

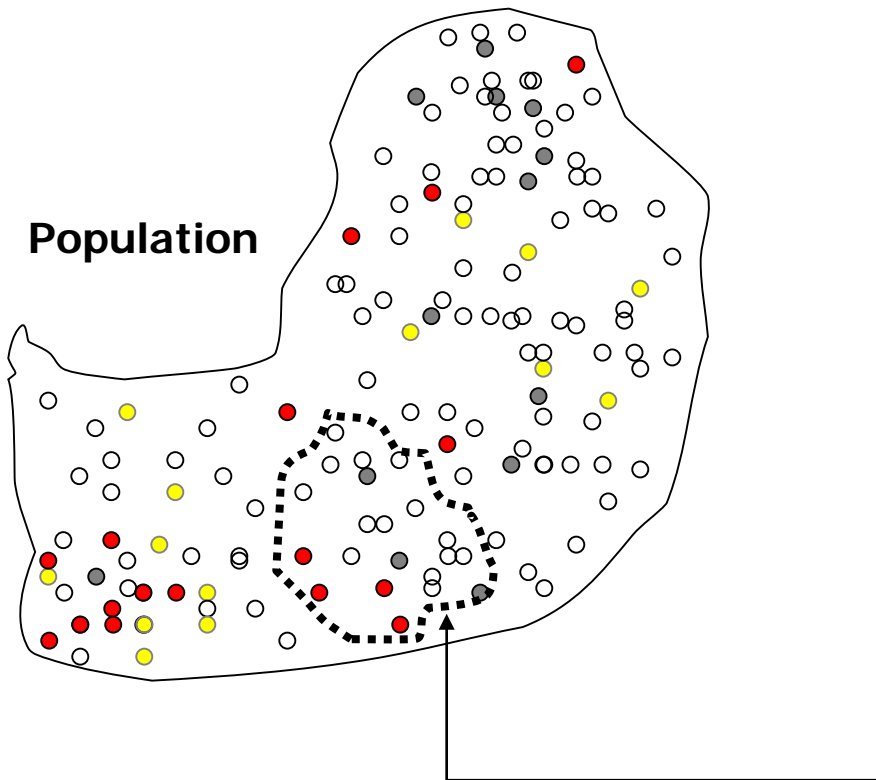
Population, Sample and Inference

Matthias RAUTERBERG

Eindhoven University of Technology

2017

Population



Definition

A **population** consists of all elements – individuals, items, or objects – whose characteristics are being studied.

The population that is being studied is also called the **target population**.

A portion of the population selected for study is referred to as a **sample**.

The random sample

Definition

A sample drawn in such a way that each element of the population has a chance of being selected is called a **random sample**.

If the chance of being selected is the *same* for each element of the population, it is called a **simple random sample**.

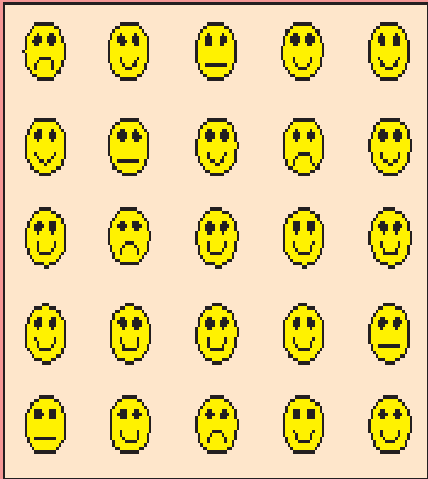

An **element** or **member** of a sample or population is a specific concept, subject or object (for example, a person, firm, item, state, or country) about which the information is collected.

A **variable** is a characteristic under study that assumes different values for different elements.

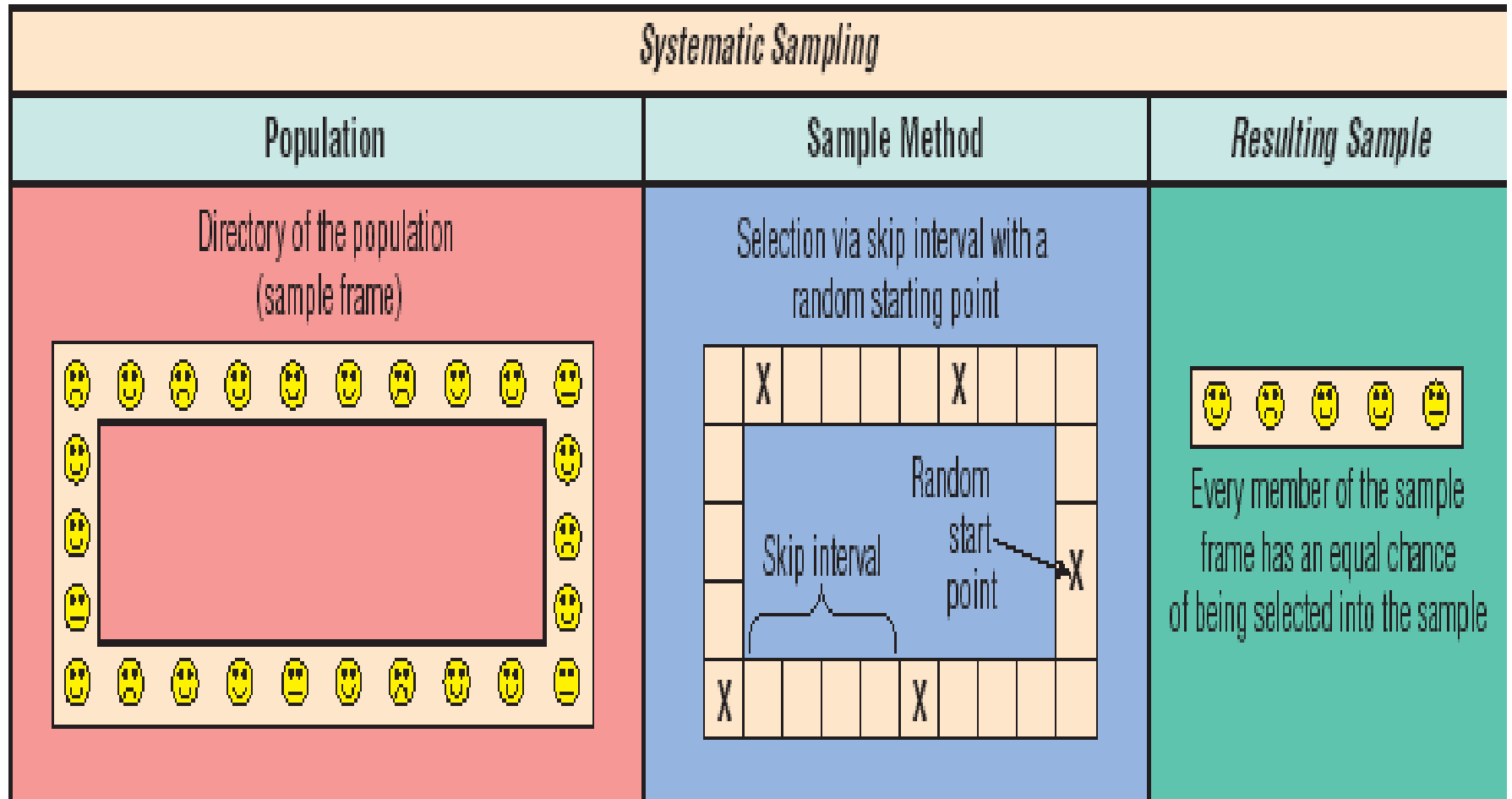
In contrast to a variable, the value of a **constant** is fixed.

The value of a variable for an element is called an **observation** or **measurement**.

Random sampling method

<i>Simple Random Sampling</i>																											
Population	Sample Method	Resulting Sample																									
<p>The population identified uniquely by number</p> 	<p>Selection by random number</p> <table border="1" data-bbox="774 715 1211 1189"><tbody><tr><td></td><td></td><td></td><td></td><td></td></tr><tr><td>X</td><td></td><td></td><td>X</td><td></td></tr><tr><td></td><td>X</td><td></td><td></td><td></td></tr><tr><td></td><td></td><td>X</td><td>X</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td></tr></tbody></table>						X			X			X						X	X							 <p>Every member of the population has an equal chance of being selected into the sample</p>
X			X																								
	X																										
		X	X																								

Systematic sampling method



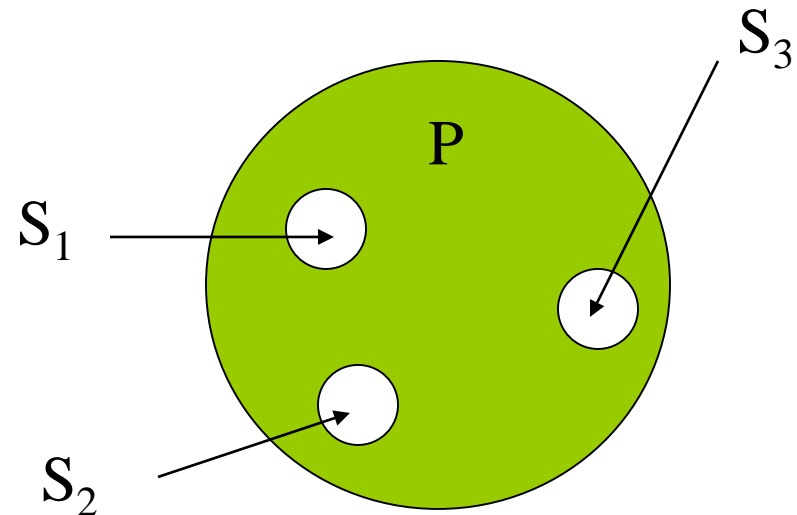
Cluster sampling method

Cluster Sampling		
Population	Sample Method	Resulting Sample
<p>The population in groups (clusters)</p> <p>A 😊 😞 😊 😊 😞</p> <p>B 😊 😞 😊 😊 😞</p> <p>C 😊 😞 😞 😊 😞</p> <p>D 😊 😞 😊 😊 😞</p> <p>E 😊 😞 😞 😊 😞</p>	<p>Random selection of 2 clusters with random selection of members of these clusters (2-stage)</p> <p>A <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/></p> <p>▼</p> <p>X x <input type="checkbox"/> x <input type="checkbox"/> <input type="checkbox"/></p> <p>C <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/></p> <p>D <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/></p> <p>▼</p> <p>X <input type="checkbox"/> x <input type="checkbox"/> x x</p>	<p>😊 😞 😊 😊 😞</p> <p>Every cluster (A, B, C, D, or E) in the population has an equal chance of being selected into the sample, and every cluster member has an equal chance of being selected from that cluster</p>

From sample to population

- Here is the problem: different samples (S_x) drawn from the same population (P) can have different properties.
- When you take a sample from a population, you only have a subset of the population--a piece of what you're trying to understand.

The **solution** to this problem is called *statistics*, in particular *inferential statistics*!



What is statistics?

Definition

Statistics is a group of methods used to collect, analyze, present, and interpret data and to make decisions.

Types of Statistics:

Descriptive Statistics consists of methods for organizing, displaying, and describing data by using tables, graphs, and summary measures.

Inferential Statistics consists of methods that use sample results to help make decisions or predictions about a population.

What is a hypothesis?

We like to think of statistical hypothesis testing as the data analysis stage of an **experiment**, in which the scientist is interested, for example, in comparing the means of a population to a specified value (e.g. mean 'usability').

A **statistical hypothesis** is a statement about the parameters of one or more populations.

One-sided and two-sided hypotheses

Two-Sided Test:

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

One-Sided Tests:

$$H_0: \mu = \mu_0$$

$$H_1: \mu > \mu_0$$

or

$$H_0: \mu = \mu_0$$

$$H_1: \mu < \mu_0$$

Testing statistical hypotheses

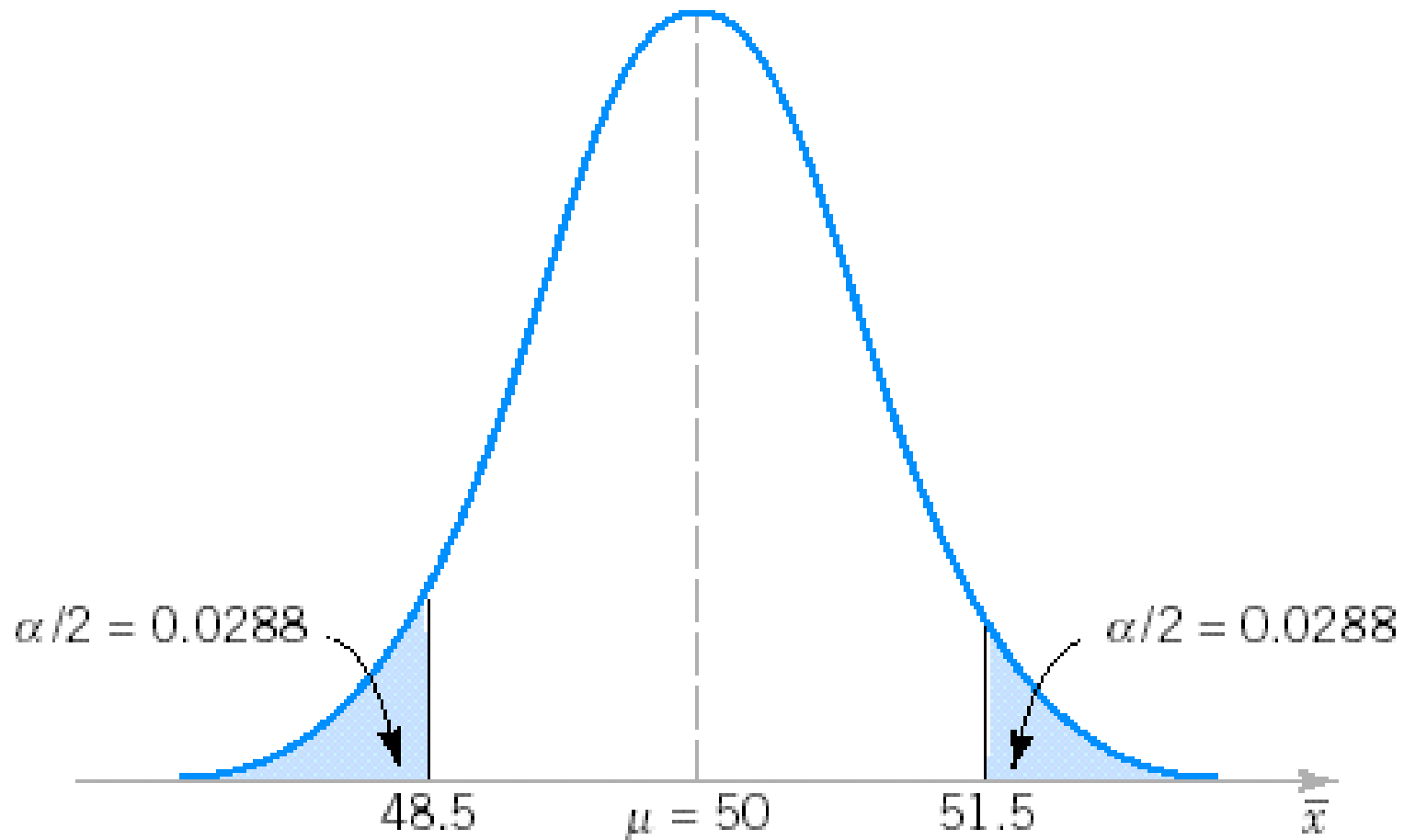


Figure 4-4 The critical region for $H_0: \mu = 50$ versus $H_1: \mu \neq 50$ and $n = 10$.

Errors in inference

- **Type I error:** Erroneously rejecting the null hypothesis. Your result is significant ($p < .05$), so you reject the null hypothesis, but the null hypothesis is actually true.
- **Type II error:** Erroneously accepting the null hypothesis. Your result is not significant ($p > .05$), so you don't reject the null hypothesis, but it is actually false.

Outcomes of a statistical analysis

H₀ True
(no correlation)

H₁ True
(correlation)

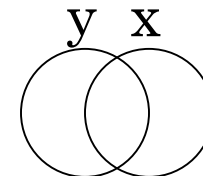
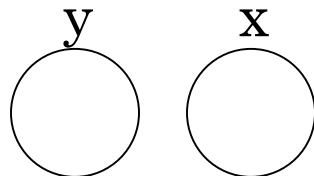
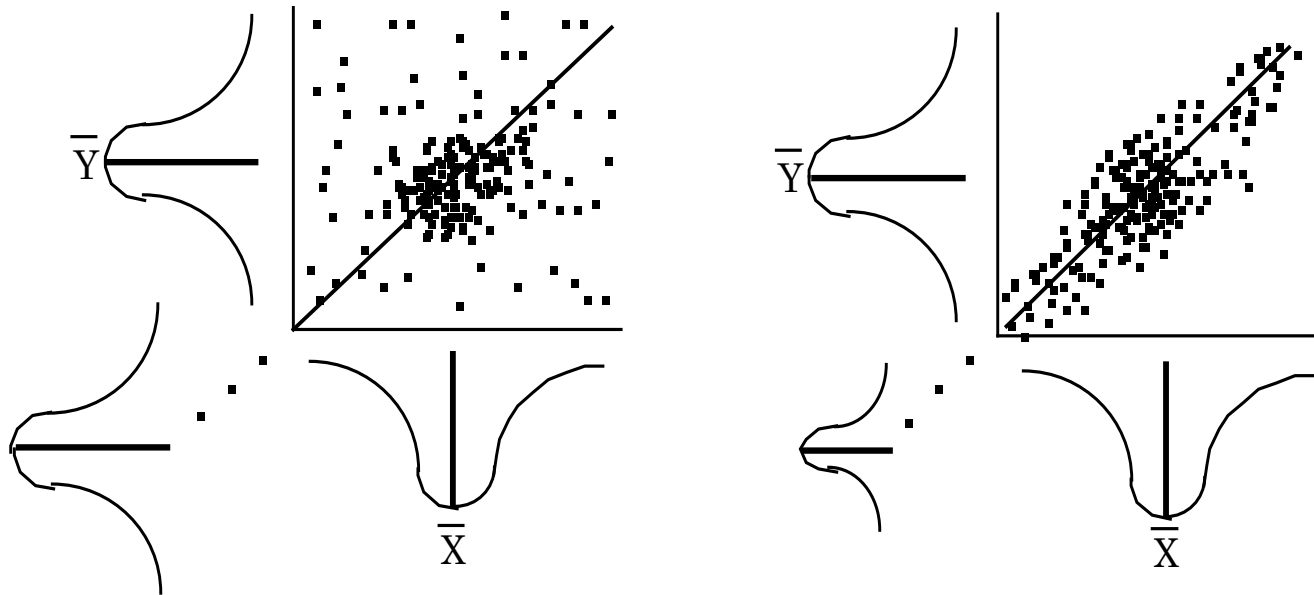
Do not reject H₀
(not stat. sig.)

Reject H₀
(stat. sig.)

Correct decision	Type II (beta error)
Type I (alpha error)	Correct decision

The analysis of variability...

Correlation is knowledge



Scale combinations leads to inference methods

Scales	Appropriate Inference Method
$N^2 * N^2$	Fisher's exact test; Odds Ratio
$N^c * N^d$	CHI ² (with $c > 2$ and/or $d > 2$)
$N^2 * O$	Mann-Whitney-U-test
$N^2 * I$	T-test (with Constant instead of N use one sample T-test)
$N_x^c * I$	[M]Anova (with $x > 1$ and/or $c > 2$)
$I_x * N^c$	Discriminant analysis (with $x > 1$ and $c > 1$)
$O * O$	Spearman/Kendall rank correlation
$I * I$	Pearson correlation
N_x	Cluster analysis (with $x > 2$)
O_x	Multi-Dimensional Scaling (with $x > 2$)
I_x	Factor analysis (with $x > 2$)

Choosing a significance level

- In general
 - Pilot program and intervention evaluations use liberal significance levels (.2 - .1) to avoid discarding effective interventions.
 - Generally accepted is a significance level of .05
 - Pure research uses conservative significance levels (.01-.001) to avoid wide dissemination of erroneous results.

Assignment-3

- Select **one** research topic with a concept to be investigated, formulate questions, identify relevant variables, define their scale type, and gather measurements (TIP: the more the better).
- Put all these measurements into one SPSS data file (case by case), specify variable names, value labels, and scale type.
- Produce with SPSS the most informative inferential statistics of all measures according their scale type.
- Prepare a powerpoint presentation.