

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 - Outdated data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

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

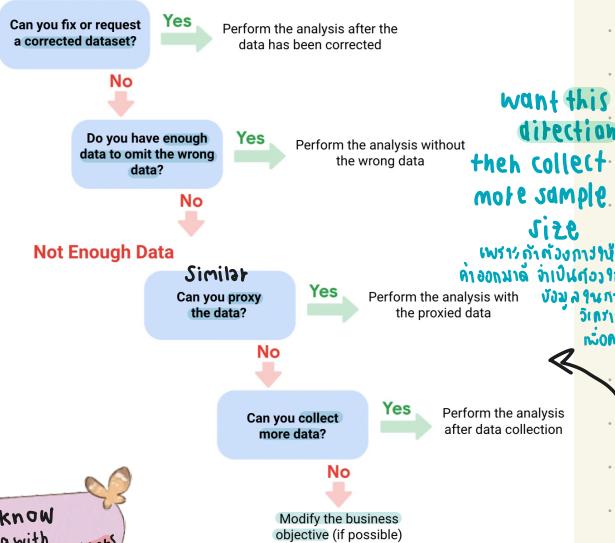
The smallest sample size for which CLT (Central Limit Theorem) is still valid is 30 (or lower than that)

Sample Size terms: Good when $\frac{S.E.}{\text{maximum amount that the}} \approx 1$ (i.e. 75% 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.

* * not have to add up to 100%.

Data Errors



Testing your data

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

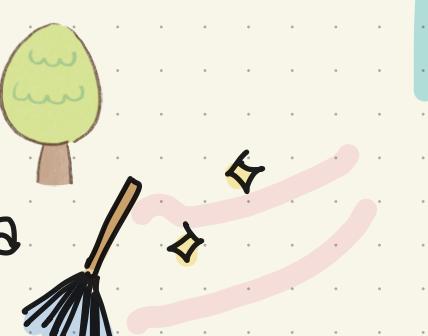
- **Hypothesis testing** → a way to see if a survey or experiment has meaningful results.

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.

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



Week 2: Clean the data in Spreadsheets