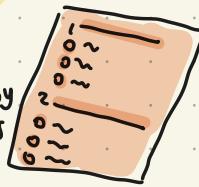


Week 1: Data Exploration Data Formats & Data Structures

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

using their
second-party
data collection
its audience



Data Collection Considerations

- ③ {
 - How the data will be collected?
 - choose data sources
 - ④ {
 - Decide what data to use
 - How much data to collect?
 - ① {
 - Select the right data type
 - ② {
 - Determine the time frame

Must know!

lives inside a computer's own systems

outside company (in case that internal data isn't enough / need more info to describe data)

collected by a researcher

from **first-hand sources** Primary

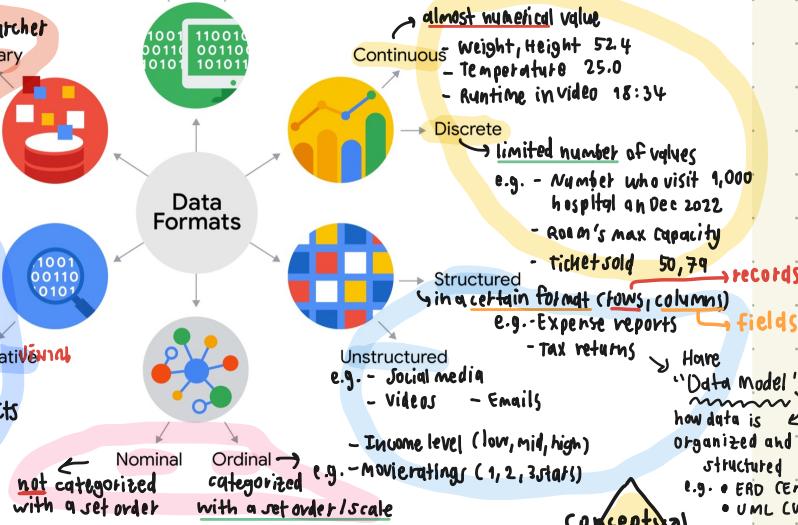
by **other people** Secondary

* subjective and explanatory *

measures of qualities and characteristics

e.g. - Exercise activity most enjoyed

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



A hand-drawn diagram of a grid with 12 horizontal rows and 4 vertical columns. The word "Fields" is written above the first column, and the word "records" is written to the right of the last row.

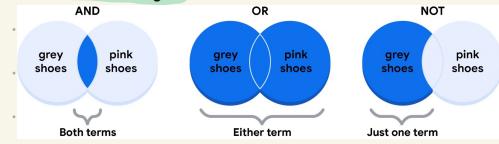
Data Types → a specific kind of data attribute that tells what kind of value data

In: ↗ Number (1, 3.14, 3.50, 1.1)

spread → text listing ("Chancery")

sheet ↴ Boolean (True/False)

Boolean logic



calculation

easy to understand
- I. After 45

Wide Da

every data subject

Lang Data format

each row is one film

paint per subject, e.g.

id	frame	name
1	Ariya	PP
2	Dohald	TRUMP

what to use depends on your work | can be
as a single row with multiple columns

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

Week 2: Data Bias (4)
Ensuring Data Integrity < Good data 🎉
Data Ethics (6) 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

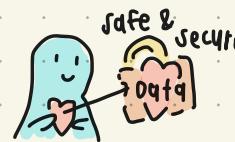
Comprehensive: research all aspects

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
ตรวจสอบได้

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

Open Data
only if...

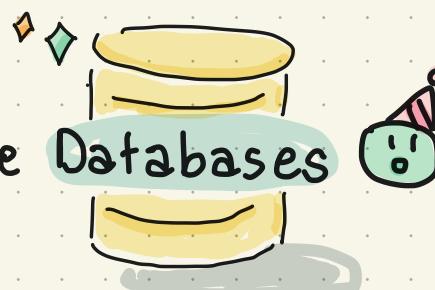
For data to be considered open, it has to: Data can be considered open when meet 3 of these standards,

- ① Be available and accessible to the public as a complete dataset.
- ② Be provided under terms that allow it to be reused and distributed.
- ③ Allow universal participation, so that everyone can use, reuse, and redistribute the data.

和个人可识别信息 (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

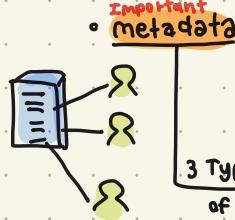


Week 3: Databases, Metadata

Working on large datasets in SQL

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

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



Important metadata

- "data about data"
- tells what data is about. ↳ a database specifically created to store metadata to make it easier to use, help track who uses metadata.
- such as Title and Description, Tags and Categories, who created it and when from Email, Photo, Spreadsheet, Websites, Books, ...
- 3 Types of metadata:
 - ① Descriptive: metadata that describes a piece of data such as ISBN behind book.
 - ② Structural: metadata that indicates how a piece of data is organized such as table of contents.
 - ③ Administrative: metadata that indicates the technical source of a digital asset such date & time photo taken.
- delineated by character / space / tab
- Comma-Separated values (CSV) ↳ helps examining a small subset of a large dataset
↳ distinguishes values from one another.



Google BigQuery

Sandbox (up to 12 projects, free)

Free trial (with Google Cloud, may have cost for upgrade)



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

- consisted of Table name, Attributes -
 - Primary key: identifier that references a column in which each value is unique ↳ can't be NULL / blank, only 1/table *
 - 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.

Note:

Complete W3 CS SQL on BigQuery Workshop

Slow on Notebook, fast on Comp then in

Normalization

is a process of organizing data in a relational database.

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

Data Governance → Helps company manage data is a process to ensure the formal management of a company's data assets by metadata Analyst/Specialist.