

1 Producing Data

Population

Sample Frame

Sample

SAMPLING

It is essential that the sample is representative of the larger population.

Types of samples:

- **Volunteer Sample:** almost always biased
- **Convenience Sample:** at the right moment and time
- **Sampling Frame:** only a selected sub-group
- **Systematic Sampling:** all pop., but not completely random
- **Probability Sampling:**

Simple Random Samp.

Custer Sampling

Stratified Sampling

Multistage sampling: consecutive forms of sampling

STUDY DESIGN

Observational:

Prospective or Retrospective:

- Potentially unlimited number of **lurking variables** make it difficult to prove **causation**

Survey:

- Open Q: too broad
- Close Q: too narrow
- Unbalanced response options
- Planting ideas in Q
- Complicated or sensitive Q

mitigate with anonymity and/or Random responses

Experiment:

Vocabulary:

- **X:** Factor
- **X1, X2, ...:** Treatments (ttt)
- **Individuals:** Subjects

Hawthorne effect: people behave differently in an experiment than in real life.

How to reduce bias:

- **randomly assign ttt:** randomized controlled experiment
- **blind or double-blind**
- Measure outside variables with **blocking** (~stratification)
- **Matched pairs** to compare responses for similar subjects

2 Exploratory Data Analysis

ONE VARIABLE

Variables

Individuals

Data

Dataset

Variable:

1. Nominal: no ranking, just differences
2. Ordinal: ordered differences
3. Interval: ordinal + equal intervals
4. Ratio: Interval + meaningful 0

Interpreting Histograms:

Shape: symmetry / skewness and peakedness (/modality)

Center:

- Mode: most frequently occurring value
- Mean: $x = \sum_{i=1}^n x_i / n$
- Median: M midpoint of the distribution

Spread:

- Standard Deviation: $s = \sqrt{\text{variance}} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$
- Range: Max - min
- Inter-Quartile Range (IQR): $Q3 - Q1$

BoxPlot

Standard Deviation Rule (aka Empirical Rule)

ONE VARIABLE

Point Estimation

We estimate an unknown parameter using a single number calculated from the sample data.

\hat{p} vs \hat{x}

Interval Estimation

We estimate an unknown parameter using an interval of values likely to contain the true value, and state how confident we are this interval is correct.

Estimating the mean μ

Estimating the proportion p

Normal Distribution

t-distribution

Central Limit Theorem:

Hypothesis Testing

Population

Sampling Distribution

Statistics vary from sample to sample due to **sampling variability**, and therefore can be regarded as **random variables** whose distribution is called **sampling distribution** (very theoretical).

PARAMETER: numbers that describe a **population**

- p (for categorical variable)
- μ (for quantitative variable)
- σ (for quantitative variable)

STATISTIC: numbers computed from a **sample**

- \hat{p}
- \bar{x}
- s

Those are actually random variables

1

2

3

Probability

TWO VARIABLES

RESPONSE (/dependent) VARIABLE

CATEGORICAL

QUALITATIVE

Linear Relationship Strength (/correlation coefficient)

Linear Regression with method of the least square criteria

Interpretation:

Scatter Plot

Hypothesis Testing

We have some **claim** about the **population**, and we check whether or not the **data** obtained from the **sample** provides **evidence against** this claim. Significance level of the test α , usually 5%

Testing for p

Testing for μ with α

Testing for μ without α

3 - Find p-value probability of \hat{p} when we assume H_0 is true

4 - Conclusion based on p-value (or with null value vs alternative value interval estimation)

TWO VARIABLES

RESPONSE (/dependent) VARIABLE

CATEGORICAL

QUALITATIVE

Linear Relationship Strength (/correlation coefficient)

Linear Regression with method of the least square criteria

Interpretation:

Scatter Plot

Probability

Notation:

Venn Diagram:

Probability Tree:

Basic Rules:

Independence Checks:

Random Variable: assigns a unique numerical value to the outcome of a random experiment

Discrete Random Variable

Continuous Random Variable

Normal Approximation of the Binomial

3 Inference