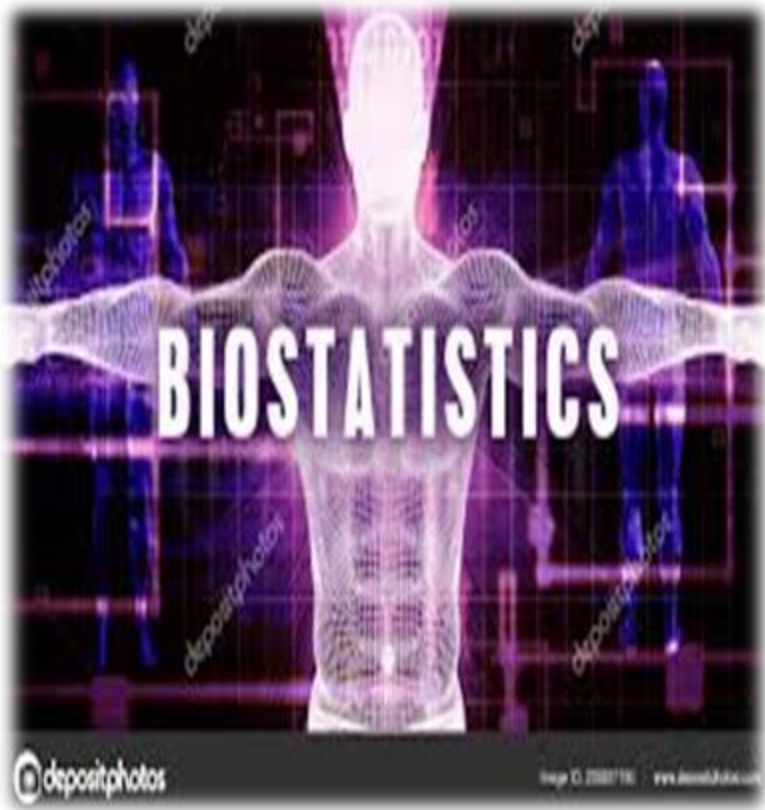


Unit One

INTRODUCTION



BIOSTATISTICS

BIOSTATISTICS

Bio: Life

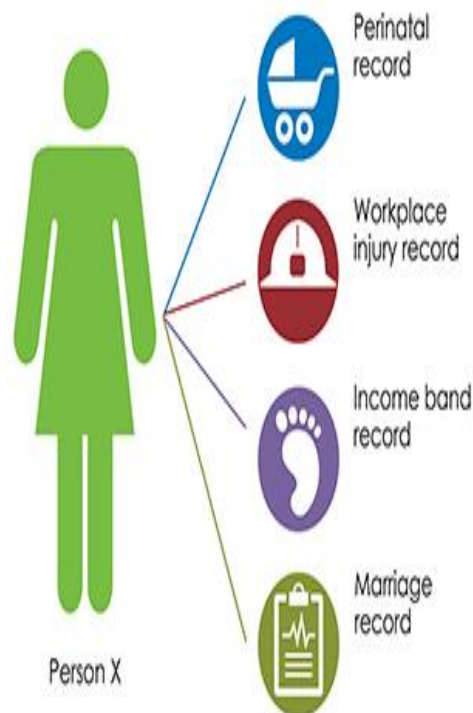
Stat: Science of facts and figures

- Provide nature and extent of illness.

Identify causation

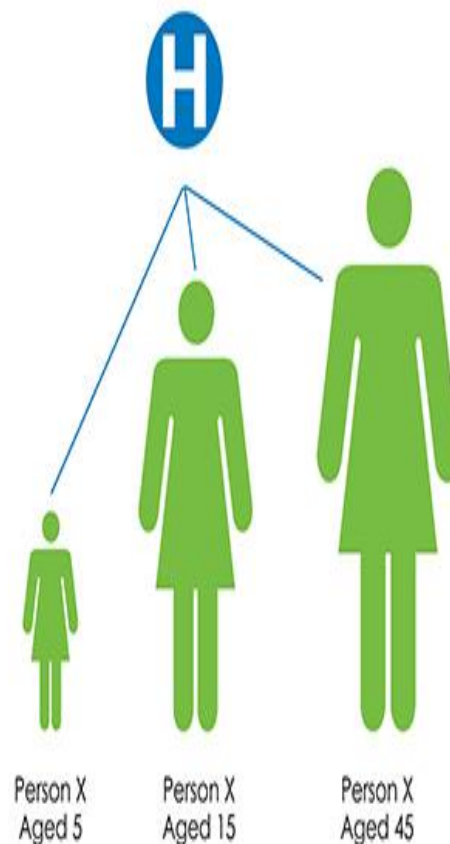
- Identify success or failure of health measures

- Research purposes/



Linkage allows information on an individual from **one data source** to be **linked** to information on the **same** individual from **another data source**.

Hospital records



Linkage can also be used for **longitudinal analyses** within the **same data source**, even when identifying information is recorded inconsistently, incompletely or has changed over time.

WHAT'S FREAKING US OUT HERE IS THAT WE'VE
FOUND A CORRELATION BETWEEN OWNING CATS
AND BEING STRUCK BY LIGHTNING



So what is Biostatics?

The application of the mathematical tools used in statistics to the fields of biological sciences and medicine.

It is a growing field with applications in many areas of biology including epidemiology, medical sciences, health sciences, educational research and environmental sciences.

Concerns of Biostatistics

Biostatistics is concerned with collection, organization, summarization, and analysis of data.

We seek to draw inferences about a body of data when only a part of the data is observed.



Why we need Statistics ?



All generalizations are false,
including this one.
Marc Twain.

To describe and summarize information
thereby reducing it to smaller, more
meaningful sets of data.

To make predictions or to generalize
about occurrences based on
observations.

To identify associations, relationships
or differences between the sets of
observations.

What is Data?

Data are **numbers** which can be measurements or can be obtained by counting.

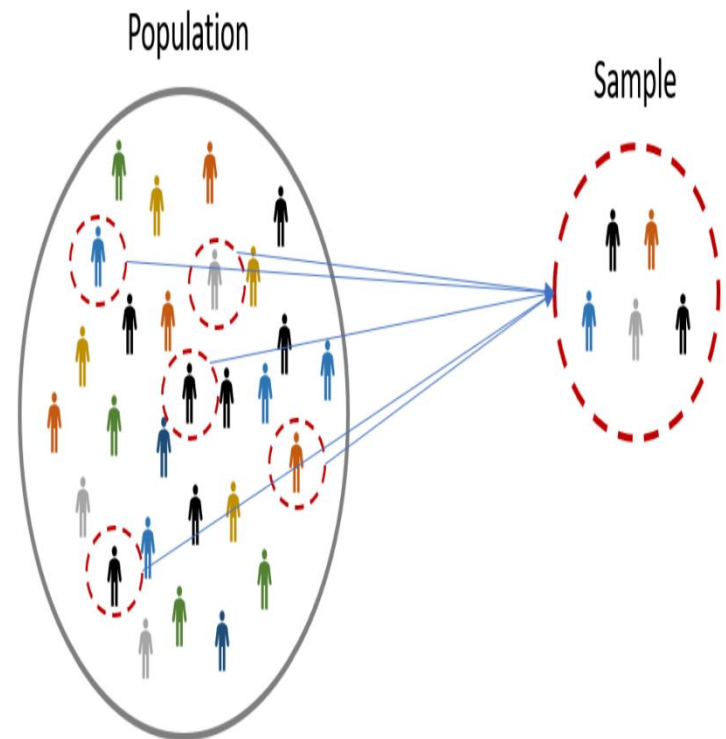
Information is the context (e.g., 12 is data, age is the information or context of data)

Biostatistics is concerned with the **interpretation of the data** and the communication of information about the data.



Populations vs. Samples

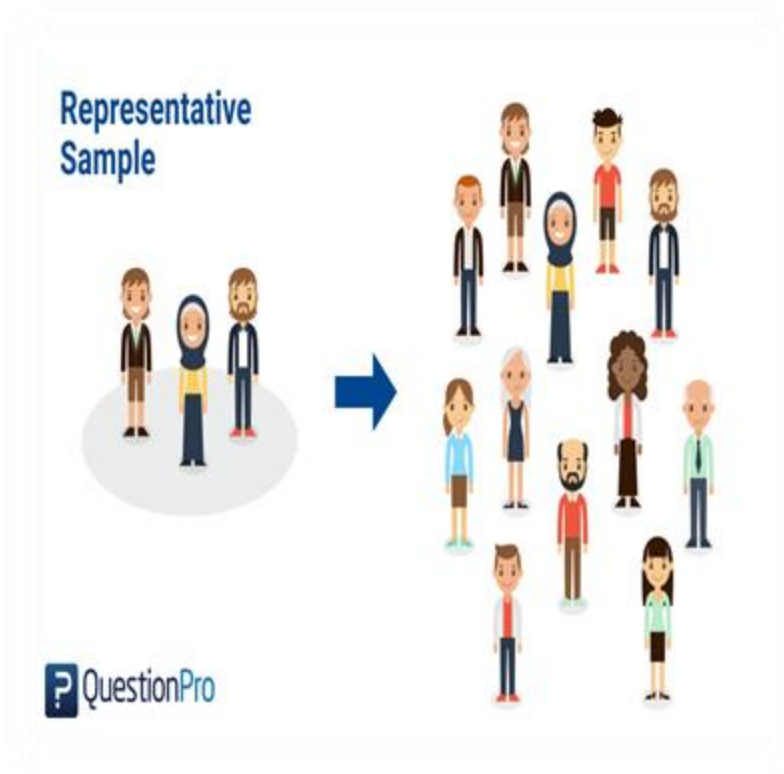
1. A **population** is the collection or set of all of the values that a variable may have. The entire category under consideration.
2. A **sample** is a part of a population. The portion of the population that is available, or to be made available, for analysis.



Population and Sampling

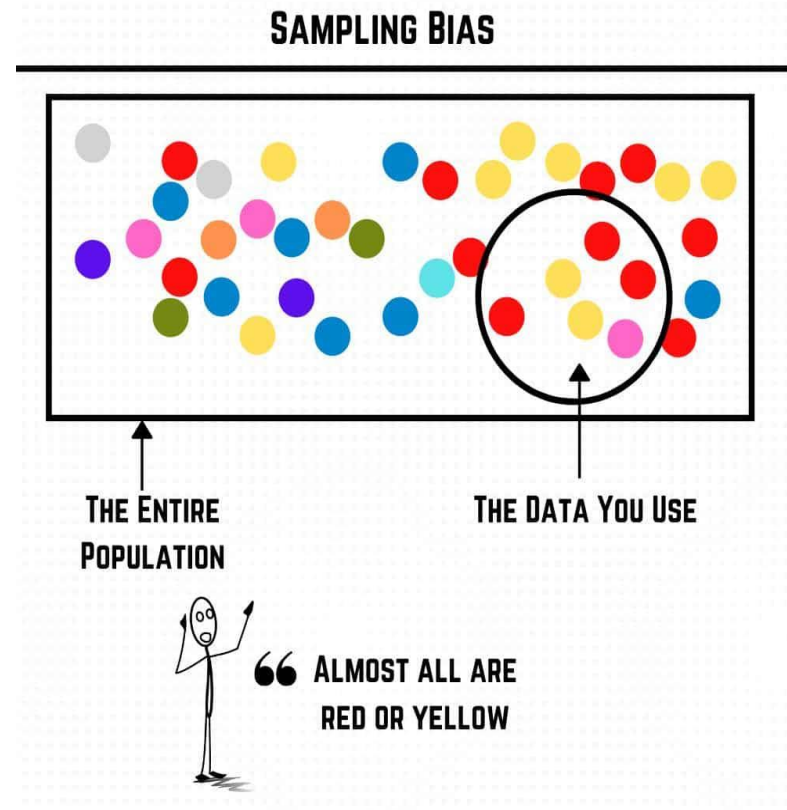
Sampling: the process of selecting portion of the population.

Representativeness: the key characteristic of the sample is close to the population.



Population and Sampling

Sampling bias: excluding any subject without any scientific rational. Or not based on the major inclusion and exclusion criteria.

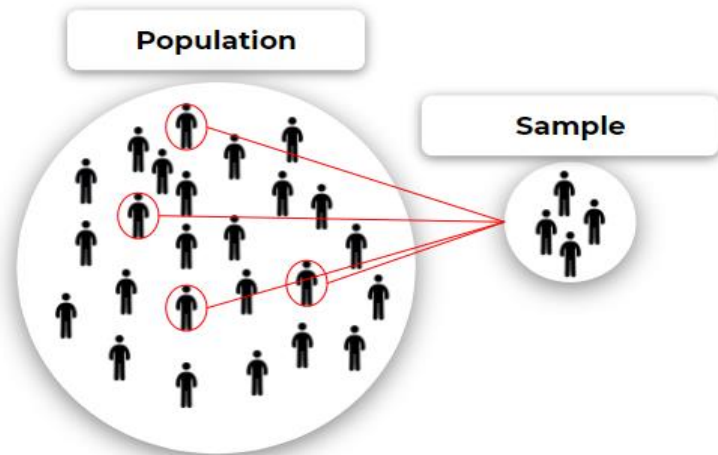


Example

Studying the self esteem and academi achievement among college students.

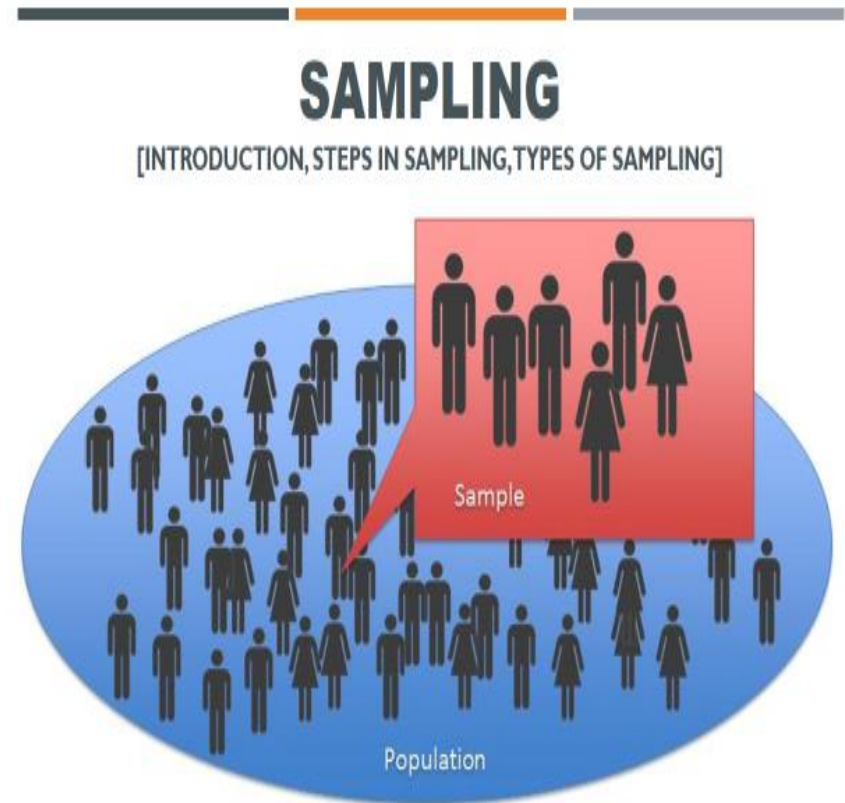
Population: all student who are enrolled in any college level.

Sample: students' college at the University of Jordan.



What is sampling?

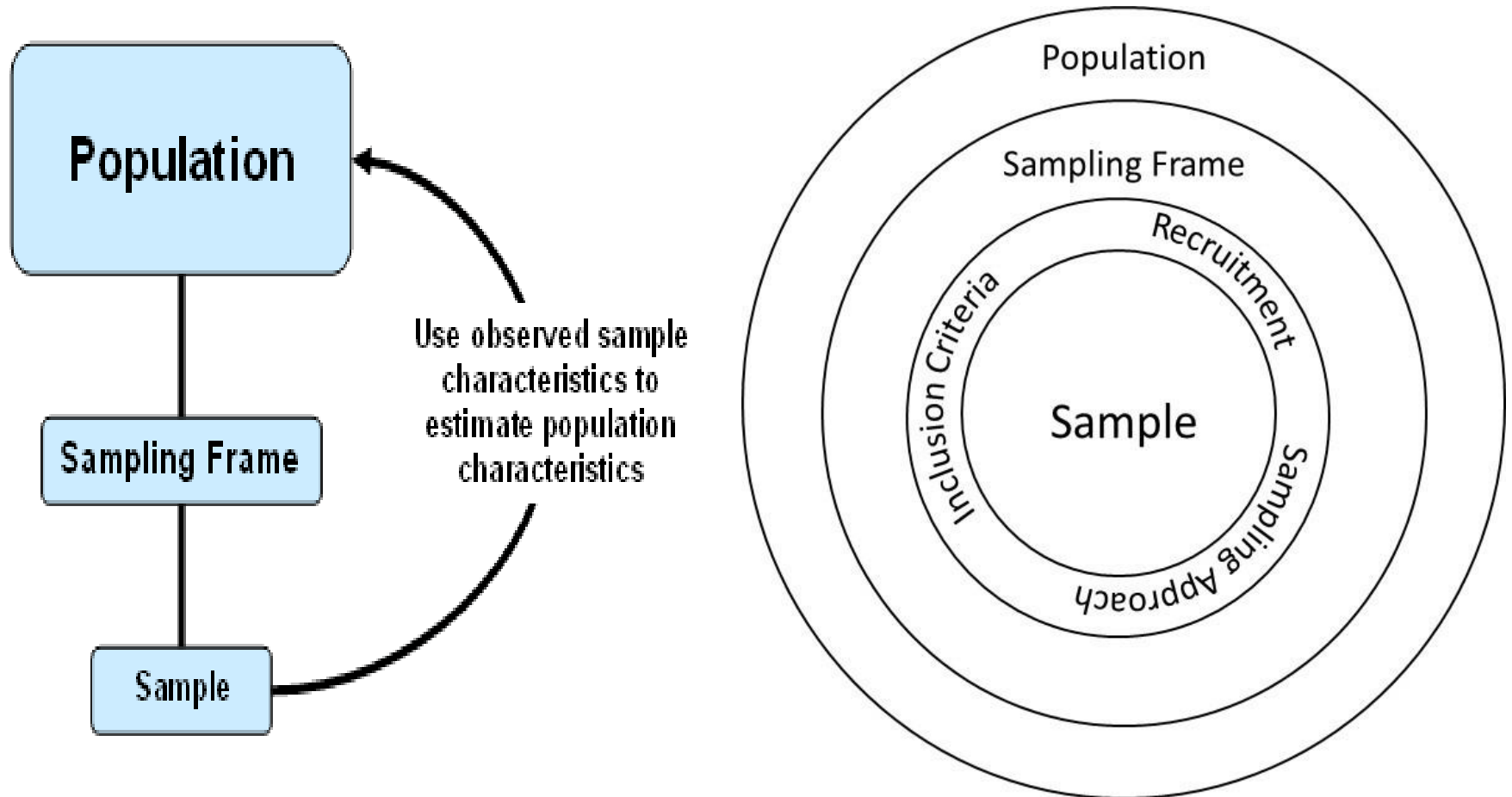
Sampling is the selection of a number of study units/subjects from a defined population.



Questions to Consider

- ❑ **Reference population** – to whom are the results going to be applied?
- ❑ What is the group of people from which we want to draw a sample (**study population**)?
- ❑ How many people do we need in our sample (**Sample Size**) ?
- ❑ How will these people be **selected**(**Sampling Method**)?

Sampling - Populations



Sampling

Element: The single member of the population (population element or population member are used interchangeably)

Sampling frame is the listing of all elements of a population, i.e., a list of all medical students at the university of Jordan, 2014-2016.

Sampling Methods

Sampling depends on the sampling frame.

Sampling frame:

is a listing of all the units that compose the study population.

Example 1:

- **Population:** Adults (18+) in Cook County
- **Possible Frame:** list of phone numbers, list of block maps, list of addresses

Example 2:

- **Population:** Females age 40–60 in Chicago
- **Possible Frame:** list of phone numbers, list of block maps

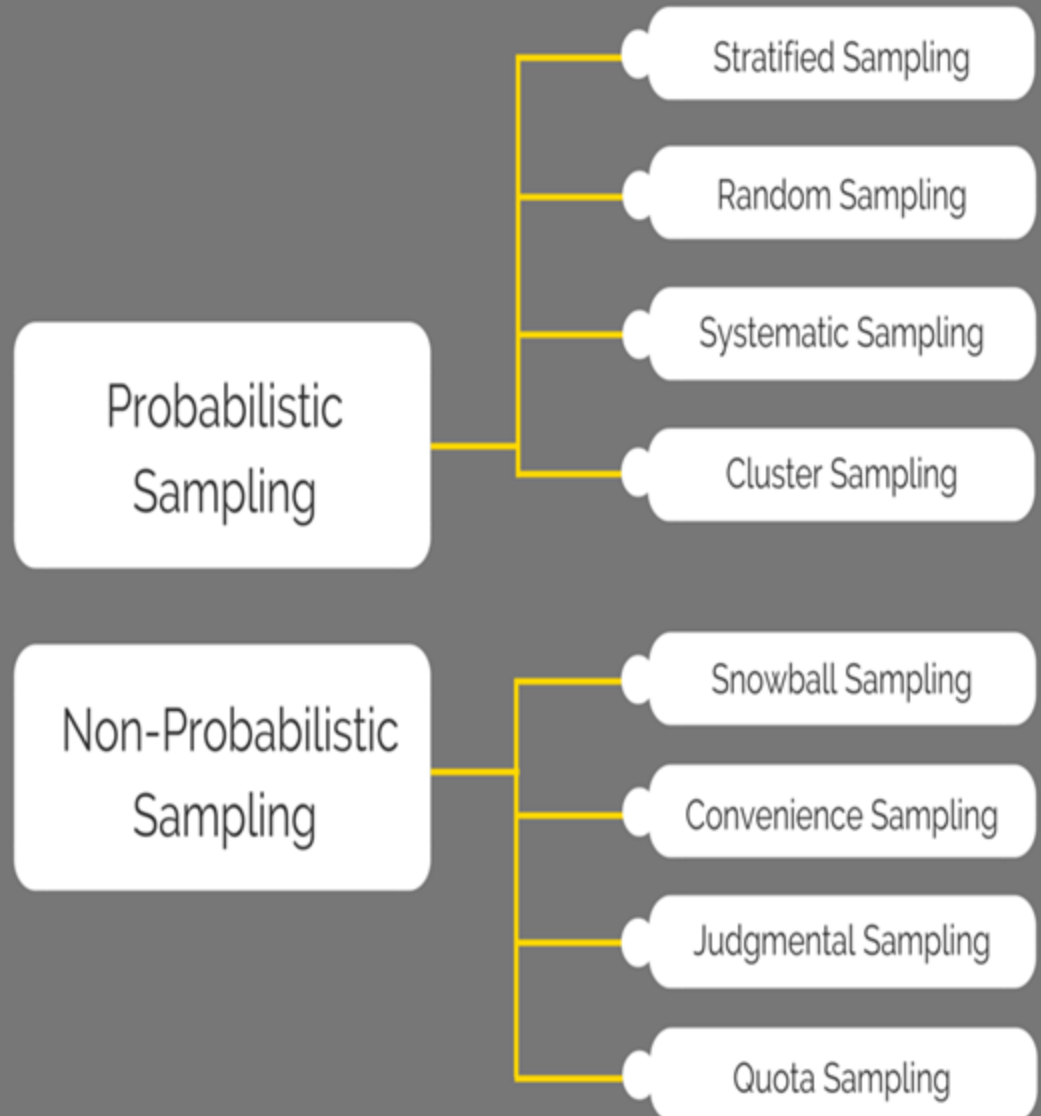
Example 3:

- **Population:** Youth age 5 to 18 in Cook County
- **Possible Frame:** List of schools



Type of sampling methods

SAMPLING TECHNIQUES

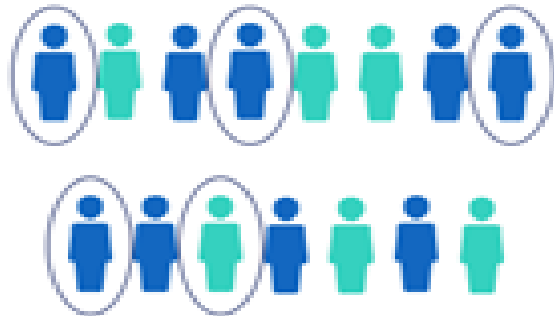


Types of Sampling Methods

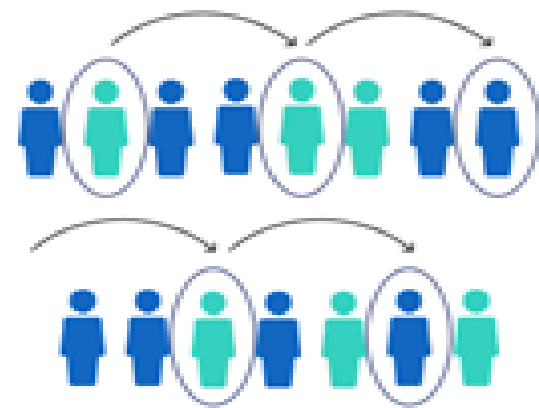
Probability Sampling Methods. Involves the use of random selection process to select a sample from members or elements of a populations.

- Simple Random Sampling
- Systematic sampling.
- Stratified sampling.
- Cluster sampling.
- Multistage sampling.

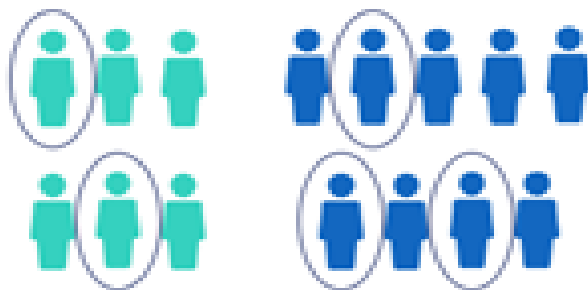
Simple random sample



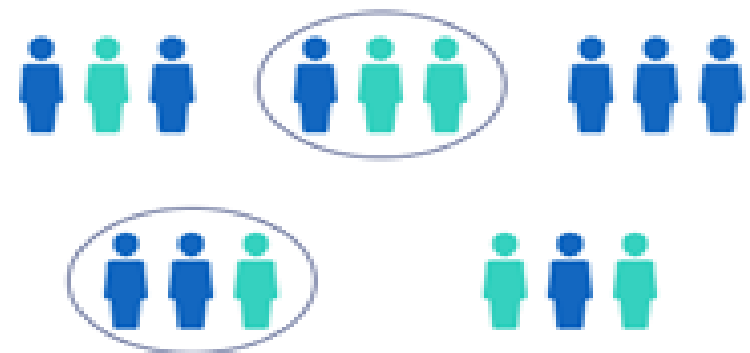
Systematic sample



Stratified sample



Cluster sample



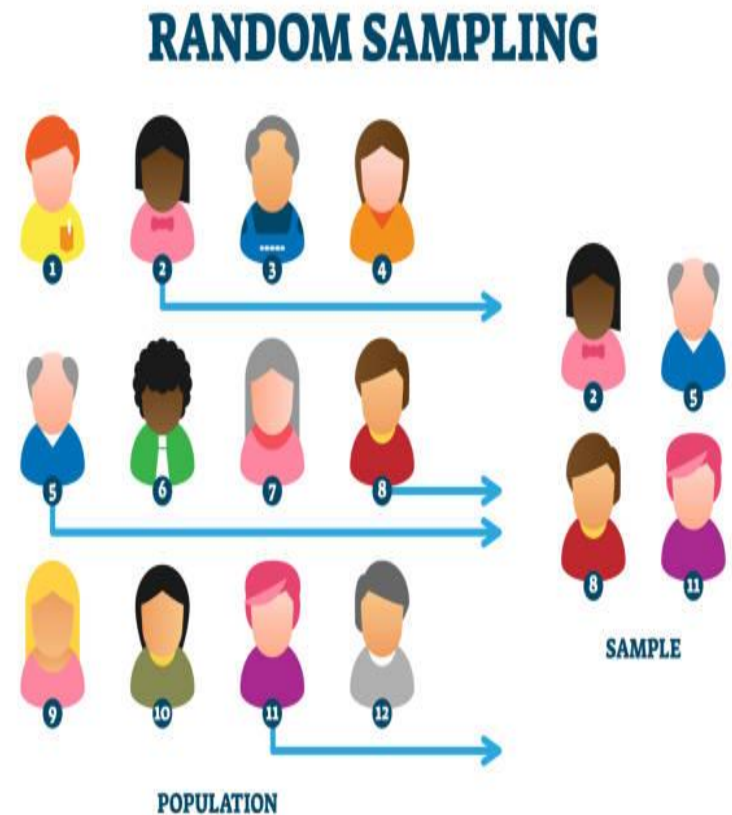
Probability Sampling Methods

Involves random selection procedures to ensure that each unit of the sample is chosen on **the basis of chance**

All units of the study population should have **an equal or at least a known chance of being included** in the sample

Requires a sampling frame

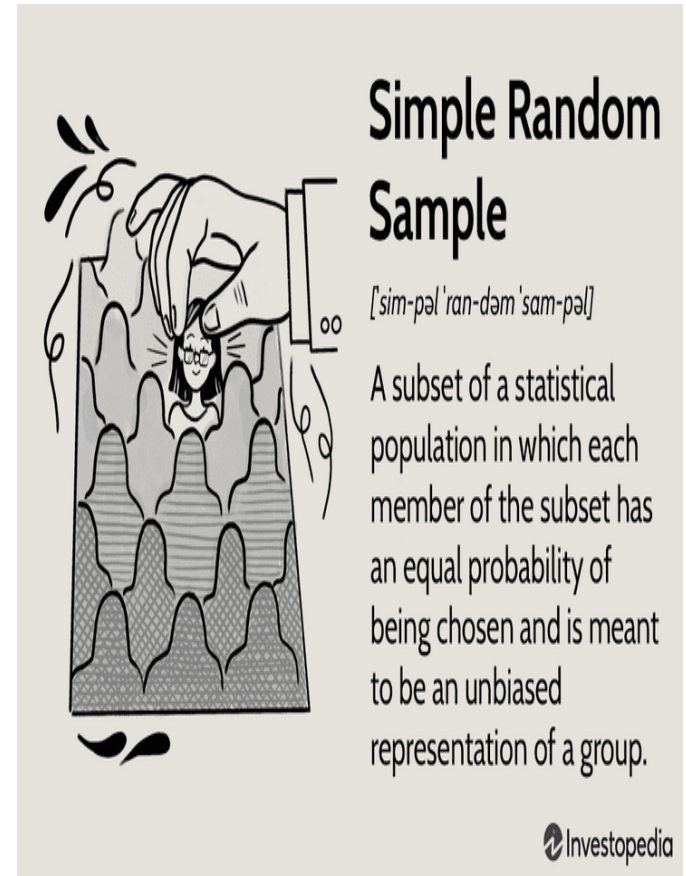
- Listing of all study units



Simple Random Sampling

This is the simplest of probability sampling

- Make a numbered list of all units in the population
- Decide on the sample size
- Select the required number of sampling units using the lottery method or a random number table



| | | | | | |
|-------|-------|-------|-------|-------|-------|
| 27798 | 12511 | 31487 | 23968 | 92052 | 69972 |
| 20741 | 65030 | 79887 | 60896 | 34880 | 80647 |
| 72357 | 22780 | 68845 | 07401 | 55229 | 40009 |
| 04912 | 18963 | 98752 | 62310 | 56615 | 37512 |
| 65671 | 11784 | 94998 | 49452 | 67552 | 87610 |
| 26900 | 84444 | 07973 | 17655 | 41558 | 29754 |
| 85865 | 45148 | 13069 | 40746 | 57146 | 28399 |
| 64831 | 27098 | 10675 | 17555 | 17833 | 72997 |
| 11063 | 88498 | 26184 | 71909 | 52135 | 78002 |

USING A RANDOM NUMBER TABLE TO MAKE A
FAIR DECISION

Systematic Sampling

Individuals are chosen at regular intervals from the sampling frame

Ideally we randomly select a number to tell us the starting point

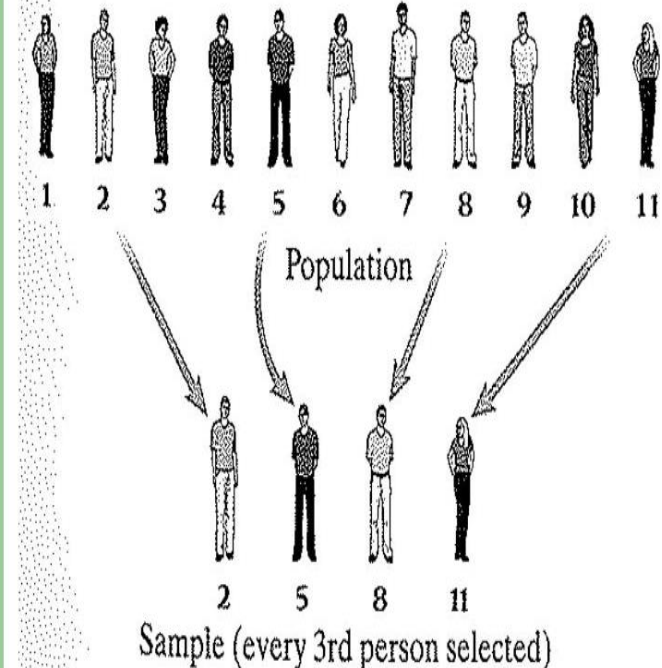
- every 5th household
- every 10th women attending ANC

Sampling fraction = $\frac{\text{Sample size}}{\text{Study population}}$

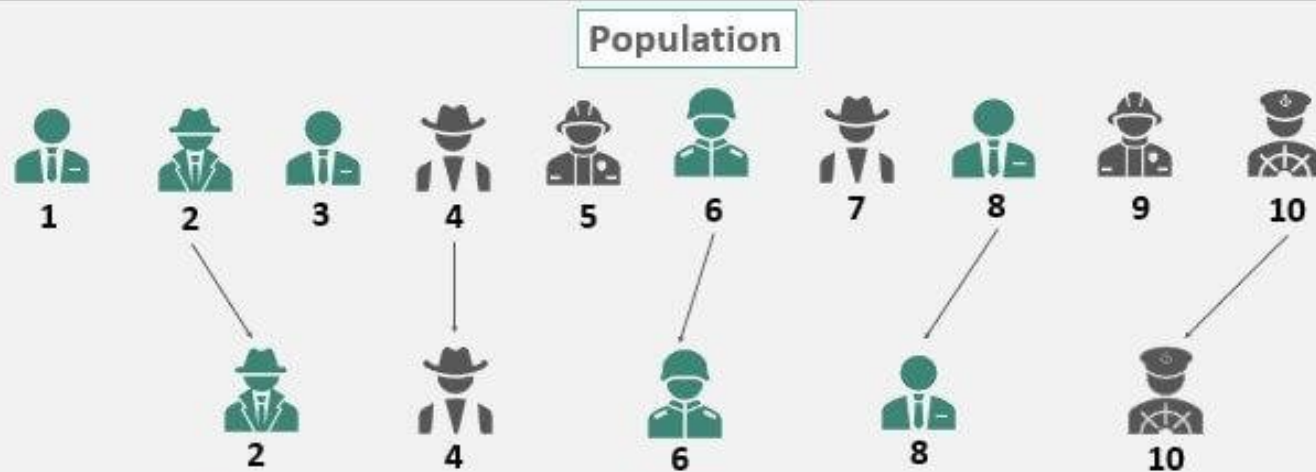
Interval size = $\frac{\text{Study population}}{\text{Sample size}}$

Interval size = $\frac{\text{study population}}{\text{Sample size}}$

Systematic Sampling



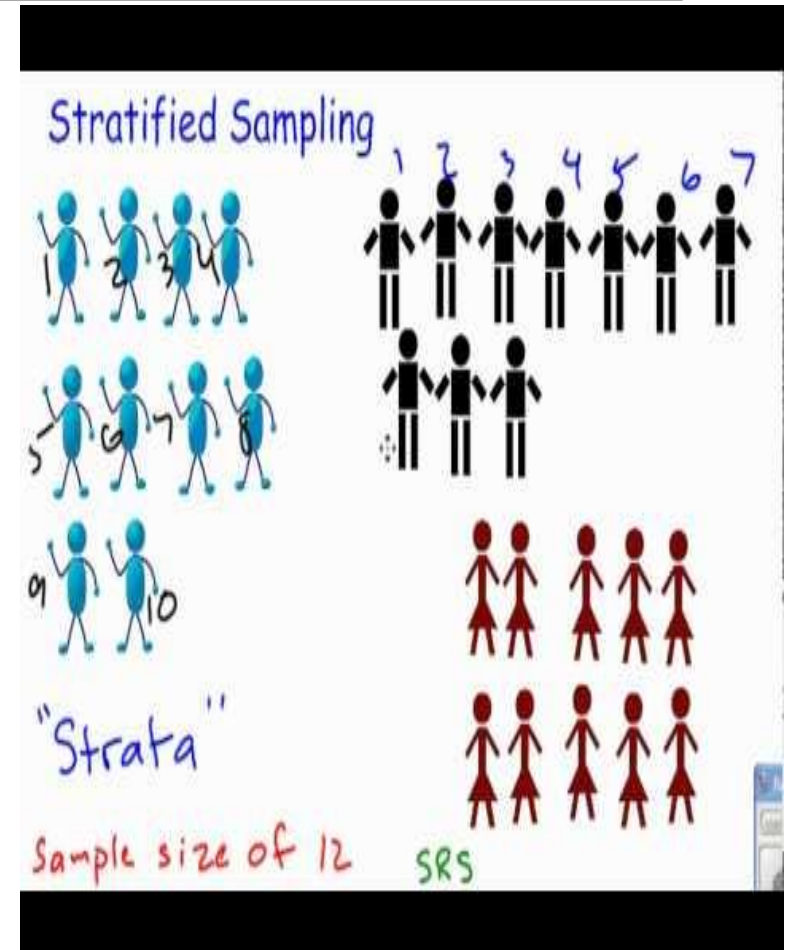
Systematic Sampling



In this case, every second person is systematically selected.

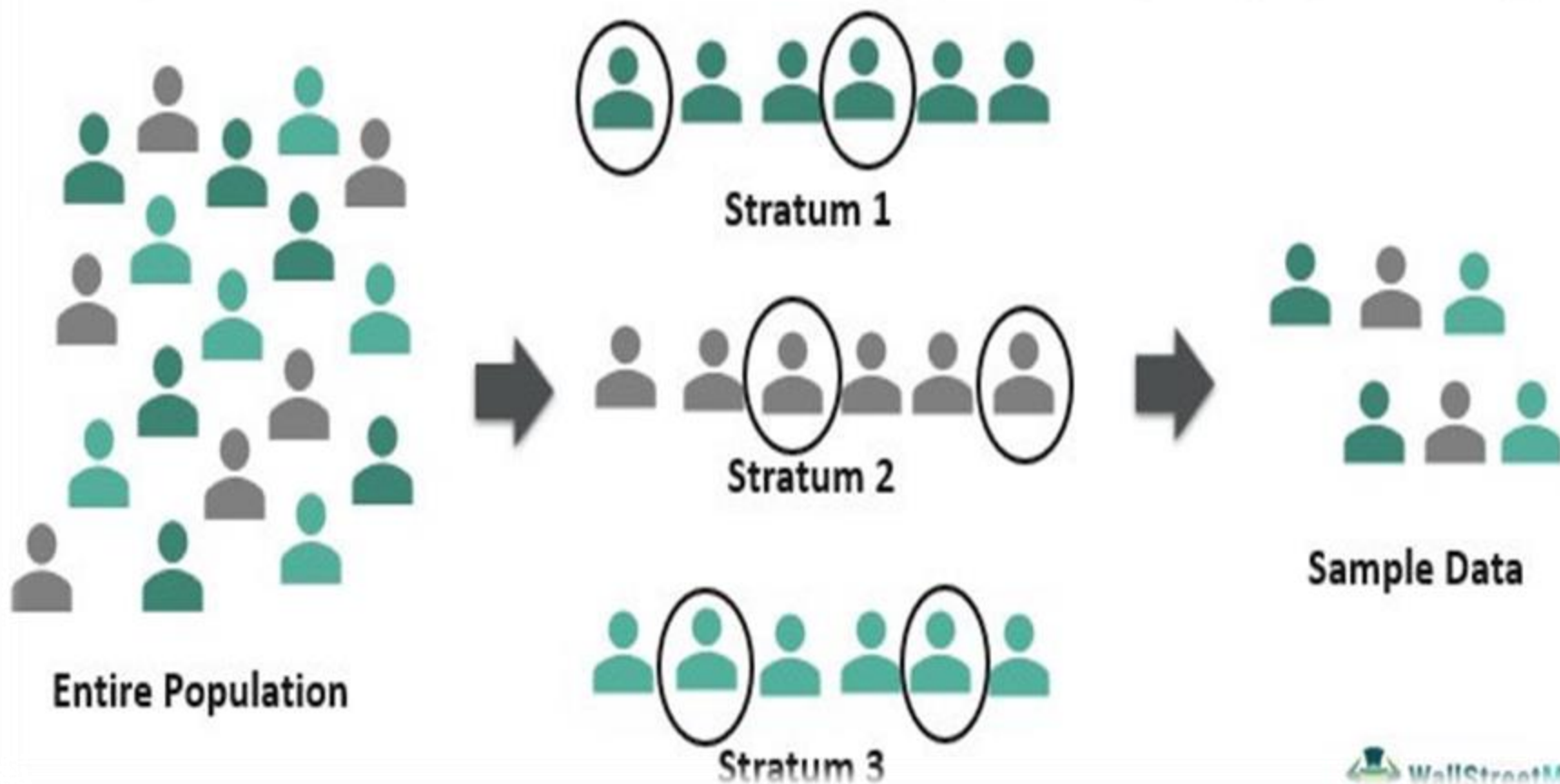
Stratified Sampling

- If we have study units with different characteristics which we want to include in the study then the sampling frame needs to be divided into strata according to these characteristics
- Ensures that proportions of individuals with certain characteristics in the sample will be the same as those in the whole study population
- Random or systematic samples of predetermined sample size will have to be obtained from each stratum based on a sampling fraction for each stratum



Stratified Sampling

Stratified sampling, refers to random sampling techniques that clubs items of whole population into different groups called strata, based on their similar characteristics. Then, samples from each stratum are taken, whether proportionately or disproportionately.



Cluster Sampling

Selection of study units (clusters) instead of the selection of individuals

All subjects/units in the cluster who meet the criteria will be sampled.

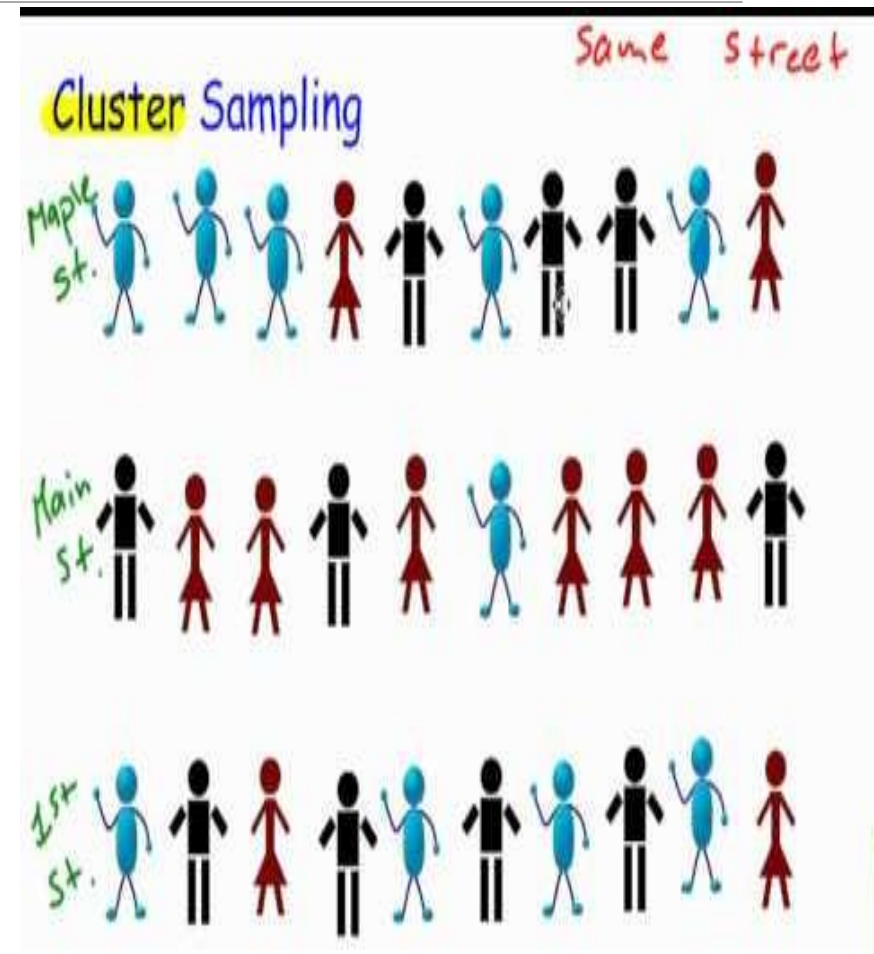
- Clusters often geographic units
- e.g. schools, villages etc

Usually used in interventional studies

- E.g. assessing immunization coverage

Advantages

- sampling frame is not required in this case
- Sampling study population scattered over a large area



Cluster sampling vs stratified sampling

Cluster Sampling

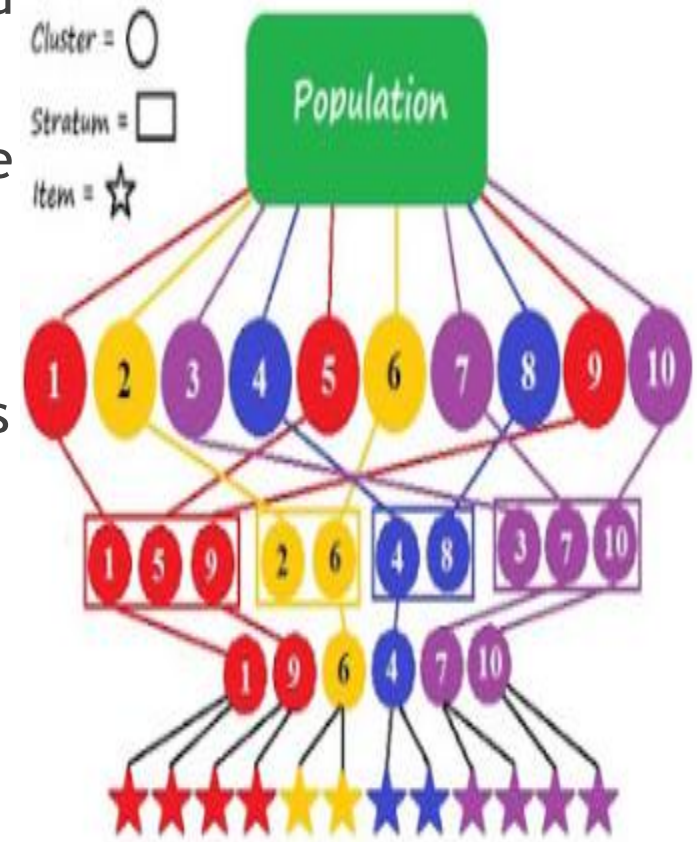
1. Cost reduction
2. Randomly selected clusters
3. Division naturally formed
4. More errors
5. Homogeneity externally

Stratified Sampling

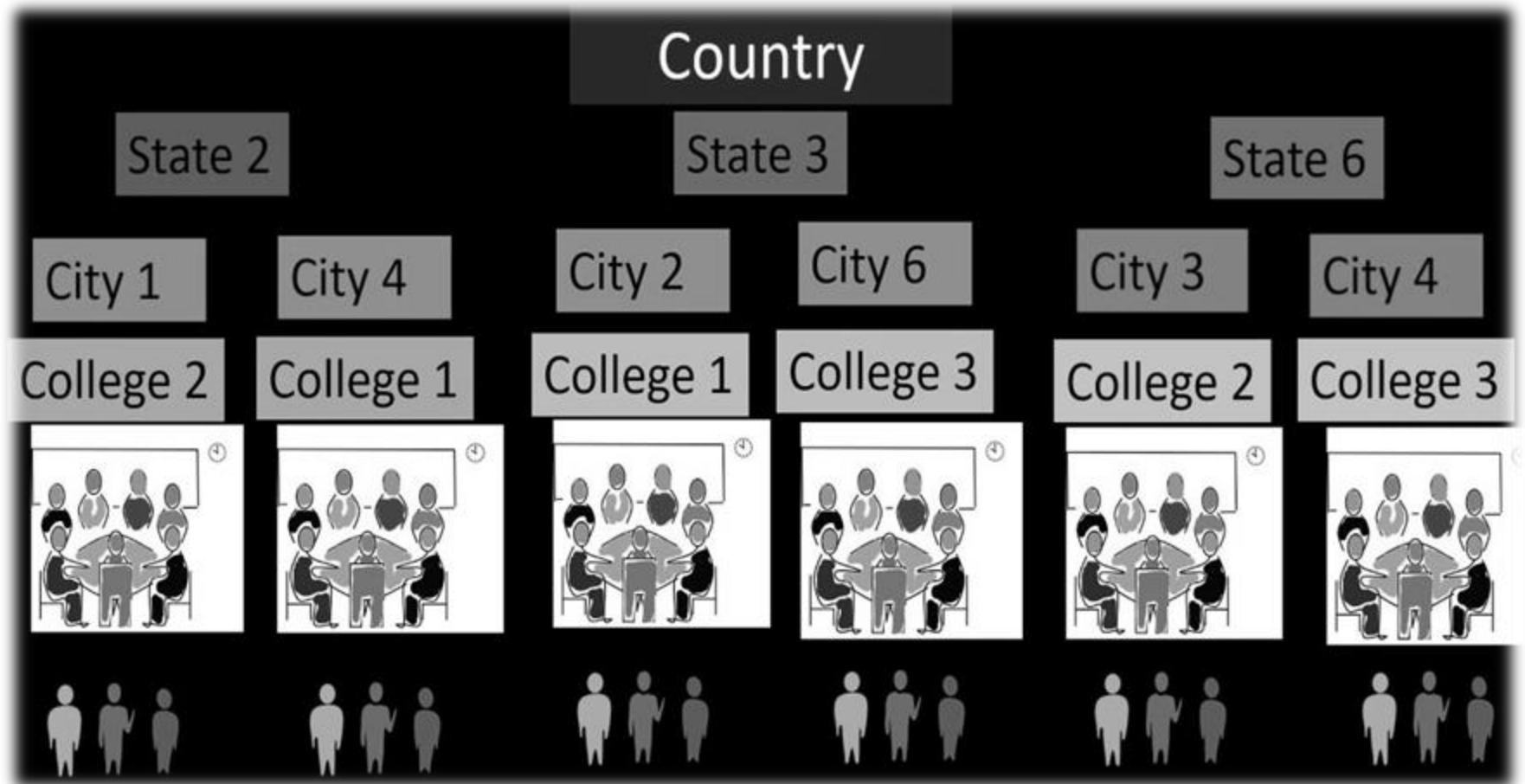
1. Enhanced precision
2. Randomly selected members from strata
3. Depends on the researcher
4. Reduced errors
5. Homogeneity internally

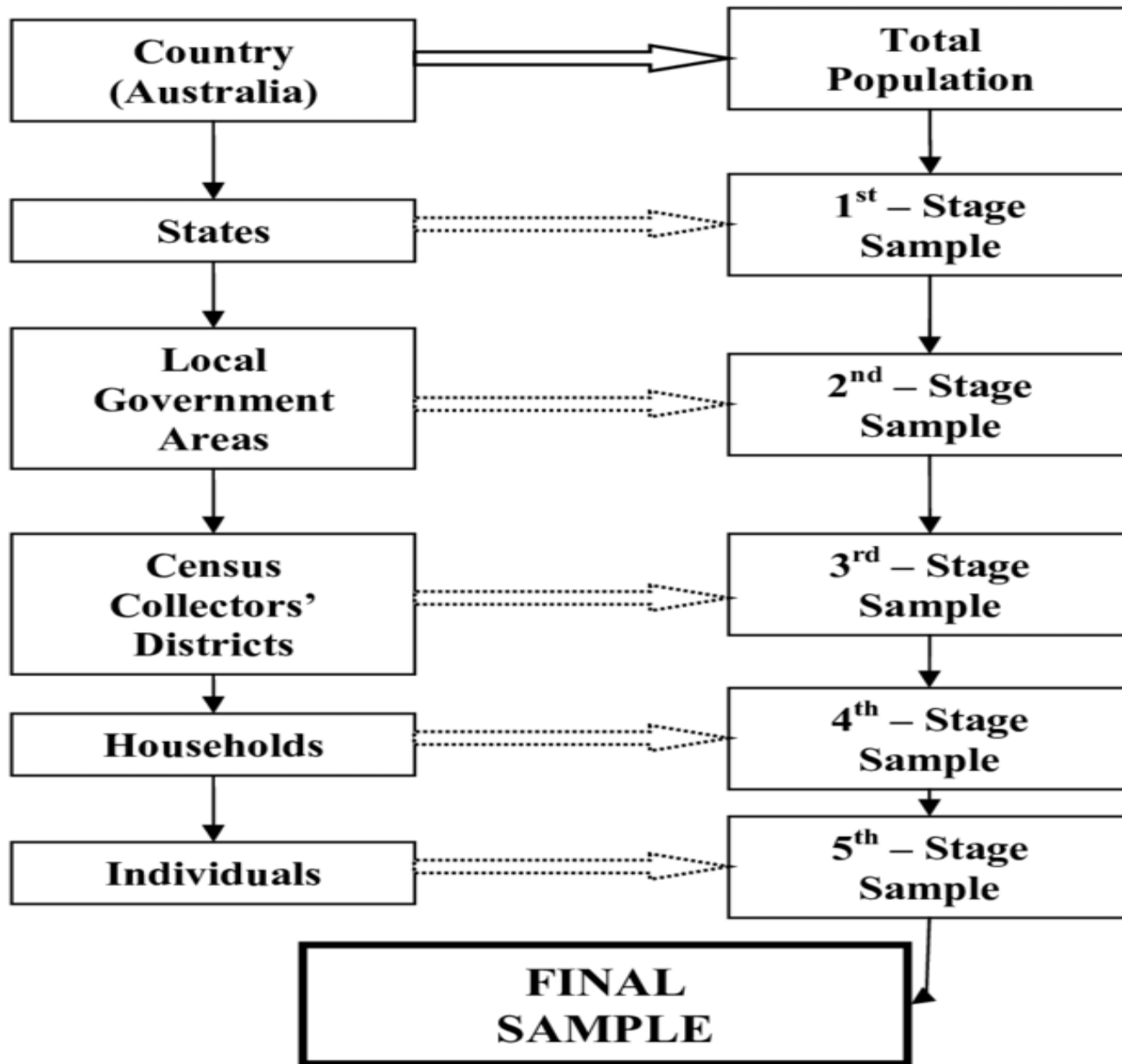
Multistage Sampling

- Involves more than one sampling method
- Is therefore carried out in phases
- Does not require a initial sampling frame of whole population
- NEED TO KNOW SAMPLING FRAME OF CLUSTERS E.G. PROVINCES
- Require sampling frames of final clusters
- Applicable to community based studies e.g. interviewing people from different villages selected from different areas, selected from different districts, provinces



Multi-stage sampling



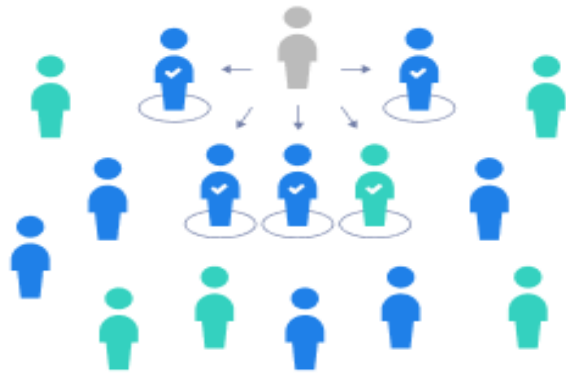


Nonprobability Sampling Methods

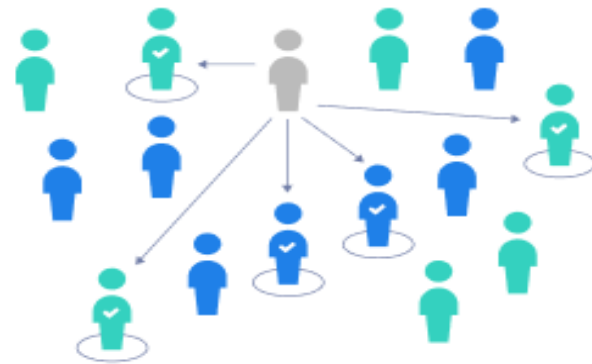
Nonprobability sampling: the sample elements are chosen from the population by nonrandom methods.

- More likely to produce a biased sample than the random sampling.
- This restricts the generalization of the study findings.
- Most frequent reasons for use of nonprobability samples involve convenience and the desire to use available subjects.

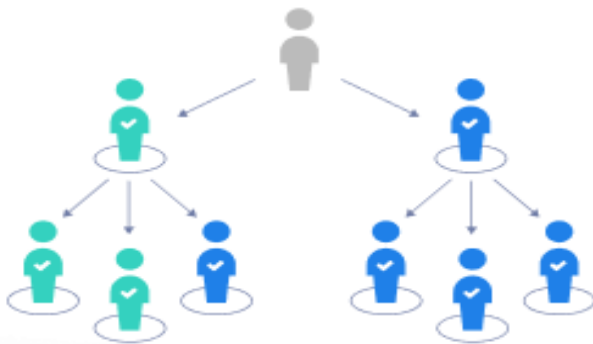
Convenience sample



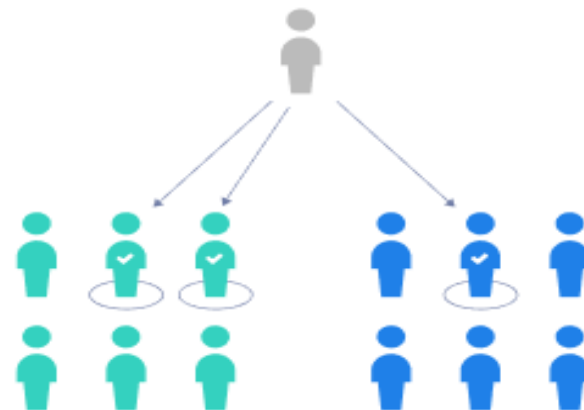
Purposive sample



Snowball sample



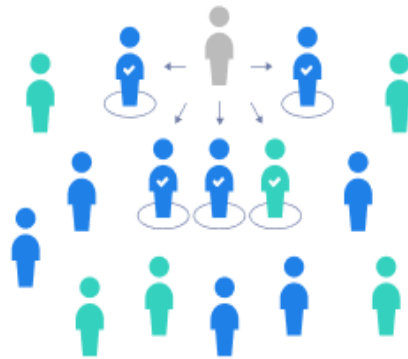
Quota sample



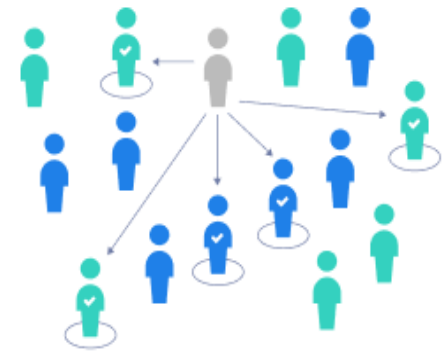
Types of Sampling Methods

Nonprobability Sampling Methods:

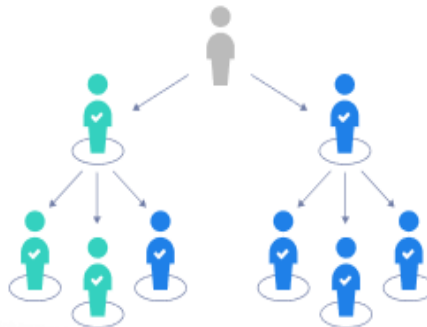
Convenience sample



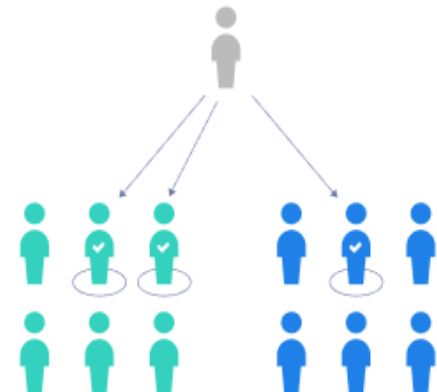
Purposive sample



Snowball sample



Quota sample



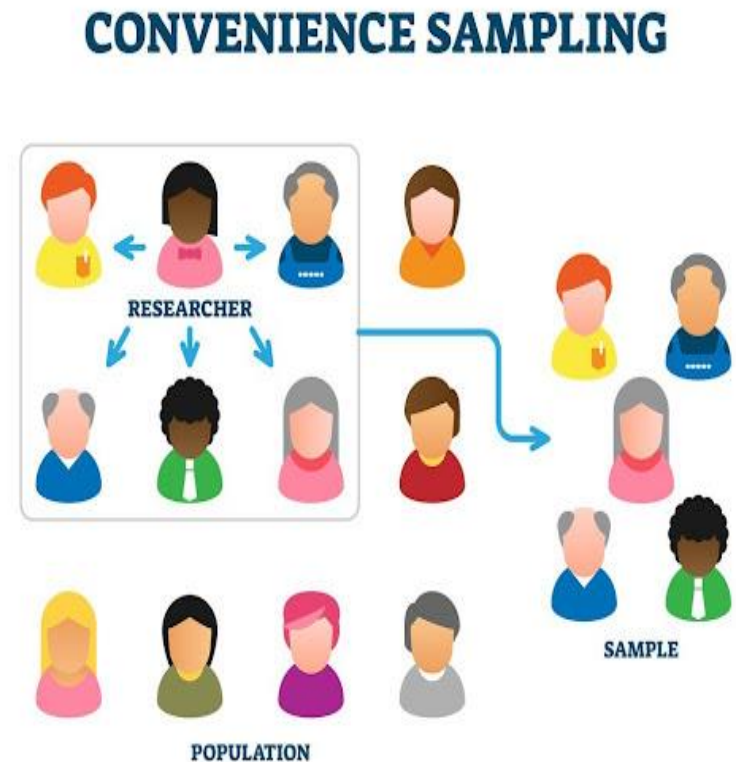
Convenience sampling

(Accidental or incidental sampling):

- People may or may not be typical of the population, no accurate way to determine their representativeness
- Most frequently used in health research

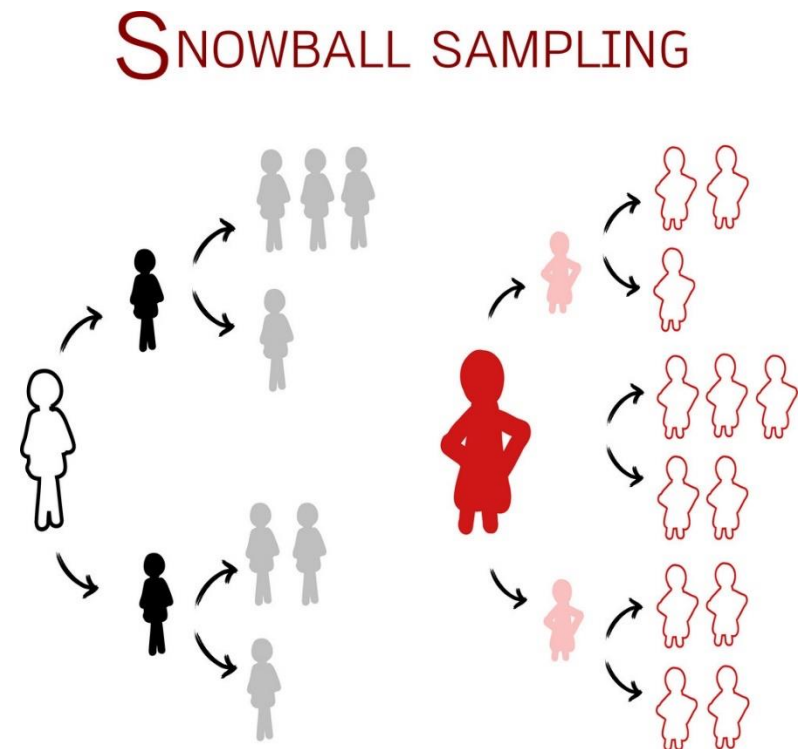
Advantages:

- Saves time and money



Snowball sampling

- A method by which the study subjects assist in obtaining other potential subjects (networking)
- Useful in topics of research where the subjects are reluctant to make their identity known, Drug users, Aids patients, etc.



Quota sampling

- In quota sampling, the sample is selected by convenience (e.g. the first 50% of males and 50% of females)
- A mean for securing potential subjects from these strata.
- In a quota sampling variables of interest to the researcher (include subject attributes), such as age, gender, educational background are included in the sample

Quota Sampling

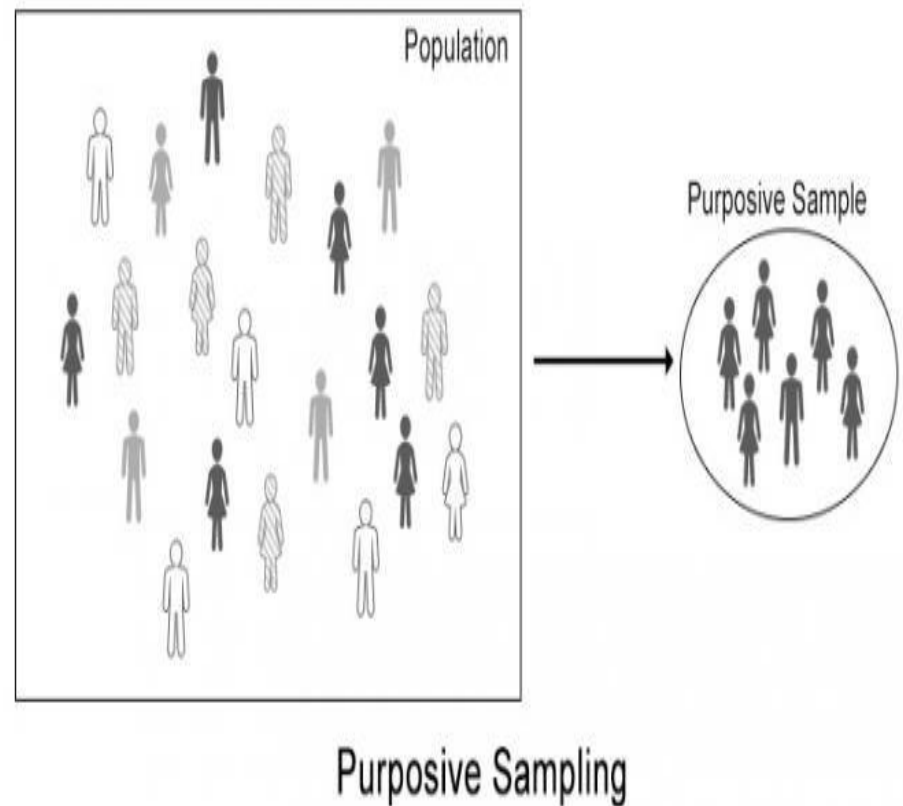
- One type of *Non-Probability Sampling*
- **IMPORTANT!** Must begin with a matrix of the target population's characteristics (i.e., % males, females, race, SES, etc).
- The set quota or % of each characteristic is fulfilled in the sampling



Purposive sampling

(handpicking, judgmental)

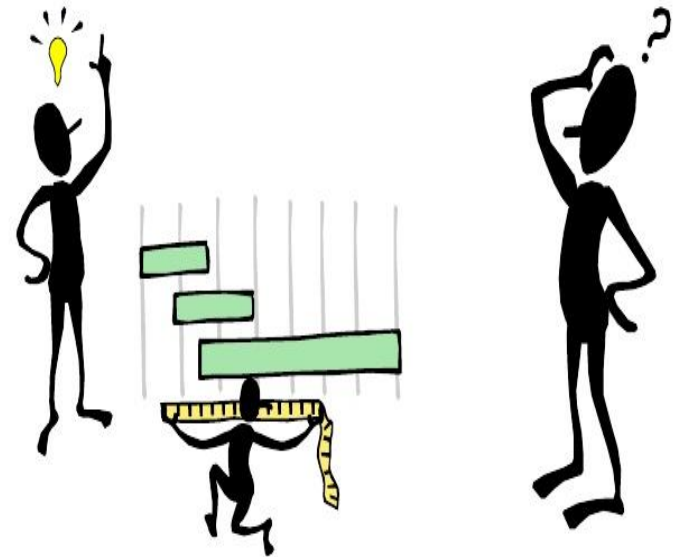
- Subjects are chosen because they are typical or representative of the accessible population, or because they are experts (more knowledgeable) in the field of research topic
- Qualitative researchers use Purposive sampling

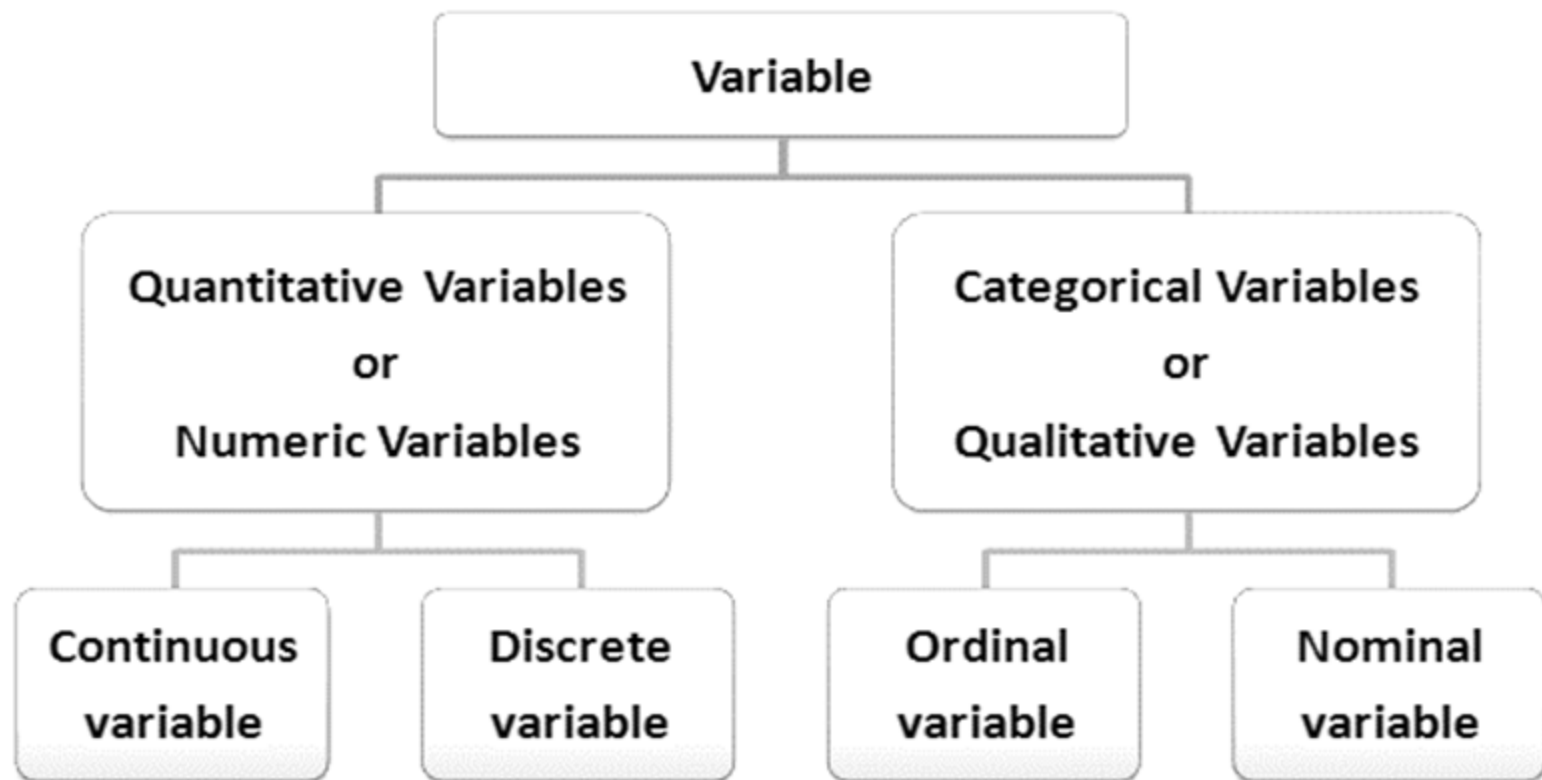


Variables

1. A variable is an object, characteristic, or property that can have different values.
2. A quantitative variable can be measured in some way.
3. A qualitative variable is characterized by its inability to be measured but it can be sorted into categories.

What is a Variable?



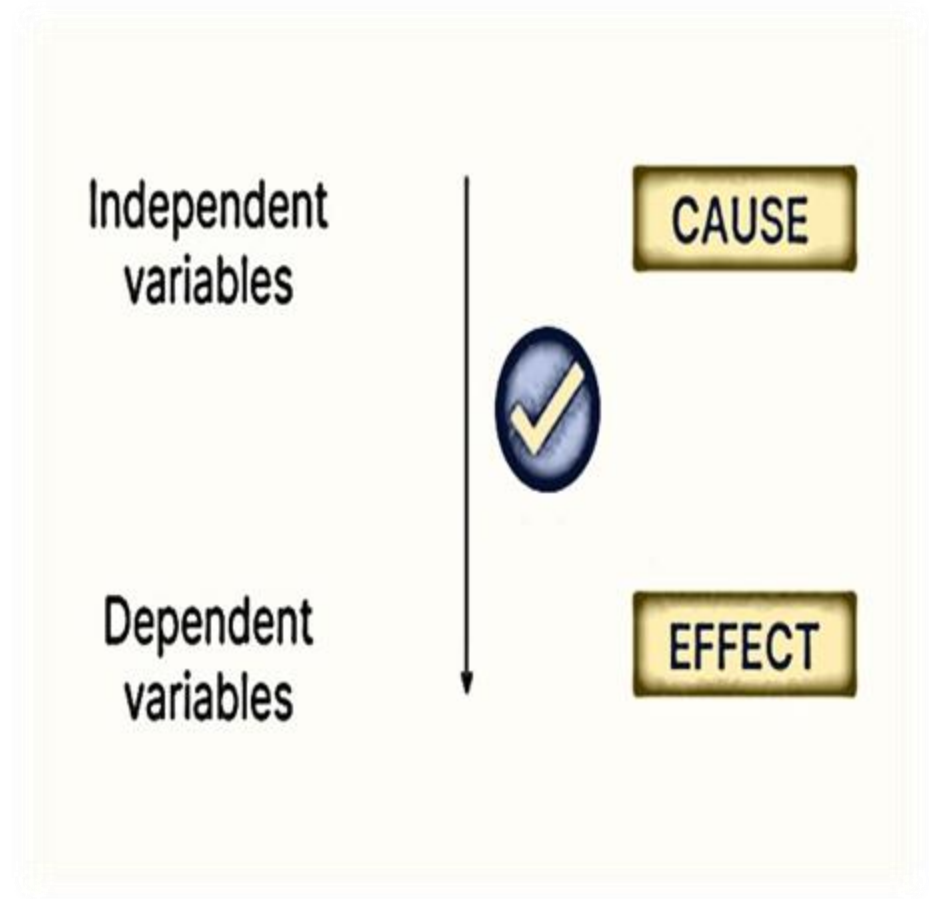


Types of Variables

Independent variable — the presumed cause (of a dependent variable)

Dependent variable — the presumed effect (of an independent variable)

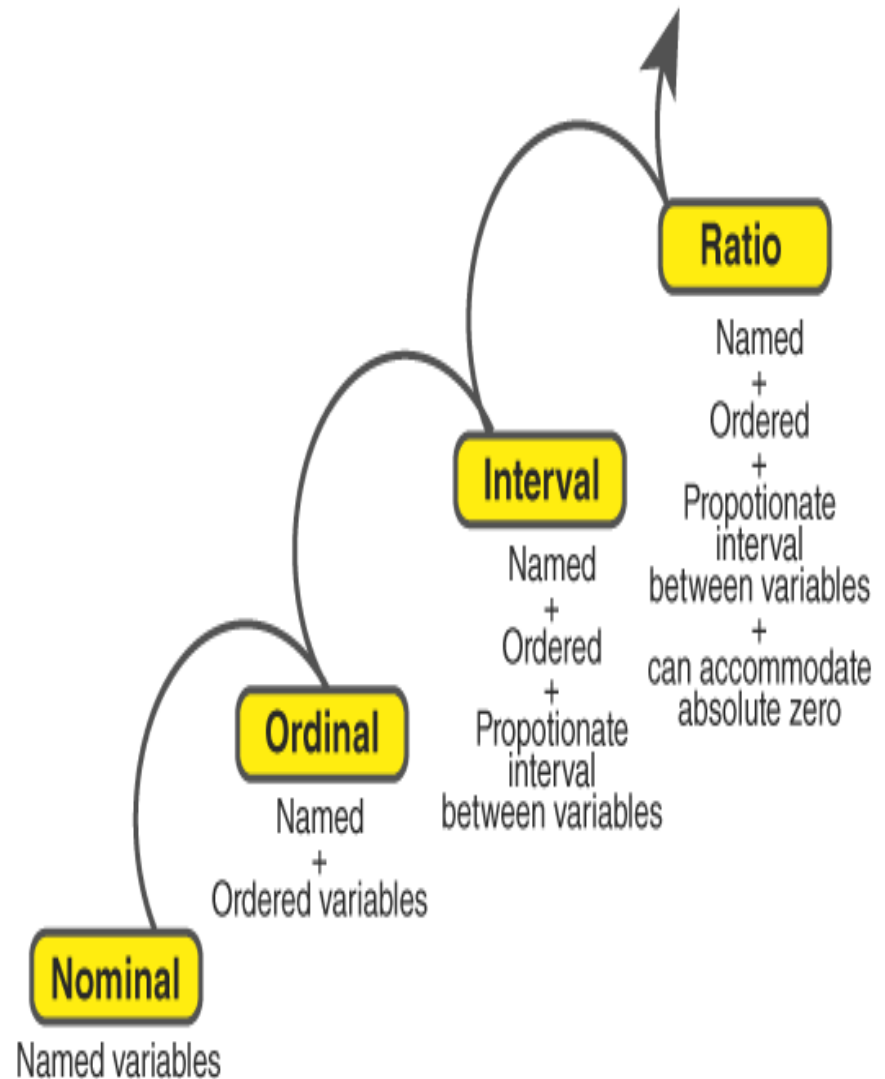
Example: Smoking (IV) → Lung cancer (DV)



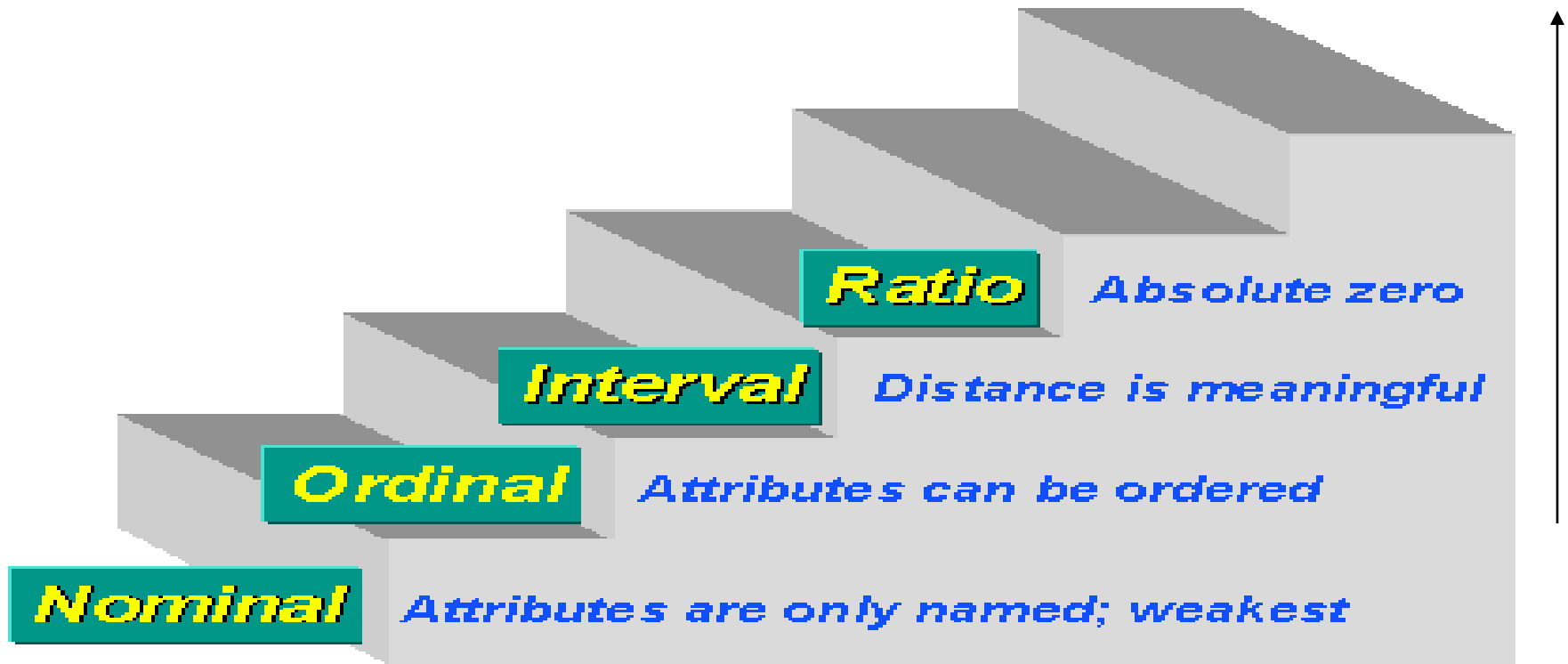
LEVELS OF MEASUREMENT

Levels of Measurement

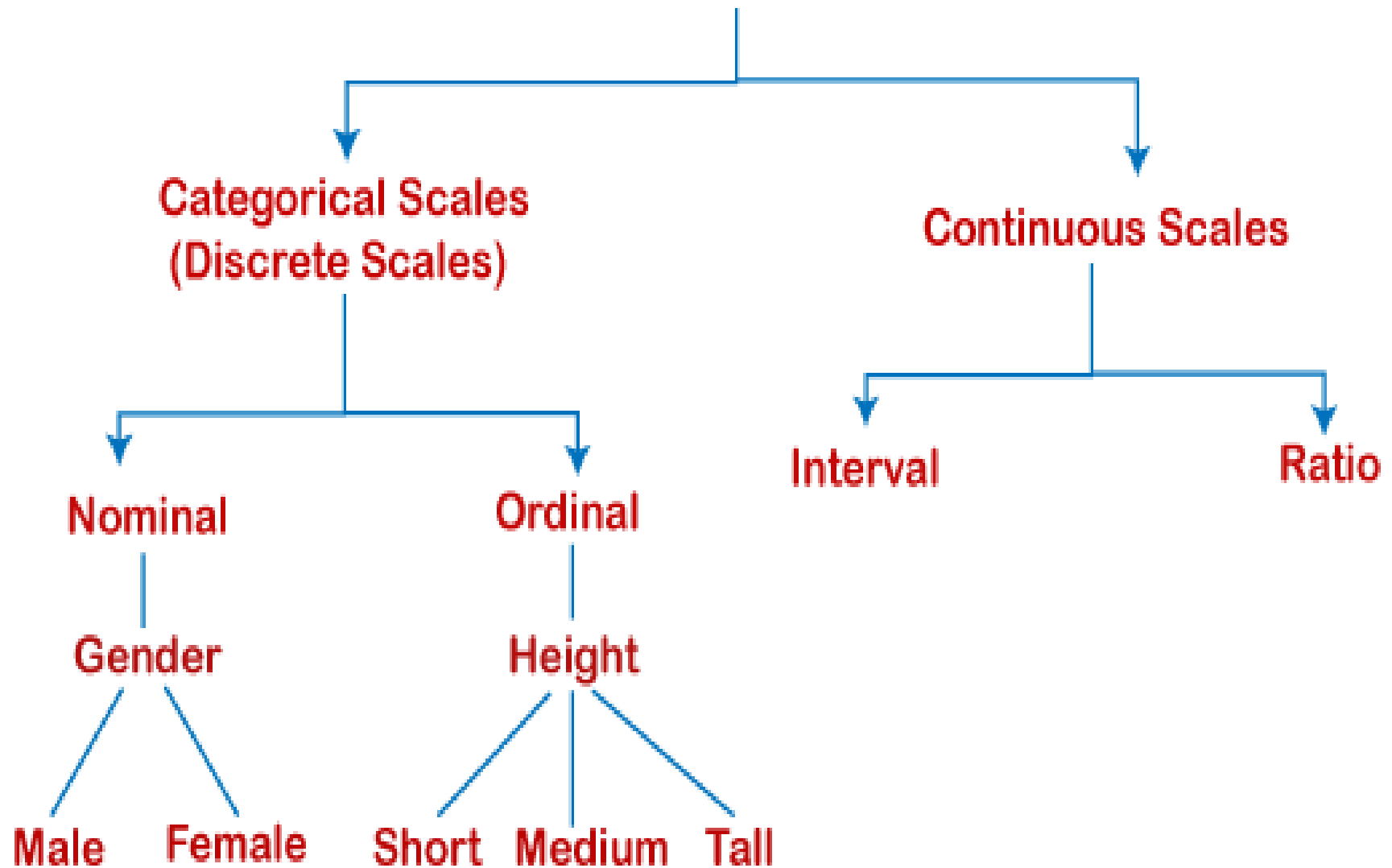
Nominal
Ordinal
Interval
Ratio



Levels of Measurement



Scales of Measurement



Nominal Level of Measurement

- Categories that are distinct from each other such as gender, religion, marital status.
- They are symbols that have no quantitative value.
- Lowest level of measurement.
- Many characteristics can be measured on a nominal scale: race, marital status, and blood type.
- Dichotomous.
- Appropriate statistics: mode, frequency
- We cannot use an average. It would be meaningless here.

| Examples of Nominal Scales | |
|--|---|
| <p>What is your gender?</p> <p><input checked="" type="radio"/> M- Male</p> <p><input type="radio"/> F- Female</p> | <p>What is your hair colour?</p> <p><input checked="" type="radio"/> 1- Brown</p> <p><input type="radio"/> 2- Black</p> <p><input type="radio"/> 3- Blonde</p> <p><input type="radio"/> 4- Gray</p> <p><input type="radio"/> 5- Other</p> |

Examples of Nominal Data



Ordinal Level of Measurement

➤ The exact differences between the ranks cannot be specified such as it indicates order rather than exact quantity.

➤ Involves using numbers to designate ordering on an attribute.

How do you feel today?

- ☒ 1 - Very Unhappy
- ☐ 2 - Unhappy
- ☐ 3 - OK
- ☐ 4 - Happy
- ☐ 5 - Very Happy

How satisfied are you with our service?

- ☒ 1 - Very Unsatisfied
- ☐ 2 - Somewhat Unsatisfied
- ☐ 3 - Neutral
- ☐ 4 - Somewhat Satisfied
- ☐ 5 - Very Satisfied

➤ Example: anxiety level: mild, moderate, severe. Statistics used involve frequency distributions and percentages.

➤ Appropriate statistics: same as those for nominal data, plus the median; but not the mean.

Interval level of Measurement

They are real numbers and the difference between the ranks can be specified.

Equal intervals, but no “true” zero.

Involves assigning numbers that indicate both the ordering on an attribute, and the distance between score values on the attribute

They are actual numbers on a scale of measurement.

Example: body temperature on the Celsius thermometer as in 36.2, 37.2 etc. means there is a difference of 1.0 degree in body temperature.

Appropriate statistics

- same as for nominal
- same as for ordinal plus,
- the mean

Examples of Interval Scales

Example 1:

How likely are you to recommend the Santa Fe Grill to a friend?

Definitely Will Not Recommend

Definitely Will Recommend

1 2 3 4 5 6 7

Example 2:

Using a scale of 0–10, with “10” being Highly Satisfied and “0” being Not Satisfied At All, how satisfied are you with the banking services you currently receive from (read name of primary bank)?
Answer: _____

Example 3:

Please indicate how frequently you use different banking methods. For each of the banking methods listed below, circle the number that best describes the frequency you typically use each method.

| Banking Methods | Never Use | | | | | | | | | | Use Very Often | | | | | | | | | |
|------------------|-----------|---|---|---|---|---|---|---|---|---|----------------|--|--|--|--|--|--|--|--|--|
| Inside the bank | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | | | | | | |
| Drive-up window | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | | | | | | |
| 24-hour ATM | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | | | | | | |
| Debit card | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | | | | | | |
| Bank by mail | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | | | | | | |
| Bank by phone | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | | | | | | |
| Bank by Internet | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | | | | | | | | |

Ratio level of Measurement

- Is the highest level of data where data can be categorized, ranked, difference between ranks can be specified and a true or natural zero point can be identified.
- A zero point means that there is a total absence of the quantity being measured.
- All scales, whether they measure weight in kilograms or pounds, start at 0. The 0 means something and is not arbitrary (SUBJECTIVE).
- Example: total amount of money.



- Is the highest level of data where data can be categorized, ranked, difference between ranks can be specified and a true or natural zero point can be identified.
- A zero point means that there is a total absence of the quantity being measured.

Ratio level of Measurement

Question 3

Which age group do you fall into?

| | | | | | |
|---------------|---|-----------------|---|---------------|---|
| 15 - 25 years | A | 26 - 35 years | B | 36 - 50 years | C |
| 51 - 65 years | D | Greater than 66 | E | | |

- All scales, whether they measure weight in kilograms or pounds, start at 0. The 0 means something and is not arbitrary (SUBJECTIVE).

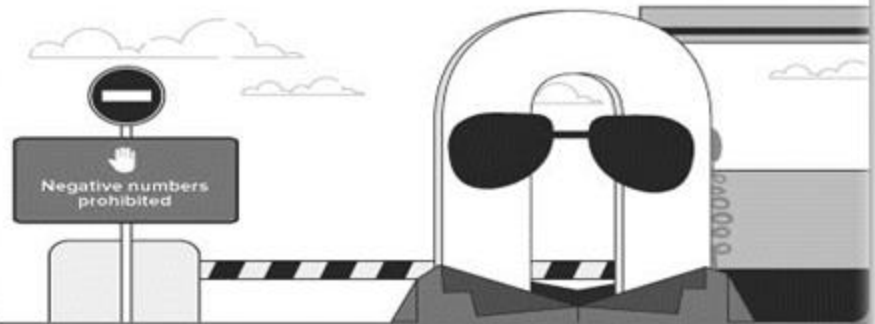
Example: total amount of money.

Characteristics of Ratio scale

Has an absolute zero characteristic. It has orders and equally distanced value between units.



Ratio scale data does not have any negative numerical value. For example, weight cannot be negative, -20 Kgs does not exist.



The data values of a ratio scale can be added, subtracted, multiplied and divided. All statistical analysis can be calculated using ratio scale.



It has ratio scale units which have several unique and useful properties. One of them is they allow unit conversion.



Examples of primary scales of measurement

Scale Nominal

Numbers
Assigned
to Runners



Finish

Ordinal

Rank Order
of Winners



Third
place



Second
place



First
place

Finish

Interval

Performance
Rating on a

8.2

9.1

9.6

0 to 10 Scale

Ratio

Time to
Finish, in

15.2

14.1

13.4

What Type of Data To collect?

The goal of the researcher is to use the highest level of measurement possible.

(A) Is nominal, so the best we can get from this data are frequencies.

- **Example: Two ways of asking about Smoking behavior. Which is better, A or B?**

(B) is ratio, so we can compute: mean, median, mode, frequencies.

(A) Do you smoke? ☐ Yes ☐ No

(B) How many cigarettes did you smoke in the last 3 days (72 hours)?

Parameter and Statistic

STATISTICS

Statistic is a descriptive measure computed from the data of the sample.

- For example, the sample mean, \bar{x} , and the standard deviation, s , are statistics.
- They are used to estimate the population parameters.

PARAMETERS

Parameter is a descriptive measure computed from the data of the population.

Parameter **VS** Statistic

When studying statistics, you might come across many terms that, if you aren't fully concentrated, will create a lot of confusion. For example, take parameter and statistic.

DEFINITION

If you are talking about a **PARAMETER**, you are talking about the **whole population**.

DEFINITION

A **STATISTIC** describes only **a sample of the population**.

Examples of Parameters

Population of
interest

Adults in
the USA

Parameter

% that are married

average age

how many are older
than 70 years old

Schools in
Africa

avg number of students

proportion that high
schools

total number of teachers

Parameters

- The population mean, μ , and the population standard deviation, σ , are two examples of population parameters.
- If you want to determine the population parameters, you have to take a census of the entire population.
- Taking a census is very costly.
- Parameters are numerical descriptive measures corresponding to populations.
- Since the population is not actually observed, the parameters are considered unknown constants.

Statistics: It is a branch of applied mathematics that deals with collecting, organizing, & interpreting data using well-defined procedures in order to make decisions.

The term **parameter** is used when describing the characteristics of the population.

The term **statistics** is used to describe the characteristics of the sample.

Types of Statistics:

- **Descriptive Statistics.** It involves organizing, summarizing & displaying data to make them more understandable.
- **Inferential Statistics.** It reports the degree of confidence of the sample statistic that predicts the value of the population parameter

Descriptive Statistics

Measures of Location

- Measures of Central Tendency:
 - Mean
 - Median
 - Mode
- Measures of noncentral Tendency-Quantiles:
 - Quartiles.
 - Quintiles.
 - Percentiles.

Measure of Dispersion (Variability):

- Range
- Interquartile range
- Variance
- Standard Deviation
- Coefficient of variation

Measures of Shape:

- Mean > Median-positive or right Skewness
- Mean = Median- symmetric or zero Skewness
- Mean < Median-Negative of left Skewness

Statistical Inference

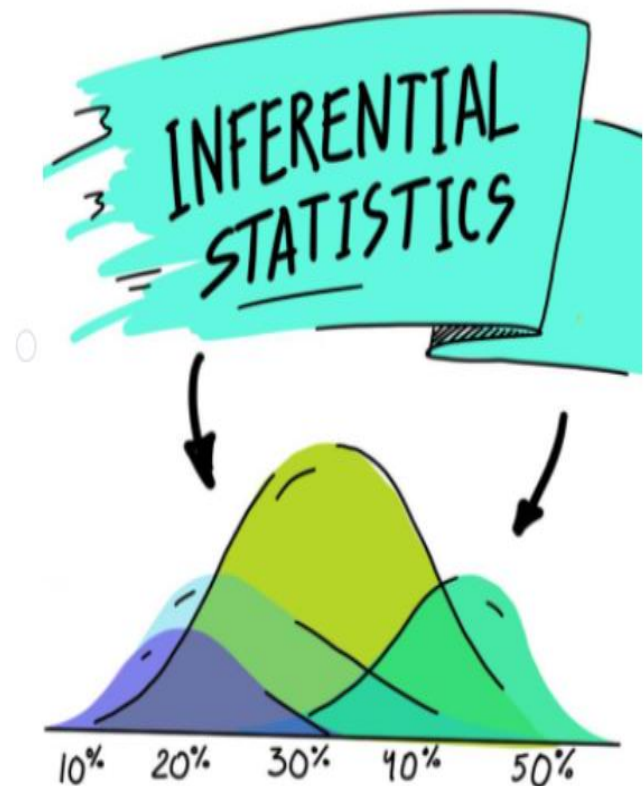
Is the procedure used to reach a conclusion about a population based on the information derived from a sample that has been drawn from that population.

Inferential Statistics

Are used to test hypotheses (prediction) about relationship between variables in the population. A relationship is a bond or association between variables.

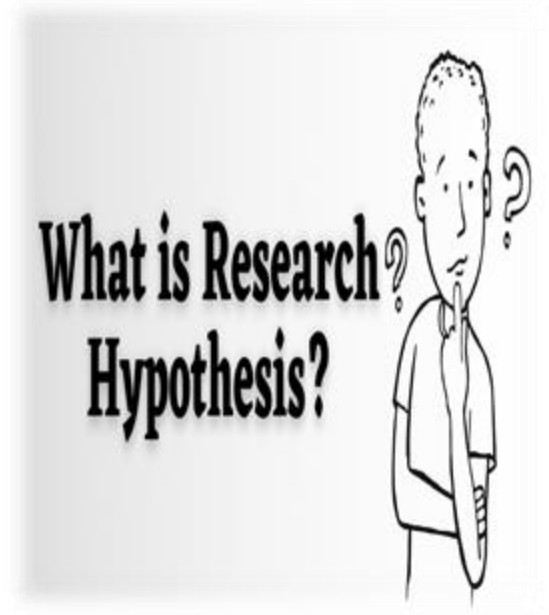
It consists of a set of statistical techniques that provide prediction about population characteristics based on information in a sample from population. An important aspect of statistical inference involves reporting the likely accuracy, or of confidence of the sample statistic that predicts the value of the population parameter.

- **Bivariate Parametric Tests:**
- One Sample t test (t)
- Two Sample t test (t)
- Analysis of Variance/ANOVA (F).
- Pearson's Product Moment Correlations (r).
- **Nonparametric statistical tests: Nominal Data:**
- Chi-Square Goodness-of-Fit Test
- Chi-Square Test of Independence
- **Nonparametric statistical tests: Ordinal Data:**
- Mann Whitney U Test (U)
- Kruskal Wallis Test (H)



Research Hypothesis

- A tentative prediction or explanation of the relationship between two or more variables.
- It's a translation of research question into a precise prediction of the expected outcomes
- In some way it's a proposal for solution/s
- In qualitative research, there is NO hypothesis

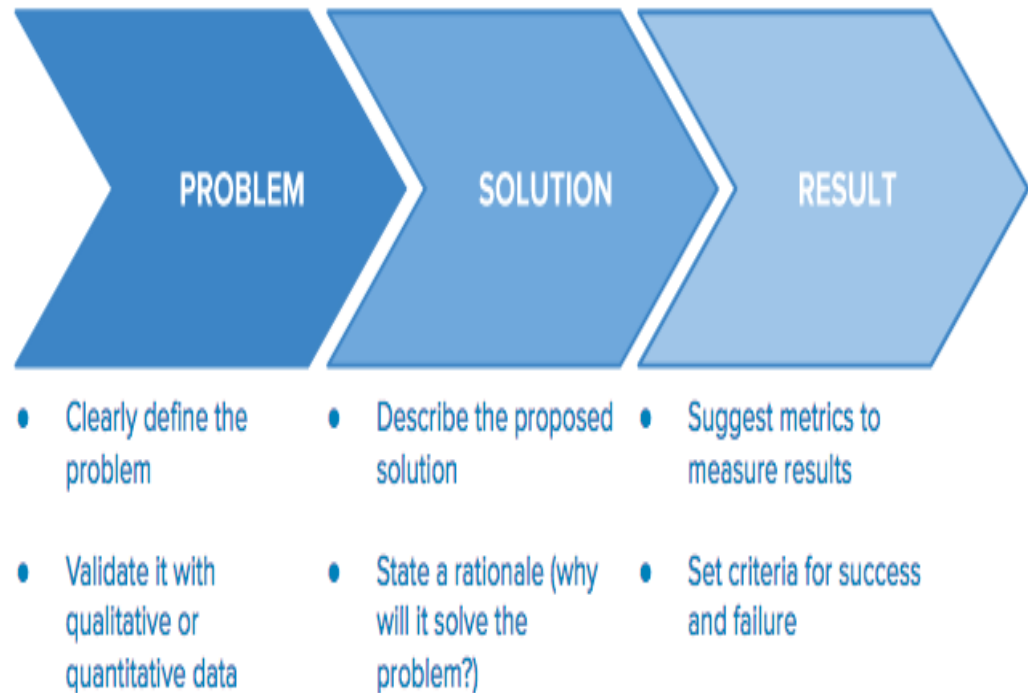


Research Hypothesis

- States a prediction
- Must always involve at least two variables
- Must suggest a predicted relationship between the independent variable and the dependent variable
- Must contain terms that indicate a relationship (e.g., more than, different from, associated with)

Hypotheses Criteria

- Written in a declarative form.
- Written in present tense.
- Contain the population
- Contain variables.
- Reflects problem statement or purpose statement.
- Empirically testable.



Hypothesis Testing

A hypothesis is made about the value of a parameter, but the only facts available to estimate the true parameter are those provided by the sample. If the statistic differs (and of course it will) from the hypothesis stated about the parameter, a decision must be made as to whether or not this difference is *significant*. If it is, the hypothesis is rejected. If not, it cannot be rejected.

H_0 : The null hypothesis. This contains the hypothesized parameter value which will be compared with the sample value.

H_1 : The alternative hypothesis. This will be “accepted” only if H_0 is rejected.

Technically speaking, we never accept H_0 . What we actually say is that we do not have the evidence to reject it.

Two Types of Errors: Alpha and Beta

Two types of errors may occur: α (alpha) and β (beta). The α error is often referred to as a Type I error and β error as a Type II error.

- You are guilty of an alpha error if you reject H_0 when it really is true.
- You commit a beta error if you “accept” H_0 when it is false.

| | | STATE OF NATURE | |
|----------|---------------------|----------------------------------|----------------------------------|
| | | H_0 Is True | H_0 Is False |
| DECISION | Do Not Reject H_0 | GOOD | β Error (Type II Error) |
| | Reject H_0 | α Error (Type I Error) | GOOD |

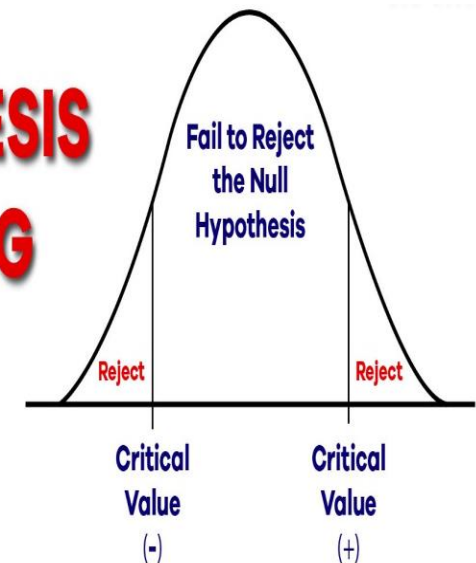
1. Formulate H_0 and H_1 . H_0 is the null hypothesis, a hypothesis about the value of a parameter, and H_1 is an alternative hypothesis.

e.g., $H_0: \mu = 12.7$ years; $H_1: \mu \neq 12.7$ years

2. Specify the level of significance (α) to be used. This level of significance tells you the probability of rejecting H_0 when it is, in fact, true. (Normally, significance level of 0.05 or 0.01 are used)
3. Select the test statistic: e.g., Z, t, F, etc.
4. Establish the critical value or values of the test statistic needed to reject H_0 . DRAW A PICTURE!
5. Determine the actual value (computed value) of the test statistic.
6. Make a decision: **Reject H_0** or **Do Not Reject H_0** .

Steps in Hypothesis Testing

HYPOTHESIS TESTING



End of Unit 1

