

Data Analytics for Data Scientists

Design of Experiments (DoE)

Lecture 01: Introduction

2025

Prof. Dr. Jürg Schwarz

Program: 16.15 until 17.55

16.15	Begin of the lesson
	Lecture: Jürg Schwarz <ul style="list-style-type: none">◦ An introductory example<ul style="list-style-type: none">◦ Historical field trial in the US◦ Details on the field trial◦ Information about the module<ul style="list-style-type: none">◦ Content / structure / assessment◦ Preview of Lecture 02
17.00	<ul style="list-style-type: none">◦ What happened so far and what follows in Lecture 02
17.10	Tutorial: Students / Jürg Schwarz / Assistants <ul style="list-style-type: none">◦ Working on the exercise<ul style="list-style-type: none">◦ Support by Jürg Schwarz / Assistants
17.55	End of the lesson

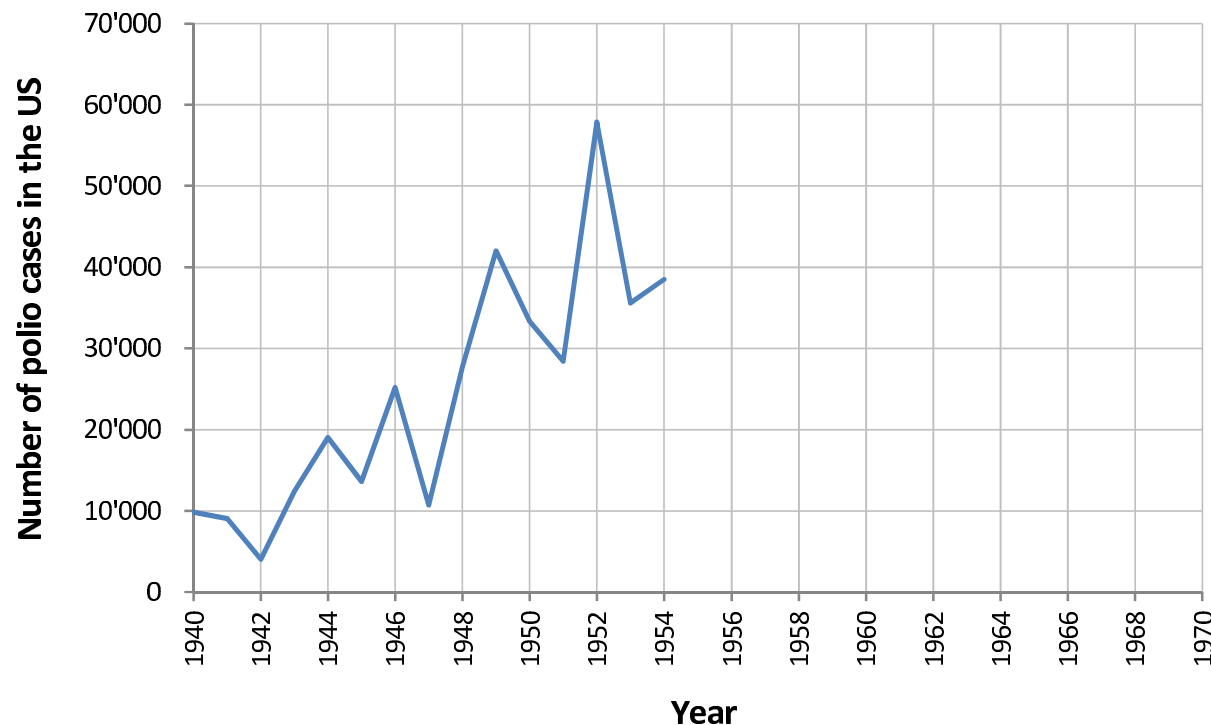
An introductory example

Historical field trial on polio in the US

Polio ...

- is a viral disease that occurs mainly in children.
- can lead to severe, permanent forms of paralysis.

Hundreds of thousands of polio cases in the US in the first half of the 20th century.



Vaccines and field trial

- Vaccines were developed at around **1950**, the most promising one by **Jonas Salk**.
- Lab tests showed that the vaccines were effective.
- In **1954**, the public health service and the NFIP* planned a **field trial** of the vaccine in children.

Design of Experiments (DoE)

- 33 states of the US were involved.
- The **population** of the trial consisted of about **1 million children** between the ages of 6 and 9.
- More than 750,000 children participated in the study.
- Procedure ...
 - some of the children were vaccinated (**Treatment**)
 - others were not (**Control**).
 - for some of the children, the parents did not give their consent (**No Consent**).

Result

- Analysis of polio cases in 1954

The results are rounded.

*The rate of polio cases ("Number of children with polio") is calculated **per 100,000 cases**.

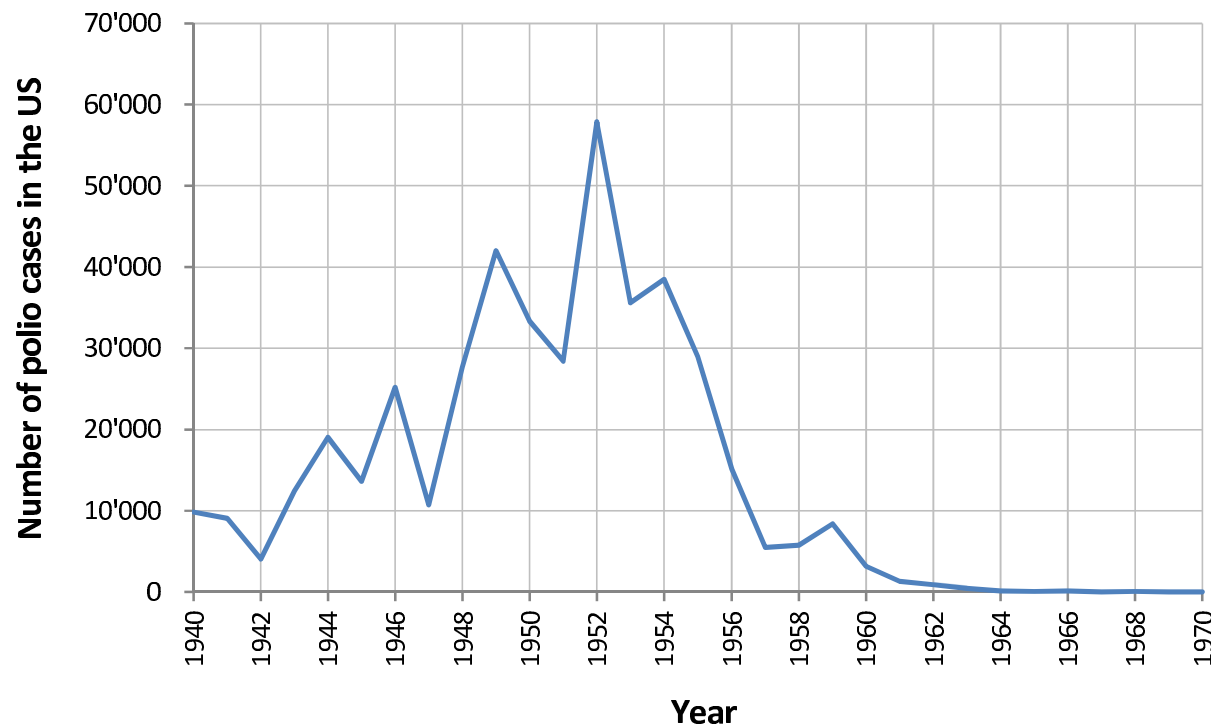
	Total number of children	Number of children with polio*
Treatment	200'000	25
Control	200'000	71
Total	400'000	96

Vaccinated children have a lower rate of polio cases (25) than non-vaccinated children (71).
The difference in the numbers between the treatment group and the control group is significant.

The vaccine that Jonas Salk developed can successfully prevent polio.

Further development

- After the field trial, the new vaccine was systematically introduced in the US.
- In the US, the health system supported polio vaccinations starting in 1960.
- The last case of polio in the US caused by a wild virus occurred in 1979. The last case in Switzerland was in 1982.



Details on the field trial on polio in the US

Reflections on the field trial

Is this statement scientifically valid?

Vaccinated children have a lower rate of polio cases (25) than unvaccinated children (71).

How can the effectiveness of the vaccine be determined in real life?

What can / must be done to conduct **scientific research**?

→ **Design of Experiments (DoE)**

The following section looks at some typical questions and problems in detail.

Research design: Observational vs. experimental

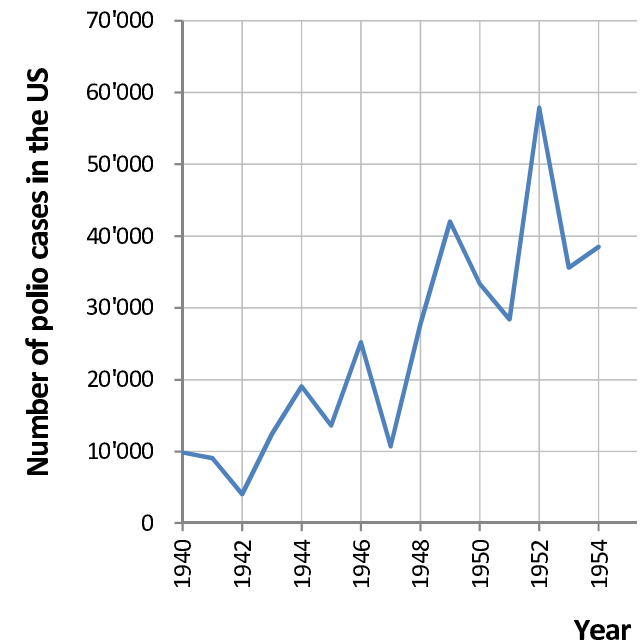
Why were not **all the children** vaccinated?

The rate of polio cases in 1954 could have been compared with the rates in previous years.

Problem concerning the observational design

Polio occurs as an epidemic, and the rate of polio cases varies greatly from year to year.

A low rate of polio cases in 1954 could have meant that 1954 was not an epidemic year – regardless of the vaccine's effect.



Solution with experimental design

Introduction of a comparison group (**Control**) consisting of children who are not vaccinated.

Comparison of the rate of polio cases among vaccinated children (**Treatment**) with the rate of polio cases among unvaccinated children (**Control**).

Research design: Randomization

Problem concerning the selection bias brought on by the parental consent

Parents with higher incomes / higher education are more likely to agree to having their child vaccinated because they are better informed about the experiment.

Therefore, the proportion of children in the vaccinated group (***Treatment***) whose parents have a higher income / higher education would be larger.

Problem of confounding* due to vulnerability to polio

Children of parents with a higher income / higher education are more vulnerable to polio.

Therefore, the vaccination would be more effective for these children on average.

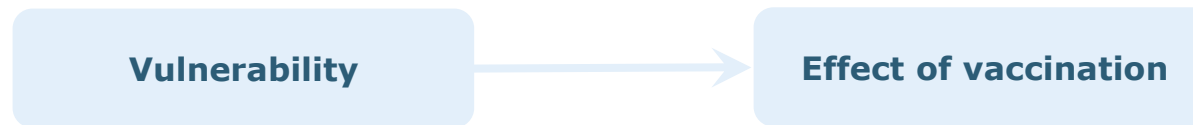
Note: This seems paradoxical because diseases are generally more common among low-income people. However, children who live in a less hygienic environment tend to become infected with mild cases of polio at an early age while still protected by their mother's antibodies. After infection, children produce their own antibodies, which later protect them against a more serious infection. Children who live in a more hygienic environment do not develop such antibodies.

Keyword: Confounding

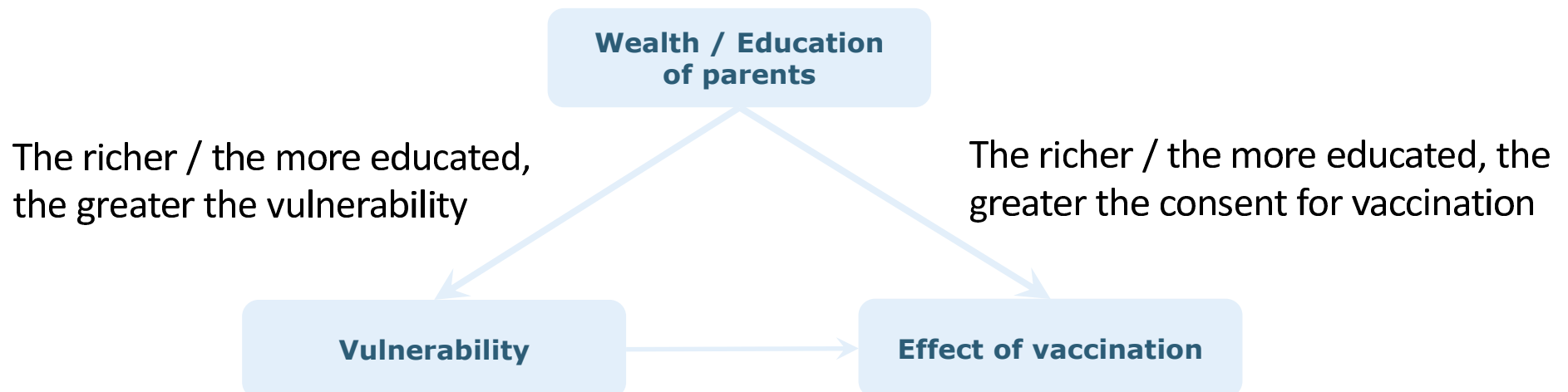
Confounding means that a factor (**Confounder**) that is not directly investigated is associated with both the independent variable and the dependent variable and accordingly causes the relationship between the two variables (**Spurious Correlation**).

Spurious Correlation

The greater the vulnerability, the more effective the vaccination.



Relationship caused by a confounder



Bias & confounding: Solution through randomization

The children are allocated **randomly** to the vaccination group (**Treatment**) and the control group (**Control**).

Ideally, each child is randomly assigned (with 50% probability) to the vaccination group or the control group.

How was randomization carried out?

The historical documents show that although randomization could not be carried out optimally, it was still adequate in terms of quality.

Liza Dawson sums it up in her article:

The design of the trial itself was an element of considerable controversy: initial plans forged by the NFIP and endorsed by Jonas Salk included vaccination of large numbers of children in a specified age group (second grade) and “observed controls” in the first and third grades of the participating schools. Leading virologists of the time challenged the lack of placebo controls and the case for placebos was won by the distinguished researcher leading the trial, [...]

However, many local public health officials had signed on to the project with the understanding that observed controls would be used and not all agreed to the placebo control design. Thus the final design was a patchwork of observed and placebo controls, although the overall trial results were unquestionably positive.

Research design: Blinding

Problem of distortion due to knowledge about the treatment

In the case of children being treated with the vaccine, knowledge about the treatment could change the effect of the vaccine.

Similarly, such knowledge could cause changes in the behavior of the doctors and nurses treating and observing the children.

Solution by blinding

Blinding is a suitable technique to avoid such distortions.

There are several stages of blinding → see [Slide 13](#)

One half [of the children] would receive vaccine; the other matching half, serving as strict controls, would receive a solution of similar appearance which should have no influence on immunity to poliomyelitis. Each child would receive the same lot of material, labeled only by code, for all three inoculations. Only the Evaluation Center would have the key.

The children in the study would be observed thereafter and all reports relating to a case of poliomyelitis would be made on a concealed or blindfold basis without knowledge of the nature of the inoculum.

Types of blinding

Open

- No blinding

Single-blind

- Test persons have **no knowledge** of group membership.

Blinding of the test persons is especially important when their attitude about the intervention affects their cooperation and compliance with regulations, or even the response to the treatment, e.g. **placebo effect**.

Double-blind

- Also, the persons carrying out the tests have **no knowledge** of group membership.

The intervention used in a study must be externally identical, with no differences in taste, smell, color, or method of application.

The "double-dummy technique," whereby an identical placebo is available for each intervention, can be applied when comparing two interventions.

Triple-blind

- Also, people doing the analysis have no knowledge about group membership.

Summary

The polio field trial is a classic **experimental design**:

Randomized double blind placebo control trial

Randomized controlled trials are generally abbreviated as **RCT**.

Characteristics and aim of an experimental designs

Treatment / exposure is planned and its effects are investigated

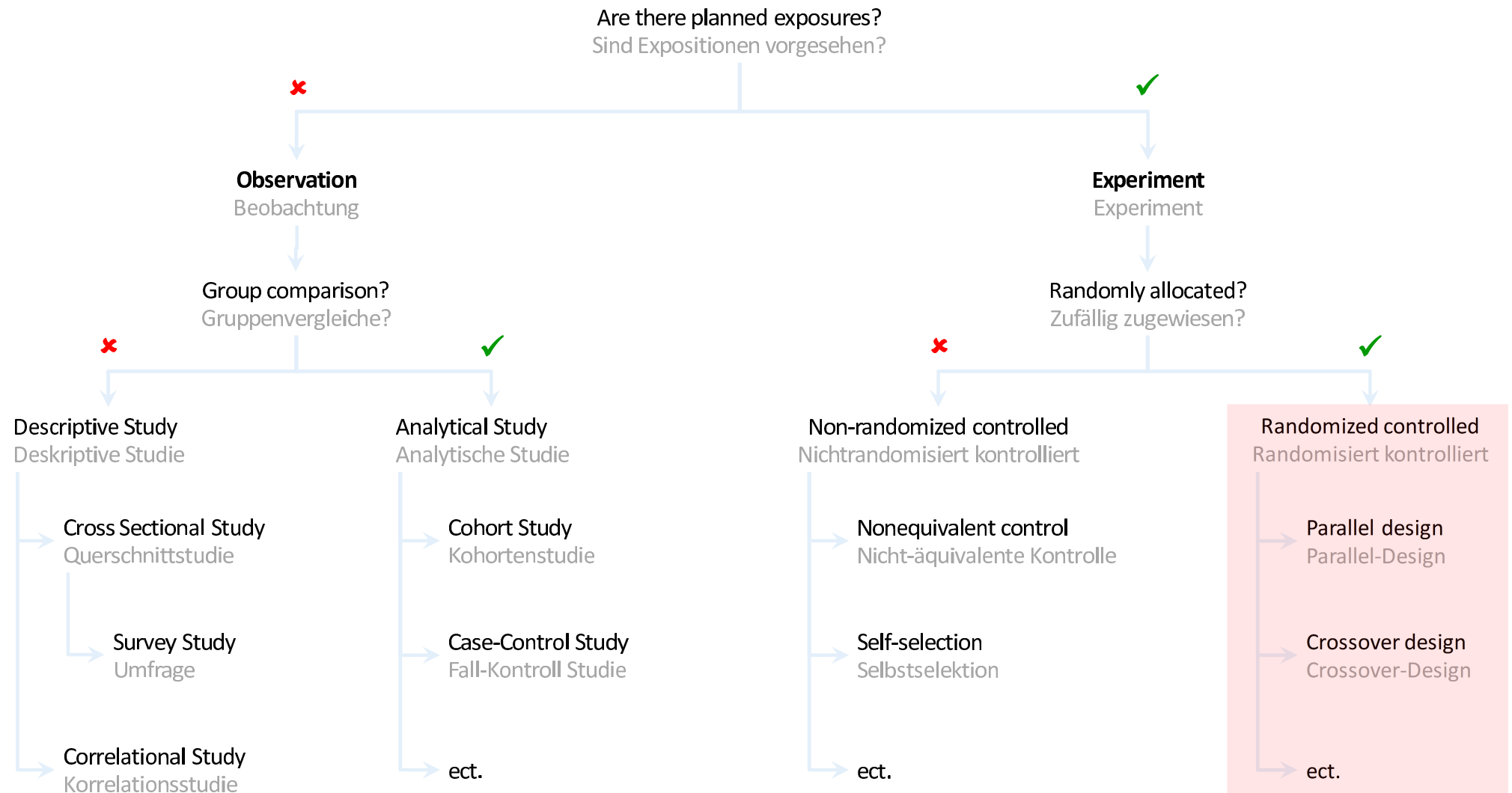
Effect and aim of randomization

- Elimination of **selection bias**
- Elimination of confounding
- Ensuring the comparability of the groups at study initiation (***Baseline Data***)

Effect and aim of blindness

- Elimination of conscious and unconscious influences on the treatment result
- Ensuring the comparability of the groups with regard to ...
 - Treatment in the course of the study
 - Evaluation of the effect of the treatment / exposure

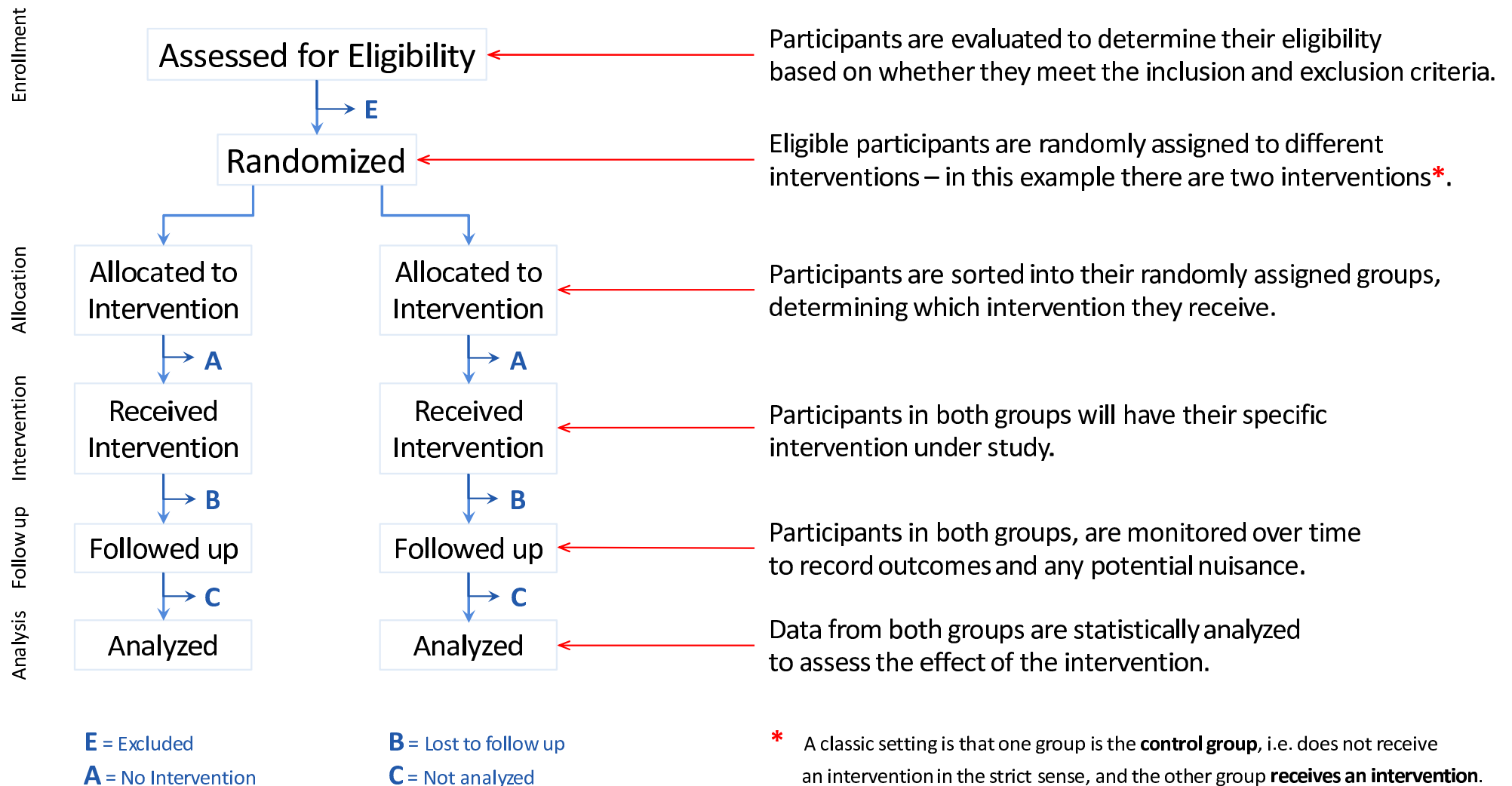
Decision Tree for Research Design – Part 1



Decision Tree for Research Design – Part 2

Type: Randomized controlled

Flowchart of the phases of a parallel randomized trial of two groups



Information about the module

Syllabus – what you always wanted to know

MSc Applied Information and Data Science	HSLU Lucerne University of Applied Sciences and Arts
Data Analytics for Data Scientists	
Design of Experiments (DoE)	
Syllabus	
February 2025	
Prof. Dr. Jürg Schwarz	
Overview and learning concept _____	2
Materials _____	3
Course assessments _____	3

Overview and learning concept _____

Materials _____

Course assessments _____

Topics and dates

Lecture	Place	Topic	Date	Exercise
01	Classroom	Introduction	20.02.2025	Introduction
02	Online	Principles	27.02.2025	02
03	Online	Introduction to design of experiments	06.03.2025	03
04	Classroom	Properties of design of experiments	13.03.2025	04
05	Online	Sampling	20.03.2025	05
	Online	Reflection / Self-study	27.03.2025	
06	Classroom	Effect size / Power analysis	03.04.2025	06
07	Classroom	Paradigms / Instead of an exercise Introduction to ANOVA	10.04.2025	07
08	Online	A/B Testing	17.04.2025	08
09	Online	Factorial designs	24.04.2025	09
10	Classroom	Large data quantities	01.05.2025	10
11	Online	Experiments in social media	08.05.2025	11
	Online	Reflection / Self-study	15.05.2025	
	Classroom	Guest Lecture	22.05.2025	

"Reflection / Self-study" allows you to process the material covered so far at your leisure.

Does not take place on site.

Is not accompanied by us either on site or online.

Combination of lectures and tutorials

Lectures on **Design of Experiments (DoE)** include highly structured content, which is why frontal teaching is used as the main method.

To make the lessons interactive nevertheless, you will participate in tutorials after the lecture during which you can apply the material that was covered.

In the **tutorials**, you will work on solving Exercises and have sufficient opportunity to ask questions and discuss possible solutions.

You will be supported by lecturers and assistants mainly in the Exercise belonging to the current course unit, and only in earlier Exercises if there is enough time.

We recommend that you organize yourselves into groups in the tutorials to discuss questions and support each other. ILIAS → Distance Learning → Zoom Links

Combination of lectures and tutorials

Timeline	Unit	Duration
16.15	Lecture in Class Schwarz	45 Minutes
17.00		
17.10	Tutorial in Groups Schwarz / Assistants	45 Minutes
17.55		

Classroom tutorials take place in the HSLU building (Z9) in reserved rooms.

Online tutorials are held via Zoom in your groups.

→ See the welcome e-mail for more details.

Course assessments

Course assessments

Exercise – Course assessment 1

You have to submit **your own solution to one** of the Exercises 02 to 06 or 08 to 11.

You are free to choose which of the Exercises you want to submit your solution to.

The deadline for submission is 8.00 a.m. on Thursday, May 15, 2025 via ILIAS.

Your solution must have been graded as "Pass" in order to take the final exam.

Details on submitting

File format PDF

File name Lastname_Firstname.pdf (e.g. Duck_Donald.pdf)

Location In the folder "Submission → Exercise" on ILIAS

Correspondence There will be no correspondence by e-mail.

Your Exercise will only be graded if it is correctly named, saved as a PDF, and uploaded to the correct location on ILIAS by the deadline.

Final exam – Course assessment 2

Date	This will be specified in the exam invitation, available by the end of the semester.
Location	This will be specified in the exam invitation, available by the end of the semester.
Duration	60 minutes
Rough duration	The exam invitation mentions more than 60 minutes duration. This is because of the additional time needed for the preparation and submission.
Admission	Your solution of an Exercise must have been graded as "Pass".
Weighting	The exam counts for 100% of the module grade.
Topics	Lectures 01 to 11 / Exercises 02 to 06 and 08 to 11, including the suggested solutions.
Level	The level of the exam corresponds to the level of the Lectures and Exercises.
Form	Paper & Pencil
Tools	No laptops or smartphones are allowed. All types of pocket calculators are allowed.
Materials	Your summary A maximum of eight A4 pages for your summary: 4 sheets double-sided or 8 sheets single-sided. The summary can contain any content, come from any source and be generated by hand or electronically.

Preview of Lecture 02

What has happened so far

Introductory example

Introduction to the topic of Design of Experiments (DoE)
using the example of polio, a classical study design

→ Randomized, controlled double-blind study

What follows in Lecture 02

Allocation

Placement of "Design of Experiments (DoE)" in the context of ...

- Scientific theory
- Research process

Research design

Reflexions on study designs

Appendix

Terminology

Research design is a strategic plan for a research project or research program, setting out the broad outline and key features of the work to be undertaken, including the methods of data collection and analysis to be employed

An **observation** is a result of study in which a variable is measured.

An **experiment** is a procedure carried out to support, refute, or validate a hypothesis. Experiments provide insight into cause-and-effect by demonstrating what outcome occurs when a particular factor is manipulated.

The **experimental group** is exposed to a condition that is hypothesized to have some causal effect. The **control group** is not, so that they act as a baseline with which to compare the results from the experimental group.

Treatment is a term used in the context of an experimental design to refer to any prescribed combination of values of explanatory variables

Randomized Controlled Trial (**RCT**) is an experiment in which subjects are randomly allocated into groups, usually called treatment and control groups, to receive or not an **intervention**.

Table of contents

An introductory example	3
Historical field trial on polio in the US	3
Details on the field trial on polio in the US.....	7
Information about the module	17
Syllabus – what you always wanted to know	17
Topics and dates	18
Combination of lectures and tutorials	19
Combination of lectures and tutorials	20
Course assessments	21
Preview of Lecture 02.....	23
What has happened so far	23
What follows in Lecture 02	23
Appendix	24
Terminology.....	24

