

Lecture

Empirical Model Building and Methods (Empirische Modellbildung und Methoden)

Prof. Dr. Dr. h.c. Dieter Rombach
Dr. Andreas Jedlitschka

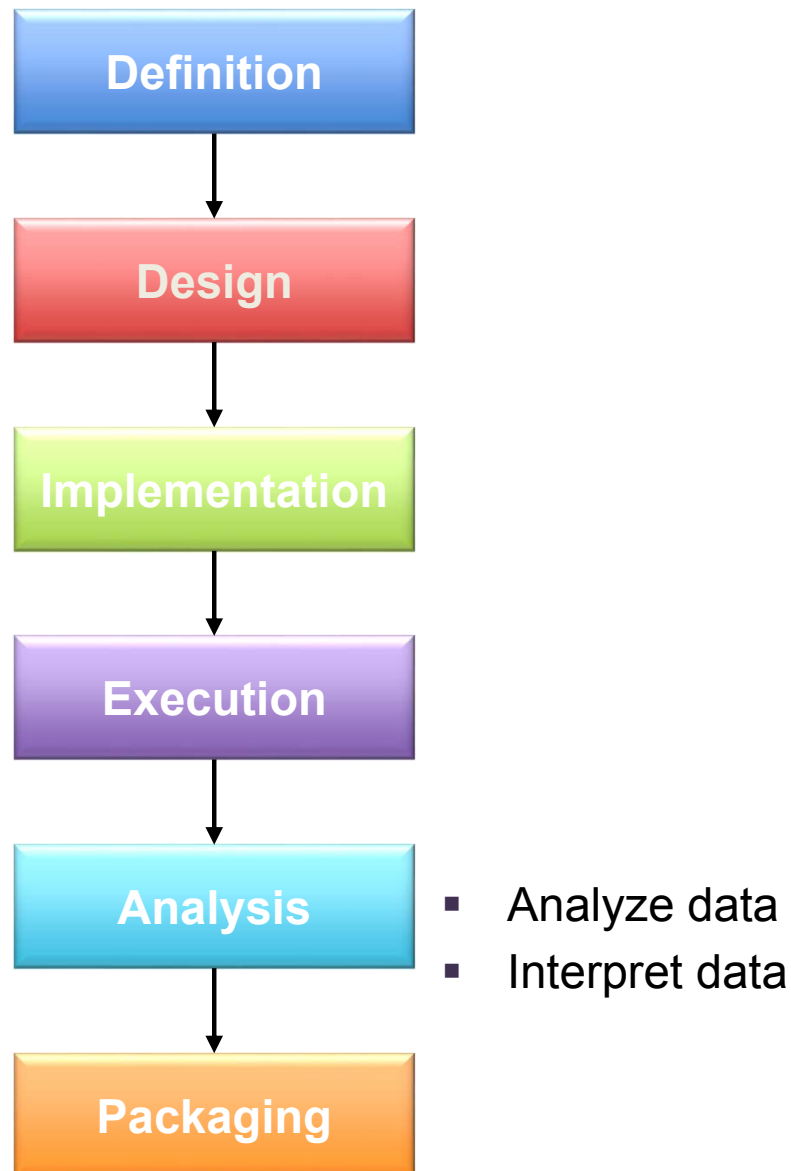
SS 2015

Chapter 3.6 – Analysis

Chapter objectives

At the end of this chapter, you should

- know the steps for analyzing results from empirical studies.
- understand practical issues concerning the analysis.



3.5.1 Overview

3.5.2 Descriptive Statistics

3.5.3 Data Reduction

3.5.4 Hypothesis Testing

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- 3.6 Data analysis**
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Purpose

- Analyze collected data for
 - Describing sample
 - Answering research questions
 - Testing hypothesis

Steps

- Explore data
 - Are they sensible? (Descriptive Statistics)
- Find unusual values
 - Missing values, outliers and inconsistent values
 - Data reduction
- Conduct statistical analysis
- Interpret results

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Exploring and describing data

- Frequencies
 - Absolute
 - Relative

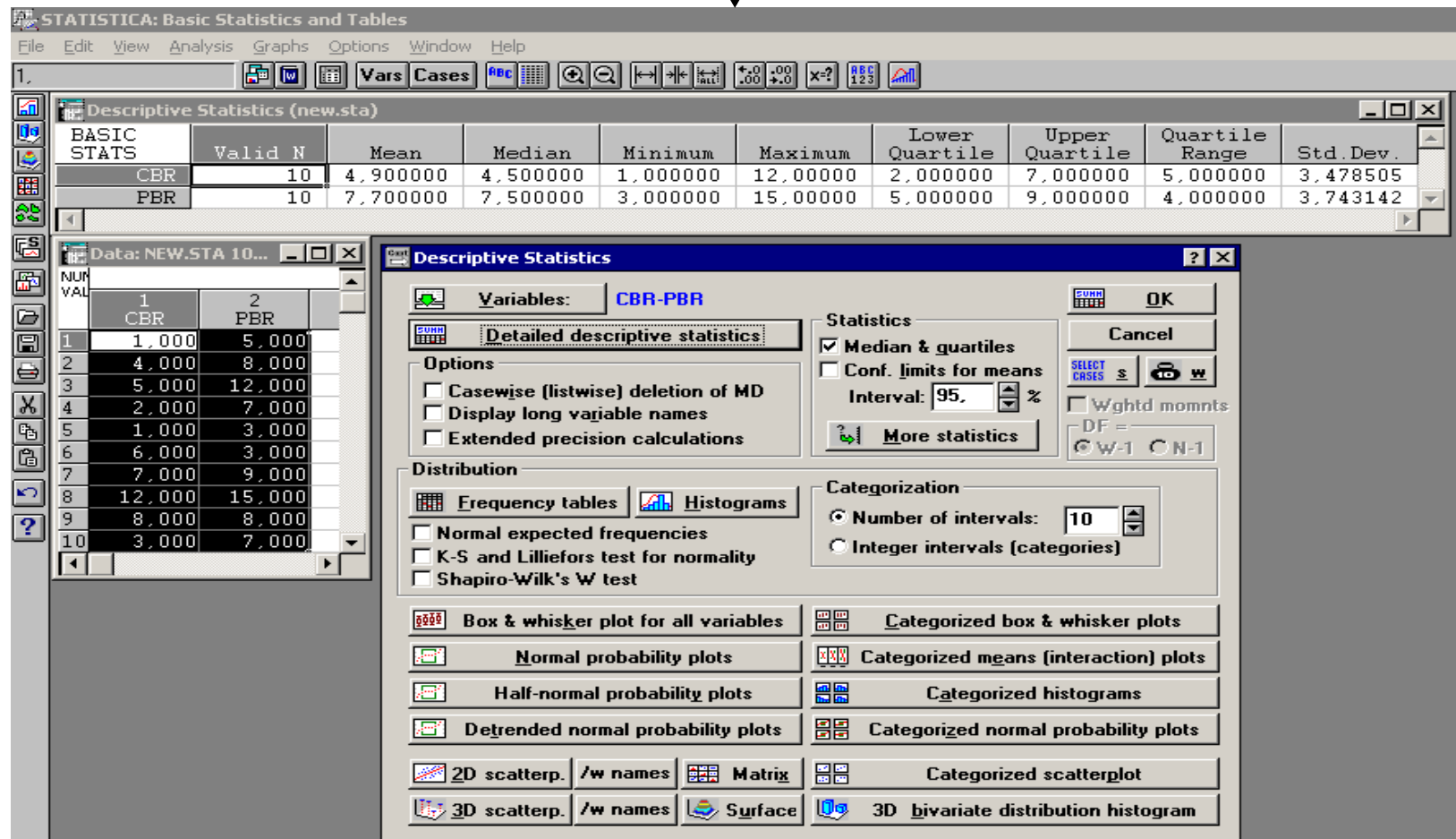
- Measure of central tendency
 - Mode: “most often”
 - Median: “middlest value”
 - Mean: “arithmetic average”

- Measures of distribution
 - Minimum/Maximum: “smallest/largest value”
 - Interquartile range: “range containing the middle 50% of data points”
 - Standard deviation

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Tool support

Statistica



The screenshot displays the STATISTICA software interface. The main window shows a 'Descriptive Statistics (new.sta)' table with the following data:

BASIC STATS	Valid N	Mean	Median	Minimum	Maximum	Lower Quartile	Upper Quartile	Quartile Range	Std.Dev.
CBR	10	4,900000	4,500000	1,000000	12,000000	2,000000	7,000000	5,000000	3,478505
PBR	10	7,700000	7,500000	3,000000	15,000000	5,000000	9,000000	4,000000	3,743142

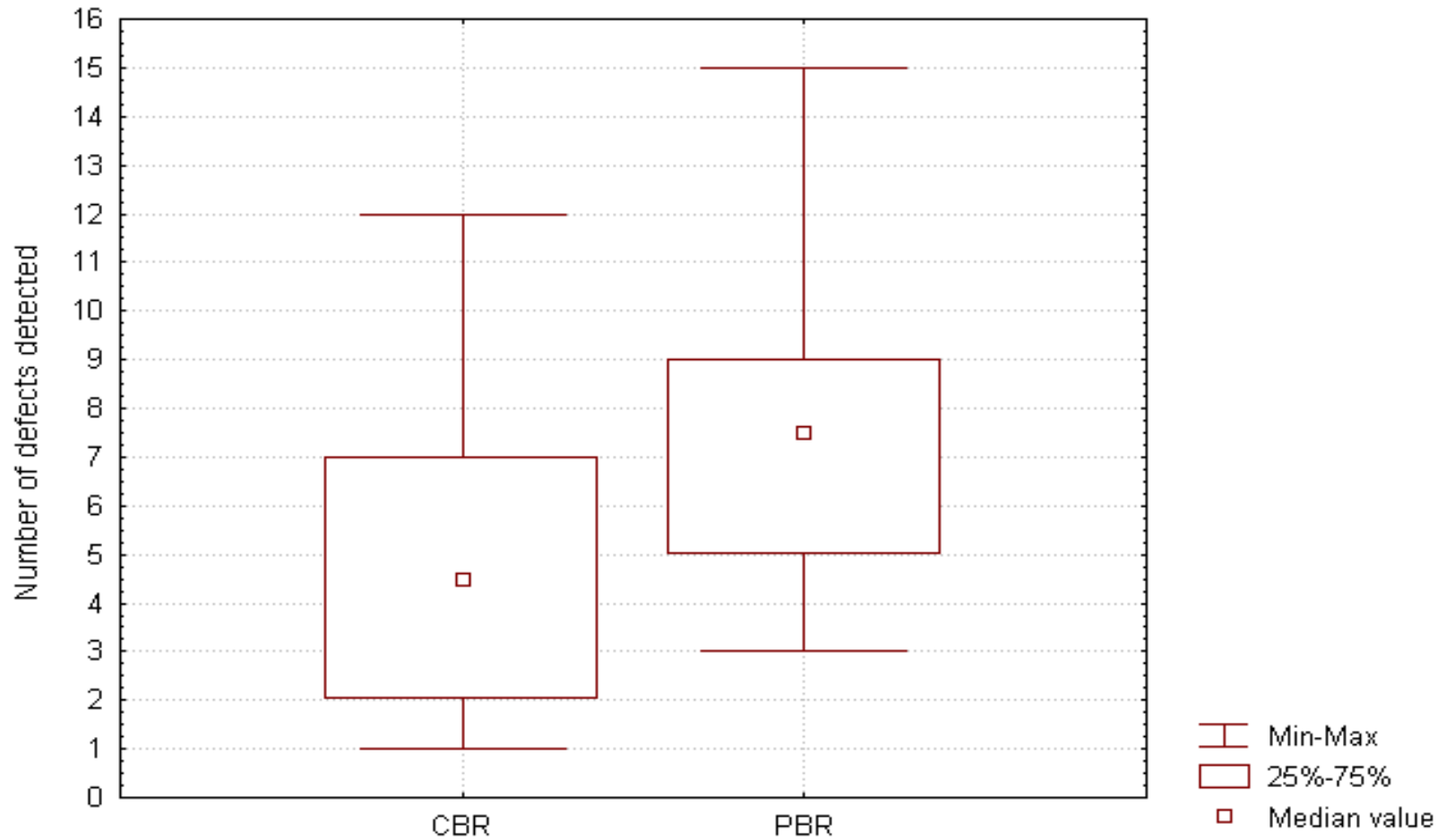
The 'Descriptive Statistics' dialog box is open, showing the following settings:

- Variables:** CBR-PBR
- Statistics:**
 - Median & quartiles
 - Conf. limits for means
 - Interval: 95. %
 - Wghtd momnts
 - DF = GW-1 CN-1
- Options:**
 - Casewise (listwise) deletion of MD
 - Display long variable names
 - Extended precision calculations
- Distribution:**
 - Normal expected frequencies
 - K-S and Lilliefors test for normality
 - Shapiro-Wilk's W test
- Categorization:**
 - Number of intervals: 10
 - Integer intervals (categories)
- Plots:**
 - Box & whisker plot for all variables
 - Normal probability plots
 - Half-normal probability plots
 - Detrended normal probability plots
 - Categorized box & whisker plots
 - Categorized means (interaction) plots
 - Categorized histograms
 - Categorized normal probability plots
 - Categorized scatterplot
 - 2D scatterp. /w names
 - 3D scatterp. /w names
 - Matrix
 - Surface
 - 3D bivariate distribution histogram

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Visualization

e.g., Boxplot



- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Purpose of data validation

- Checking consistency and credibility
- Ensuring the correctness and completeness of the collected data
- Identifying missing variables
- Identifying outliers

# Defects Subject 1	Effort Subject 1 (minutes)	# Defects Subject 2	Effort Subject 2 (minutes)	Total Effort
13333	0	12	120	120
23	60	10	65	125
	20	23	80	80

Exceptionally high/low values

Null values reasonable?

Missing values

Missing records

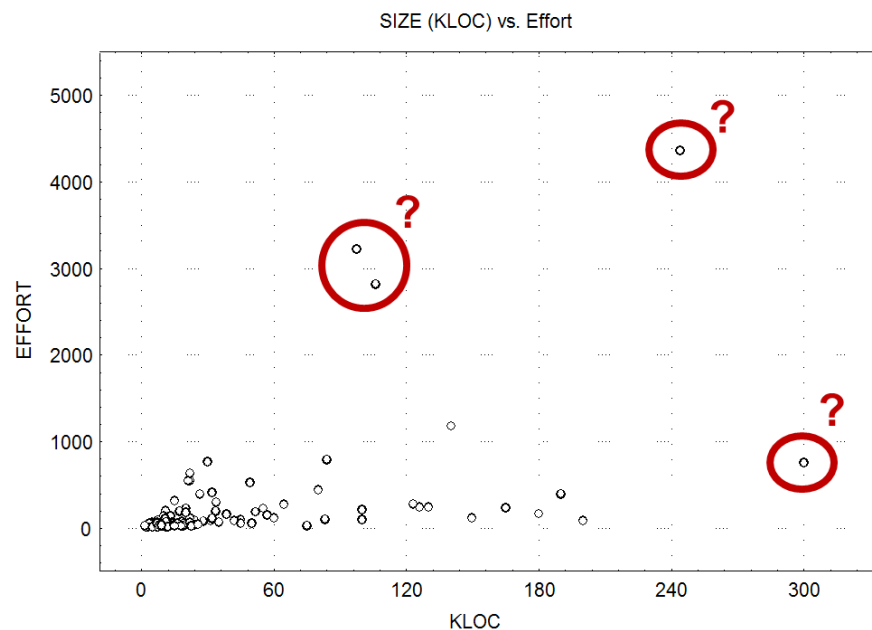
Inconsistent values

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Outliers

- Data points that are much larger or much smaller than one could expect looking at the other points

- Problems
 - They may be hinting at hidden influences
 - Comparison with the rest of data may make no sense
 - Most statistical techniques are not robust against outliers.



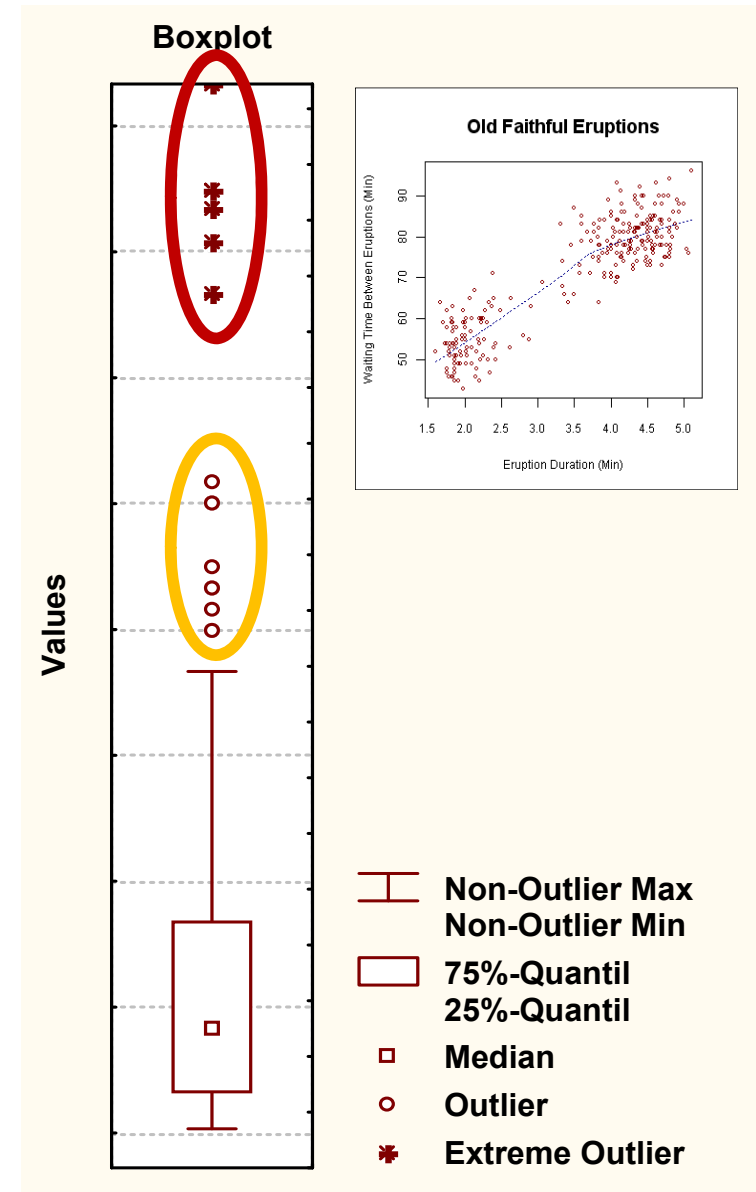
- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
 - Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Identifying outliers

- Create box-plots
- Analyse if there are systematic explanations for these values
- Perform sensitivity analysis with respect to measures of central tendency and distribution

Options for dealing with outliers

- Exclude the outlying observations
- Report two analysis (with and without outliers)
- (Transform the variable)
- Gather more observations



- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Definition

- Systematic procedure for determining whether experimental results (of a sample) provide support for the validity of a particular theory/hypothesis in the overall population.
- It implies decision making based on probabilities!

Basic questions

- Is the relationship between two or more variables due to chance?
 - Mathematical Statistics
- If chance can be excluded, what does it mean?
 - Interpretation

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Approach (Falsification principle)

- “Opposite-of-what-you-predict” reasoning
- Drawing conclusions by
 - “identifying the probability of rejecting the inverse hypothesis” (e.g., > 0.95) or
 - “of not rejecting the inverse hypothesis although the hypothesis is true” (e.g., < 0.05)

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Logic of significance test: Example (1/2)

- Research hypotheses
 - H2: Individuals applying each PBR perspective (PERSP) respectively perform better than individuals applying Ad-hoc reading (N) with respect to their mean defect detection rate (DDR).
 - $PERSP = \{\text{tester perspective (T), user perspective (U), designer perspective (D)}\}$
 - $DDR = \text{number of defects found by individual} / \text{total number of defects in the document}$

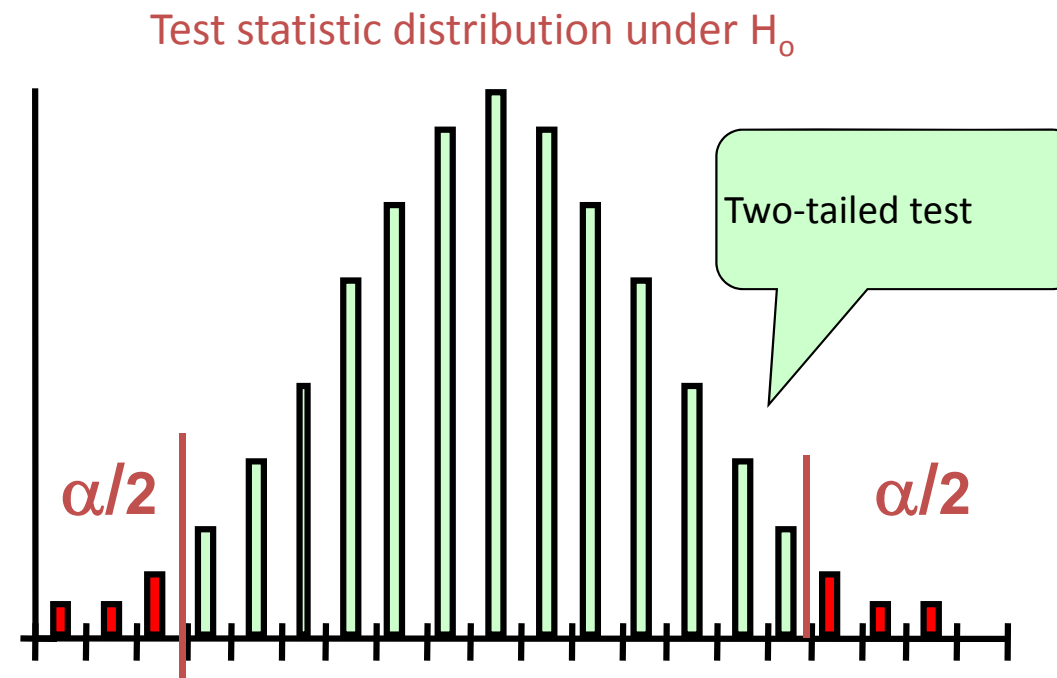
- Statistical hypotheses
 - Null hypothesis
 - Individuals applying ad-hoc reading perform better than individuals applying each PBR perspective
 - $H_{02} = DDR(PBR\ PERSP) \leq DDR(N)$

 - Alternative hypothesis
 - $H_{12} = DDR(PBR\ PERSP) > DDR(N)$

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Logic of significance test: Example (2/2)

- Statistical hypotheses



- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Procedure (1/2)

1. Formulate the alternative and null hypothesis
2. Select statistical test considering data distribution
 - Normal distribution → Parametric tests
 - Non-normal or ordinal/nominal distribution → Non-parametric tests
3. Select significance level (α -value) and perform power analysis
 - α conventionally 0.05 or 0.01
 - Power = $1 - \beta$ (β conventionally 0.2)
 - Determine optimal sample size based on α , effect size and power
 - Determine α based on sample size, effect size and power

		In the population ...	
		H_0 is true	H_0 is false
Decision	H_0 is not rejected	Correct outcome True negative	Type II error (β) False negative
	H_0 is rejected	Type I error (α) False positive	Correct outcome True positive

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- 3.6 Data analysis**
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Purpose

- (a) compare treatments or (b) correlate variables
 - If compare treatments
 - How many independent variables did you manipulate
 - How many do you wish to compare at any one time
 - Related or unrelated sample

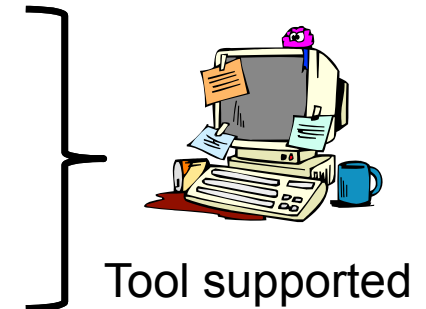
Type of Data

- (a) nominal, (b) ordinal or (c) at least interval
 - If interval
 - Are the scores normally distributed

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Procedure (1/2)

1. Determine test statistic of your sample
 - e.g., t-value for student t-test
2. Compare test statistic of your sample with comparison distribution



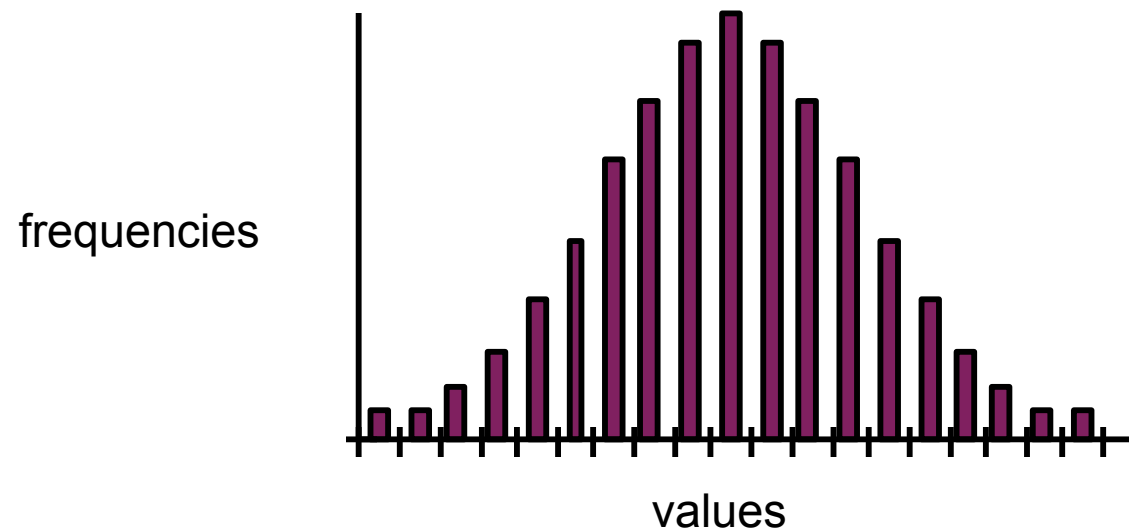
WS95/96	N	df. Effect	MS effect	df. Error	MS error	F	P-value
D vs. N	9	1	355.5	16	63.26	5.620	.301
U vs. N	8	1	6.250	14	66.36	.094	.736
T vs. N	8	1	115.5	14	99.81	1.158	.300

WS96/97	N	df. Effect	MS effect	df. Error	MS error	F	P-value
D vs. N	8	1	246.7	14	55.27	4.464	.053
U vs. N	9	1	547.0	16	71.49	7.651	.014
T vs. N	9	1	101.0	16	70.07	1.442	.247

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Assumption: Normal Distribution

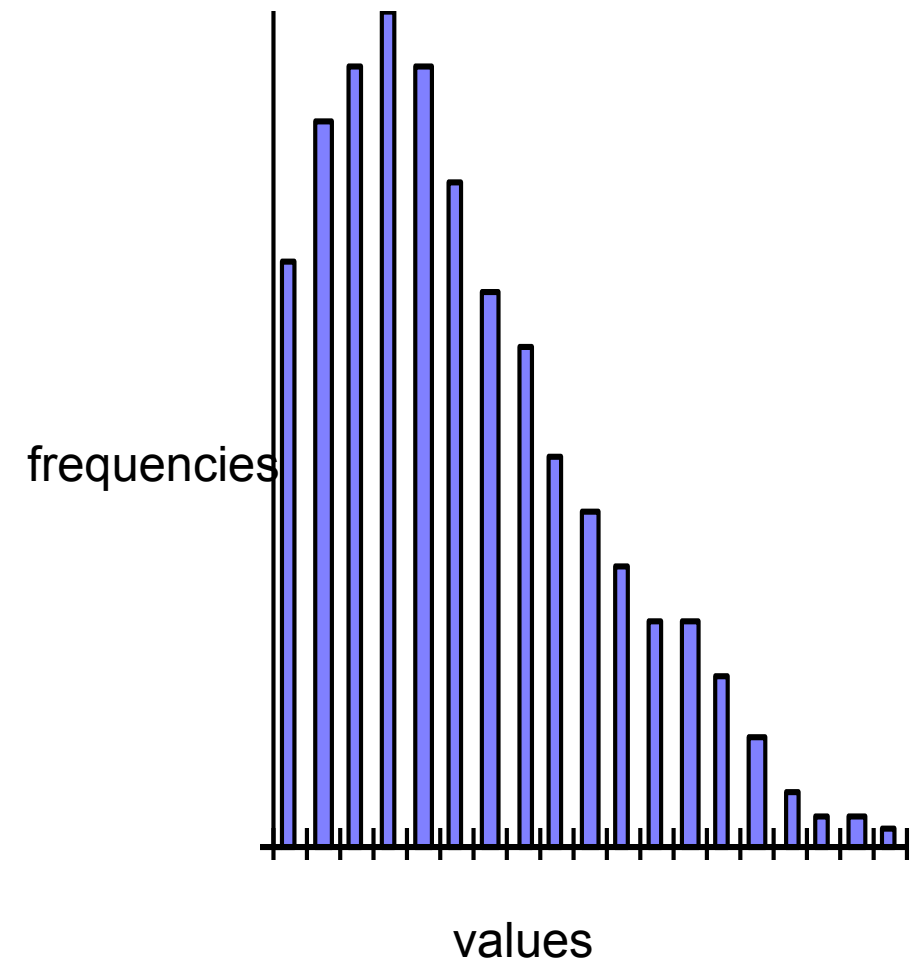
- Symmetrical, no beginning or ending point, 68.26% of the area within one standard deviation from the mean.
- Rare in software engineering
- Many variables have properties that are close enough, but frequently some variables are skewed to the right.



- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Other distributions:

- e.g., Skewed distribution
 - Typical distribution in SE
 - Examples
 - Productivity
 - Change density



- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

What do we do when the populations sampled are not normally distributed?

Central Limit Theorem (simplified):

- Irrespective of the distribution of the parent population - given that its mean m and a variance s^2 , and so long as the sample size n is large, the distribution of sample means is approximately normal with mean m and variance s^2 / n .

Beware of small sample sizes!

- Consider nonparametric tests.

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Statistical tests exist for different purposes, with different preconditions, and different power!

Statistical Test Selection (1/2)

depends on further assumptions

What do you want to do?	Number of variables	Subjects in condition	Parametric Test	Non parametric Test
Differences between conditions	One variable: two treatments	Independent	Independent t-test	Mann-Whitney U test
		Dependent	Paired t-test	Wilcoxon matched pairs test
	One variable: > 2 treatments	Independent	One factor independent ANOVA	Kruskal-Wallis-One way ANOVA
		Dependent	One factor repeated measures ANOVA	Friedman ANOVA
	Two or more treatments	Independent/Dependent	Variation of ANOVA-Analysis	
Compare frequency counts (in categories)				Chi-square
Correlate variables	Two variables		Pearson's r	Spearman's r
	More than two variables		Multiple correlation R	

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Statistical Test Selection (2/2)

- Each test makes a number of assumptions
 - Experimental design
 - Distribution of data
 - Outliers
- Be aware that your data set fulfills the assumptions.
- Statistical testing is just the means to an end not an end in itself.
- More difficult than running statistical tests
 - Interpretation of the results.
 - What does the results mean?

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Statistical Tools

- SPSS (<http://www.spss.com/>)
- Statistica (<http://www.statsoftinc.com>)
- R Package (<http://www.r-project.org/>, Open source)
- JMP (<http://www.jmpdiscovery.com/>)
- S-Plus (<http://www.splus.mathsoft.com>)
- Minitab (www.minitab.com/)
- Excel
- ...

- 3.1 Introduction
- 3.2 Definition
- 3.3 Design
- 3.4 Implementation
- 3.5 Execution
- Data analysis
 - 3.6.1 Overview
 - 3.6.2 Descriptive Statistics
 - 3.6.3 Data Reduction
 - 3.6.4 Hypothesis Testing
- 3.7 Packaging

Interpretation

- If null-hypothesis is rejected
 - There is an effect.
 - Calculate „real“ effect size.
 - Is effect in concordance with theory?
 - Do you need to modify the theory?

- If null-hypothesis is not rejected
 - It is not possible to conclude there is no effect!
 - There is not sufficient evidence to accept there is an effect.
 - Discuss descriptive statistics.
 - Calculate „real“ effect size. Determine sample size required to reject null hypothesis, if effect is in „right“ direction.
 - **Formulate new hypothesis**

		In the population ...	
		H_0 is true	H_0 is false
Decision	H_0 is not rejected	Correct outcome True negative	Type II error (β) False negative
	H_0 is rejected	Type I error (α) False positive	Correct outcome True positive

Discuss and interpret the results in the context of your theory!

References

- Bortz, J. and Döring, N. (2006). **Forschungsmethoden und Evaluation für Human- und Sozialwissenschaftler (4 Auflage)**. Berlin: Springer Verlag.
- Juristo, N., and Moreno, A. (2001) **Basics of Software Engineering Experimentation**, Kluwer Academic Publishers.
- Jedlitschka, A., Ciolkowski, M., Pfahl, D. (2008). **Reporting Controlled Experiments in Software Engineering**. In: Shull, F., Singer, J., Sjöberg, D.I. (Eds.). **Guide to Advanced Empirical Software Engineering**. Springer.
- Wohlin, Runeson, Höst, Ohlsson, Regnell, Wesslén (2000). **Experimentation in Software Engineering: An Introduction**, Kluwer Academic Publishers.