

## Week 1: Data Integrity

e.g. change all of the dates to the same format

However, in case that the data only partially aligns with an objective,  
↳ Alignment to business objectives + newly discovered variables + constraints  
= "Accurate Conclusions"

### Data Integrity

→ the accuracy, completeness, consistency, and trustworthiness of data throughout its lifecycle.

these prevent data integrity

Data can be compromised when:

- Data Replication → the process of **storing** data in multiple locations  
↳ lack data integrity, not consistency ✗
- Data Transfer → the process of **copying** data from a storage device → memory, or from one comp → another
- Data Manipulation → the process of **changing** data to make it more organized & easier to read

### Data Issues that may occur:

- Duplicate data
- Not enough data / Insufficient → wait / find more other sources
- \* types of Insufficient data
  - Data from only 1 source.
  - Data that keeps updating.
  - Outdated data
  - Geographically-limited data

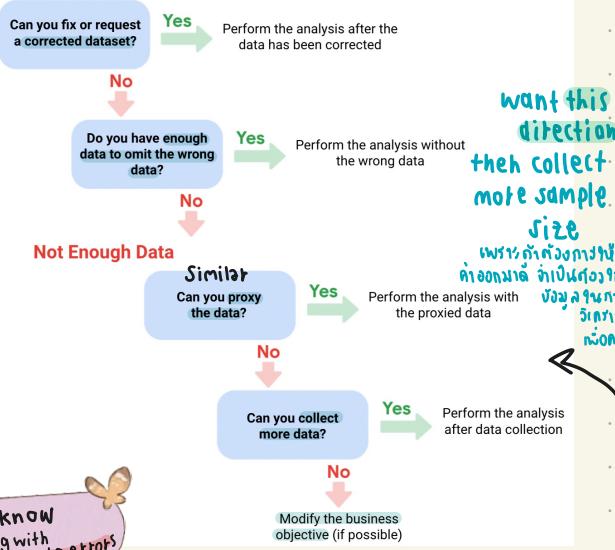


The most common workaround is when there isn't time to collect data, you can perform **proxy data** from other datasets

e.g.

Use data from another city with a similar size and demographic in case that you don't have the data for a particular city

### Data Errors



### Testing your data

**statistical power** → the probability of getting meaningful results from a test.

- **Hypothesis testing** → a way to see if a survey or experiment has meaningful results.

e.g. test whether a milkshake campaign satisfies it

- Want statistically significant result (can believe) about 0.8 or 80%  
this means that 80% → results are reliable

want this direction  
then collect more sample.

size  
minimum number  
necessary  
to get  
significance  
margin of error

### Pre-cleaning steps

1. Determine data integrity by assessing the overall accuracy, consistency, and completeness of the data.
2. Connect objectives to data by investigating how your business objectives can be served into the data.
3. Know when to stop collecting data.



# Next, let's clean the data



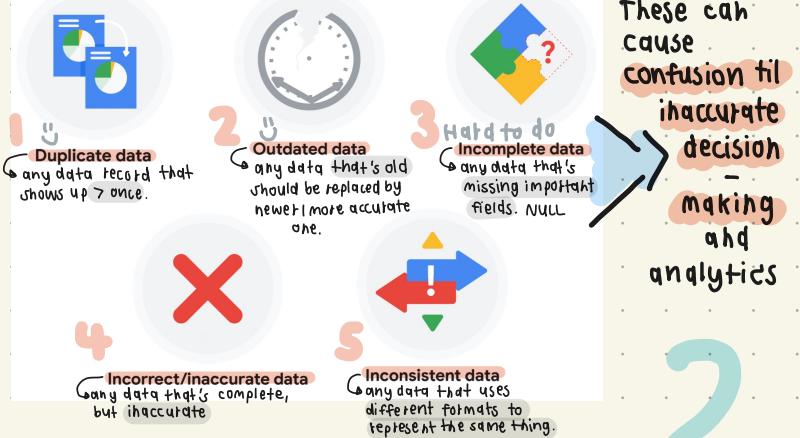
1

## Week 2: Clean the data in Spreadsheets

want to solve

**Dirty Data** → is incomplete, incorrect, or irrelevant to the problem you're trying to solve.

→ Types of Dirty Data:



Roles working with data:

คุณครุ่งข้อมูล

Data Engineers

transform data into a useful format for analysis and give it a reliable infrastructure.  
(Internal data)

Data Warehousing Specialists

→ develop processes and procedures to effectively store and organize data.  
→ make sure data is secure, available, and not lost.

→ Spreadsheet Workshop Commands → to clean up your data

Used:

- conditional formatting
- clear formatting
- change format
- Data > Split text to columns
- Data > Data cleanup > Remove Duplicates / Trim whitespace (e.g. extra spaces).  
( $=TRIM()$  or  $=TRIMC()$ )
- PivotTable
- chart plotting to see something unusual

Function:

- COUNTIF()
- LEN() count string length
- LEFT(), RIGHT(), MID()  
to extract substring from left, right, and middle in order.
- CONCATENATE() concat ≥ 2 strings  
vertical
- VLOOKUP() like PK, FK on the DB, link between 2 tables.

\* **Data Mapping** → the process of matching fields from one source to another.

- o Primary unique value in column.
- o Foreign key link to PK in another table.

Workflow Automation → the process of automating parts of your work.

Some parts can't be automated such as **Communicating with your team & stakeholders**,

**Present your findings** because there is no replacement of person-to-person communications,

Some can be partially automated: **Preparing data/Cleaning data**, **Data Exploration**,

But **Modeling data** can be automated.

validity (conform to defined business rules/constraints)  
accuracy (conform to true value)  
completeness (all is known)  
consistency (all is equivalent)

\* **Clean Data** → is complete, correct, and relevant to the problem you're trying to solve.

→ very important for external data (however internal data is also required to clean data too)

o **Data validation**: a tool for checking the accuracy and quality of data before adding/importing it such as provide a fixed length/format when receiving inputs.

→ tends to make data misaligned, inconsistent

**Data merging** → the process of combining two or more datasets into a single dataset.

\* **Want compatibility** → how well two/more datasets are able to work together

Some errors you might come across while cleaning your data include:

- Not checking for spelling errors "John" → "Jon"
- Forgetting to document errors write errors you solve, can fix later when it doesn't work.
- Not checking for misaligned values wrong cell
- Overlooking missing values
- Silence collecting data
- Not analyzing the system prior to data cleaning understand where this bad data comes from (e.g. not getting a spell check, lack of formats)
- Looking at a subset of data and not the whole picture e.g. using diff sources, you need to look over all data, in case some might be repeated.
- Losing track of the business objectives losing track on what you want to solve.
- Not fixing the source of the error Root cause of error
- Not backing up your data prior to data cleansing na spate
- Not accounting for data cleaning in your deadlines/process not forget deadlines

Next, clean your data on SQL ❤️

