STATISTICAL DATA ANALYSIS FOR THE PHYSICAL SCIENCES

Data analysis lies at the heart of every experimental science. Providing a modern introduction to statistics, this book is ideal for undergraduates in physics. It introduces the necessary tools required to analyse data from experiments across a range of areas, making it a valuable resource for students.

In addition to covering the basic topics, the book also takes in advanced and modern subjects, such as neural networks, decision trees, fitting techniques and issues concerning limit or interval setting. Worked examples and case studies illustrate the techniques presented, and end-of-chapter exercises help test the reader's understanding of the material.

ADRIAN BEVAN is a Reader in Particle Physics in the School of Physics and Astronomy, Queen Mary, University of London. He is an expert in quark flavour physics and has been analysing experimental data for over 15 years. Cambridge University Press 978-1-107-03001-5 - Statistical Data Analysis for the Physical Sciences Adrian Bevan Frontmatter More information Cambridge University Press 978-1-107-03001-5 - Statistical Data Analysis for the Physical Sciences Adrian Bevan Frontmatter <u>More information</u>

STATISTICAL DATA ANALYSIS FOR THE PHYSICAL SCIENCES

ADRIAN BEVAN Queen Mary, University of London



 $\ensuremath{\mathbb{C}}$ in this web service Cambridge University Press

CAMBRIDGE UNIVERSITY PRESS Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore, São Paulo, Delhi, Mexico City

Cambridge University Press The Edinburgh Building, Cambridge CB2 8RU, UK

Published in the United States of America by Cambridge University Press, New York

www.cambridge.org Information on this title: www.cambridge.org/9781107670341

© A. Bevan 2013

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2013

Printed and bound in the United Kingdom by the MPG Books Group

A catalogue record for this publication is available from the British Library

ISBN 978-1-107-03001-5 Hardback ISBN 978-1-107-67034-1 Paperback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication, and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

Contents

	Preface		
1	 Introduction 1.1 Measuring g, the coefficient of acceleration due to gravity 1.2 Verification of Ohm's law 1.3 Measuring the half-life of an isotope 1.4 Summary 	1 1 5 7 10	
2	Sets2.1Relationships between sets2.2SummaryExercises	12 13 17 18	
3	Probability3.1Elementary rules3.2Bayesian probability3.3Classic approach3.4Frequentist probability3.5Probability density functions3.6Likelihood3.7Case studies3.8SummaryExercises	20 21 21 24 25 26 27 27 32 33	
4	 Visualising and quantifying the properties of data 4.1 Visual representation of data 4.2 Mode, median, mean 4.3 Quantifying the spread of data 4.4 Presenting a measurement 4.5 Skew 	35 35 37 39 41 43	

vi	Contents			
	 4.6 Measurements of more than one observable 4.7 Case study 4.8 Summary Exercises 	44 52 53 53		
		55		
5	Useful distributions5.1Expectation values of probability density functions5.2Binomial distribution5.3Poisson distribution5.4Gaussian distribution5.5 χ^2 distribution5.6Computational issues5.7SummaryExercises	56 57 57 62 65 67 68 70 70		
	Exercises	70		
6	 Uncertainty and errors 6.1 The nature of errors 6.2 Combination of errors 6.3 Binomial error 6.4 Averaging results 6.5 Systematic errors and systematic bias 6.6 Blind analysis technique 6.7 Case studies 6.8 Summary Exercises 	72 72 75 79 81 82 84 85 90 91		
7	 Confidence intervals 7.1 Two-sided intervals 7.2 Upper and lower limit calculations 7.3 Limits for a Gaussian distribution 7.4 Limits for a Poisson distribution 7.5 Limits for a binomial distribution 7.6 Unified approach to analysis of small signals 7.7 Monte Carlo method 7.8 Case studies 7.9 Summary Exercises 	93 93 94 96 98 100 101 105 106 111 112		
8	 Hypothesis testing 8.1 Formulating a hypothesis 8.2 Testing if the hypothesis agrees with data 8.3 Testing if the hypothesis disagrees with data 	114 114 115 117		

	Contents		
	8.4 8.5 8.6 8.7 8.8 Exerc	Hypothesis comparison Testing the compatibility of results Establishing evidence for, or observing a new effect Case studies Summary ises	117 119 120 124 125 126
9	Fitting 9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Exerc	Optimisation The least squares or χ^2 fit Linear least-squares fit Maximum-likelihood fit Combination of results Template fitting Case studies Summary ises	128 128 131 134 136 140 142 142 150 151
10	Multi 10.1 10.2 10.3 10.4 10.5 10.6 10.7 10.8 Exerc	variate analysis Cutting on variables Bayesian classifier Fisher discriminant Artificial neural networks Decision trees Choosing an MVA technique Case studies Summary ises	153 154 157 158 162 169 171 174 177
	Apper Apper Apper Apper Apper Refere	adix A Glossary adix B Probability density functions adix C Numerical integration methods adix D Solutions adix E Reference tables	181 186 198 201 207 216
	Index		218

Cambridge University Press 978-1-107-03001-5 - Statistical Data Analysis for the Physical Sciences Adrian Bevan Frontmatter More information

Preface

The foundations of science are built upon centuries of careful observation. These constitute measurements that are interpreted in terms of hypotheses, models, and ultimately well-tested theories that may stand the test of time for only a few years or for centuries. In order to understand what a single measurement means we need to appreciate a diverse range of statistical methods. Without such an appreciation it would be impossible for scientific method to turn observations of nature into theories that describe the behaviour of the Universe from sub-atomic to cosmic scales. In other words science would be impracticable without statistical data analysis. The data analysis principles underpinning scientific method pervade our everyday lives, from the use of statistics we are subjected to through advertising to the smooth operation of SPAM filters that we take for granted as we read our e-mail. These methods also impact upon the wider economy, as some areas of the financial industry use data mining and other statistical techniques to predict trading performance or to perform risk analysis for insurance purposes.

This book evolved from a one-semester advanced undergraduate course on statistical data analysis for physics students at Queen Mary, University of London with the aim of covering the rudimentary techniques required for many disciplines, as well as some of the more advanced topics that can be employed when dealing with limited data samples. This has been written by a physicist with a non-specialist audience in mind. This is not a statistics book for statisticians, and references have been provided for the interested reader to refer to for more rigorous treatment of the techniques discussed here. As a result this book provides an up-to-date introduction to a wide range of methods and concepts that are needed in order to analyse data. Thus this book is a mixture of a traditional text book approach and a teach by example approach. By providing these opposing viewpoints it is hoped that the reader will find the material more accessible. Throughout the book, a number of case studies are presented with possible solutions discussed in detail. The purpose of these sections is to consolidate the more abstract notions discussed in the book and

х

Preface

apply them to an example. In some instances the case study may appear somewhat abstract and specific to scientific research; however, where possible more widely applicable problems have been included. At the end of each chapter there is a summary of the main issues raised, followed by a number of example questions to help the reader practise and gain a deeper understanding of the material included. Solutions to questions are presented at the end of the book.

The Introduction motivates the importance of studying statistical methods when analysing data by looking at three common problems encountered early within the life of a physicist: measuring *g*, testing Ohm's law and studying the law of radioactive decay. Following this motivational introduction the book is divided into two parts: (i) the foundations of statistical data analysis from set notation through to confidence intervals, and (ii) discussion of more advanced topics in the form of optimisation, fitting, and data mining. The material in the first part of the book is ordered logically so that successive sections build on material discussed in the earlier ones, while the second part of the book contains stand alone chapters that depend on concepts developed in the first part. These later chapters can be read in any order.

The first part of this book starts with an introduction to sets and Venn diagrams that provide some of the language that we use to discuss data. Having developed this language, the concept of probability is formally introduced in Chapter 3. Readers who are familiar with these concepts already may wish to skip over the first two chapters and proceed straight to the discussion in Chapter 4 on how to visualise and quantify data. Distributions of data are often described by simple functions that are used to represent the probability of observing data with a certain value. A number of useful distributions are described in Chapter 5, and Appendix B builds on this topic by discussing a number of additional functions that may be of use. Measurements are based on the determination of some central value of an observable quantity, with an uncertainty or error on that observable. Issues surrounding uncertainties and errors are introduced in Chapter 6, and this topic is further developed in Chapter 7. Chapter 8 discusses hypothesis testing and brings together many of the earlier concepts in the book.

The second part of the book presents more advanced topics. Chapter 9 discusses fitting data given some assumed model using χ^2 and likelihood methods. This relies heavily on concepts developed in Chapters 5 and 6, and Appendix B. Chapter 10 discusses data mining, or how to efficiently separate two classes of data, for example signal from background using numerical methods. The methods discussed include the use of 'cut-based' selection, the Bayesian classifier, Fisher's linear discriminant, artificial neural networks, and decision trees.

To avoid interrupting the flow of the text, a number of detailed appendices have been prepared. The most important of these appendices is a collection of probability

Preface

tables, which is conveniently located at the end of the book in order to provide a quick reference to the reader. There is also a glossary of terms intended to help the reader when referring back to the book some time after an initial reading. Appendices listing a number of commonly used probability density functions, and elementary numerical integration techniques have also been provided. While these are not strictly required in order to understand the concepts introduced in the book, they have been included in order to make this a more complete resource for readers who wish to study this topic beyond an undergraduate course.

There are a number of technical terms introduced throughout this book. When a new term is introduced, that term is highlighted in *bold-italic text* to help the reader refer back to this description at a later time.

I would like to thank colleagues who have provided me with feedback on the draft of this book, and in particular Peter Crew.