

V. BARNETT
T. LEWIS

OUTLIERS IN STATISTICAL DATA

3RD EDITION

WILEY SERIES IN PROBABILITY
AND MATHEMATICAL STATISTICS



WILEY

0212

B261

E-3

9462991

Outliers in Statistical Data

Third Edition



E9462991

VIC BARNETT

Rothamsted Experimental Station, UK

and

TOBY LEWIS

University of East Anglia, UK



JOHN WILEY & SONS

Chichester • New York • Brisbane • Toronto • Singapore

Copyright © 1978, 1984, 1994 by John Wiley & Sons Ltd,
Baffins Lane, Chichester,
West Sussex P019 1UD, England
National Chichester (0243) 779777
International +44 243 779777

All rights reserved.

No part of this book may be reproduced by any means,
or transmitted, or translated into a machine language
without the written permission of the publisher.

Other Wiley Editorial Offices

John Wiley & Sons, Inc., 605 Third Avenue,
New York, NY 10158-0012, USA

Jacaranda Wiley Ltd, 33 Park Road, Milton,
Queensland 4064, Australia

John Wiley & Sons (Canada) Ltd, 22 Worcester Road,
Rexdale, Ontario M9W 1L1, Canada

John Wiley & Sons (SEA) Pte Ltd, 37 Jalan Pemimpin #05-04,
Block B, Union Industrial Building, Singapore 2057

Library of Congress Cataloging-in-Publication Data

Barnett, Vic.

Outliers in statistical data / Vic Barnett and Toby Lewis. — 3rd ed.
p. cm.

Includes bibliographical references and index.

ISBN 0 471 93094 6

I. Outliers (Statistics) I. Lewis, Toby. II. Title.

QA276.B2849 1994

519.5'2—dc20

93-29289

CIP

British Library Cataloguing in Publication Data

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

ISBN 0 471 93094 6

Typeset in 10/12pt Times by Pure Tech Corporation, Pondicherry, India
Printed and bound in Great Britain by Biddles Ltd,
Guildford and King's Lynn

Outliers in Statistical Data

Third Edition

Preface to First Edition

The concept of an outlier has fascinated experimentalists since the earliest attempts to interpret data. Even before the formal development of statistical method, argument raged over whether, and on what basis, we should discard observations from a set of data on the grounds that they are 'unrepresentative', 'spurious', or 'mavericks' or 'rogues'. The early emphasis stressed the contamination of the data by unanticipated and unwelcome errors or mistakes affecting some of the observations. Attitudes varied from one extreme to another: from the view that we should never sully the sanctity of the data by daring to adjudge its propriety, to an ultimate pragmatism expressing 'if in doubt, throw it out'.

The present views are more sophisticated. A wider variety of aims are recognized in the handling of outliers, outlier-generating models have been proposed, and there is now available a vast array of specific statistical techniques for processing outliers. The work is scattered throughout the literature of the present century, shows no sign of any abatement, but has not previously been drawn together in a comprehensive review. Our purpose in writing this book is to attempt to provide such a review, at two levels. On the one hand we seek to survey the existing state of knowledge in the outlier field and to present the details of selected procedures for different situations. On the other hand we attempt to categorize differences in attitude, aim, and model in the study of outliers, and to follow the implications of such distinctions for the development of new research approaches. In offering such a comprehensive overview of the principles and methods associated with outliers we hope that we may help the practitioner in the analysis of data and the researcher in opening up possible new avenues of enquiry.

Early work on outliers was (inevitably) characterized by lack of attention to the modelling of the outlier-generating mechanism, by informality of technique with no backing in terms of a study of the statistical properties of proposed procedures, and by a leaning towards the hardline view that outliers should be either rejected or retained with full import. Even today sufficient attention is not always paid to the form of the outlier model, or to the practical purpose of investigating outliers, in the presentation of methods for processing outliers. Many procedures have an *ad hoc*,

intuitively justified, basis with little external reference in the sense of the relative statistical merits of different possibilities. In reviewing such techniques we will attempt to set them, as far as possible, within a wider framework of model, statistical principle, and practical aim, and we shall also consider the extent to which such basic considerations have begun to formally permeate outlier study over recent years.

Such an emphasis is reflected in the structure of the book.* The opening two chapters are designed respectively to motivate examination of outliers and to pose basic questions about the nature of an outlier. Chapter 1 gives a general survey of the field. In Chapter 2 we consider the various ways in which we can model the presence of outliers in a set of data. We examine the different interests (from *rejection* of unacceptable contamination, through the *accommodation* of outliers with reduced influence in robust procedures applied to the whole set of data, to specific *identification* of outliers as the facets of principal interest in the data). We discuss the statistical respectability of distinct methods of study, and the special problems that arise from the dimensionality of the data set or from the purpose of its analysis (single-sample estimation or testing, regression, analysis of data from designed experiments, examination of slippage in multisample data, and so on).

Chapter 3 examines at length the assessment of discordancy of outliers in single univariate samples. It discusses basic considerations and also presents a battery of techniques for practical use with comment on the circumstances supporting one method rather than another.

Chapter 4, on the accommodation of outliers in single univariate samples, deals with inference procedures which are robust in the sense of providing protection against the effect of outliers. Chapter 5 is concerned with processing several univariate samples both with regard to the relative slippage of the distributions from which they arise and (to a lesser extent) in relation to the accommodation of outliers in robust analysis of the whole set of data.

Chapters 6 and 7 extend the ideas and methods (in relation to the three interests: rejection, accommodation, identification) to single multivariate samples and to the analysis of data in regression, designed experiments, or time series situations. Chapter 8 gives fuller and more specific attention to the implications of adopting a Bayesian, or a non-parametric, approach to the study of outliers. The concluding Chapter 9 poses a few issues for further consideration or investigation.

The book aims to bring together in a logical framework the vast amount of work on outliers which has been scattered over the years in the various professional journals and texts, and which appears to have acquired a new

* Chapter numbers have been changed in the Second and Third Editions (see *Preface* and *Contents* to each Edition).

lease of life over the last decade or so. It is directed to more than one kind of reader: to the student (to inform him of the range of ideas and techniques), to the experimentalist (to assist him in the judicious choice of methods for handling outliers), and to the professional statistician (as a guide to the present state of knowledge and a springboard for further research).

The level of treatment assumes a knowledge of elementary probability theory and statistical method such as would be acquired in an introductory university-level course. The methodological exposition leans on an understanding of the principles and practical implications of testing and estimation. Where basic modelling and demonstration of statistical propriety are discussed, a more mathematical appreciation of basic principles is assumed, including some familiarity with optimality properties of methods of constructing tests and estimators and some knowledge of the properties of order statistics. Proofs of results are formally presented where appropriate, but at a heuristic rather than highly mathematical level.

Extensive tables of appropriate statistical functions are presented in an Appendix, to aid the practical worker in the use of the different procedures. Many of these tables are extracted from existing published tables; we are grateful to all the authors and publishers concerned, and have made individual acknowledgement at the appropriate places in our text. Other tables have been specially produced by us. The whole set of tables has been presented in as compact and consistent a style as possible. This has involved a good deal of selection and re-ordering of the previously published material; we have aimed as far as possible to standardize the ranges of tabulated values of sample size, percentage point, etc.

Copious references are given throughout the text to source material and to further work on the various topics. They are gathered together in the section entitled 'References and Bibliography' with appropriate page references to places of principal relevance in the text. Additional references augment those which have been discussed in the text. These will of course appear without any page reference, but will carry an indication of the main area to which they are relevant.

It is, of course, a privilege and pleasure to acknowledge the help of others. We thank Dave Collett, Nick Fieller, Agnes Herzberg, and David Kendall for helpful comments on early drafts of some of the material. We are particularly grateful to Kim Malafant who carried out the extensive calculations of the new statistical tables in Chapter 3. Our grateful thanks go also to Hazel Howard who coped nobly with the typing of a difficult manuscript.

We are solely responsible for any imperfections in the book and should be glad to be informed of them.

VIC BARNETT
TOBY LEWIS
July, 1977

Preface to Second Edition

The first edition of this book took the opportunity of drawing together for the first time the vast and assorted literature of over a century on the topic of outliers. In offering a combination of specific methods for the experimentalist and a 'state-of-the-art' review for the professional statistician it was designed to provide a logical framework for, and comprehensive coverage of, the variety of topics and emphases in outlier study. About four hundred references were used as basic source material, in a climate where the subject seemed clearly 'to have acquired a new lease of life over the last decade or so'. What could not be anticipated was the sheer force and momentum of this revitalization of interest. In the six years since the appearance of the first edition about three hundred more published articles have appeared: in crude numerical terms, a seventy-five per cent expansion over the whole previous history of the subject. It is inevitable therefore that we should seek to reassess the current situation in the light of so many new ideas and refinements. The second edition aims to do this.

The modifications in the revised edition are of three types. Firstly, there are areas of enquiry which represent new topics of outlier study or which now need to be described within the context of outlier methodology. These are reflected in the new chapters (7 and 11) on outliers in directional data, and in time series, respectively. Secondly, there are the many contributions which refine, reassess, or extend our knowledge on specific aspects of outlier investigation. Such developments are incorporated by means of substantial expansion, and some judicious pruning, of the discussion of almost all the topics of the earlier edition of the book. Particular attention has been given to discordancy tests for univariate and multivariate samples and for data from structured models (linear models generally, and specific aspects of regression and designed experiment situations). At a more general level, there has been a welcome growth of emphasis on methods of accommodation (robust inference in the face of outliers) and on informal (often graphical) descriptive procedures. These latter areas of development have contributed to the stimulus for a final type of modification—a reordering or reemphasis of some of the basic ideas and principles. In particular this has prompted a separation of general approach from specific results for the study of univariate samples with regard both to tests of

discordancy and methods of accommodation, and a reversal of order in the treatment of these two aspects. Thus we now have Chapter 3 and 4 on accommodation and Chapters 5 and 6 on discordancy testing, for univariate samples, distinguishing between general principle (Chapters 3 and 5) and specific method (Chapter 4 and 6) in each case. Accordingly the chapters dealing with slippage, multivariate outliers, outliers in linear models (regression, designed experiments), and Bayesian methods are now renumbered as Chapters 8, 9, 10, and 12, respectively.

There has been much activity in recent years on the theme of 'influential observations': as sample values which have disproportionate effect on estimates or tests of parameters in a model (particularly in the field of regression). Whilst not entirely coincident concepts, outliers and influential observations are clearly related in important aspects. Similarities and distinctions of aim and method are highlighted at appropriate stages in our revised study of outlier methods, especially in Chapter 10.

We are grateful to those who have pointed out misprints and ambiguities in the first edition; particularly to Dr N. A. Campbell. We have taken the opportunities of remedying such matters.

Whilst the new edition reflects the many developments and extensions mentioned above, it is important to stress that the basic emphasis and aim is unchanged. We continue to eschew the two extremes of pre-digested recipes for instant outlier management on the one hand and indulgence in mathematical formality or sophistication for its own sake on the other. Our aim remains one of explaining basic principle and developing associated method to a level where outlier techniques can be soundly and sanely applied (with relevant illustration and tabulation) and the researcher can be provided with a springboard for further exploration of a fascinating statistical topic.

VIC BARNETT
TOBY LEWIS
December, 1983

Preface to Third Edition

Since the publication of the second edition of this book, work in outlier methodology has continued unabated. Indeed, over 1000 new refereed publications have appeared in the literature. These not only extend previous knowledge in existing fields, but have opened up whole new areas of enquiry. Clearly, a new edition was needed to do justice to these many developments and to maintain our objective of providing a comprehensive and up to date coverage of the subject.

We have thoroughly revised and updated the material on the range of topics that were covered in the earlier editions. At the same time, extra material has been introduced to reflect new themes and changing emphases.

In contrast, a small amount of material has been omitted, in particular the chapter on the peripheral topic of *outlying subsamples*.

Topics on which the coverage *is new or has been substantially changed or extended* include:

- **basic principles**; distribution theory under contamination models, measures of efficiency and performance for multiple outliers, assessment of masking and swamping, allocation of outliers.
- **univariate data**; new tests (including extreme value and Weibull), wider study of robustness and accommodation (including logistic and double exponential distributions), additional tables.
- **multivariate and structured data**; estimation of individual components of vector parameters, use of correlation estimators, deletion methods, elliptically symmetric distributions, graphical methods, least median of squares and L_1 -norm methods, multivariate linear model, multiple outliers, non-linear regression (including logistic and generalized linear models).
- **special topics**; Bayesian methods, time series (ARIMA model, distinction of AO and IO outliers, model specification, new accommodation methods, diagnostics, multiple time series), directional data (methods for axial and vectorial spherical data, accommodation).

Further special topics are dealt with in *new chapters*—these include new methods for outliers in contingency tables, problems of sample surveys, statistical software and international standards and regulations. Practical

illustrations remain important for reinforcement of ideas—new examples are included as also is discussion of data studies in the literature.

To cope with the vast expansion of material, it has been advantageous to restructure the book, which is now divided into the following four distinct parts: *Basic Principles*, *Univariate Data*, *Multivariate and Structured Data*, *Special Topics*. These reflect the headings used above to categorize the major areas of change. The *References and Bibliography* section has been substantially expanded to cover the new material.

We would like to thank Eileen Stoydin, Marian Joyce and Sharon Wilson for their most efficient assistance in the preparation of the typescript.

We hope that this new edition will prove to be valuable both to research statisticians and to practitioners in various fields.

VIC BARNETT
TOBY LEWIS
April, 1993

Contents

PART I BASIC PRINCIPLES	1
CHAPTER 1 INTRODUCTION	3
1.1 <i>Human error and ignorance</i>	10
1.2 <i>Outliers in relation to probability models</i>	11
1.3 <i>Outliers in more structured situations</i>	16
1.4 <i>Bayesian methods</i>	23
1.5 <i>Statistical computing</i>	24
1.6 <i>Survey of outlier problems</i>	24
CHAPTER 2 WHY DO OUTLYING OBSERVATIONS ARISE AND WHAT SHOULD ONE DO ABOUT THEM?	27
2.1 <i>Early informal approaches</i>	27
2.2 <i>Origin of outliers, statistical methods and aims</i>	31
2.2.1 <i>The nature and origin of an outlier</i>	32
2.2.2 <i>Relevant statistical procedures for handling outliers</i>	34
2.2.3 <i>Different aims in examining outliers</i>	39
2.3 <i>Models for contamination</i>	43
2.4 <i>Outlier proneness, outlier resistance, outlier labelling</i>	52
CHAPTER 3 THE ACCOMMODATION APPROACH: ROBUST ESTIMATION AND TESTING	55
3.1 <i>Performance criteria</i>	60
3.1.1 <i>Efficiency measures for estimators</i>	60
3.1.2 <i>Distributional properties of contaminated distributions</i>	66
3.1.3 <i>The qualitative approach: influence curves</i>	68
3.1.4 <i>Robustness of confidence intervals</i>	74
3.1.5 <i>Robustness of significance tests</i>	75
3.2 <i>Methods of accommodation—general and outlier-specific</i>	78
3.2.1 <i>Estimation of location</i>	78
3.2.2 <i>Estimation of dispersion</i>	83
3.2.3 <i>Hypothesis tests and confidence intervals</i>	85
3.2.4 <i>Adaptive robust procedures</i>	87

CHAPTER 4 TESTING FOR DISCORDANCY: PRINCIPLES AND CRITERIA	89
4.1 <i>Construction of discordancy tests</i>	94
4.1.1 <i>Test statistics</i>	94
4.1.2 <i>Statistical bases for construction of tests</i>	98
4.1.3 <i>Inclusive and exclusive measures</i>	108
4.1.4 <i>Masking and swamping</i>	109
4.1.5 <i>Assessment of significance</i>	115
4.2 <i>The assessment of test performance</i>	121
4.2.1 <i>Measures of performance</i>	121
4.2.2 <i>Distributional properties under the alternative hypothesis</i>	125
4.3 <i>The multiple outlier problem</i>	125
4.3.1 <i>Block procedures for multiple outliers in univariate samples</i>	133
4.3.2 <i>Consecutive procedures for multiple outliers in univariate samples</i>	136
PART II UNIVARIATE DATA	141
CHAPTER 5 ACCOMMODATION PROCEDURES FOR UNIVARIATE SAMPLES	143
5.1 <i>Estimation of location</i>	143
5.1.1 <i>Estimators based on trimming or Winsorization</i>	143
5.1.2 <i>L-estimators (linear order statistics estimators)</i>	146
5.1.3 <i>M-estimators (maximum likelihood type estimators)</i>	148
5.1.4 <i>R-estimators (rank test estimators)</i>	152
5.1.5 <i>Other estimators</i>	153
5.2 <i>Estimation of scale or dispersion</i>	155
5.3 <i>Hypothesis tests and confidence intervals</i>	157
5.4 <i>Accommodation of outliers in univariate normal samples</i>	158
5.5 <i>Accommodation of outliers in gamma (including exponential) samples</i>	174
5.6 <i>Accommodation of outliers in logistic and double exponential samples</i>	185
5.6.1 <i>Estimators of the mean μ</i>	186
5.6.2 <i>Estimators of the scale parameter σ</i>	187
5.6.3 <i>Robust estimation for logistic samples</i>	188
5.6.4 <i>Robust estimation for double exponential samples</i>	189

CHAPTER 6 SPECIFIC DISCORDANCY TESTS FOR OUTLIERS IN UNIVARIATE SAMPLES	191
6.1 <i>Guide to use of the tests</i>	191
6.2 <i>Discordancy tests for gamma (including exponential) samples</i>	193
6.2.1 <i>Gamma samples: contents list and details of tests</i>	195
6.3 <i>Discordancy tests for normal samples</i>	216
6.3.1 <i>Normal samples: contents list and details of tests</i>	217
6.4 <i>Discordancy tests for samples from other distributions</i>	250
6.4.1 <i>Log-normal samples</i>	251
6.4.2 <i>Truncated exponential samples</i>	251
6.4.3 <i>Uniform samples</i>	252
6.4.4 <i>Gumbel, Fréchet, and Weibull samples</i>	254
6.4.5 <i>Pareto samples</i>	260
6.4.6 <i>Poisson samples</i>	261
6.4.7 <i>Binomial samples</i>	264
 PART III MULTIVARIATE AND STRUCTURED DATA	 267
CHAPTER 7 OUTLIERS IN MULTIVARIATE DATA	269
7.1 <i>Principles for outlier detection in multivariate samples</i>	269
7.2 <i>Accommodation of multivariate outliers</i>	273
7.2.1 <i>Estimation of parameters in multivariate distributions</i>	273
7.2.2 <i>Estimation of individual components of vector parameters</i>	279
7.2.3 <i>Outliers in multivariate analyses</i>	282
7.3 <i>Discordancy tests</i>	283
7.3.1 <i>Multivariate normal samples</i>	284
7.3.2 <i>Multivariate exponential samples</i>	293
7.3.3 <i>Multivariate Pareto samples</i>	295
7.3.4 <i>Discordancy tests for more general distributions</i>	296
7.3.5 <i>A transformation approach</i>	296
7.4 <i>Informal methods for multivariate outliers</i>	297
7.4.1 <i>Marginal outliers and linear constraints</i>	300
7.4.2 <i>Graphical and pictorial methods</i>	301
7.4.3 <i>Principal component analysis method</i>	303
7.4.4 <i>Use of reduction measures in the form of generalized distances</i>	306
7.4.5 <i>Function plots</i>	308
7.4.6 <i>Correlation methods and influential observations</i>	309
7.4.7 <i>A 'gap test' for multivariate outliers</i>	311

CHAPTER 8 THE OUTLIER PROBLEM FOR STRUCTURED DATA: REGRESSION, THE LINEAR MODEL AND DESIGNED EXPERIMENTS	315
8.1 <i>Outliers in simple linear regression</i>	320
8.1.1 <i>Discordancy tests for simple linear regression</i>	321
8.1.2 <i>Outlier accommodation in simple linear regression</i>	323
8.1.3 <i>Outliers in linear structural and functional models</i>	325
8.2 <i>Outliers with general linear models</i>	328
8.2.1 <i>Residual-based methods</i>	328
8.2.2 <i>Augmented residual-based, and non-residual-based, methods</i>	340
8.2.3 <i>Accommodation of outliers for the linear model</i>	346
8.3 <i>Outliers in non-linear regression</i>	351
8.4 <i>Outliers in designed experiments</i>	354
8.5 <i>Graphical methods and diagnostics for linear model outliers</i>	365
8.6 <i>Some further comments on influential observations</i>	371
PART IV SPECIAL TOPICS	375
CHAPTER 9 BAYESIAN APPROACHES TO OUTLIERS	377
9.1 <i>Basic Bayesian considerations</i>	377
9.2 <i>Bayesian accommodation of outliers</i>	380
9.3 <i>Bayesian assessment of contamination</i>	388
CHAPTER 10 OUTLIERS IN TIME SERIES: AN IMPORTANT AREA OF OUTLIER STUDY	395
10.1 <i>Detection and testing of outliers in time series</i>	396
10.2 <i>Accommodation of outliers in time series</i>	401
10.2.1 <i>Time-domain characteristics</i>	401
10.2.2 <i>Frequency-domain characteristics</i>	407
10.2.3 <i>Specific inference problems</i>	409
10.3 <i>Bayesian methods</i>	413
10.4 <i>Comment</i>	414
CHAPTER 11 OUTLIERS IN DIRECTIONAL DATA	417
11.1 <i>Outliers on the circle</i>	420
11.2 <i>Outliers on the sphere</i>	424
11.2.1 <i>Vectorial data</i>	424
11.2.2 <i>Axial data</i>	428

CHAPTER 12	SOME LITTLE-EXPLORED AREAS: CONTINGENCY TABLES AND SAMPLE SURVEYS	431
12.1	<i>Outliers in contingency tables</i>	431
12.2	<i>Outliers in sample surveys</i>	440
CHAPTER 13	IMPORTANT STRANDS: COMPUTER SOFTWARE, DATA STUDIES, STANDARDS AND REGULATIONS	449
13.1	<i>Outliers and statistical computing</i>	449
13.2	<i>Outliers in detailed data studies</i>	454
13.3	<i>Outliers in standards requirements</i>	455
CHAPTER 14	PERSPECTIVE	459
APPENDIX:	STATISTICAL TABLES	463
	<i>Contents list</i>	465
REFERENCES AND BIBLIOGRAPHY		527
INDEX		574

PART I

Basic Principles