

Week 1: Data Exploration

protecting data
ethical data analysis

Third-party data (not reliable)
outsources not directly collected

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



First-party data

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

Second-party data

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

Data Collection Considerations

1. Select the right data type
2. Determine the time frame
3. How the data will be collected?
4. Choose data sources
5. Decide what data to use
6. How much data to collect?

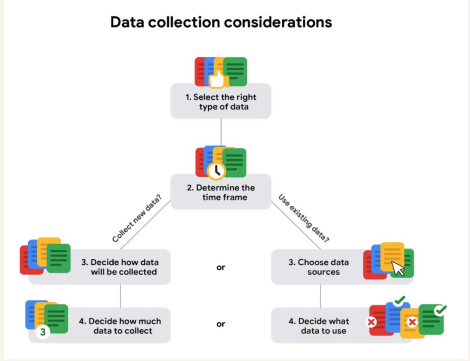
Population → entire
Sample → some groups

Need answer immediately → Historical Data

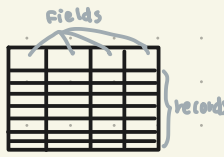
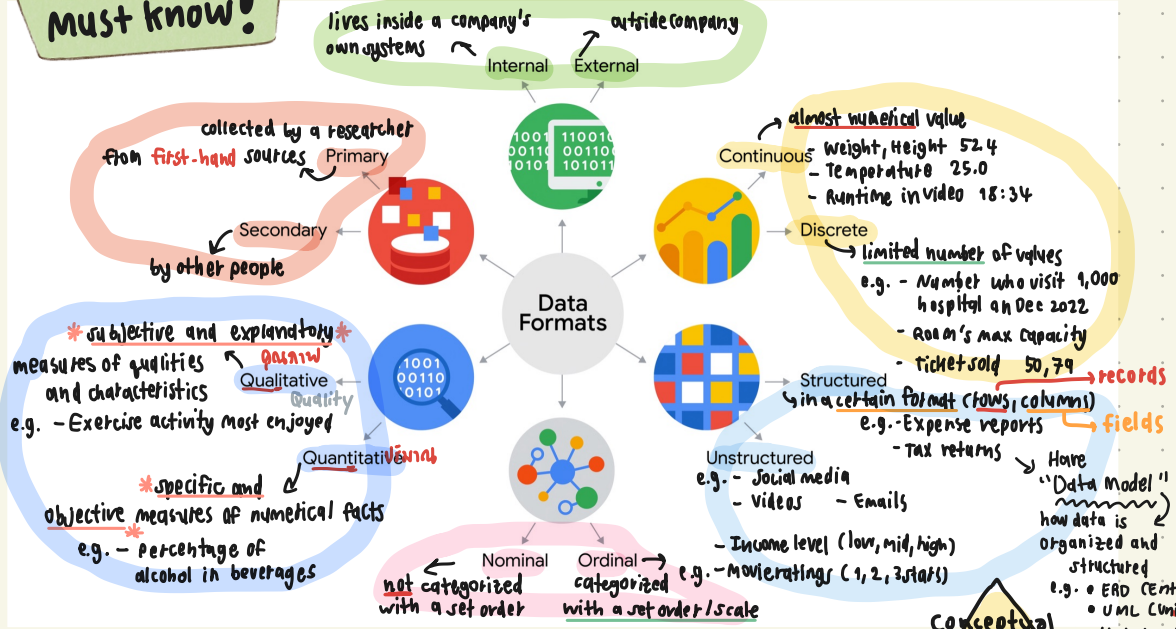
Data Formats

Quantitative data → e.g. numbers, range, price
Qualitative data → cannot be counted, measured
e.g. text, description, heading

- Discrete → counted & limited number
e.g. price, rating stars (no decimal accepted version)
- Continuous → unspecified number of possible points
e.g. weight, height, temperature
- Nominal → categorized without a set order
(no sequence)
e.g. Yes/No
- Ordinal → within an order
e.g. 1st, 2nd, 3rd



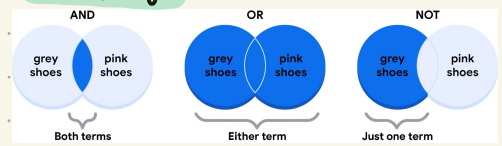
Must know!



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

In spread sheet
Number (1, 3.14, 3.50, 11)
Text (string ("Chan")) → not used in calculation
Boolean (True/False)

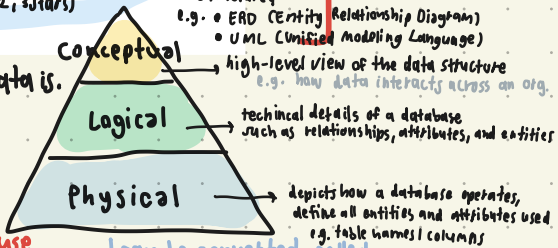
Boolean logic



id	fname	lname
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



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