

## Research Targets

Population: Entire group we wish to know something about

Sample: A proportion of the population selected in the study

Sampling frame: "Source Material" from which sample is drawn

Census: An attempt to reach out to the entire population of interest

## Major Biases

- Selection bias refers to the researcher's biased selection of participants
- Non-response bias refers to participants' non-participation in the research

*every unit in a sampling frame has a known non-zero probability of being selected*

## Probability Sampling Methods

- Simple random sampling: A sample of size n is chosen from the sampling frame such that every unit has an equal chance to be selected
- Systematic sampling: The xth unit is chosen from every n/k units where x,k are chosen integers and n is the size of the sampling frame
- Stratified sampling: The population is divided into groups (strata) and SRS is applied to each strata to form the sample
- Cluster sampling: The population is divided into clusters and a fixed number of clusters are chosen using SRS

## Non-Probability Sampling Methods

Convenience sampling: Subjects are chosen based on ease of availability

Volunteer sampling: Subjects volunteer themselves into a sample

## Generalisability Criteria

- Sampling frame  $\geq$  population
- Probability sampling method implemented (selection bias  $\downarrow$ )
- Large sample size (variability and random error  $\downarrow$ )
- Minimise non-response rate

## Variable Types

Categorical: Variables that take on mutually exclusive categories

Numerical: Variables with numerical values where arithmetic can be performed meaningfully

## Variable Sub-Types

Ordinal: Categorical variables where there is some natural ordering

Nominal: Categorical variable where there is no intrinsic ordering

Discrete: Numerical variable with gaps in the set of possible numbers

Continuous: Numerical variable that can be all values in a given range

Random: Numerical variable with probabilities assigned to each value

## Properties of Mean ( $\bar{x}$ ) and Median (r)

- Adding c to all data points changes  $\bar{x}$  to  $\bar{x} + c$  and r to  $r + c$
- Multiplying c to all data points changes  $\bar{x}$  to  $c\bar{x}$  and r to  $cr$

## Properties of Standard Deviation and IQR

- $s_x$  and IQR are positive and 0 only when all data points are identical
- Adding c to all data points does not change  $s_x$  and IQR
- Multiplying c to all data points changes  $s_x$  to  $|c|s_x$  and IQR to  $|c|IQR$

## Study Designs *cause and effect*

Experimental study: The independent variable is intentionally manipulated to observe its effect on the dependent variable

Observational study: Individuals are observed and variables are measured without any manipulation

*Method to allocate subj into treatment and control group is Random Assignment*

## Blinding

- Single blinding is achieved when subjects do not know what group they belong to
- Double blinding is achieved when neither the subjects nor the assessors are aware of the assignment

$$\text{Sample Variance, } \text{Var} = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1};$$

$$\text{Standard Deviation, } s_x = \sqrt{\text{Var}}.$$

$$\text{coefficient of variation} = \frac{s_x}{\bar{x}}.$$

Establishing association	
Positive association between A and B: (any of the following)	Negative association between A and B: (any of the following)
$\text{rate}(A   B) > \text{rate}(A   NB)$ $\text{rate}(B   A) > \text{rate}(B   NA)$ $\text{rate}(NA   NB) > \text{rate}(NA   B)$ $\text{rate}(NB   NA) > \text{rate}(NB   A)$	$\text{rate}(A   B) < \text{rate}(A   NB)$ $\text{rate}(B   A) < \text{rate}(B   NA)$ $\text{rate}(NA   NB) < \text{rate}(NA   B)$ $\text{rate}(NB   NA) < \text{rate}(NB   A)$

## Symmetry Rules

- $\text{rate}(A | B) > \text{rate}(A | NB) \iff \text{rate}(B | A) > \text{rate}(B | NA)$
- $\text{rate}(A | B) < \text{rate}(A | NB) \iff \text{rate}(B | A) < \text{rate}(B | NA)$
- $\text{rate}(A | B) = \text{rate}(A | NB) \iff \text{rate}(B | A) = \text{rate}(B | NA)$

### Base Rule on Rates

The overall rate(A) will always lie between  $\text{rate}(A | B)$  and  $\text{rate}(A | NB)$ .

#### Consequence 1:

The closer  $\text{rate}(B)$  is to 100%, the closer  $\text{rate}(A)$  is to  $\text{rate}(A | B)$ .

#### Consequence 2:

If  $\text{rate}(B) = 50\%$ , then  $\text{rate}(A) = \frac{1}{2}[\text{rate}(A | B) + \text{rate}(A | NB)]$ .

#### Consequence 3:

If  $\text{rate}(A | B) = \text{rate}(A | NB)$ , then  $\text{rate}(A) = \text{rate}(A | B) = \text{rate}(A | NB)$ .

## Simpson's Paradox *confounder exists (one way)*

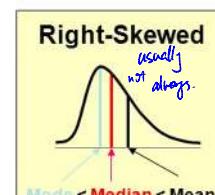
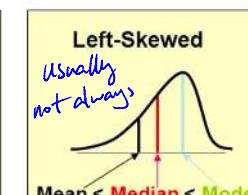
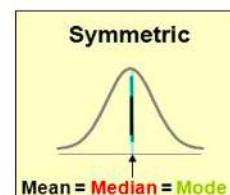
A phenomenon in which a trend appears in more than half of the groups of data but changes when the groups are combined

## Confounders *Random assignments is a preferred soln $\Rightarrow$ equal proportion*

- A third variable that is associated with both the independent and dependent variables *Prove this to show confounder exists*
- When a confounder is present, segregate the data by the confounding variable. This method is called slicing

## Outliers

- An outlier is an observation that falls well above or below the overall bulk of the data
- A general rule is that outliers should not be removed unnecessarily
- x is an outlier if  $x > Q3 + 1.5 \cdot \text{IQR}$  or  $x < Q1 - 1.5 \cdot \text{IQR}$



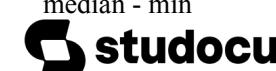
## Analysing Histograms

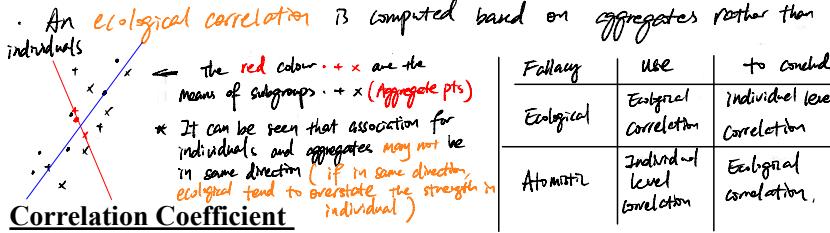
- The peak/s show the mode
- A more spread out histogram shows higher variability

*median and mode robust statistics because not affected by outlier by too much*

## Analysing Box Plots

- The center is the median
- The smallest and largest non-outlier values are the whiskers
- Skewness can be observed by comparing max - median and median - min





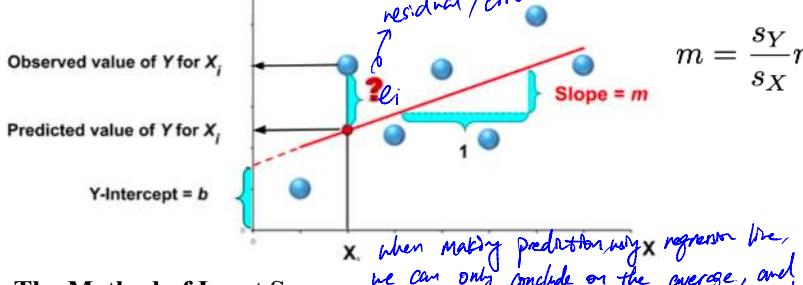
## Correlation Coefficient

- Measure of the linear association between two variables
- $-1 \leq r \leq 1$
- $0 \text{ to } \pm 0.3 = \text{weak}$ ,  $\pm 0.3 \text{ to } \pm 0.7 = \text{moderate}$ ,  $\pm 0.7 \text{ to } \pm 1 = \text{strong}$
- Removing outliers can increase, decrease, or cause no change to  $r$

Properties of  $r$   $r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x}, \frac{y_i - \bar{y}}{s_y} \right)$   $\Rightarrow$  if  $x$  and  $y$  are evenly distributed,  $r=0$ .

- $r$  is not affected by interchanging the  $x$  and  $y$  variables
- $r$  is not affected by adding a number to all values of a variable
- $r$  is not affected by multiplying a number to all values of a variable

Association  $\neq$  Causation,  $r$  value is a purely statistical relationship



## The Method of Least Squares

- For a line fit through a set of data points, the distance between the observed value and predicted outcome is the error ( $e$ )
- The sum of squares of errors is given by  $e_1^2 + e_2^2 + \dots + e_n^2$
- Our target is to minimize the sum of squares of errors
- Say we are finding linear relation between A and B. The regression line formed when A is independent and B dependent is Not interchangeable with A as dependent and B as independent.

- For an event E,  $0 \leq P(E) \leq 1$
- For a sample space S,  $P(S) = 1$
- The sum of squares of errors is given by  $e_1^2 + e_2^2 + \dots + e_n^2$
- For mutually exclusive events E and F,  $P(E \cup F) = P(E) + P(F)$

## Conditional Probability

$$P(E | F) = \frac{P(E \cap F)}{P(F)} \Leftrightarrow P(E|F)P(F) = P(E \cap F).$$

Conditionally independent:

## Probability in Independent Events

For independent events A and B: *Conjunction Fallacy*  
 $P(A) = P(A | B)$   
 $P(A) \times P(B) = P(A \cap B)$

## Sensitivity and Specificity

*True positive*  
 $\text{Sensitivity} = P(\text{Test Positive} | \text{Individual is infected})$   
*True negative*  
 $\text{Specificity} = P(\text{Test Negative} | \text{Individual is not infected})$

## Law of Total Probability

Formally, the law of total probability states that if  $E$ ,  $F$  and  $G$  are events from the same sample space  $S$  such that

- (1)  $E$  and  $F$  are mutually exclusive; and
- (2)  $E \cup F = S$ .

Then,

$$P(G) = P(G | E) \times P(E) + P(G | F) \times P(F).$$

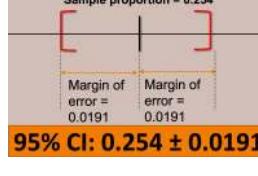
Sample statistic = population parameter + bias + random error	Confidence interval for population proportion:	Confidence interval for population mean $\mu$ :
$P \hat{\pm} z^* \sqrt{\frac{P(1-P)}{n}}$	$\bar{x} \pm t^* \times \frac{s}{\sqrt{n}}$	
$\hat{p}$ = sample proportion		event: subcollection of sample space.
$z^*$ = z-value for n-distribution		outcome: one result in the sample space.
n = sample size.		

## Normal Distribution

- A class of continuous random variables denoted as  $N(x,y)$  where  $x = \text{mean}$  and  $y = \text{variance}$
- The density curve is usually bell-shaped
- The peak of the curve occurs at the mean and the curve is symmetrical about the mean
- For a normal distribution, mean = median = mode

## Confidence Intervals

A confidence interval is a range of values likely to contain a population parameter based on a certain degree of confidence



We are 95% confident that the population parameter lies within the confidence interval

Another interpretation is that 95% of the researchers who repeat the experiment will have intervals that contain the population parameter

It is a common mistake to say that there is 95% chance that the population parameter lies within the confidence interval

## Properties of Confidence Intervals

- The larger the sample size, the smaller the random error and narrower the confidence interval
- The higher the confidence level, the wider the confidence interval

## Null and Alternative Hypothesis

- The null hypothesis asserts the stand of no effect, meaning that the variances in the sample are not inherent in the population and occurred by random chance when choosing sample
- The alternative hypothesis is what we wish to confirm and pit against the null hypothesis
- Through hypothesis testing, we wish to reject the null hypothesis in favour of the alternative hypothesis

## Significance Level

- How convincing the evidence should be before rejecting the null hypothesis. ( $0 \leq SL \leq 1$ )
- The lower the significance level, the greater the evidence required

## p-Value

- The probability of obtaining a test result at least as extreme as the result observed, assuming the null hypothesis is true
- Alternatively, the probability of observing a test result that favours the alternative hypothesis at least as much as what is observed, assuming the null hypothesis is true
- If  $p\text{-value} \geq SL$ , do not reject null hypothesis
- If  $p\text{-value} < SL$ , reject the null hypothesis

Never accept null hypothesis or reject alternate hypothesis

One-sample t-test	Chi-squared test
Mainly used to test difference between sample mean and a known or hypothesised mean.	Mainly used to test for association between two categorical variables.
Population distribution should be approximately normal if sample size is small.	Data required for the test is the count for the categories of a categorical variable.
Data used should be acquired via random sampling.	Data used should be acquired via random sampling.