

SCIENTIFIC INFERENCE

Providing the knowledge and practical experience to begin analysing scientific data, this book is ideal for physical sciences students wishing to improve their data handling skills.

The book focuses on explaining and developing the practice and understanding of basic statistical analysis, concentrating on a few core ideas, such as the visual display of information, modelling using the likelihood function, and simulating random data.

Key concepts are developed through a combination of graphical explanations, worked examples, example computer code and case studies using real data. Students will develop an understanding of the ideas behind statistical methods and gain experience in applying them in practice. Further resources are available at www.cambridge.org/9781107607590, including data files for the case studies so students can practice analysing data, and exercises to test students' understanding.

SIMON VAUGHAN is a Reader in the Department of Physics and Astronomy, University of Leicester, where he has developed and runs a highly regarded course for final year physics students on the subject of statistics and data analysis.

SCIENTIFIC INFERENCE

Learning from data

SIMON VAUGHAN

University of Leicester





CAMBRIDGE
 UNIVERSITY PRESS

Shaftesbury Road, Cambridge CB2 8EA, United Kingdom
 One Liberty Plaza, 20th Floor, New York, NY 10006, USA
 477 Williamstown Road, Port Melbourne, VIC 3207, Australia
 314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi – 110025, India
 103 Penang Road, #05–06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of Cambridge University Press & Assessment,
 a department of the University of Cambridge.

We share the University's mission to contribute to society through the pursuit of
 education, learning and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781107024823

© S. Vaughan 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 & Assessment.

First published 2013

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

Library of Congress Cataloging-in-Publication data

Vaughan, Simon, 1976– author.

Scientific inference : learning from data / Simon Vaughan.

pages cm

Includes bibliographical references and index.

ISBN 978-1-107-02482-3 (hardback) – ISBN 978-1-107-60759-0 (paperback)

1. Mathematical statistics – Textbooks. I. Title.

QA276.V34 2013

519.5 – dc23 2013021427

ISBN 978-1-107-02482-3 Hardback

ISBN 978-1-107-60759-0 Paperback

Cambridge University Press & Assessment 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.

For my family

Contents

| | <i>page x</i> |
|---|---------------|
| <i>For the student</i> | xii |
| <i>For the instructor</i> | xii |
| 1 Science and statistical data analysis | 1 |
| 1.1 Scientific method | 1 |
| 1.2 Inference | 3 |
| 1.3 Scientific inference | 6 |
| 1.4 Data analysis in a nutshell | 7 |
| 1.5 Random samples | 8 |
| 1.6 Know your data | 10 |
| 1.7 Language | 11 |
| 1.8 Statistical computing using R | 12 |
| 1.9 How to use this book | 12 |
| 2 Statistical summaries of data | 14 |
| 2.1 Plotting data | 14 |
| 2.2 Plotting univariate data | 16 |
| 2.3 Centre of data: sample mean, median and mode | 18 |
| 2.4 Dispersion in data: variance and standard deviation | 21 |
| 2.5 Min, max, quantiles and the five-number summary | 24 |
| 2.6 Error bars, standard errors and precision | 25 |
| 2.7 Plots of bivariate data | 28 |
| 2.8 The sample correlation coefficient | 36 |
| 2.9 Plotting multivariate data | 38 |
| 2.10 Good practice in statistical graphics | 43 |
| 2.11 Chapter summary | 44 |
| 3 Simple statistical inferences | 46 |
| 3.1 Inference about the mean of a sample | 46 |
| | vii |

| | | |
|------|---|-----|
| viii | <i>Contents</i> | |
| | 3.2 Difference in means from two samples | 49 |
| | 3.3 Straight line fits | 51 |
| | 3.4 Linear regression in practice | 56 |
| | 3.5 Residuals: what lies beneath | 58 |
| | 3.6 Case study: regression of Reynolds' data | 59 |
| | 3.7 Chapter summary | 63 |
| 4 | Probability theory | 64 |
| | 4.1 Experiments, outcomes and events | 64 |
| | 4.2 Probability | 69 |
| | 4.3 The rules of the probability calculus | 72 |
| | 4.4 Random variables | 82 |
| | 4.5 The visual perception of randomness | 89 |
| | 4.6 The meaning of 'probability' and 'random' | 89 |
| | 4.7 Chapter summary | 92 |
| 5 | Random variables | 94 |
| | 5.1 Properties of random variables | 94 |
| | 5.2 Discrete random variables | 100 |
| | 5.3 Continuous random variables | 110 |
| | 5.4 Change of variables | 116 |
| | 5.5 Approximate variance relations (or the propagation of errors) | 120 |
| | 5.6 Chapter summary | 122 |
| 6 | Estimation and maximum likelihood | 124 |
| | 6.1 Models | 124 |
| | 6.2 Case study: Rutherford & Geiger data | 125 |
| | 6.3 Maximum likelihood estimation | 129 |
| | 6.4 Weighted least squares | 133 |
| | 6.5 Case study: pion scattering data | 139 |
| | 6.6 Chapter summary | 140 |
| 7 | Significance tests and confidence intervals | 142 |
| | 7.1 A thought experiment | 142 |
| | 7.2 Significance testing and test statistics | 143 |
| | 7.3 Pearson's χ^2 test | 146 |
| | 7.4 Fixed-level tests and decisions | 153 |
| | 7.5 Interpreting test results | 156 |
| | 7.6 Confidence intervals on MLEs | 159 |
| | 7.7 Chapter summary | 166 |

Contents

ix

| | | |
|---|--|-----|
| 8 | Monte Carlo methods | 169 |
| | 8.1 Generating pseudo-random numbers | 169 |
| | 8.2 Estimating sampling distributions by Monte Carlo | 175 |
| | 8.3 Computing confidence by <i>bootstrap</i> | 181 |
| | 8.4 The power of Monte Carlo | 183 |
| | 8.5 Further reading | 184 |
| | 8.6 Chapter summary | 184 |
| Appendix A Getting started with statistical computation | | 185 |
| | A.1 What is R? | 185 |
| | A.2 A first R session | 185 |
| | A.3 Entering data | 187 |
| | A.4 Quitting R | 188 |
| | A.5 More mathematics | 188 |
| | A.6 Writing your own R scripts | 189 |
| | A.7 Producing graphics in R | 190 |
| | A.8 Saving graphics in R | 192 |
| | A.9 Good practice with R | 193 |
| Appendix B Data case studies | | 195 |
| | B.1 Michelson's speed of light data | 195 |
| | B.2 Rutherford–Geiger radioactive decay | 196 |
| | B.3 A study of fluid flow | 198 |
| | B.4 The HR diagram | 199 |
| | B.5 A particle physics experiment | 202 |
| | B.6 Atmospheric conditions in New York City | 205 |
| Appendix C Combinations and permutations | | 207 |
| | C.1 Permutations | 207 |
| | C.2 Combinations | 208 |
| | C.3 Probability of combinations | 209 |
| Appendix D More on confidence intervals | | 210 |
| Appendix E Glossary | | 214 |
| Appendix F Notation | | 219 |
| | <i>References</i> | 221 |
| | <i>Index</i> | 223 |

For the student

Science is not about certainty, it is about dealing rigorously with uncertainty. The tools for this are statistical. Statistics and data analysis are therefore an essential part of the scientific method and modern scientific practice, yet most students of physical science get little explicit training in statistical practice beyond basic error handling. The aim of this book is to provide the student with both the knowledge and the practical experience to begin analysing new scientific data, to allow progress to more advanced methods and to gain a more statistically literate approach to interpreting the constant flow of data provided by modern life.

More specifically, if you work through the book you should be able to accomplish the following.

- Explain aspects of the scientific method, types of logical reasoning and data analysis, and be able to critically analyse statistical and scientific arguments.
- Calculate and interpret common quantitative and graphical statistical summaries.
- Use and interpret the results of common statistical tests for difference and association, and straight line fitting.
- Use the calculus of probability to manipulate basic probability functions.
- Apply and interpret model fitting, using e.g. least squares, maximum likelihood.
- Evaluate and interpret confidence intervals and significance tests.

Students have asked me whether this is a book about statistics or data analysis or statistical computing. My answer is that they are so closely connected it is difficult to untangle them, and so this book covers areas of all three.

The skills and arguments discussed in the book are highly transferable: statistical presentations of data are used throughout science, business, medicine, politics and the news media. An awareness of the basic methods involved will better enable you to use and critically analyse such presentations – this is sometimes called *statistical literacy*.

For the student

xi

In order to understand the book, you need to be familiar with the mathematical methods usually taught in the first year of a physics, engineering or chemistry degree (differential and integral calculus, basic matrix algebra), but this book is designed so that the probability and statistics content is entirely self-contained.

For the instructor

This book was written because I could not find a suitable textbook to use as the basis of an undergraduate course on scientific inference, statistics and data analysis. Although there are good books on different aspects of introductory statistics, those intended for physicists seem to target a post-graduate audience and cover either too much material or too much detail for an undergraduate-level first course. By contrast, the ‘Intro to stats’ books aimed at a broader audience (e.g. biologists, social scientists, medics) tend to cover topics that are not so directly applicable for physical scientists. And the books aimed at mathematics students are usually written in a style that is inaccessible to most physics students, or in a recipe-book style (aimed at science students) that provides ready-made solutions to common problems but develops little understanding along the way.

This book is different. It focuses on explaining and developing the *practice* and *understanding* of basic statistical analysis, concentrating on a few core ideas that underpin statistical and data analysis, such as the visual display of information, modelling using the likelihood function, and simulating random data. Key concepts are developed using several approaches: verbal exposition in the main text, graphical explanations, case studies drawn from some of history’s great physics experiments, and example computer code to perform the necessary calculations.¹ The result is that, after following all these approaches, the student should both understand the ideas behind statistical methods and have experience in applying them in practice.

The book is intended for use as a textbook for an introductory course on data analysis and statistics (with a bias towards students in physics) or as self-study companion for professionals and graduate students. The book assumes familiarity with calculus and linear algebra, but no previous exposure to probability or statistics

¹ These are based on R, a freely available software package for data analysis and statistics and used in many statistics textbooks.

is assumed. It is suitable for a wide range of undergraduate and postgraduate science students.

The book has been designed with several special features to improve its value and effectiveness with students:

- several complete data analysis case studies using real data from some of history's great experiments
- 'example boxes' – approximately 20 boxes throughout the text that give specific, worked examples for concepts as they are discussed
- 'computer practice boxes' – approximately 90 boxes throughout the text that give working R code to perform the calculations discussed in the text or produce the plots shown
- graphical explanations of important concepts
- appendices that provide technical details supplementary to the main text
- a well-populated glossary of terms and list of notational conventions.

The emphasis on a few core ideas and their practical applications means that some subjects usually covered in introductory statistics texts are given little or no treatment here. Rigorous mathematical proofs are not covered – the interested reader can easily consult any good reference work on probability theory or mathematical statistics to check these. In addition, we do not cover some topics of 'classical' statistics that are dealt with in other introductory works. These topics include

- more advanced distribution functions (beta, gamma, multinomial, . . .)
- ANOVA and the generalised linear model
- characteristic functions and the theory of moments
- decision and information theories
- non-parametric tests
- experimental design
- time series analysis
- multivariate analysis (principal components, clustering, . . .)
- survival analysis
- spatial data analysis.

Upon completion of this book the student should be in a much better position to understand any of these topics from any number of more advanced or comprehensive texts.

Perhaps the 'elephant in the room' question is: what about Bayesian methods? Unfortunately, owing to practical limitations there was not room to include full chapters developing Bayesian methods. I hope I have designed the book in such a way that it is not wholly frequentist or Bayesian. The emphasis on model fitting

using the likelihood function (Chapter 6) could be seen as the first step towards a Bayesian analysis (i.e. implicitly using flat priors and working towards the posterior mode). Fortunately, there are many good books on Bayesian data analysis that can then be used to develop Bayesian ideas explicitly. I would recommend Gelman *et al.* (2003) generally and Sivia and Skilling (2006) or Gregory (2005) for physicists in particular. Albert (2007) also gives a nice ‘learn as you compute’ introduction to Bayesian methods using R.