Chapman & Hall/CRC Interdisciplinary Statistics Series

Power Analysis of Trials with Multilevel Data



Mirjam Moerbeek Steven Teerenstra



A CHAPMAN & HALL BOOK

Power Analysis of Trials with Multilevel Data

CHAPMAN & HALL/CRC

Interdisciplinary Statistics Series

Series editors: N. Keiding, B.J.T. Morgan, C.K. Wikle, P. van der Heijden

Published titles

AGE-PERIOD-COHORT ANALYSIS: NEW MODELS, METHODS, AND EMPIRICAL APPLICATIONS Y. Yang and K. C. Land

ANALYSIS OF CAPTURE-RECAPTURE DATA R.S. McCrea and B.J.T. Morgan

AN INVARIANT APPROACH TO STATISTICAL ANALYSIS OF SHAPES S. Lele and J. Richtsmeier

ASTROSTATISTICS G. Babu and E. Feigelson

BAYESIAN ANALYSIS FOR POPULATION ECOLOGY R. King, B. J.T. Morgan, O. Gimenez, and S. P. Brooks

BAYESIAN DISEASE MAPPING: HIERARCHICAL MODELING IN SPATIAL EPIDEMIOLOGY, SECOND EDITION A. B. Lawson

BIOEQUIVALENCE AND STATISTICS IN CLINICAL PHARMACOLOGY S. Patterson and B. Jones

CLINICAL TRIALS IN ONCOLOGY, THIRD EDITION S. Green, J. Benedetti, A. Smith, and J. Crowley

CLUSTER RANDOMISED TRIALS R.J. Hayes and L.H. Moulton

CORRESPONDENCE ANALYSIS IN PRACTICE, SECOND EDITION M. Greenacre

DESIGN AND ANALYSIS OF QUALITY OF LIFE STUDIES IN CLINICAL TRIALS, SECOND EDITION D.L. Fairclough

DYNAMICAL SEARCH L. Pronzato, H. Wynn, and A. Zhigljavsky

FLEXIBLE IMPUTATION OF MISSING DATA S. van Buuren

GENERALIZED LATENT VARIABLE MODELING: MULTILEVEL, LONGITUDI-NAL, AND STRUCTURAL EQUATION MODELS A. Skrondal and S. Rabe-Hesketh

GRAPHICAL ANALYSIS OF MULTI-RESPONSE DATA K. Basford and J. Tukey

INTRODUCTION TO COMPUTATIONAL BIOLOGY: MAPS, SEQUENCES, AND GENOMES M. Waterman

MARKOV CHAIN MONTE CARLO IN PRACTICE W. Gilks, S. Richardson, and D. Spiegelhalter

MEASUREMENT ERROR ANDMISCLASSIFICATION IN STATISTICS AND EPIDE-MIOLOGY: IMPACTS AND BAYESIAN ADJUSTMENTS P. Gustafson

MEASUREMENT ERROR: MODELS, METHODS, AND APPLICATIONS

J. P. Buonaccorsi

MEASUREMENT ERROR: MODELS, METHODS, AND APPLICATIONS

J. P. Buonaccorsi

Published titles

MENDELIAN RANDOMIZATION: METHODS FOR USING GENETIC VARIANTS IN CAUSAL ESTIMATION S.Burgess and S.G.Thompson

META-ANALYSIS OF BINARY DATA USINGPROFILE LIKELIHOOD D. Böhning, R. Kuhnert, and S. Rattanasiri

POWER ANALYSIS OF TRIALS WITH MULTILEVEL DATA M. Moerbeek and S. Teerenstra

STATISTICAL ANALYSIS OF GENE EXPRESSION MICROARRAY DATA T. Speed

STATISTICAL AND COMPUTATIONAL PHARMACOGENOMICS R. Wu and M. Lin

STATISTICS IN MUSICOLOGY J. Beran

STATISTICS OF MEDICAL IMAGING T. Lei

STATISTICAL CONCEPTS AND APPLICATIONS IN CLINICAL MEDICINE J. Aitchison, J.W. Kay, and I.J. Lauder

STATISTICAL AND PROBABILISTIC METHODS IN ACTUARIAL SCIENCE P.J. Boland

STATISTICAL DETECTION AND SURVEILLANCE OF GEOGRAPHIC CLUSTERS P. Rogerson and I. Yamada

STATISTICS FOR ENVIRONMENTAL BIOLOGY AND TOXICOLOGY A. Bailer and W. Piegorsch

STATISTICS FOR FISSION TRACK ANALYSIS R.F. Galbraith

VISUALIZING DATA PATTERNS WITH MICROMAPS D.B. Carr and L.W. Pickle

Chapman & Hall/CRC Interdisciplinary Statistics Series

Power Analysis of Trials with Multilevel Data

Mirjam Moerbeek

Utrecht University, The Netherlands

Steven Teerenstra

Radboud University Medical Center, The Netherlands



CRC Press is an imprint of the Taylor & Francis Group, an **informa** business A CHAPMAN & HALL BOOK CRC Press Taylor & Francis Group 6000 Broken Sound Parkway NW, Suite 300 Boca Raton, FL 33487-2742

© 2016 by Taylor & Francis Group, LLC CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works Version Date: 20150417

International Standard Book Number-13: 978-1-4987-2990-1 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright. com (http://www.copyright.com/) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at http://www.taylorandfrancis.com

and the CRC Press Web site at http://www.crcpress.com

Contents

Li	st of	figures	xi
Li	st of	tables	xv
Pr	Preface x		
1	Intr	oduction	1
	1.1	Experimentation	2
		1.1.1 Problems with random assignment	5
	1.2	Hierarchical data structures	6
	1.3	Research design	9
		1.3.1 Cluster randomized trial	10
		1.3.2 Multisite trial	11
		1.3.3 Pseudo cluster randomized trial	12
		1.3.4 Individually randomized group treatment trial	12
		1.3.5 Longitudinal intervention study	13
		1.3.6 Some guidance to design choice	14
	1.4	Power analysis for experimental research	15
	1.5	Aim and contents of the book	18
		1.5.1 Aim	18
		1.5.2 Contents	18
2	Mu	tilevel statistical models	21
	2.1	The basic two-level model	21
	2.2	Estimation and hypothesis test	26
	2.3	Intraclass correlation coefficient	29
	2.4	Multilevel models for dichotomous outcomes	32
	2.5	More than two levels of nesting	35
	2.6	Software for multilevel analysis	37
3	Con	cepts of statistical power analysis	39
	3.1	Background of power analysis	39
		3.1.1 Hypotheses testing	39
		3.1.2 Power calculations for continuous outcomes	41
		3.1.3 Power calculations for dichotomous outcomes	45
		3.1.3.1 Risk difference	45
		3.1.3.2 Odds ratio	46

	3.2	Types of power analysis
	3.3	Timing of power analysis 49
	3.4	Methods for power analysis
	3.5	Robustness of power and sample size calculations 52
	3.6	Procedure for a priori power analysis
		3.6.1 An example
	3.7	The optimal design of experiments
		3.7.1 An example (continued) $\ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots 59$
	3.8	Sample size and precision analysis
	3.9	Sample size and accuracy of parameter estimates 61
4	Clu	ster randomized trials 63
	4.1	Introduction
	4.2	Multilevel model
	4.3	Sample size calculations for continuous outcomes
		4.3.1 Factors that influence power
		$4.3.2 \text{Design effect} \dots \dots \dots \dots \dots \dots \dots \dots \dots $
		4.3.3 Sample size formulae for fixed cluster size or fixed
		number of clusters
		4.3.4 Including budgetary constraints
	4.4	Sample size calculations for dichotomous outcomes 78
		4.4.1 Risk difference
		4.4.2 Odds ratio
	4.5	An example
5	Imp	proving statistical power in cluster randomized trials 83
	5.1^{-1}	Inclusion of covariates
	5.2	Minimization, matching, pre-stratification
	5.3	Taking repeated measurements
	5.4	Crossover in cluster randomized trials
	5.5	Stepped wedge designs 101
6	Mu	ltisite trials 107
	6.1	Introduction 107
	6.2	Multilevel model
	6.3	Sample size calculations for continuous outcomes 115
		6.3.1 Factors that influence power
		$6.3.2$ Design effect \ldots 118
		6.3.3 Sample size formulae for fixed cluster size or fixed
		number of clusters
		6.3.4 Including budgetary constraints
		6.3.5 Constant treatment effect
	6.4	Sample size calculations for dichotomous outcomes
		6.4.1 Odds ratio
	6.5	An example
		•

7	Pse	eudo cluster randomized trials	129
	7.1	Introduction	129
	7.2	Multilevel model	132
	7.3	Sample size calculations for continuous outcomes	134
		7.3.1 Factors that influence power	134
		7.3.2 Design effect	136
		7.3.3 Sample size formulae for fixed cluster size or fixed	
		number of clusters	137
	7.4	Sample size calculations for binary outcomes	138
	7.5	An example	140
8	Ind	ividually randomized group treatment trials	141
	8.1	Introduction	141
	8.2	Multilevel model	143
		8.2.1 Clustering in both treatment arms	143
		8.2.2 Clustering in one treatment arm	145
	8.3	Sample size calculations for continuous outcomes	146
		8.3.1 Clustering in both treatment arms	146
		8.3.1.1 Factors that influence power	146
		8.3.1.2 Sample size formulae for fixed cluster sizes .	147
		8.3.1.3 Including budgetary constraints	148
		8.3.2 Clustering in one treatment arm	150
		8.3.2.1 Factors that influence power	150
		8.3.2.2 Sample size formulae for fixed cluster sizes .	151
		8.3.2.3 Including budgetary constraints	151
	8.4	Sample size calculations for dichotomous outcomes	153
		8.4.1 Clustering in both treatment arms	153
		8.4.2 Clustering in one treatment arm	154
	8.5	An example	155
9	Lon	ngitudinal intervention studies	159
	9.1	Introduction	159
	9.2	Multilevel model	161
	9.3	Sample size calculations for continuous outcomes	165
		9.3.1 Factors that influence power	165
		9.3.2 Sample size formula for fixed number of measurements	168
		9.3.3 Including budgetary constraints	169
	9.4	Sample size calculations for dichotomous outcomes	170
		9.4.1 Odds ratio	171
	9.5	The effect of drop-out on statistical power	172
		9.5.1 The effects of different drop-out patterns	173
		9.5.2 Including budgetary constraints	179
	9.6	An example	180

10 Extensions: three levels of nesting and factorial designs	183	
10.1 Introduction \ldots	183	
10.2 Three-level cluster randomized trials	184	
10.3 Multisite cluster randomized trials	188	
10.4 Repeated measures in cluster randomized trials and multisite		
trials	193	
10.5 Factorial designs	198	
10.5.1 Continuous outcome	198	
10.5.2 Binary outcome	199	
10.5.3 Sample size calculation for factorial designs	200	
11 The problem of unknown intraclass correlation coefficients11.1 Estimates from previous research11.2 Sample size re-estimation11.3 Bayesian sample size calculation11.4 Maximin optimal designs	 203 204 205 211 214 	
12 Computer software for power calculations	217	
12.1 Introduction	217	
12.2 Computer program SPA-ML	218	
References		

List of figures

$11 \\ 12$
12
13
14
14
23
24
32
34
40
42
58
60
70
72
78
87
93
98

5.4	Efficiency of the cluster randomized individual crossover design relative to the cluster randomized design as a function of ρ_2 and for different values of ρ_1	101
$5.5 \\ 5.6$	Graphical representation of a stepped wedge design The ratio of the number of clusters needed of the stepped wedge versus the posttest design	102 106
$6.1 \\ 6.2 \\ 6.3$	Power as a function of the number of clusters Power as a function of cluster size	117 118
7.1	Relative efficiency of pseudo cluster and individual randomiza- tion compared to that of cluster randomization as a function of the ratio of the between- and within-variance of cluster average n_1q .	119
8.1 8.2	Power as a function of the number of clusters in both treatment conditions and for common cluster sizes 20 and 40 Power as a function of the number of control subjects for dif	149
8.3	Fower as a function of the number of control subjects, for dif- ferent values of the number of clusters in the experimental con- dition and for cluster sizes 20 and 40	152 157
$9.1 \\ 9.2$	Linear relation between time and response Power as a function of number of subjects and number of measurement occasions	161 167
$9.3 \\ 9.4$	Power as a function of study duration and variance ratio Number of subjects and relative costs as a function of number	167
$9.5 \\ 9.6 \\ 9.7$	of measurement occasions in a trial without drop-out Effect of drop-out on statistical power	$170 \\ 173 \\ 175$
9.8	and power	176
9.9	ment occasions per subject and power	177
9.10	Number of subjects and relative costs as a function of the number of measurement occasions.	178
10.1	Graphical representation of $s.e.(\hat{\beta}_1)$ for a three-level cluster ran- domized trial as a function of n_1 and n_2	188

10.2	Graphical representation of <i>s.e.</i> $(\hat{\beta}_1)$ for a multisite cluster randomized trial as a function of n_1 and n_2	191
11.1	The effect of an incorrect a priori intraclass correlation coefficient estimate on power for the test on treatment effect.	205
11.2	Densities of prior distribution of intraclass correlation coefficient, required number of clusters to receive a power level 0.8,	
	and power at 48 clusters.	213
11.3	Graphical derivation of maximin optimal designs	215
12.1	Main menu window of SPA-ML program	220
12.2	Power protocol window of SPA-ML program	221
12.3	Power graph window of SPA-ML program	222

List of tables

$\begin{array}{c} 4.1 \\ 4.2 \\ 4.3 \\ 4.4 \\ 4.5 \end{array}$	Example data for cluster randomized trial \ldots Results based on the multilevel model \ldots Results based on the mixed effects ANOVA model \ldots Results based on the mixed effects ANOVA model \ldots Results based on the mixed effects using optimal allocation ratio r as a function of costs ratio c_E/c_C \ldots Results based on the mixed effects allocation ratio r as a function of costs ratio c_E/c_C \ldots Results based on the mixed effects allocation ratio r as a function of costs ratio c_E/c_C \ldots Results based on the mixed effects allocation ratio r as a function of costs ratio c_E/c_C \ldots Results based on the mixed effects allocation ratio r and	66 67 68 69 75
$ \begin{array}{r} 6.1 \\ 6.2 \\ 6.3 \\ 6.4 \\ 6.5 \end{array} $	Example data for multisite trial	$110 \\ 111 \\ 113 \\ 114 \\ 115$
$9.1 \\ 9.2$	Orthogonal polynomial contrast coefficients Results of simulation study for a longitudinal intervention study with a dichotomous outcome and 160 subjects	164 172
10.1 10.2 10.3	The required number of sites, number of clusters per site or cluster size when the other two are known for a three-level cluster randomized trial	185 190 193
$\begin{array}{c} 11.1\\ 11.2 \end{array}$	Overview of literature that presents ICC estimates \ldots . Empirical type I error rate α and power $1 - \beta$ for standard design and re-estimation design.	206 211
12.1 12.2	Overview of designs for individually randomized trials available in SPA-ML program	223 224

12.3	Overview of designs for multisite trials available in SPA-ML	
	program	225
12.4	Overview of designs for cluster randomized trials available in	
	SPA-ML program	226
12.5	Overview of designs for pseudo cluster trials are available in	
	SPA-ML program	226
12.6	Overview of designs for longitudinal trials available in SPA-ML	
	program	227

Preface

What is the most effective method to lose body weight: a diet, a diet combined with physical exercise or a diet combined with a peer pressure group? Do special training programs increase unemployed people's chances of finding jobs on the labor market? Does a hormone treatment have an effect on the psychosocial functioning of children with growth retardation? Such questions are asked daily, not only by scientists but also by the general public. The similarity of these questions is that they focus on effective methods to delay or prevent disease or some unwanted and undesirable behavior, opinion or condition.

To compare various alternative methods in an experimental setting, it is important to seek a design that results in the highest statistical power at the lowest costs. It is common knowledge that the power level of a statistical test increases with increasing sample size. A sample size that is too small may result in not being able to detect an effect of a treatment, while an excessive sample size may be a waste of time and money of both the researchers and the trial's participants. It is therefore important to carefully calculate the required sample size before a trial is conducted.

Many textbooks and computer programs are of aid to the researcher who is planning a trial. These books and programs mainly apply to trials where the outcome measurements on disease, behavior or condition of a given subject are uncorrelated to those of other subjects in the trial. This assumption is very likely to be violated in trials where subjects are nested within groups, such as in multicenter clinical trials with patients nested within clinics and schoolbased smoking prevention interventions with pupils nested within schools. In the latter example, a pupil's smoking behavior will be influenced not only by the treatment condition he or she is assigned to, but also by the smoking behavior of other pupils in the same school, the teachers' smoking behavior and the school's policy toward smoking.

Ignoring the nested data structure may result in a study that is either under- or overpowered. Moreover, while the sample size affects a trial's power level, the allocation of units is also of importance. In the school-based smoking intervention, for instance, one has to decide whether to sample many schools and include just a few pupils per school, or to sample few schools and include many pupils per school.

In the past decades many journal papers that focus on statistical power analysis and optimal design for trials with nested data structures have appeared. It is now time to compile the published findings in these papers in one single book, and to present a related computer program to perform the power calculations. This is the book that you currently hold in your hands.

The two authors have traveled a long way before we were able to write this book. During this journey various people have been of invaluable support to our careers in the field of applied statistics.

Mirjam Moerbeek. My journey started in spring 1996, when I was appointed PhD researcher at the Department of Methodology and Statistics at Maastricht University, the Netherlands. Under the supervision of Martijn Berger and Gerard van Breukelen I wrote my PhD thesis on the design and analysis of multilevel intervention studies. Part of this research was done in collaboration with Weng Kee Wong and conducted at the Department of Biostatistics at the University of California at Los Angeles. During my PhD research, I frequented meetings of the Netherlands' Multilevel Modeling group, which were organized by Tom Snijders and Cora Maas. Both of them may be considered pioneers of multilevel modeling in the Netherlands.

After receiving my PhD, I continued my research on optimal design and power analysis at the Department of Methodology and Statistics at Utrecht University, the Netherlands. Peter van der Heijden encouraged me to apply for a prestigious research grant from the Netherlands Organization of Scientific Research (NWO). I was awarded a Veni grant in 2003 to extend my research on optimal design for trials with multilevel and longitudinal data, and a Vidi and Aspasia grant in 2008 to study optimal designs for discrete-time survival analysis and to write this book.

The Department of Methodology and Statistics at Utrecht University is the best place to conduct research like this because it has an enthusiastic group of researchers in the field of multilevel analysis, among whom are Joop Hox, Rens van de Schoot and Leoniek Wijngaards-de Meij. Together we organize the biennial International Conference on Multilevel Analysis.

Part of the research in this book was done by former PhD students of mine: Elly Korendijk studied the effects of a misspecification of the intraclass correlation coefficient in the design phase of a cluster randomized trial and Esther Oomen-de Hoop studied the design and analysis of stepped-wedge cluster randomized trials. Charlotte Rietbergen, a former master student of mine, studied crossover in cluster randomized trials. Sander van Schie, also a former master student, studied sample size re-estimation in cluster randomized trials and prepared an overview of papers with estimates of the intraclass correlation coefficient. Katarzyna Jóźwiak wrote the SPA-ML computer program to perform many of the sample size and power calculations that are presented in this book.

Steven Teerenstra. What I like about statisticians is their diversity of origins. My interest in statistics arose when I was looking for a more practically oriented job after finishing my appointment as a PhD student in mathematics in Nijmegen. At that time (2002) I had a freelance job in the Radboud

Preface

University Nijmegen Medical Center to build a research database and I was asked to also think about the statistical analysis of the data. I contacted a former colleague of mine, Wiebe Pestman, who had already travelled the road from mathematics to statistics. He gave me not only statistical advice but also suggested I become a statistician. Without formal statistical training or working experience, my first application outside Nijmegen was not successful. However, opportunities were closer by than I thought. After attending a talk at the department of Epidemiology, Biostatistics and HTA in Nijmegen (now the Department for Health Evidence), I had a chat with Gerhard Zielhuis, who told me that the biostatistics group was looking for a candidate. Under guidance of George Borm, my journey in biostatistics began: hands-on and steep. Because much of my consultancy involved cluster randomized trials, George encouraged me to do research in this topic. A choice that I enjoy till today, because the techniques also gave me insight in other fields of statistics such as longitudinal data and meta-analysis.

While becoming acquainted with the literature in the field, I came across Mirjam's papers. I saw that we had common interests and contacted her to collaborate, which led to our joint endeavors. A research grant by NWO allowed me to work with Anouk Spijker to investigate the efficiency of ANCOVA analysis of cluster randomized trials. With a grant from the Radboud University Nijmegen Medical Center, I engaged Esther Oomen-de Hoop as a PhD student to investigate practical solutions for cluster randomized trials with few clusters.

> Mirjam Moerbeek, Utrecht, the Netherlands Steven Teerenstra, Nijmegen, the Netherlands

1

Introduction

The scientific community performs tens of thousands of experiments each year to improve the treatment of diseases and psychological and social conditions with the aim to increase quality of life and life expectancy. Examples include substance use prevention and cessation interventions, trials that compare cognitive behavioral group therapy to pharmacotherapy for the treatment of social phobia, and trials that evaluate the effects of growth hormones on the psychosocial functioning of children with growth retardation.

The process of data collection can take up to several decades in the case of long-term interventions. In addition, efforts should be paid to good data analysis, including the use of an appropriate statistical model and proper treatment of missing data arising from drop-out and non-response. It is therefore of utmost importance to plan an experiment in such a way that the best treatments are selected with highest probability.

In the design phase of a trial, the treatments to be compared are selected, and choices are made on eligibility criteria, the duration of the study and the best training of professionals who are delivering the treatments to their patients or clients. Another important choice in the design phase is the required number of subjects, as the probability to detect a significant difference between treatment conditions depends on the chosen sample size. This probability is called the statistical power and the calculation of the required sample size to achieve a desired power level is called a power analysis.

In many trials in the social and biomedical sciences, the data have a socalled hierarchical structure, meaning that subjects are nested within groups. For such trials, not only the total number of subjects needs to be determined, but also their allocation over the groups. Should few groups with many subjects each be sampled, or many groups with few subjects each?

Similar questions may be asked for the design of longitudinal trials with repeated measurements over time that are nested within subjects. Should few subjects be sampled and should they be measured often, or should many subjects be sampled and they be measured just a few times?

Guidelines for sample sizes for trials with hierarchical data structures have been studied extensively in the last two decades. The aim of this book is to provide formulae to perform power calculations to social and biomedical researchers and to statisticians working in these fields of science. This first chapter serves as an introduction to basic concepts with respect to experimentation, hierarchical data structures, and study design and further elucidates the need for power calculations. A further description of the aim and contents of the book is given at the end of this chapter.

1.1 Experimentation

This section aims to provide an overview of the basic concepts of experimentation: controlling, randomization, stratification or blocking and matching, replication and blinding. Although these concepts appear obvious nowadays, their formal introduction was initiated less than a hundred years ago by the work of Fisher in the field of agriculture in the 1920s and 1930s (Fisher, 1926, 1935). The first modern randomized controlled trial in health care research is generally considered to be that of the Medical Research Council (1948), that compared the effect of streptomycin and bed rest to that of bed rest only on the treatment of pulmonary tuberculosis. The first randomized controlled trial in social research appears to have been conducted in the late 1920s (Forsetlund, Chalmers, & Bjørndal, 2007). For a more thorough introduction to experimentation the reader is referred to Jadad and Enkin (2007) and Torgerson and Torgerson (2008).

To evaluate the effect of a new type of therapy, surgery, teaching method or any other type of treatment, the measurements on outcome variables of subjects in the new treatment group have to be compared to those receiving a standard treatment, no treatment at all, or those who are assigned to a waiting list. The latter groups of subjects are often called control groups, and the corresponding treatment condition is called the control condition. Including a control is important since it enables the researcher to compare the performance of the new treatment to the performance of a standard treatment or no treatment at all. Without the control group, a researcher is unable to tell whether the performance of the new treatment. In other words, the control serves as a benchmark and enables the estimation of the treatment effect, which is the difference in performance of the two treatment conditions on the outcome variable.

If one wishes to generalize the findings to all subjects within the population, a random draw should be made from this population. To achieve a fair comparison of treatment conditions, the treatment groups should be as similar as possible at baseline with respect to all measured and unmeasured covariates, which are variables that are supposed to have an effect on the outcome. For instance, as parents' smoking behavior has an effect on adolescent smoking behavior, it is important that the baseline percentages of adolescents whose parents smoke are equal over the treatment conditions in a smoking prevention intervention that targets adolescents.

For large sample sizes, this can be achieved by randomly assigning adoles-

cents to treatment conditions and the design is called a completely randomized design. For small sample sizes, random assignment does not guarantee balance with respect to parent smoking behavior. This can be resolved by stratifying adolescents with respect to their parents' smoking behavior: one stratum consisting of adolescents who have one or two parents who smoke, while adolescents in the other stratum do not have a parent who smokes. Randomization to treatment conditions is then done within strata and the design is called a stratified randomized design. Alternative terminology is to use blocking, blocks and randomized block design wording instead of stratifying, strata and stratified randomized block design is more often used for factors that are generally determinable and often controllable in the experiment and for which the sample size can be determined upfront, and called blocking factors (e.g., temperature).

The terms stratifying, strata and stratified randomized design apply more often to factors that cannot be determined upfront (e.g., concomitant medication use). Often the distinction is not clear-cut. An option that is more useful when multiple covariates are present is the matched-pair design. With this design, pairs of subjects who are as similar as possible with respect to all covariates are formed, and randomization is done within each pair. Special algorithms are available to balance treatment groups with respect to multiple covariates; see for instance G. F. Borm, Hoogendoorn, Den Heijer, and Zielhuis (2005). If such balancing methods have not performed well, or if these methods were not used at all, then the method of statistical control can be used to correct for between-group differences. This is achieved by including important covariates as predictor variables in a regression or analysis of covariance model.

Whichever strategy is chosen to prevent imbalance or to correct for it, it is important that relevant covariates are defined in the planning of a trial to be able to measure them in the data collection phase. It is therefore necessary to search the literature for relevant covariates prior to data collection.

Several replications should be made to estimate the variability of the error terms in a statistical model. That is, more than one subject should be available within each combination of treatment condition and covariate values. In a trial that randomizes to treatment conditions within gender, for instance, at least two males and two females are required in the control condition, and another two males and two females in the experimental condition. Replication should not be confused with taking repeat measurements on the same subjects. If repeated measurements over time are taken on each subject, then time should be included as a predictor in the statistical model.

In an ideal case, the trial is double-blind, meaning that neither the researchers nor the subjects know to which treatment condition each subject is assigned. Only a third party has access to the key that identifies to which condition each subject is randomized. With single blinding only the study participants are unaware to which treatment condition they are assigned. Double blinding is a means to lessen observer bias that occurs when the results are influenced by conscious or unconscious bias on the part of the observer. Such bias occurs when observers tend to adjust their scores on a subject's outcome variables to their expectations or judgments of the (known) treatment condition to which a subject is assigned.

Another reason for double blinding is to reduce selection bias which may occur when recruiters know beforehand the treatment to which a subject is assigned. The risk of both biases is illustrated by the finding that the outcomes and characteristics of included participants may be different in blinded trials than in unblinded trials (Veerus, Fischer, Hakama, & Hemminki, 2012) and hence blinding may influence both internal and external validity.

Unfortunately, blinding is not always an option. Double blinding is often used in pharmaceutical trials where different substances are compared. The pills or injections containing these substances should be as similar with respect to size, color, smell, weight and so forth as possible. In social and behavioral trials, treatment often relies on interpersonal interactions, such as risk-reduction sessions, peer pressure groups and training programs and blinding is often not possible.

The degree to which a researcher has control over the environment in which the experiment is conducted determines whether the experiment is a laboratory or field experiment. In a laboratory (pure) experiment, the control over the environment is maximal and all subjects are exposed to the same influences expect for the treatment condition to which they are assigned. Consider as an example a trial that investigates the effect of the opponent's ethnicity on the aggressiveness of players of dual-player computer games. Participants are assigned to an opponent of the same or different ethnicity and all other sources of variation are kept under the experimenter's control. Thus, the participants play the same computer game in rooms with the same background music, temperature, lighting, accessibility to food and drinks and so forth.

Laboratory experiments are often artificial and doubts on generalization of results to the real world exist. In a field experiment, the daily life environment is treated as the laboratory. Randomization to treatments is under experimental control but the experimenter does not have an influence on uncontrolled events. This is not considered a problem as long as both treatment groups are equally exposed to uncontrolled events. However, a problem arises when the one treatment group is exposed to uncontrolled events to a higher degree than the other group and/or when the uncontrolled events interact with treatment condition. In both situations, the external factors influence the one group more than the other, thereby threatening the validity of the study. Consider as an example the effects of exercise on the reduction of body weight in a controlled trial. A difference in body weight between both treatment groups can only be attributed to the effects of exercise if the calorie intake is equal across groups.

1.1.1 Problems with random assignment

Although randomization is a strong tool to achieve treatment groups that are equal with respect to relevant covariates, it also has some drawbacks. In the ideal case, randomization is done by the researcher or an independent statistician, but in practice it is often done by those who control access to potential participants, such as general practitioners, therapists and school principals. They often tend to assign the participants they consider most deserving or most likely to benefit to the new treatment group, a process which is called selection bias.

People who are recruited to participate in an experiment are often only willing to do so if they are assigned to the interesting and promising new treatment. Some of them will even try to force the person performing the randomization to assign them to the new treatment, resulting in a treatment group that consists of more demanding or more assertive subjects than the control condition, a situation that is most undesirable when the outcome variable is related to these qualities of character.

When participants cannot influence the randomization procedure, they can still undermine the study. For instance, those assigned to the control condition can try to benefit as much from the new treatment as possible by contacting participants in the new treatment group. This is especially a threat if the new treatment condition consists of oral or written information to improve lifestyles, and less so when the new condition consists of medication or surgery. This process is called control group contamination since information on the new treatment somehow leaks to participants in the control group. In addition, drop-out rates among those in the control group are often higher than those in the new treatment group and disappointment may be common after allocation to the control condition (Lindström, Sundberg-Petersson, Adami, & Tönnesen, 2010).

Resentment of those in the control condition can be lessened by using a placebo, which is an irrelevant treatment in the control group, such as a sugar pill in pharmaceutical trials. As an alternative, one can inform the participants in the control group on the necessity of random assignment and the use of a control group and assure them of having one of the more desirable treatments at the end of the trial, provided these treatments have a relevant effect and do not show harmful side effects.

A related problem occurs when participants in the new treatment condition do not wish to participate if they find this new treatment undesirable, if they expect its effects are negligible or if they expect unwanted or harmful side effects. For instance, parents who smoke might not be interested in their children participating in a smoking prevention intervention. This problem may be solved to only include those participants who have expressed their willingness to participate. As is obvious, this may induce problems with respect of generalization to the general population.

In some cases randomization is not possible due to ethical aspects. For

instance, a general practitioner might not be willing to let only half of his patients benefit from the new and promising treatment. In other cases, one cannot force subjects to adhere to the treatment assigned. For instance, one cannot force subjects to smoke or not in a trial on the effects of smoking on lung cancer. In such cases, treatments can only be compared on basis of a quasi-experimental design. With this design, the experimenter cannot randomly assign subjects to treatment conditions and must rely on existing differences between people with respect to the variable of interest. The pitfall is that this variable may be correlated with covariates that also have an effect on the outcome and an exhaustive literature review to detect all such covariates is therefore required.

The sample size formulae that are presented in the next chapters are valid for experimental designs where randomization is done in such a way that the treatment groups are comparable with respect to covariates. In Section 5.2 the effects of matching and pre-stratification on power are treated in further detail.

1.2 Hierarchical data structures

In many experiments in the social and biomedical sciences, the data have a so-called hierarchical structure. This means that subjects are nested within groups that may themselves be nested within higher order groups, and so on. The words *nested*, *clustered* and *multilevel* are often used as synonyms to the word *hierarchical*. Examples are clinical trials with patients nested within clinics, school-based smoking and substance abuse prevention interventions with pupils nested within clinics nested within schools, and studies in the field of clinical psychology with clients nested within therapists.

Hierarchical data structures arise from multistage sampling where sampling is done in a number of successive steps. For instance, therapists are sampled from a population of therapists in the first step. In the second step, a sample of clients is taken from therapists selected in the first step. Only those patients whose therapists were selected in the first step have a chance of enrolling into the trial; the changes of the other patients decline to zero once the first step is executed. Thus, selection probabilities are not constant when multistage sampling is used. The sampling scheme for two-stage sampling is illustrated in Figure 1.1. The large circles represent the therapists; the small circles represent their clients. The black large circles are therapists who are not selected in the first step, and neither are any of their patients in the second step, as represented by the small black circles. The white large circles represent therapists selected in the first step. For each of these, some of his or her patients are selected in the second step and some of them are not, as represented by the small white and black circles, respectively.



Figure 1.1: Graphical representation of a two-stage sampling scheme.

The advantage of the two-stage sampling scheme is that it is often less expensive than taking a simple random sample where each patient has equal changes of being selected. This is explained by the fact that the number of therapists is often lower if multistage sampling is used, meaning that fewer therapists have to be communicated with and fewer therapists have to be trained to deliver the new type of therapy to their clients. In addition, the patients of the same therapists are geographically organized, which may reduce travel costs.

Similar arguments hold for longitudinal studies where repeated measurements are nested within subjects. With such studies it may be more convenient to follow the same subjects over time than to take a new sample of subjects at each time point. In other words, a cohort design may be more convenient than a cross-sectional design.

The other side of the coin is that data analysis and power calculations for trials with hierarchical data are generally more complicated than those for trials without nesting because measurements of subjects within the same group are likely to be correlated. In other words, a subject's measurements on health, behavior, attitude, opinion and other relevant outcomes cannot be regarded as independent of those from other subjects within the same group. Such dependent data may arise from therapist effects, which occur when the success of treatment depends on the skill and training of the person delivering the treatment, or on the quality of the equipment that is used by this person. This is often the case in trials that rely on surgery, interviewing or physical or