

## 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**
- what makes data insufficient
- 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  - Outdated data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

### Way to address Insufficient data

solved by

- ✓ Identify trends with the available data
- ✓ Wait for more data (if time allows)
- ✓ Talk with stakeholders and adjust your objective
- ✓ Look for a new dataset

The smallest sample size for which CLT (Central Limit Theorem) is still valid is 30 (no lower than that)

**Sample Size terms:** Good when  $\frac{S.E.}{\text{margin of error}} \geq 5$  (i.e. if 5% of respondents report they would buy the product again, margin of error is 1%, then the population size is 90-80%)

1. **Margin error** → result of sample is expected to differ from result of entire population, the smaller, the closer result differences. (commonly used is 95%)
2. **Confidence level** → is how confident you are in the survey results. (How much you can rely on your sample size)
3. **Statistical Significance** → is the determination of whether your result could be done due to random chance or not. The greater, the less due to chance.

### 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.

### Testing your data

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

o **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



# Next, let's clean the data



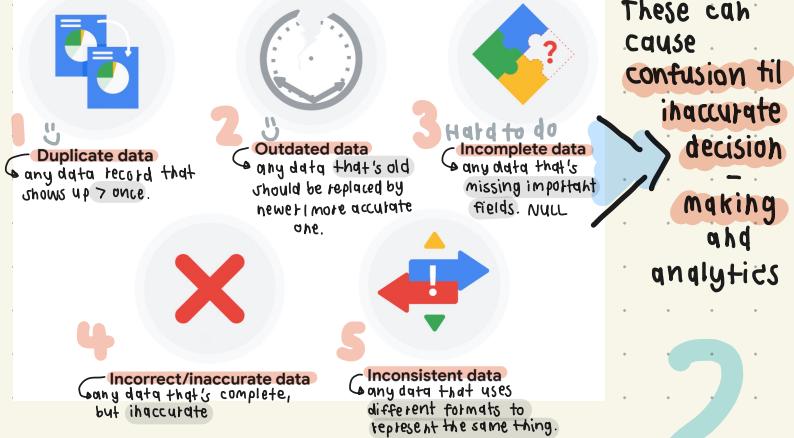
1

## 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:



# 2

Roles working with data:

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:



Next, clean your data on SQL ❤️



# Week 3: Clean the data in SQL (SEQUEL)

→ used to interact with database programs

What tool to use, depends on where the data lives

Spreadsheets	SQL Databases
Smaller datasets	Larger datasets
Create graphs & visualizations in the same program.	Prepare data for further analysis in another software.
Best when working solo on a project.	Great for collaborative work and tracking queries run by all users.
Store locally	Stored across a database
Built-in functionality	Useful across multiple programs

## Advanced cleaning functions →

### ① CAST()

\* Beware of number that's string, when searching data it isn't the same, you need to do TYPECASTING to convert data to another type e.g. `x CAST(purchase_price AS FLOAT64)`

BigQuery stores number in bits 64 system ↑

`x CAST(date AS date)` // change datetime to date format

### ② CONCAT()

to concatenate between text strings

e.g. `CONCAT(product_code, product_color)`

### ③ COALESCE()

to return non-null values in a list.

e.g. `COALESCE(product, product_code)`

↑ look for this first, if it's null then look & return from this instead

if non null record = null then it returns product\_code instead

## → SQL Workshop Commands →

- `SELECT` data
- `INSERT INTO` customer\_data.customer\_addr (customer\_id, address, country) column name  
`VALUES (2645, '333 SQL Road', 'USA');`
- `UPDATE` customer\_data.customer\_addr SET country = "US" new value WHERE customer\_id = 2645;
- `CREATE TABLE` IF NOT EXISTS
- `DROP TABLE` IF EXISTS
- `SELECT DISTINCT` customer\_id FROM   ;
- `IS NULL`
- `LENGTH (~)` → return length of text string
- `TRIM (~)` → trim in case it has whitespaces
- `SUBSTR (~, # start at, length)`
- `ASC, DESC` when ordering records

3

no duplicate  
customer\_id  
(instead of remove  
duplicates)