**Core Statistics**

Based on two courses for new graduate students, Core Statistics provides concise coverage of the fundamentals of inference for parametric statistical models, including both theory and practical numerical computation. The book considers both frequentist maximum likelihood estimation and Bayesian stochastic simulation, focusing on general methods applicable to a wide range of models and emphasizing the common questions addressed by the two approaches.

This compact book aims to cover the core knowledge needed by beginning graduate students in statistical subjects, at a level suitable for those going on to develop new methods or undertaking novel applications of statistical modelling: Bayesian and frequentist approaches to modelling and inference; practical computational implementation, including numerical issues; brief coverage of some essential probability; and the essentials of R as a statistical programming language. Aimed also at any quantitative scientist who uses statistical methods, this book will deepen readers' understanding of why and when methods work and explain how to develop suitable methods for non-standard situations, such as in ecology, big data analysis and genomics.

SIMON N. WOOD works as a professor of statistics at the University of Bath, with particular interests in statistical computing, methodology of smoothing and environmental statistics.

INSTITUTE OF MATHEMATICAL STATISTICS
TEXTBOOKS

*Editorial Board*
D. R. Cox (University of Oxford)
B. Hambly (University of Oxford)
S. Holmes (Stanford University)
X.-L. Meng (Harvard University)

IMS Textbooks give introductory accounts of topics of current concern suitable for advanced courses at master's level, for doctoral students and for individual study. They are typically shorter than a fully developed textbook, often arising from material created for a topical course. Lengths of 100–290 pages are envisaged. The books typically contain exercises.

Other Books in the Series

1. *Probability on Graphs*, by Geoffrey Grimmett
2. *Stochastic Networks*, by Frank Kelly and Elena Yudovina
3. *Bayesian Filtering and Smoothing*, by Simo Särkkä
4. *The Surprising Mathematics of Longest Increasing Subsequences*, by Dan Romik
5. *Noise Sensitivity of Boolean Functions and Percolation*, by Christophe Garban and Jeffrey E. Steif

# Core Statistics

SIMON N. WOOD
*University of Bath*

CAMBRIDGE
UNIVERSITY PRESS

# Contents

v

*Contents*

*Contents* vii

# Preface

This book is aimed at the numerate reader who has probably taken an introductory statistics and probability course at some stage and would like a brief introduction to the core methods of statistics and how they are applied, not necessarily in the context of standard models. The first chapter is a brief review of some basic probability theory needed for what follows. Chapter 2 discusses statistical models and the questions addressed by statistical inference and introduces the maximum likelihood and Bayesian approaches to answering them. Chapter 3 is a short overview of the R programming language. Chapter 4 provides a concise coverage of the large sample theory of maximum likelihood estimation, and Chapter 5 discusses the numerical methods required to use this theory. Chapter 6 covers the numerical methods useful for Bayesian computation, in particular Markov chain Monte Carlo. Chapter 7 provides a brief tour of the theory and practice of linear modelling. Appendices then cover some useful information on common distributions, matrix computation and random number generation. The book is neither an encyclopedia nor a cookbook, and the bibliography aims to provide a compact list of the most useful sources for further reading, rather than being extensive. The aim is to offer a concise coverage of the core knowledge needed to understand and use parametric statistical methods and to build new methods for analysing data. Modern statistics exists at the interface between computation and theory, and this book reflects that fact. I am grateful to Nicole Augustin, Finn Lindgren, the editors at Cambridge University Press and the students on the Bath course 'Applied Statistical Inference' and the Academy for PhD Training in Statistics course 'Statistical Computing' for many useful comments.