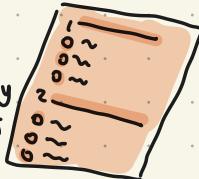


Week 1: Data Exploration

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.

outside org.

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?

④ choose data sources

④ Decide what data to use

④ How much data to collect?

① Select the right data type

② Determine the time frame

Data formats

Quantitative data → e.g. numbers, range, price

Qualitative data → cannot be counted, measured
e.g. text, description, heading

Population → entire sample → some groups
Need answer immediately → Historical Data

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

Continuous → unspecified number of possible points
have decimal limited

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

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

from External source → 3 Spreadsheet
• Google Sheets → IMPORTRANGE()
• Excel → IMPORTHTML()

Must know!

lives inside a company's own systems

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

Secondary data

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

Internal

External

External / Secondary data

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

Continuous
limited number of values
e.g. - Number who visit a hospital in Dec 2022

- Room's max capacity
- Tickets sold 50, 79

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

Structured
fields

e.g. - Tax returns

Have "Data Model"
how data is organized and structured

e.g. - Entity Relationship Diagram (ERD)
• UML (Unified Modeling Language)

high-level view of the data structure
e.g. how data interacts across an org.

Technical details of a database such as relationships, attributes, and entities

depicts how a database operates, defines all entities and attributes used e.g. table names / columns

can be converted called "Data Transformation"

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

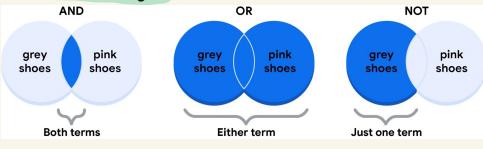
Data Types → a specific kind of data attribute that tells what kind of value data is.

In Number (1, 3.14, 3.50, 11.)

Spreadsheet → Text / String ("Chan") → not used in calculation

Sheet → Boolean (True / False)

Boolean logic



easy to understand

& often use

Wide Data format

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

what to use depends on your work

depends on your work

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 |

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

• Data Bias → type of error that systematically skews results.

→ Sampling Bias → a sample isn't representative of the population (small group sampling)

Prevent: Random sampling

- Randomly sample it with fairness e.g. not all random is all female/male

• Visualization easy to recognize it's

Types of Data Bias

→ 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.

Sometimes "Good" can be subjective, therefore to define what's good with "ROCCC"

• "Good" Data Sources

Reliable: ☺

Original: with original sources

Comprehensive: research all aspects

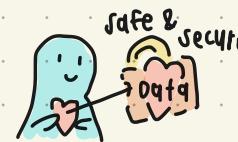
Current: to present up-to-date

Cited: more credible e.g. Academic Research, Government data, vetted Public Datasets



All affects way we make sense of data

"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

สิ่งที่ต้องคำนึงถึง
ตรวจสอบ, แก้ไข,
รักษาความปลอดภัย
และตรวจสอบต่อไป

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

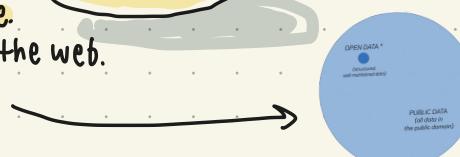
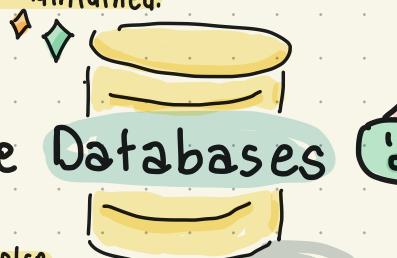
Open Data
only if...

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.

Open data is the information that has been published on government-sanctioned portals the data is structured, open-licensed, and well-maintained.

Next week,
we'll play with the Databases



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:

• Person's address
• Credit card information
• Social security number

Public data is the data that exists everywhere else.
It's freely available, but not really accessible on the web.
The data is unstructured.



Week 3: Databases, Metadata

Working on large datasets in SQL

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

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

Important metadata

→ "data about data" → **Metadata Repository**

→ 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.



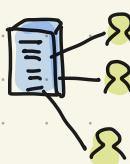
Sandbox (up to 12 projects, free)

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

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.



3 Relational Database

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

→ consisted of Table name, Attributes -

"us"
"US"
"us"
"Us"
country_code
↑
not the same

Primary key: identifier that references a column in which each value is unique → can't be null / blank, only 1/1 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.



Naming Conventions

- **SQL Dialects**, some SQL dialects are case-sensitive such as Big Query, however most of them is case-insensitive like MySQL, PostgreSQL, SQL Server.
- Column names should be all lowercase. → when create a new column use "snake_case" format e.g. as table_name
- Table names — " — in CamelCase.
- To treat the text as string use "" or ''

* Avoid spaces and special characters in file names *

naming conventions (meaningful), include date & version, short & sweet)
foldeing (organize files into foldets)

metadata to describe data
automatically back up your files

also separate unfinished files from finished one in the folder by creating "archives"

4

Week 4: Protecting Data



Benefits of organizing data:

- Easier to find & use.
- Avoid making mistakes during your analysis.
- protect your data.

"Access Control"

Data security

→ is protecting data from unauthorized access / corruption by adopting safety measures.

Security measures:

1) Encryption

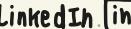
2) Tokenization: replaces the data elements you want to protect with randomly generated data or "Token". The original data is stored in a separated location and mapped to the tokens. (Tokenized data and Token Mapping)

(optional)

Week 5: Online Presence



on



Making connections

Allows recruiter to see yourself → make sure the post content is appropriate

both online & offline

Networking → professional relationship building.

→ meet people with same interests.

→ mentors

5



You did a great job!
Keep on !!

