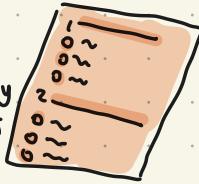


Week 1: Data Exploration

Data Formats & Data Structures

Data can be collected by

1. Interviews
2. Observations → mostly used by scientists
3. forms
4. Questionnaires
5. Surveys
6. Cookies → track people online's activities and interests



in owner's name

First-party data

↳ data collected by an individual or group using their own resources.

in a third party's name

Second-party data

↳ data collected by a group directly from its audience and then sold.

Third-party data (not reliable)

↳ outsources not directly collected doesn't have a direct relationship with data.

Data Collection Considerations

③ How the data will be collected?

choose data sources

④ Decide what data to use

How much data to collect?

Select the right data type

Determine the time frame

Population → entire sample → some groups

Need answer immediately → Historical Data

Discrete → counted & limited number
e.g. price, rating (no decimal accepted version)
have decimal limited

Continuous → unspecified number of possible points
3.41.....

Nominal → categorized without a set order (no sequence)
e.g. Yes/No

1, 2, 3, 4, 5
Ordinal → within an order
e.g. 1, 2, 3, 4, 5, stars

Must know!

lives inside a company's own systems
Internal

outside company
External

collected by a researcher from first-hand sources

Primary

Secondary
by other people

* subjective and explanatory
measures of qualities and characteristics
e.g. - Exercise activity most enjoyed

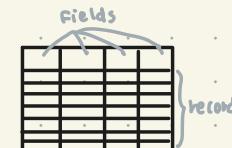
* specific and objective measures of numerical facts
e.g. - Percentage of alcohol in beverages

0101
110010
00110
101011

almost numerical value
Continuous
Weight, Height 52.4
- Temperature 25.0
- Runtime in video 18:34

Discrete
limited number of values
e.g. - Number who visit a hospital in Dec 2022
- Room's max capacity
- Tickets sold 50, 79

Structured
in a certain format (rows, columns)
records
e.g. - Expense reports
- Tax returns



records

fields

Nominal
not categorized with a set order

Ordinal → categorized with a set order/scale

Quantitative

Qualitative

Qualitative
Quality

Venue

Unstructured

e.g. - Social media

- Videos

- Emails

Conceptual

e.g. - Movie ratings (1, 2, 3 stars)

Logical

- Income level (low, mid, high)

Physical

- Movie ratings (1, 2, 3 stars)

Data Types → a specific kind of data attribute that tells what kind of value data is.

In Number (1, 3.14, 3.50, 11)

→ not used in calculation

Spread Text / String ("Chan")

Sheet Boolean (True/False)

Boolean logic

AND

grey shoes

pink shoes

Both terms

OR

grey shoes

pink shoes

Either term

NOT

grey shoes

pink shoes

Just one term

ID	frame	lane
1	Ariya	PP
2	Donald	Trump

Wide Data format

every data subject has a single row with multiple columns to hold the values

Long Data format

each row is one time point per subject, each subject will have data in multiple rows

country	countryname	Year
ARG	Argentina	2010
ARG	Argentina	2012
ARG	Argentina	2011

what to use depends on your work

can be converted called "Data Transformation"

process of changing data's format, structure, or values

Data collection considerations



1. Select the right type of data



2. Determine the time frame



3. Decide how data will be collected



4. Decide how much data to collect



or



3. Choose data sources



4. Decide what data to use

1

Week 2: Data Bias (4)
Ensuring Data Integrity < Good data (😊) Bad data (😢)

- Bias → a preference in favor of or against a person, group of people / things.



2

Types of Data Bias

- **Data Bias** → type of error that systematically skews results.
 - **Sampling Bias** → a sample isn't representative of the population (small group sampling)
Prevent:
 - Randomly sample it with fairness e.g. not all random is all female/male
 - Visualization easy to recognize it's
 - **Observer / Experimenter / Research Bias** → Different peoples observe things differently.
R^A / R^B
 - **Interpretation Bias** → Interpret ambiguous things differently (they might have diff backgrounds)
 - **Confirmation Bias** → search for interpret info in a way that confirms pre-existing beliefs.

All affects way we make sense of data

• "Good" Data Sources

Reliable: ☺

Original: with original sources

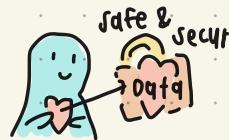
Comprehensive: research all aspects

Current: to present up-to-date

Cited: more credible e.g. Academic Research, Government data, vetted public datasets



"Bad" Data Sources don't ROCCC



• Data Ethics (Do / Don't)

↳ well-founded standards of right and wrong, and data privacy included

o GDPR (General Data Protection Regulation of the European Region)

** o Aspects of Data Ethics **

- 1) **Ownership**: individuals own the raw data they provide and have primary control over its usage.
- 2) **Transaction transparency**: all data-processing activities and algorithms should be completely explainable.
- 3) **Consent**: an individual's right to know explicitly details about how and why their data will be used.
- 4) **Currency**: individuals should be aware of financial transactions.
- 5) **Privacy (personal)**: preserving a data subject's information ("Data Protection")
- 6) **Openness**: free access, usage, sharing of their data e.g. Kaggle datasets

สิ่งที่ต้องการ
ตรวจสอบ, แก้ไข,
รับเอกสาร
และติดตาม

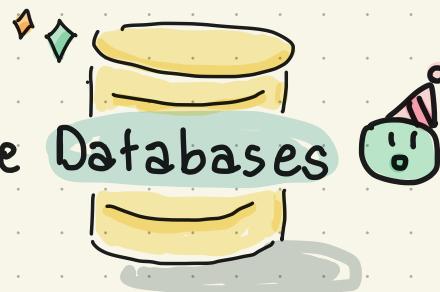
- Names
- IP addresses
- Medical Records
- Email addresses
- Photographs



Personal identifiable information (PII)
is data that likely to identify
a person and make information
known about them.
+ It's important to keep it safe!
such as:

- o Person's address
- o Credit card information
- o Social security number

Next week,
we'll play with the Databases



Week 3: Databases, Metadata

Working on large datasets in SQL

- Primary key
- Candidate / Alternate / Secondary key is a unique identifier in a row PK is a unique identifier each row in primary candidate key

Database → is a collection of data stored in a computer system.

• metadata → "data about data"

↳ tells what data is about.



Relational Database

→ is a database that contains a series of related tables that can be connected via their relationships.

↳ consisted of Table name, Attributes

Note:

* ✨ Inspecting a dataset (nhìn lướt) ôn tập

Normalization

is a process of organizing data in a relational database.

- e.g. • creating tables,
- establishing relationships between tables

→ can't be null / blank, only 1 table *

Primary key: identifier that references a column in which each value is unique

Foreign key: reference key, how one table can have ≥ 1 connected to another table

Composite key: primary key that constructed using multiple columns in a table.