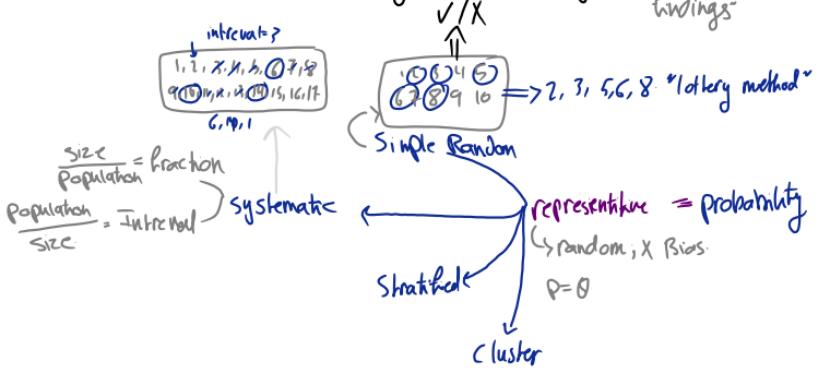


BioStatistics

unit 1 :-

- Time consuming,
- expensive
- fast & efficient
- not always possible
- strong external validity " allows for generalizing our findings"



* Data collection:-

① Primary \Rightarrow "you" collect it for a specific use
 ↳ questionnaire, survey, ...

② Secondary \Rightarrow they're already collected ; medical records, published research
 ↳ might be collected for a different reason.

Data source \Leftarrow population $\xrightarrow{\text{Sampling}}$ Sample [has to representative of the population]

Types of Sampling:

Sampling Biases:

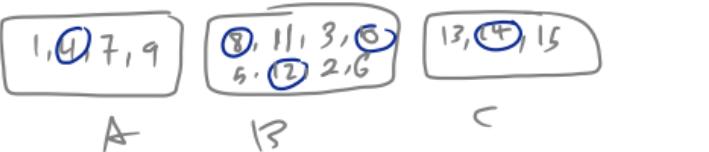
- X too small
- X too large
- exclusion w/o scientific rationale.

* process

- ① Define population \Rightarrow 1 member is called an "element"
 - ② Frame \Rightarrow all accessible members/ population elements \leadsto table
 - ③ Method
 - ④ Size \Rightarrow Cochran's formula = $\frac{Z^2}{e^2} (pq) / (1-p)$
 - ⑤ execute the process
- depends on the CI
- e^2 margin of error "precision"
-

1. external validity

2. Representation: $x/V \leftarrow$



A, B, C: characteristics needed for the study & it has to be proportional to the population.

O: selected by randomly & systematically,

1. x sampling frame $V/X \leftarrow$

2. costly

3. external validity

4. speed

5. not always possible

Cluster
↓
predefined groups:



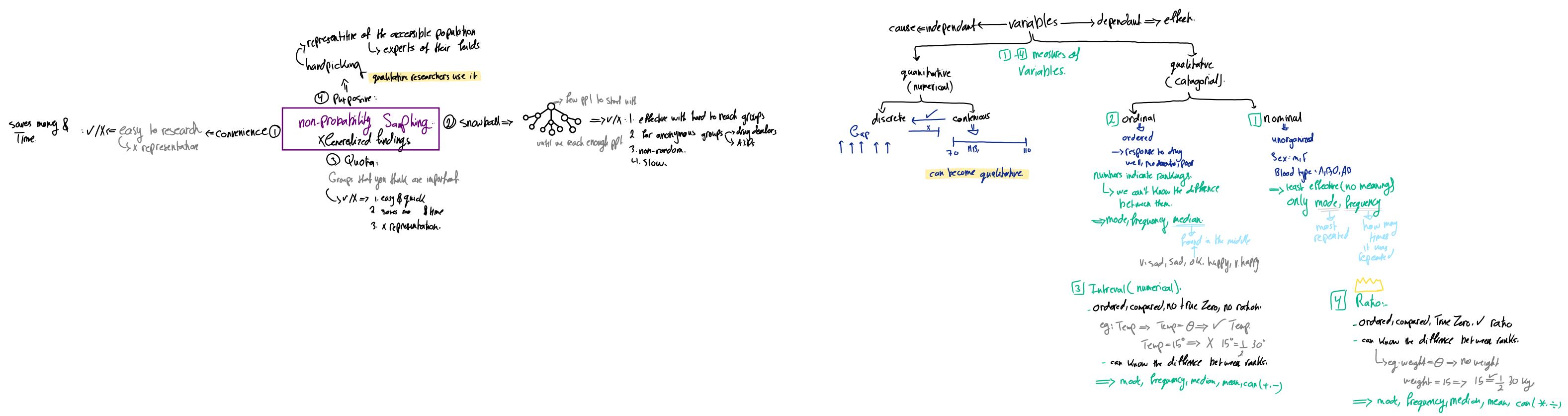
Data from selected clusters

↳ all units who meet the criteria will be sampled.

usually used in **interventional studies**

↳ e.g.: immunization coverage

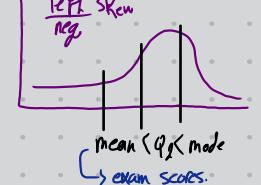
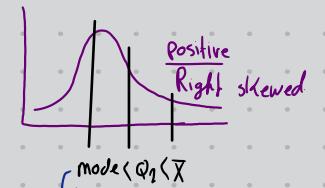




Descriptive statistics

Parameter	Sample
μ	\bar{x}
σ^2	s^2
σ	s

Parameter \leftarrow population \leftarrow Data \rightarrow organized & classified into groups \sim organized
 Sample \leftarrow Statistics \leftarrow unorganized nor summarized \sim raw data



frequency unaffected by extreme values
unique

(3) Mode:
most repeated value
not unique
could have no mode

(1) Mean (M, \bar{x})

$$\bar{x} = \frac{\sum x}{n}$$

$$\mu = \frac{\sum x}{N}$$

always present
 $+/- 0$
unique
Sensitive to extreme values

(2) Median

- middle value
odd \Rightarrow always in the middle
even \Rightarrow two middle digits

(2) Deciles (10 parts)

(3) Percentile (100 parts)

$Q_1 - Q_3 = Q_3$
 $Q_1 - Q_2 < Q_2 - Q_1$
 $Q_2 - 10 > Q_3 - Q_2$

10% 20% 30% 40% 50% 60% 70% 80% 90%

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

tendency

location

quartiles

Dispersion

(spread of values).

Range (Max-Min).

- influenced by extreme values.

IQR (Q₃ - Q₁)

Deciles (10 parts)

percentile (100 parts)

not influenced by outliers

Standard deviation ($\sigma = \sqrt{\text{var}}$)

\uparrow std \rightarrow \uparrow variability \rightarrow not a v. reliable data.

Variation ($\frac{\sum (x - \bar{x})^2}{n-1} / \frac{\sum (x - M)^2}{N}$)

\uparrow Standard error (e):
 $\frac{s}{\sqrt{n}} \rightarrow \uparrow$ error

Data classification

(1) Spacial
Geographical

(2) Temporal
over a period of time
"chronological"

(3) Qualitative

(4) Quantitative

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

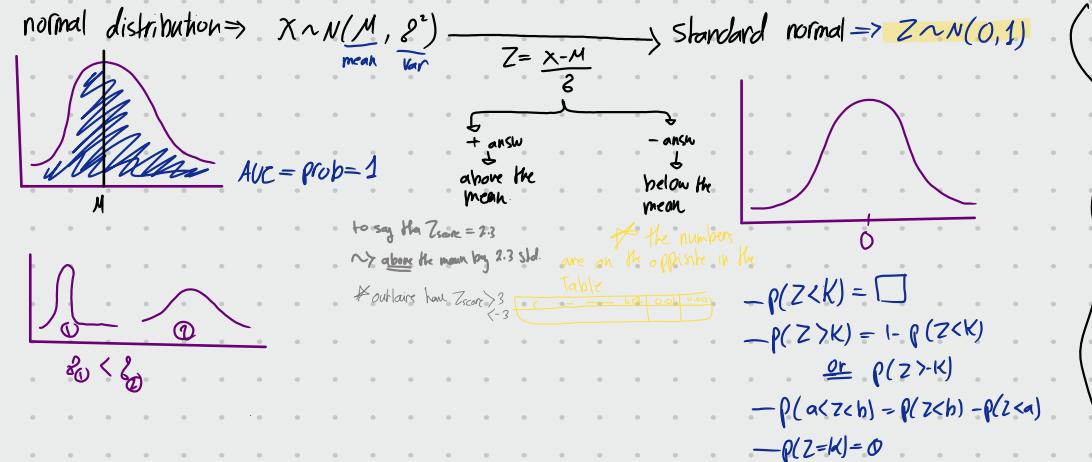
A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

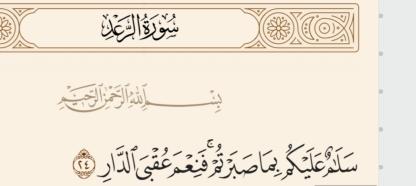
A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			20

A	5	0.15	25
B	8	0.4	40
C	7	0.35	35
Total			



Central limit theorem:-
 sample is distributed normally
 $X \sim N(\mu, \frac{\sigma^2}{n}) \Rightarrow st. dev. = \frac{\sigma}{\sqrt{n}}$

(1) from a normally distributed population.
 (2) $n > 30 \Rightarrow$ because of skewed distribution



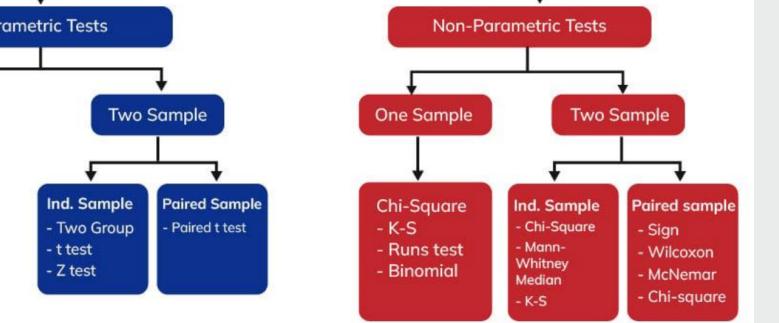
Inferential statistics
 ↗ relation between variables
 ↗ confidence interval
 ↗ Test of hypothesis

- * Research hypothesis:
 * qualitative has no hypothesis
 - predict between dependent vs. independent variable
 - contain the population
- ↗ difference \Rightarrow significant \Rightarrow we reject the null hypothesis
- H_0 : null \sim hypothesized parameter "fail to reject"
 H_1 : alternative \sim accepted if H_0 is rejected

	H_0 true	H_0 false
Rej H_0	Type I (α) error	power of test ($1 - \beta$)
Acc H_0		Type II (β) error

Type I, α : $P(\text{rej } H_0 \text{ when it's true})$
 Type II, β : $P(\text{acc } H_0 \text{ when it's false}) \Rightarrow$ more wrong
 $1 - \beta = \text{the power of test}$
 $P(\text{rej } H_0 \text{ when it's false})$

Parametric & Non-Parametric Test



Difference between Parametric & Non-parametric test

Parametric test	Non parametric test
1. Used for ratio or interval data	For ordinal or nominal data
2. Used for Normal distribution	Any distribution
3. Mean is usual central measure	Median is usual central measure
4. Information about population is completely known	No information available
5. Specific assumptions made regarding population	Assumption free test

\rightarrow T test \rightarrow One Sample
 \rightarrow 2 samples with unknown variances.

\rightarrow ANOVA $\sim 3 \leqslant$ Samples
 ↗ data as an interval/ratio (quantitative).

Chi-square:

H_0 : variables are independent (no association)
 H_1 : variables are dependent (✓ association)

observed values \Rightarrow

$$\text{expected value} = \frac{\text{Row} \times \text{Col}}{\text{Total}}$$

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

distance $O \neq E \uparrow, \chi^2 \uparrow \Rightarrow$ rej H_0
 $O \neq E \downarrow, \chi^2 \downarrow \Rightarrow$ acc H_0

* P value $< 0.05 \Rightarrow$ we reject $H_0 \Rightarrow$ ✓ association significance

P value $> 0.05 \Rightarrow$ we fail to reject $H_0 \Rightarrow X$ association