

Biostatistics

Nuno Sepúlveda, 06.10.2025

About myself

BSc and MSc in Applied Mathematics (Statistics)

PhD in Biomedical Sciences

@Gulbenkian Institute for Science (2001-2009, Portugal)

Theoretical Immunology Group / Quantitative Biology Group

@London School of Hygiene and Tropical Medicine (2010-2019, United Kingdom)

Research Fellow in Statistical Genetics and Genetic Epidemiology

Assistant Professor in Biostatistics and Statistical Genetics

@Charité Medical University of Berlin (2020-2021, Germany)

Consultant in Bioinformatics and Biostatistics

@Politechnika Warszawa (2021-Current, Poland)

Visiting Professor (ULAM Programme - NAWA)

Assistant Profesor

My research

Tell me about yourself

Syllabus

1. General review

- a. Population/Sample/Sample size
- b. Type of Data – quantitative and qualitative variables
- c. Common probability distributions/popular tests

2. Applications in Medicine

- a. Construction and analysis of diagnostic tools – Binomial distribution, ROC curve, sensitivity, specificity, Rogal-Gladen estimator
- b. Estimation of treatment effects - generalized linear models
- c. Survival analysis - Kaplan-Meier curve, log-rank test, Cox's proportional hazards model

3. Applications in Genetic and Epigenetic Data

- a. Genetic association studies – Hardy-Weinberg test, homozygosity, minor allele frequencies, additive model, multiple testing correction
- b. Methylation association studies – M versus beta values, estimation of biological age

4. Applications in Serological Data Analysis

- a. Determination of seropositivity using Gaussian mixture models
- b. Reversible catalytic models for estimating seroconversion rate
- c. Sample size calculation for estimating seroconversion rate

Course material

https://github.com/immune-stats/Biostatistics_2025_2026/

Software



Version 4.5.1

R Studio/Posit

Communication

nuno.sepulveda@pw.edu.pl

Two people should be chosen as the main contact points with me

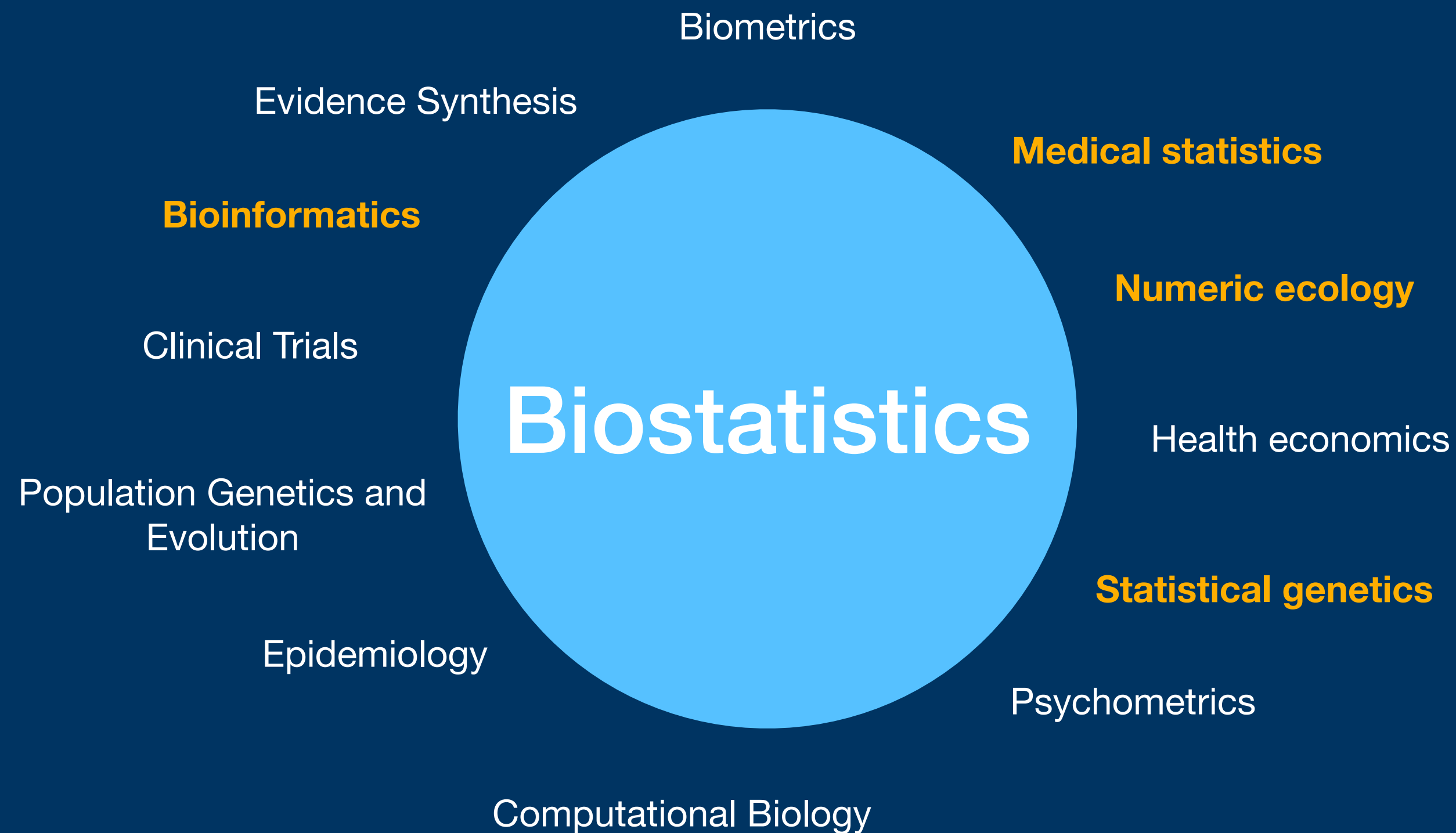
Evaluation

Group Project + Presentation (40%)

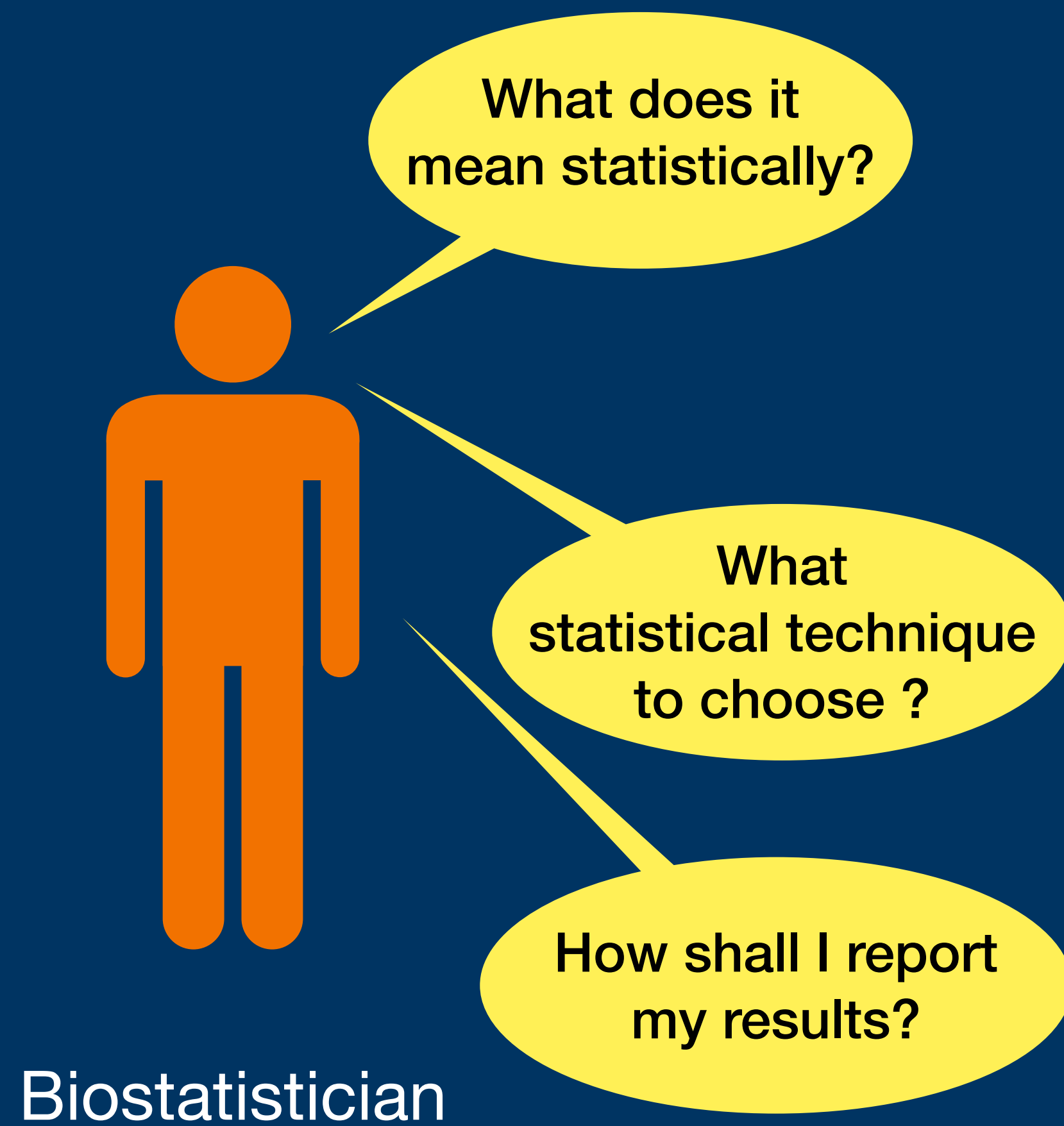
Oral examination (60%)

Biostatistics

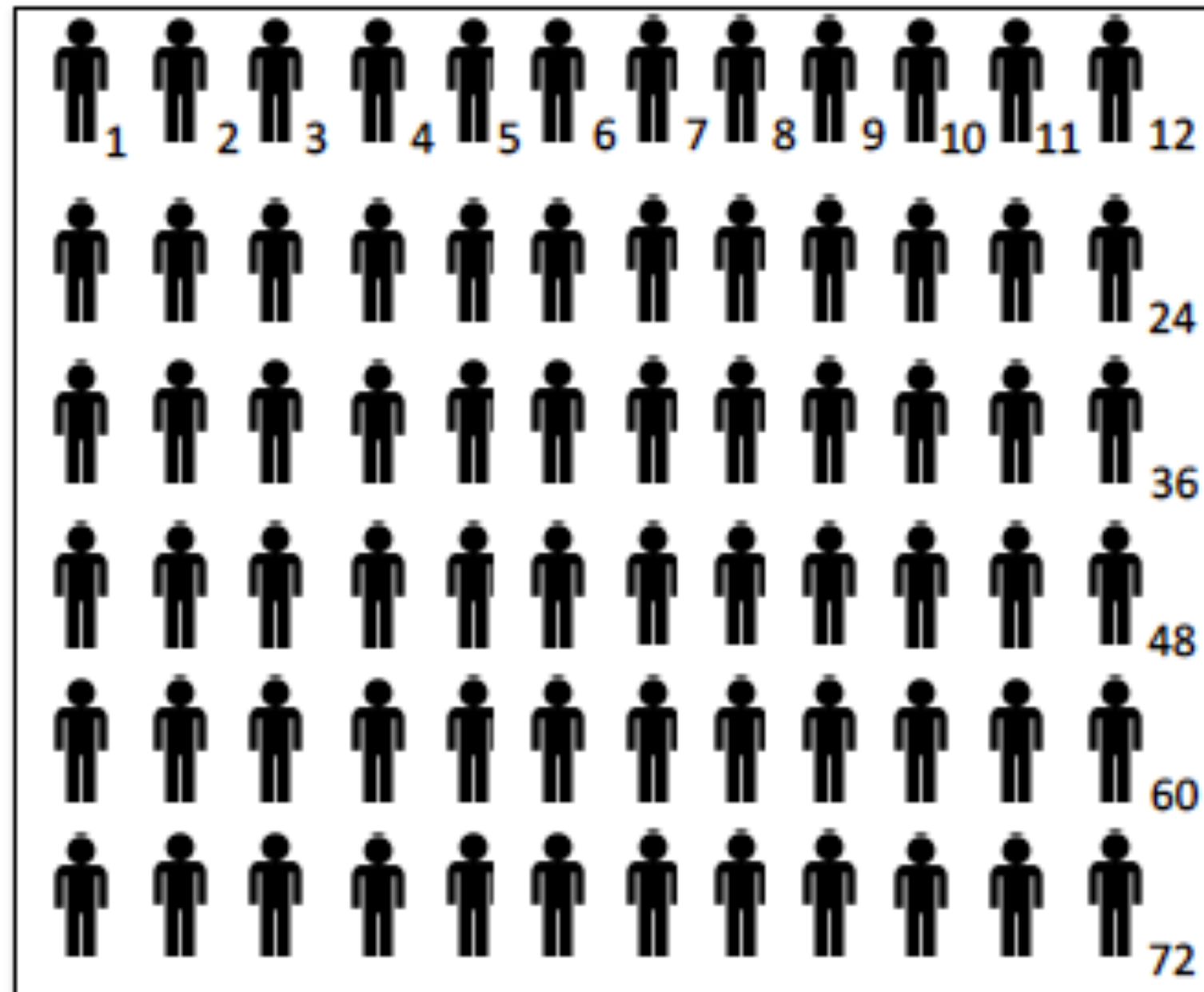
Application of statistical techniques to scientific research in health-related fields, including medicine, biology, and public health, and the development of new tools to study these areas.



Importance of communication



Population



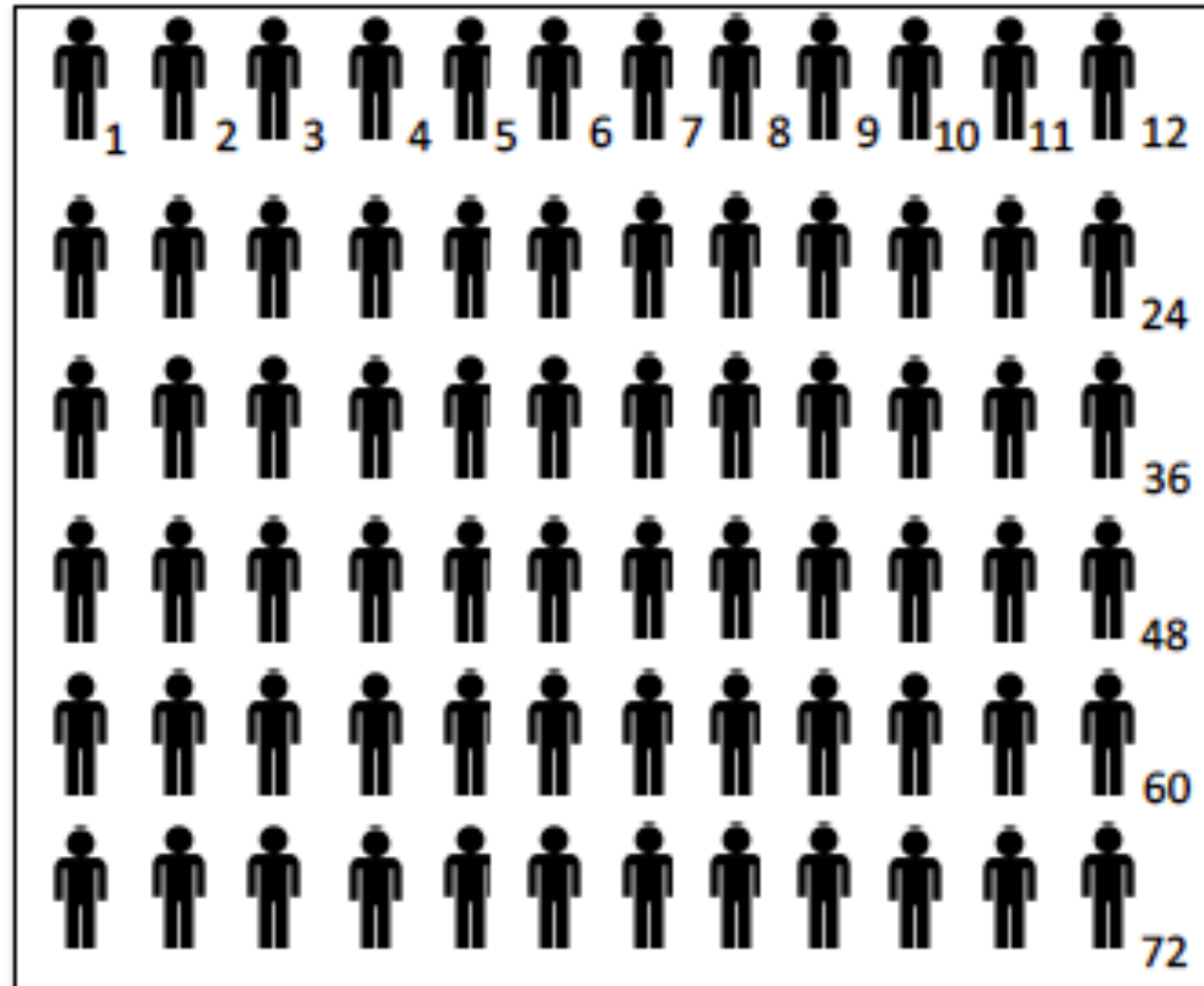
Population:

The complete set of individuals from which you want to learn something.

Population Size:

The total number of individuals in the population.

Census



Census:

It is a study conducted in the entire population.
It might require a large set of resources.

Example:


Data collected by the National Office for
Statistics.




Study of a rare disease.

Sample



 Sampled individuals

 Not-sampled individuals

Sample:

A set of individuals which it is thought to be **representative** of the whole population.

Sample size:

The total number of individuals included in the sample.

Randomisation



Sampled individuals



Not-sampled individuals

Randomized sample:


Individuals should be **randomly** selected from the population.


Exercise:

Can you randomly select a new sample of 5 individuals from this population?

I. Unstratified sampling



 Sampled individuals

 Not-sampled individuals

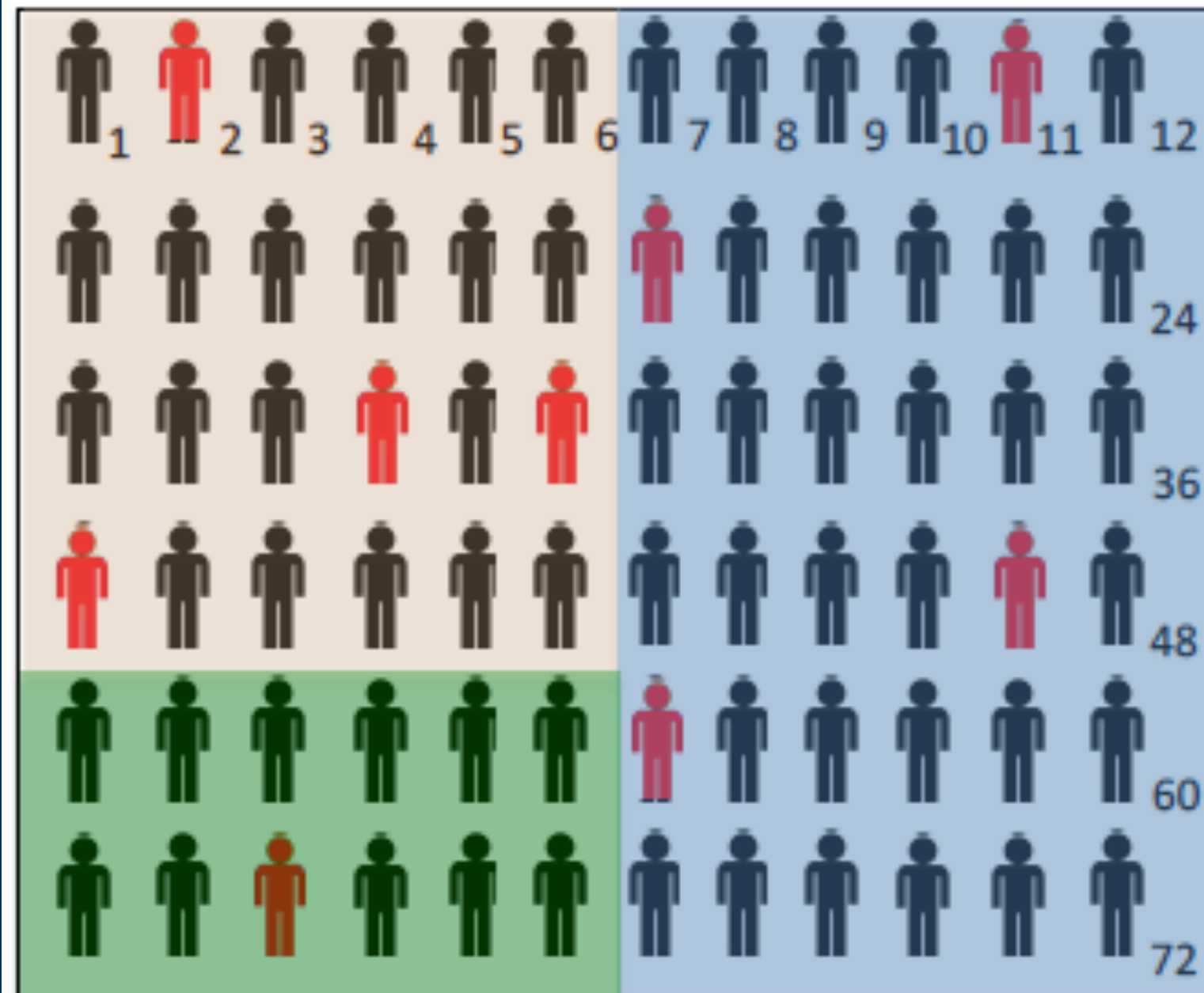
Unstratified sampling


Individuals should be **randomly** selected from the population.


It is chosen when little information is known about the population under study.

II. stratified sampling

Population is composed of distinct subpopulations (stratum) with similar pattern of response.



 Sampled individuals

 Not-sampled individuals

Stratified sampling:

Individuals should be **randomly** selected from each stratum of the population.

The total number of sampled individuals per stratum should scale with the respective total number in the population.

Main advantage:

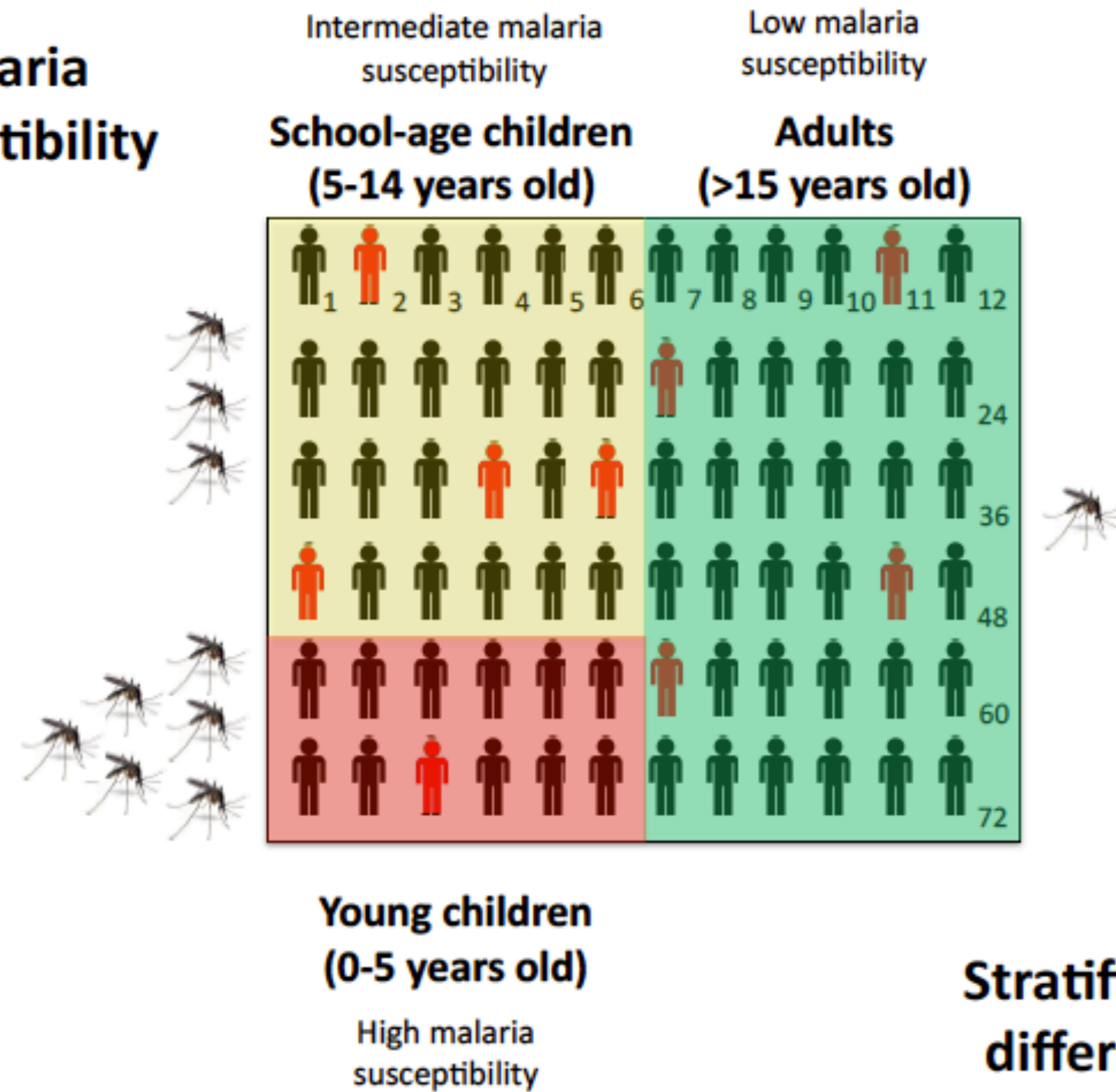
Increase precision on estimates.

Avoids confounding and/or controls confounder effects.

Requirement:

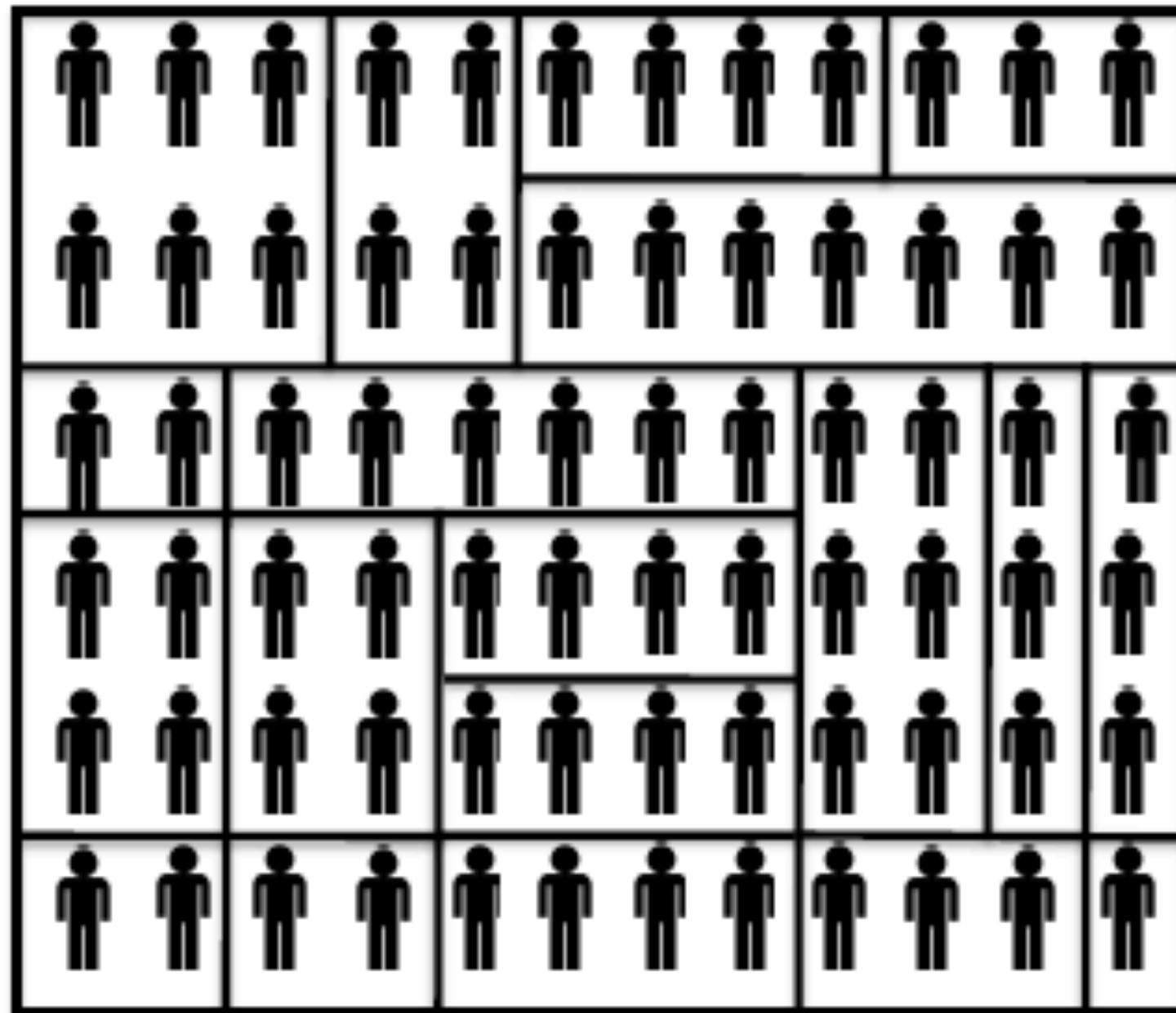
You need to know how to define the strata!

Malaria susceptibility



Stratified sampling by different age groups

III. Cluster sampling



Sampling by clusters:

Population is divided into different clusters.

Select clusters randomly.

All individuals within a cluster are measured.

Example of clusters:

Household or compound.

Common sampling strategies: cluster sampling



Sampling by clusters:

Population is divided into different clusters.

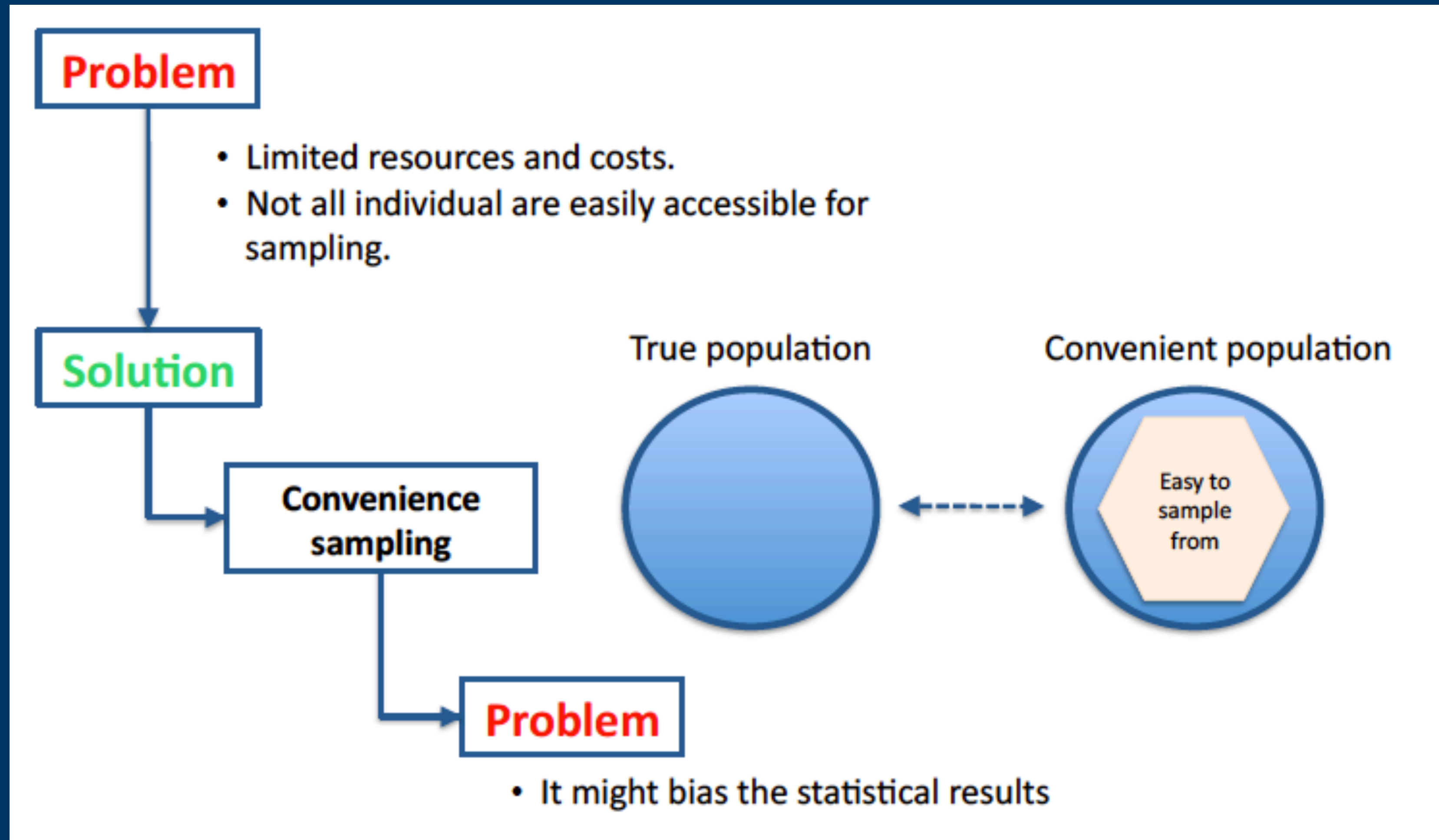
Select clusters randomly.

All individuals within a cluster are measured.

Example of clusters:

Household or compound.

In practice



Type of variables

Quantitative

Continuous – measurements with virtual infinite precision

Ex: height, weight, time until cure, etc.

Discrete – count data

Ex: number of infection episodes per person, number of treatment doses per patient.

Qualitative

Binary – two categories

Ex: cured/not cured, presence/absence, wild type/mutated allele, etc.

Polytomous – Many categories

Ex: eye colour, genotype, socioeconomic status, ethnicity, etc.

Other types of variables

Images

Medical imaging

Videos

Symbolic/distribucional data

Two cautionary notes

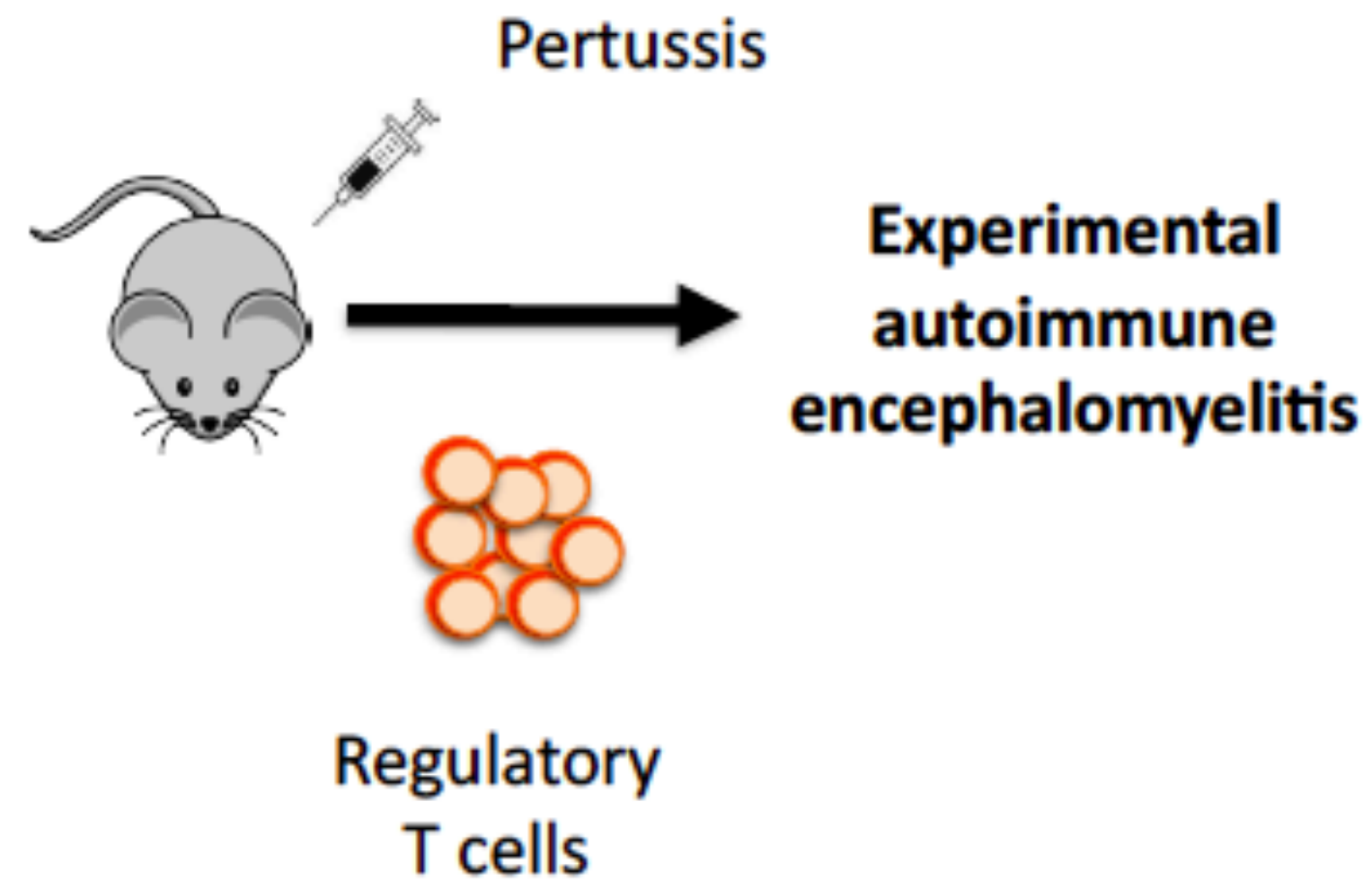
1. The true type of a variable is dependent on the unit of analysis

Define the unit of analysis as a function of the objective

2. Qualitative variables might be “hidden” in apparently quantitative variables

Always read first the data dictionary before doing any analyses

Example from the literature: EAE score



Quintana et al (2008). Nature

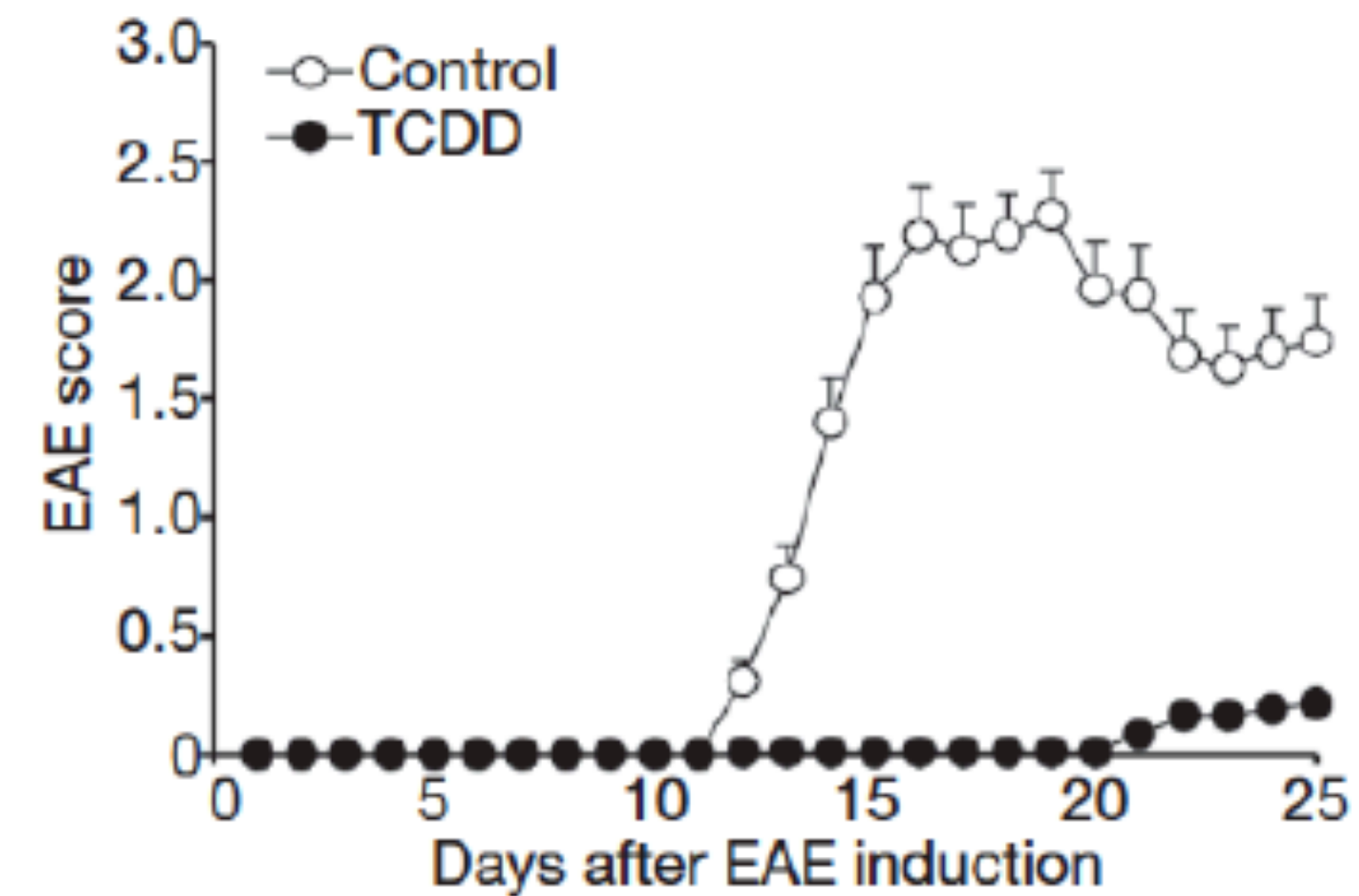
Vol 452 | 1 May 2008 | 441-452 | www.nature.com

nature

ARTICLES

Control of T_{reg} and T_H17 cell differentiation by the aryl hydrocarbon receptor

Francisco J. Quintana¹, Alexandre S. Basso¹, Antonio H. Iglesias¹, Thomas Korn¹, Mauricio F. Farez², Estelle Bettelli¹, Mario Caccamo², Mohamed Galka² & Howard L. Weiner¹



(Sigma-Aldrich) intraperitoneally on days 0 and 2. Clinical signs of EAE were assessed according to the following score: 0, no signs of disease; 1, loss of tone in the tail; 2, hindlimb paresis; 3, hindlimb paralysis; 4, tetraplegia; 5, moribund.

Summarising data

Summary statistics

Maximum
Minimum
Mean
Median
Quantiles
Quartiles
Mode
Proportion
Frequencies

Standard deviation
Variance
Variation coefficient
Interquartile range

**When to use each
one of these
summary tools?**

Visualization tools

Boxplots
Scatterplots
Histograms
ECDF plots
Density plots
Strip plots

Barplots
Piecharts

Heatmaps

Common probability distributions

Which distributions do you know?

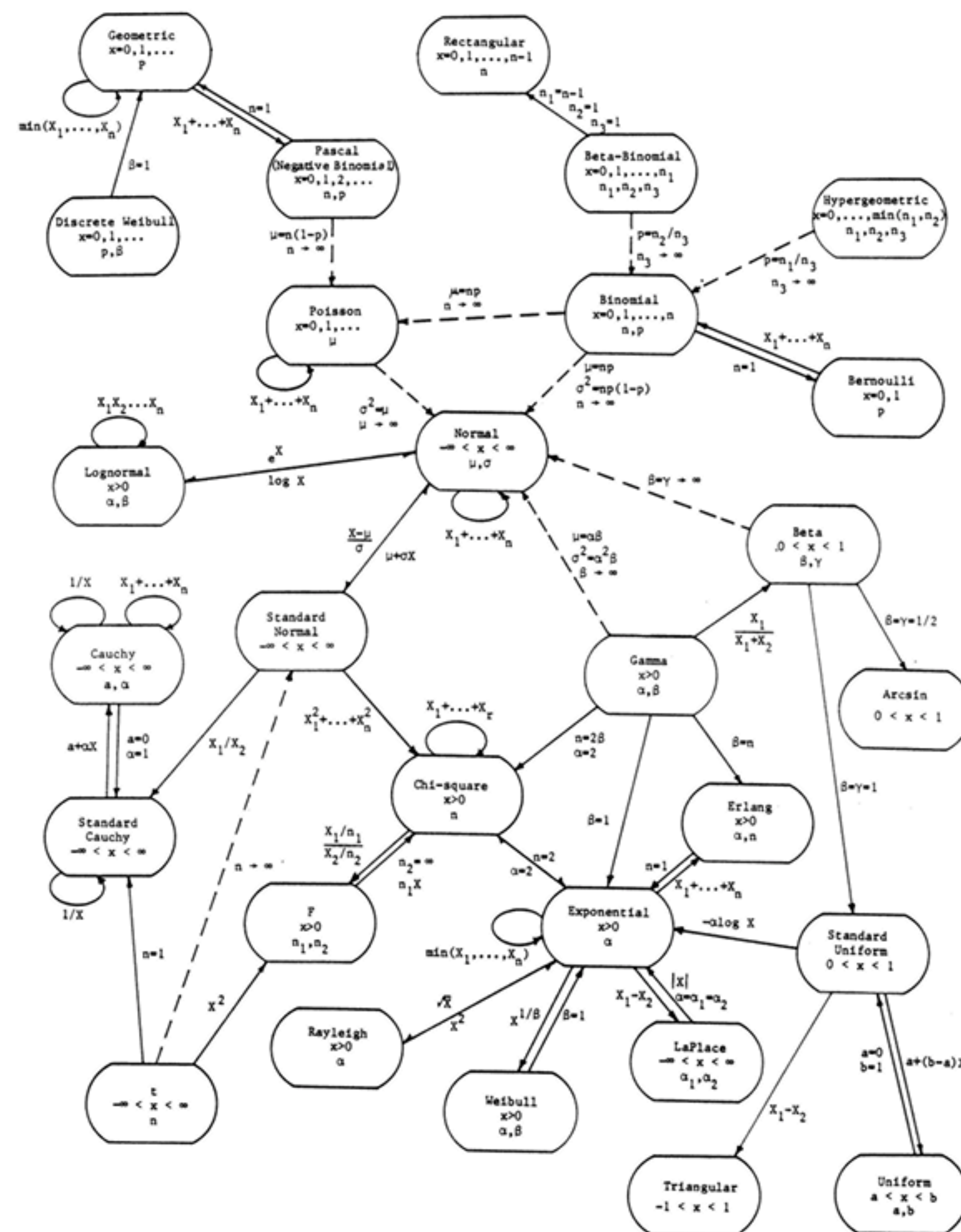


Figure 1. Relationships Among Distributions.

Leemis (1986). The American Statistician, 40 , 143

Statistical tests

Which statistical tests do you remember?

Warm-up exercise

Estimating medium- and long-term trends in malaria transmission by using serological markers of malaria exposure

C. J. Drakeley^{*†‡}, P. H. Corran^{*‡§}, P. G. Coleman^{*}, J. E. Tongren^{*}, S. L. R. McDonald^{*}, I. Carneiro^{*}, R. Malima^{†¶}, J. Lusingu^{†¶}, A. Manjurano^{†¶}, W. M. M. Nkya^{†¶}, M. M. Lemnge^{†¶}, J. Cox^{*}, H. Reyburn^{*†}, and E. M. Riley^{*.***}

^{*}Department of Infectious and Tropical Diseases, London School of Hygiene and Tropical Medicine, Keppel Street, London WC1E 7HT, United Kingdom; [†]Joint Malaria Programme, P.O. Box 2228, Moshi, Tanzania; [§]National Institute for Biological Standards and Control, South Mimms EN6 3QG, United Kingdom; [¶]Kilimanjaro Christian Medical Centre, P.O. Box 3010, Moshi, Tanzania; and [¶]Amani Medical Research Institute, National Institute for Medical Research, P.O. Box 4, Amani, Tanzania

Edited by Louis H. Miller, National Institutes of Health, Rockville, MD, and approved February 23, 2005 (received for review November 23, 2004)

Warm-up exercise

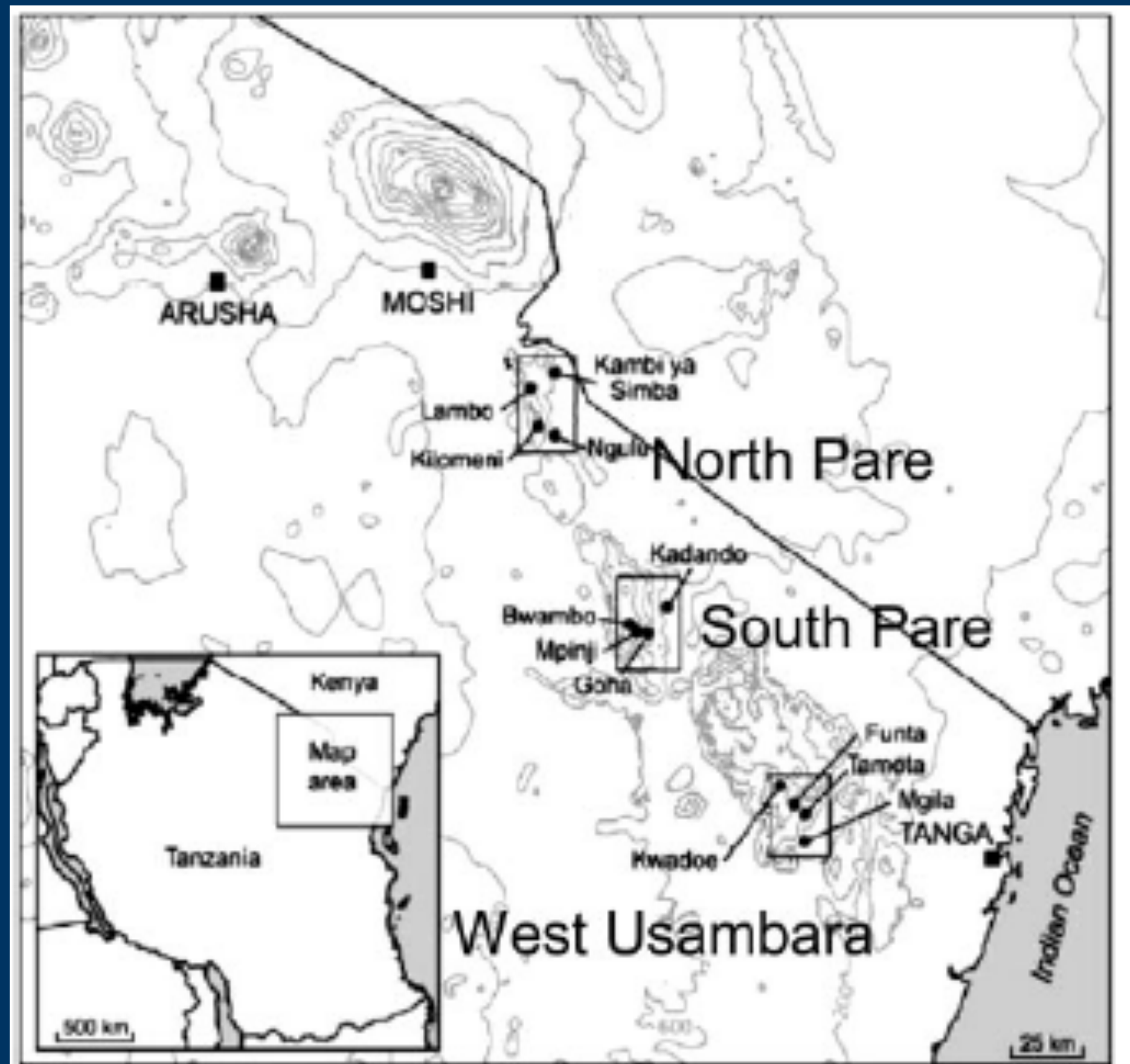


Fig. 1. Map of the study area showing the three altitude transects and 12 study villages.

Cross-sectional study

Stratified sampling (three age groups: 0-4, 5-14, 15-45)

24 villages in 6 altitude transects

~8146 individuals (6 months-45 years old)

Gender and age distributions matched across villages

Github: [Data/data_tanzania.csv](#)

Warm-up exercise



Can you check whether
gender and age are indeed matched in these
villages?

Researcher

Warm-up exercise



Can you check whether proportion of infection varies with village?

Researcher

Warm-up exercise



Can you check whether
the proportion of infection is related to
altitude?

Researcher