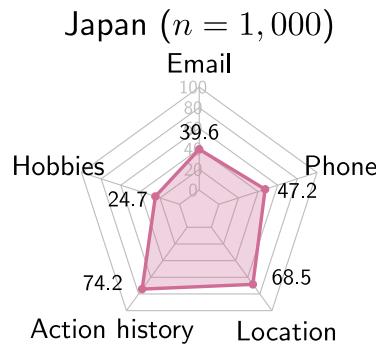
Survey: "Data I don't want to provide or disclose under any circumstances"

General personal data (not highly sensitive personal data; $n = 1,000$):



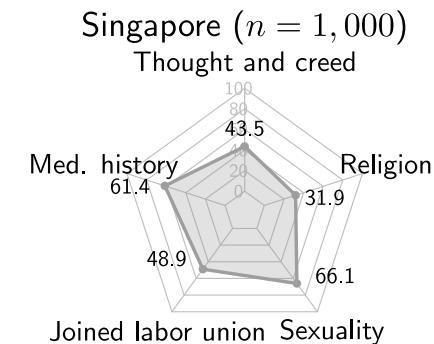
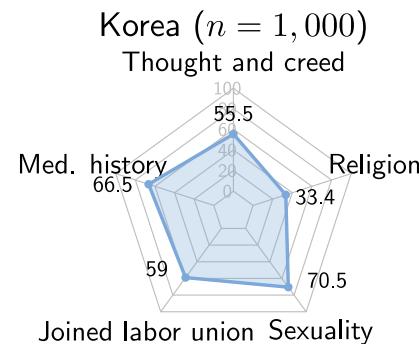
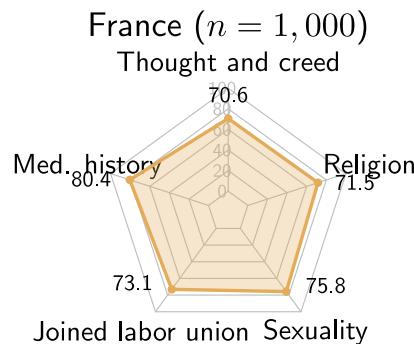
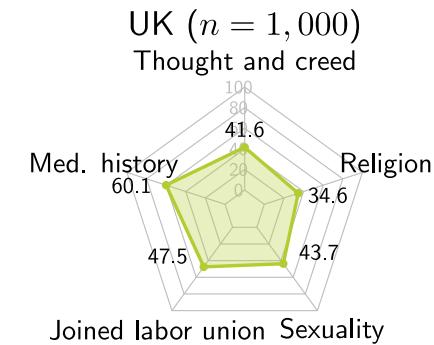
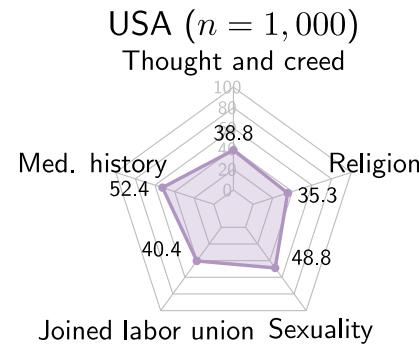
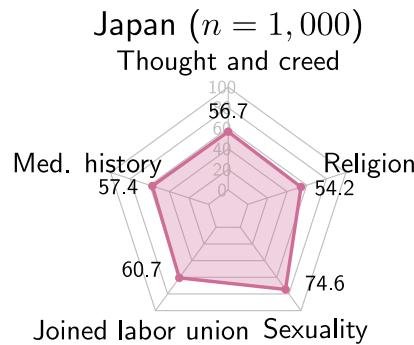
Survey: "Data I don't want to provide or disclose under any circumstances"

Sensitive personal data (data requiring careful handling; $n = 1,000$):



Survey: "Data I don't want to provide or disclose under any circumstances"

Highly sensitive personal data ($n = 1,000$):



Anonymization and pseudonymization of data

Anonymizing data:

- | *Making it impossible to identify an individual.*

Example:

Delete personal identifiable elements (name, address, etc.).

Pseudonymizing data:

Make sure that an individual is not identifiable *without additional information*.

Example:

Encrypt person-identifiable information (e.g., name) by ID.

Anonymization and pseudonymization of data

The screenshot shows a web page titled "Japanese Law Translation". On the left, there is a sidebar with a table of contents for the Act on the Protection of Personal Information (Partly unenforced). The sidebar includes sections such as Article 95 (Exceptions to the Due Date for Decisions on Correction), Article 96 (Transfer of Cases), Article 97 (Notification to A Party to which Personal Information an Administrative Entity Holds is Provided), Subsection 3 Ceasing to Use Personal Information an Administrative Entity Holds, Article 98 (Right to Request Ceasing to Use Personal Information), and Article 99. The main content area displays the title "Act on the Protection of Personal Information (Partly unenforced)" and the date "Act No. 57 of May 30, 2003". Below the title is a "Table of Contents" section listing various chapters and sections of the law.

Act on the Protection of Personal Information (Partly unenforced)

Act No. 57 of May 30, 2003

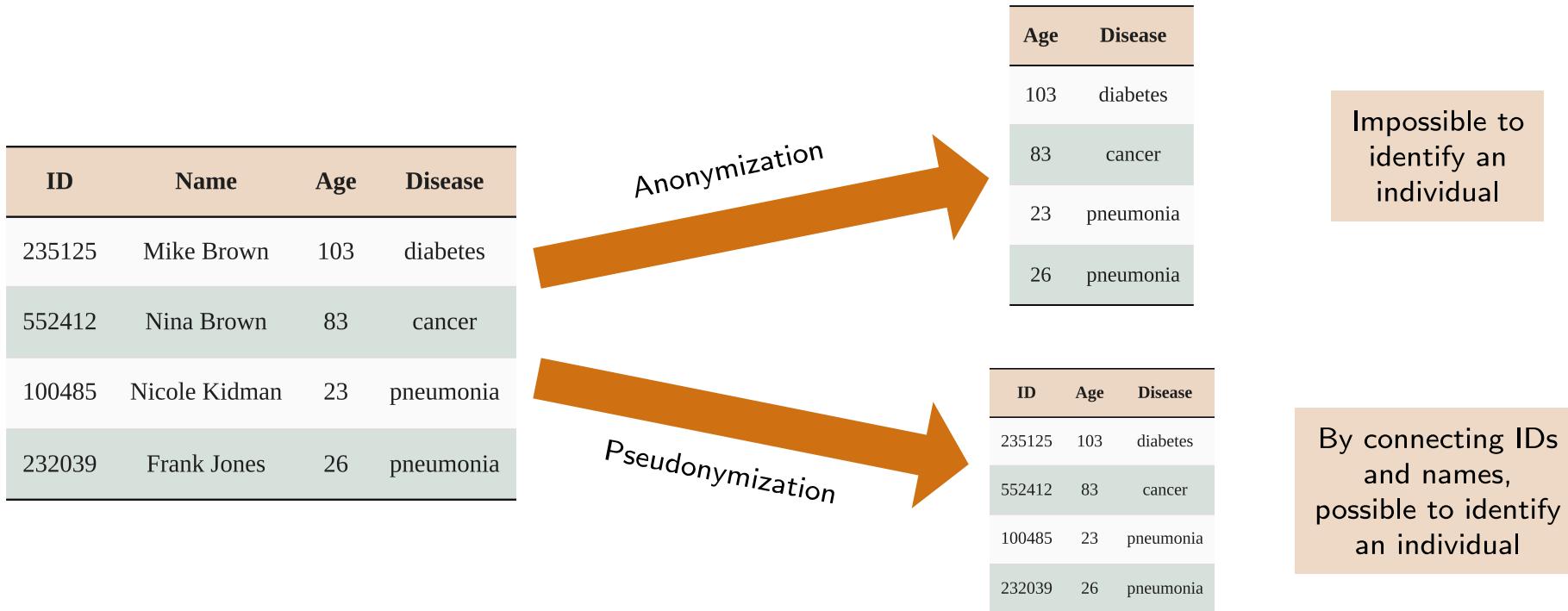
Table of Contents

Chapter I General Provisions (Articles 1 to 3)
Chapter II Responsibilities of the National and Local Governments (Articles 4 to 6)
Chapter III Measures to Protect Personal Information
Section 1 Basic Policy on the Protection of Personal Information (Article 7)
Section 2 Measures Taken by the National Government (Articles 8 to 11)
Section 3 Measures Taken by the Local Government (Articles 12 to 14)
Section 4 Cooperation between the National and the Local Governments (Article 15)
Chapter IV Obligations of Businesses Handling Personal Information
Section 1 General Provisions (Article 16)
Section 2 Obligations of Businesses Handling Personal Information and Businesses Handling Information Related to Personal Information (Articles 17 to 40)

From the Act on the Protection of Personal Information (APPI):

- **Anonymized data:**
In principle, it *can be provided* to a third party without consent.
- **Pseudonymized data:**
In principle, it is *impossible to provide* to a third party without consent.

Anonymization and pseudonymization of data: Example



Anonymization and pseudonymization of data: Example

Even if you think that you have anonymized or pseudonymized data sufficiently, it may still be possible to identify an individual.

Example:

When **only one** person has age over 100 in the data.

The diagram illustrates a transformation process. On the left, a table shows four individuals with their ages and diseases. Mike is the only one with an age over 100. A large orange arrow points from this table to the right, labeled "correct". On the right, the same four individuals are shown again, but their ages have been pseudonymized: Mike's age is now "Over 80", while Nina, Nicole, and Frank have ages in the range of 20-30. This transformation is labeled "Impossible to identify Mike" at the bottom right.

	Age	Disease
Mike	103	diabetes
Nina	83	cancer
Nicole	23	pneumonia
Frank	26	pneumonia

correct

	Age	Disease
Mike	Over 80	diabetes
Nina	Over 80	cancer
Nicole	20~30	pneumonia
Frank	20~30	pneumonia

Possible to identify Mike Impossible to identify Mike

| **Always** pay close attention to personal identification!!

So, do you now see any problems?



I want to broadcast about the approval rate of the governor of Hiroshima among Hiroshima citizens during tomorrow's news.

Let's ask people at Hiroshima station. We tell upfront that the survey is to broadcast the approval rating of the governor at my news.

(After broadcast) To prove that I didn't fake it, let's publish the collected data on our website!

ID	Name	Sex	Approve (0=no, 1=yes)
1	Mike Brown	male	1
2	Sue Allen	female	0
3	John Watson	male	1
4	Tom Hanks	male	1
5	Norah Jones	female	0
6	Scarlett Johansson	female	0
7	Anne Hathaway	female	0
8	Meryl Streep	female	1

Problems:

- **Individuals are clearly identifiable.**
Anonymization/pseudonymization is needed.
- **Insufficient informed consent.**
Persons must be carefully informed about, and informed consent must be obtained afterwards.

Other sorts of problems



I want to broadcast about the approval rate of the governor of Hiroshima among Hiroshima citizens during tomorrow's news.

Let's ask people at Hiroshima station.
We tell upfront that the survey is to broadcast the approval rating of the governor at my news.

Although you want to know about the approval rate of **all Hiroshima citizens**, data were only collected at Hiroshima station.

The estimated approval rate is **biased**:

*The results do not reflect **all** residents of Hiroshima.*

Selection bias

Selection bias:

*Bias caused by an **inappropriate** way of collecting data.
As a result, the collected data do not reflect the population at large.*

Example:

Study the *prevalence rate* of a certain disease among Japanese people.

Data collection plan:

Collect data from patients who visited hospitals.



Q: What is the problem?

A: Many people who go to the hospital are already sick.

The prevalence rate of such a sample will be **biased** upwards (inflated, that is, larger than it actually is in the population).

Information bias

Information bias:

*Problems related to confusing information or improper measurement methods.
As a result, the observed measurements may be biased.*

Example:

Ask people "Did you catch a cold last month?"

✗ Only people who recognized their symptoms as a "cold" would answer "yes" to this question.

Example:

Ask people "How much savings do you have?"

✗ People may answer more than their actual savings (e.g., to conceal financial problems).

Confounding bias

Confounding bias:

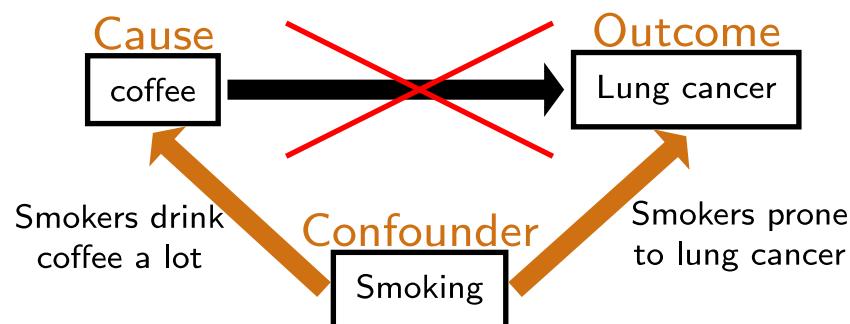
*Bias caused by factors (**confounders**) affecting both a 'cause' and an 'outcome'. This could lead to a **wrong** research outcome.*

Example:

Does drinking coffee **increase** the risk of lung cancer?

✗ Well, not necessarily.

Here is an alternative *plausible* explanation:



The relation between "coffee" and "lung cancer" no longer seems causal once "smoking" is taken into account:

Smoking is a **confounder**.

Recap: Which bias is taking place here?



I want to broadcast about the approval rate of the governor of Hiroshima among Hiroshima citizens during tomorrow's news.

Let's ask people at Hiroshima station. We tell upfront that the survey is to broadcast the approval rating of the governor at my news.

Although you want to know about the approval rate of **all Hiroshima citizens**, data were only collected at Hiroshima station.

The estimated approval rate is **biased**:

*The results do not reflect **all** residents of Hiroshima.*

Selection bias:

It is important to sample more randomly from the entire Hiroshima Prefecture.

There are **many** other biases. To reduce biases, we need to think **carefully** about how to collect data.

Summary

Data acquisition and open data:

- Data is **crucial** for data science!
- Open data is convenient, responsible, and is easy to use.

Data science ethics:

When collecting data by yourself and utilizing them, be aware of the following:

- **Ethical** considerations.
- **Anonymization** and **pseudonymization** of data.
- Various **biases**.