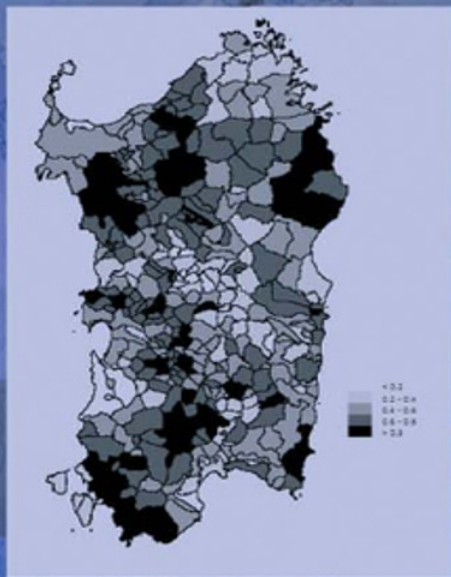


---

# STATISTICAL METHODS IN SPATIAL EPIDEMIOLOGY

## Second Edition



Andrew B. Lawson



# **Statistical Methods in Spatial Epidemiology**

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: *David J. Balding, Peter Bloomfield, Noel A. C. Cressie, Nicholas I. Fisher, Iain M. Johnstone, J. B. Kadane, Geert Molenberghs, Louise M. Ryan, David W. Scott, Adrian F. M. Smith, Jozef L. Teugels;*  
Editors Emeriti: *Vic Barnett, J. Stuart Hunter, David G. Kendall*

A complete list of the titles in this series appears at the end of this volume.

# Statistical Methods in Spatial Epidemiology

Second Edition

**Andrew B. Lawson**

*Department of Epidemiology and Biostatistics,  
University of South Carolina,  
Columbia, USA*



John Wiley & Sons, Ltd

Copyright © 2006

John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester,  
West Sussex PO19 8SQ, England

Telephone (+44) 1243 779777

Email (for orders and customer service enquiries): [cs-books@wiley.co.uk](mailto:cs-books@wiley.co.uk)

Visit our Home Page on [www.wiley.com](http://www.wiley.com)

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, electronic, mechanical, photocopying, recording, scanning or otherwise, except under the terms of the Copyright, Designs and Patents Act 1988 or under the terms of a licence issued by the Copyright Licensing Agency Ltd, 90 Tottenham Court Road, London W1T 4LP, UK, without the permission in writing of the Publisher. Requests to the Publisher should be addressed to the Permissions Department, John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex PO19 8SQ, England, or emailed to [permreq@wiley.co.uk](mailto:permreq@wiley.co.uk), or faxed to (+44) 1243 770620.

Designations used by companies to distinguish their products are often claimed as trademarks. All brand names and product names used in this book are trade names, service marks, trademarks or registered trademarks of their respective owners. The Publisher is not associated with any product or vendor mentioned in this book.

This publication is designed to provide accurate and authoritative information in regard to the subject matter covered. It is sold on the understanding that the Publisher is not engaged in rendering professional services. If professional advice or other expert assistance is required, the services of a competent professional should be sought.

#### ***Other Wiley Editorial Offices***

John Wiley & Sons Inc., 111 River Street, Hoboken, NJ 07030, USA

Jossey-Bass, 989 Market Street, San Francisco, CA 94103-1741, USA

Wiley-VCH Verlag GmbH, Boschstr. 12, D-69469 Weinheim, Germany

John Wiley & Sons Australia Ltd, 42 McDougall Street, Milton, Queensland 4064, Australia

John Wiley & Sons (Asia) Pte Ltd, 2 Clementi Loop #02-01, Jin Xing Distripark, Singapore 129809

John Wiley & Sons Canada Ltd, 22 Worcester Road, Etobicoke, Ontario, Canada M9W 1L1

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

#### ***British Library Cataloguing in Publication Data***

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

ISBN-13: 978-0-470-01484-4

ISBN-10: 0-470-01484-9

**“... a story is a letter the author writes to themselves, to tell  
themselves things that they would be unable to discover  
otherwise.”**

***after* Carlos Ruiz Zafón**





**‘to Keir, Fraser, and Hugh and all my family’**



# Contents

<b>Preface and Acknowledgements to Second Edition</b>	<b>xv</b>
<b>Preface and Acknowledgements</b>	<b>xvii</b>
<b>I The Nature of Spatial Epidemiology</b>	<b>1</b>
<b>1 Definitions, Terminology and Data Sets</b>	<b>3</b>
1.1 Map Hypotheses and Modelling Approaches . . . . .	5
1.2 Definitions and Data Examples . . . . .	7
1.2.1 Case event data . . . . .	7
1.2.2 Count data . . . . .	8
1.3 Further Definitions . . . . .	10
1.3.1 Control events and processes . . . . .	10
1.3.2 Census tract information . . . . .	10
1.3.3 Clustering definitions . . . . .	10
1.4 Some Data Examples . . . . .	11
1.4.1 Case event examples . . . . .	11
1.4.2 Count data examples . . . . .	19
<b>2 Scales of Measurement and Data Availability</b>	<b>25</b>
2.1 Small Scale . . . . .	26
2.2 Large Scale . . . . .	26
2.3 Rate Dependence . . . . .	27
2.4 Data Quality and the Ecological Fallacy . . . . .	27
2.5 Edge Effects . . . . .	28
<b>3 Geographical Representation and Mapping</b>	<b>31</b>
3.1 Introduction and Definitions . . . . .	31
3.2 Maps and Mapping . . . . .	32
3.2.1 Statistical maps and mapping . . . . .	34
3.2.2 Object process mapping . . . . .	34
3.2.3 Geostatistical mapping . . . . .	36

- 3.3 Statistical Accuracy . . . . . 37
- 3.4 Aggregation . . . . . 37
- 3.5 Mapping Issues Related to Aggregated Data . . . . . 37
- 3.6 Conclusions . . . . . 39
  
- 4 Basic Models 41**

  - 4.1 Sampling Considerations . . . . . 41
  - 4.2 Likelihood-Based and Bayesian Approaches . . . . . 42
  - 4.3 Point Event Models . . . . . 42
    - 4.3.1 Point process models and applications . . . . . 43
    - 4.3.2 The basic Poisson process model . . . . . 44
    - 4.3.3 Hybrid models and regionalisation . . . . . 49
    - 4.3.4 Bayesian models and random effects . . . . . 50
    - 4.3.5 MAP estimation, empirical Bayes and full Bayesian analysis 52
    - 4.3.6 Bivariate/multivariate models . . . . . 53
    - 4.3.7 Hidden structure and mixture models . . . . . 56
    - 4.3.8 Space-time extensions . . . . . 56
  - 4.4 Count Models . . . . . 58
    - 4.4.1 Standard models . . . . . 60
    - 4.4.2 Approximations . . . . . 63
    - 4.4.3 Random-effect extensions . . . . . 63
    - 4.4.4 Hidden structure and mixture models . . . . . 64
    - 4.4.5 Space-time extensions . . . . . 65

  
- 5 Exploratory Approaches, Parametric Estimation and Inference 67**

  - 5.1 Exploratory Methods . . . . . 68
    - 5.1.1 Cartographic issues . . . . . 69
    - 5.1.2 Case event mapping . . . . . 71
    - 5.1.3 Count mapping . . . . . 75
  - 5.2 Parameter Estimation . . . . . 80
    - 5.2.1 Case event likelihood models . . . . . 80
    - 5.2.2 Count event likelihood models . . . . . 85
    - 5.2.3 Approximations . . . . . 87
    - 5.2.4 Bayesian models . . . . . 88
  - 5.3 Residual Diagnostics . . . . . 96
  - 5.4 Hypothesis Testing . . . . . 98
  - 5.5 Edge Effects . . . . . 99
    - 5.5.1 Edge effects in case events . . . . . 101
    - 5.5.2 Edge effects in counts . . . . . 101
    - 5.5.3 Edge weighting schemes and MCMC methods . . . . . 102
    - 5.5.4 Discussion . . . . . 104
    - 5.5.5 The Tuscany example . . . . . 105

**II Important Problems in Spatial Epidemiology 109**

**6 Small Scale: Disease Clustering 111**

6.1 Definition of Clusters and Clustering . . . . . 112

6.2 Modelling Issues . . . . . 115

6.3 Hypothesis Tests for Clustering . . . . . 118

    6.3.1 General non-specific clustering . . . . . 118

    6.3.2 Specific clustering . . . . . 121

6.4 Space-Time Clustering . . . . . 123

    6.4.1 Modelling issues . . . . . 123

    6.4.2 Hypothesis testing . . . . . 126

6.5 Clustering Examples . . . . . 127

    6.5.1 Humberside example . . . . . 127

    6.5.2 Larynx cancer example . . . . . 131

    6.5.3 Count data clustering example . . . . . 133

    6.5.4 Space-time clustering examples . . . . . 136

6.6 Other Methods Related to Clustering . . . . . 138

    6.6.1 Wombling . . . . . 140

**7 Small Scale: Putative Sources of Hazard 143**

7.1 Introduction . . . . . 143

7.2 Study Design . . . . . 144

    7.2.1 Retrospective and prospective studies . . . . . 144

    7.2.2 Study region design . . . . . 145

    7.2.3 Replication and controls . . . . . 146

7.3 Problems of Inference . . . . . 147

    7.3.1 Exploratory techniques . . . . . 148

7.4 Modelling the Hazard Exposure Risk . . . . . 153

7.5 Models for Case Event Data . . . . . 162

    7.5.1 Estimation . . . . . 164

    7.5.2 Hypothesis tests . . . . . 164

    7.5.3 Diagnostic techniques . . . . . 166

7.6 A Case Event Example . . . . . 167

7.7 Models for Count Data . . . . . 169

    7.7.1 Estimation . . . . . 171

    7.7.2 Hypothesis tests . . . . . 171

7.8 A Count Data Example . . . . . 172

7.9 Other Directions . . . . . 174

    7.9.1 Multiple disease analysis . . . . . 174

    7.9.2 Space-time modelling . . . . . 184

    7.9.3 Space-time exploratory analysis . . . . . 184

    7.9.4 Space-time Bayesian analysis . . . . . 185

**8 Large Scale: Disease Mapping 189**

8.1 Introduction . . . . . 189

- 8.2 Simple Statistical Representation . . . . . 189
  - 8.2.1 Crude rates . . . . . 190
  - 8.2.2 Standardised mortality/morbidity ratios, standardisation and relative risk surfaces . . . . . 191
  - 8.2.3 Interpolation . . . . . 193
  - 8.2.4 Exploratory mapping methods . . . . . 193
- 8.3 Basic Models . . . . . 194
  - 8.3.1 Likelihood models . . . . . 194
  - 8.3.2 Random effects and Bayesian models . . . . . 197
- 8.4 Advanced Methods . . . . . 201
  - 8.4.1 Non-parametric methods . . . . . 202
  - 8.4.2 Incorporating spatially correlated heterogeneity . . . . . 203
  - 8.4.3 Case event modelling . . . . . 206
- 8.5 Model Variants and Extensions . . . . . 209
  - 8.5.1 Semiparametric modelling . . . . . 209
  - 8.5.2 Geographically weighted regression . . . . . 210
  - 8.5.3 Mixture models . . . . . 211
- 8.6 Approximate Methods . . . . . 212
- 8.7 Multivariate Methods . . . . . 213
- 8.8 Evaluation of Model Performance . . . . . 216
- 8.9 Hypothesis Testing in Disease Mapping . . . . . 219
  - 8.9.1 First-order effects . . . . . 219
  - 8.9.2 Second-order and variance effects . . . . . 221
- 8.10 Space-Time Disease Mapping . . . . . 222
- 8.11 Spatial Survival and Longitudinal Data . . . . . 229
  - 8.11.1 Spatial survival analysis . . . . . 229
  - 8.11.2 Spatial longitudinal analysis . . . . . 231
  - 8.11.3 Spatial multiple event modelling . . . . . 232
- 8.12 Disease Mapping: Case Studies . . . . . 232
  - 8.12.1 Eastern Germany . . . . . 232
  - 8.12.2 Ohio respiratory cancer . . . . . 239
- 9 Ecological Analysis and Scale Change . . . . . 247**
  - 9.1 Ecological Analysis: Introduction . . . . . 247
  - 9.2 Small-Scale Modelling Issues . . . . . 252
    - 9.2.1 Hypothesis tests . . . . . 253
    - 9.2.2 Ecological aggregation effects . . . . . 253
  - 9.3 Changes of Scale and MAUP . . . . . 255
    - 9.3.1 MAUP: the *modifiable areal unit problem* . . . . . 255
    - 9.3.2 Large-scale issues . . . . . 260
  - 9.4 A Simple Example: Sudden Infant Death in North Carolina . . . . . 261
  - 9.5 A Case Study: Malaria and IDDM . . . . . 263
- 10 Infectious Disease Modelling . . . . . 269**
  - 10.1 Introduction . . . . . 269

- 10.2 General Model Development . . . . . 270
- 10.3 Spatial Model Development . . . . . 273
  - 10.3.1 Count data . . . . . 273
  - 10.3.2 Individual-level data . . . . . 278
- 10.4 Modelling Special Cases for Individual-Level Data . . . . . 280
  - 10.4.1 Proportional hazards interpretation . . . . . 280
  - 10.4.2 Subgroup modifications . . . . . 281
  - 10.4.3 Cluster function specification . . . . . 282
- 10.5 Survival Analysis with Spatial Dependence . . . . . 283
- 10.6 Individual-Level Data Example . . . . . 284
  - 10.6.1 Distribution of susceptibles  $S(x, t)$  . . . . . 285
  - 10.6.2 The spatial distance function  $h$  . . . . . 285
  - 10.6.3 The function  $g$  . . . . . 285
  - 10.6.4 Fitting the model . . . . . 286
  - 10.6.5 Revised model . . . . . 287
- 10.7 Underascertainment and Censoring . . . . . 288
- 10.8 Conclusions . . . . . 289

**11 Large Scale: Surveillance 293**

- 11.1 Process Control Methodology . . . . . 294
- 11.2 Spatio-Temporal Modelling . . . . . 295
- 11.3 S-T Monitoring . . . . . 297
  - 11.3.1 Fixed spatial and temporal frame . . . . . 297
  - 11.3.2 Fixed spatial frame and dynamic temporal frame . . . . . 301
- 11.4 Syndromic Surveillance . . . . . 304
- 11.5 Multivariate–Multifocus Surveillance . . . . . 305
- 11.6 Bayesian Approaches . . . . . 308
  - 11.6.1 Bayesian alarm functions, Bayes factors and syndromic analyses . . . . . 308
- 11.7 Computational Considerations . . . . . 310
- 11.8 Infectious Diseases . . . . . 311
- 11.9 Conclusions . . . . . 312

**Appendix A Monte Carlo Testing, Parametric Bootstrap and Simulation**

- Envelopes 313**
- A.1 Nuisance Parameters and Test Statistics . . . . . 313
- A.2 Monte Carlo Tests . . . . . 314
- A.3 Null Hypothesis Simulation . . . . . 315
  - A.3.1 Spatial case . . . . . 316
  - A.3.2 Spatio-temporal case . . . . . 318
- A.4 Parametric Bootstrap . . . . . 319
  - A.4.1 Bayesian spatial models . . . . . 322
  - A.4.2 Spatio-temporal case . . . . . 323
- A.5 Simulation Envelopes . . . . . 324

<b>Appendix B Markov Chain Monte Carlo Methods</b>	<b>325</b>
B.1 Definitions . . . . .	325
B.2 Metropolis and Metropolis–Hastings Algorithms . . . . .	326
B.2.1 Metropolis algorithm . . . . .	326
B.2.2 Metropolis–Hastings extension . . . . .	327
B.2.3 The Gibbs sampler . . . . .	327
B.2.4 M–H versus Gibbs algorithms . . . . .	328
B.2.5 Examples . . . . .	328
<b>Appendix C Algorithms and Code</b>	<b>331</b>
C.1 Data Exploration . . . . .	331
C.2 Likelihood and Bayesian Models . . . . .	335
C.3 Likelihood Models . . . . .	336
C.3.1 Case event data . . . . .	336
C.3.2 Count data . . . . .	340
C.4 Bayesian Hierarchical Models . . . . .	341
C.4.1 Case event data . . . . .	341
C.4.2 Count data . . . . .	344
C.5 Space-Time Analysis . . . . .	346
C.5.1 Data exploration . . . . .	346
C.5.2 Likelihood models . . . . .	349
C.5.3 Bayesian models . . . . .	351
C.5.4 Infectious disease models . . . . .	357
<b>Appendix D Glossary of Estimators</b>	<b>359</b>
D.1 Case Event Estimators . . . . .	359
D.2 Tract Count Estimators . . . . .	361
<b>Appendix E Software</b>	<b>363</b>
E.1 Software . . . . .	363
E.1.1 Spatial statistical tools . . . . .	363
E.1.2 Geographical information systems . . . . .	365
<b>Bibliography</b>	<b>367</b>
<b>Index</b>	<b>389</b>



# Preface and Acknowledgements to Second Edition

Since the appearance of the first edition of this book there has been a considerable development of interest in statistical methodology in the area of spatial epidemiology. This development has seen the increased output of research papers and books marking the maturity of certain areas of concern. For example, close to that date when the edited volume by Elliott *et al.* (2000) appeared, and since special issues of the *Journal of the Royal Statistical Society, Series A* (2001), *Environmental and Ecological Statistics* (2005), *Statistical Methods in Medical Research* (2005, 2006) and *Statistics in Medicine* (2006) have all contributed to the appearance of novel methodology. The development of software has also facilitated the wider use of the more advanced methods. In particular, the availability of free packages such as R, WinBUGS and SaTScan has led to wide dissemination of the available methods.

In particular, the area of disease map modelling has seen much development with Bayesian modelling as a particular feature. The use of mixture models and variants of likelihoods has seen development, while the routine application of sophisticated random-effect models is now relatively straightforward. The areas of disease clustering, ecological analysis and infectious disease modelling have all seen advances. In addition, the area of surveillance has re-emerged due to interest in early detection of potential bioterrorism attacks and in particular syndromic surveillance has become a major focus.

I would like to take this opportunity to acknowledge the influence and support of the following: Linda Pickle (NIH), Ram Tiwari (NIH), Martin Kulldorff, Dan Wartenburg, Peter Rogerson, Andrew Moore, Sudipto Banerjee, Ken Kleinman, William Browne, Carmen Vidal Rodeiro, Monir Hossain, Allan Clark, Yang Wang, Yuan Liu, Bo Ma, Huafeng Zhou. Finally I should also like to acknowledge the helpful interactions with staff at Wiley Europe over the years: Kathryn Sharples, Sian Jones, Helen Ramsey, Sharon Clutton and Lucy Bryan.

**Andrew Lawson, Columbia, South Carolina**  
*December 2005*



# Preface and Acknowledgements

The development of statistical methods in spatial epidemiology has had a chequered career. One of the earliest examples of the analysis of geographical locations of disease in relation to a putative health hazard was John Snow's analysis of cholera cases in relation to the location of the Broad Street water pump in London (Snow, 1854). However, until recently, developments in statistical methods in this area have been sporadic. While medical geography developed in the 1960s (Howe, 1963), only a number of papers on space-time clustering (Mantel, 1967; Knox, 1964) appeared in the statistical literature. More recently, developments of methods in spatial statistics, image processing, and in particular Bayesian methods and computation, have seen parallel developments in methods for spatial epidemiology (see Marshall (1991b) for a review). It is notable that methods for the analysis of case locations around a source of hazard (such as Snow's cholera map) have only recently been developed (Diggle, 1989; Lawson, 1989). The current increased level of interest in statistical methods in spatial epidemiology is a reflection, in part, of the increased concern in society for environmental issues and their relation to the health of individuals. Hence, the 'detection' of pollution sources or sources of health hazard can be seen as the backdrop to many studies in environmental epidemiology (Diggle, 1993). The correct allocation of resources for health care in different areas by health services is also greatly enhanced by the development of statistical methods which allow more accurate depiction of 'true' disease incidence and its relation to explanatory variables. Previous work in this area has been reviewed by Lawson and Cressie (2000), and Marshall (1991b) and Elliott *et al.* (1992a) discuss the general epidemiological issues surrounding spatial epidemiological problems.

It is the purpose of this book to provide an overview of the main statistical methods currently available in the field of spatial epidemiology. Inevitably, some selectivity in choice of methods reviewed will be apparent, but it is hoped that our coverage will encompass the most important areas of development. One area which we do not examine in detail is that of space-time analysis of epidemiological data, although the modelling of infectious disease data is considered in Chapter 11.

As this book is mainly a review of recent research work, its target audience is largely confined to those with some statistical knowledge and is appropriate for

third level degree and postgraduate students in statistics, or epidemiology with a strong statistical background.

A considerable number of people have directly or indirectly contributed to the production of this book. First, I acknowledge the support of Sharon Clutton and Helen Ramsey at Wiley and Tony Johnson of Statistics in Medicine for their support from Budapest onwards. Fundamental influences in the development of my ideas in spatial epidemiology have been Richard Cormack and Peter Diggle. I also acknowledge the encouragement of Noel Cressie, who has supported my work through visits to Iowa State and Ohio State Universities, and important collaborations with Martin Kulldorff, Annibale Biggeri, Dankmar Boehning, Peter Schlattmann, Emmanuel Lesaffre, Jean-Francois Viel, Adrian Baddeley, Niels Becker and Andrew Cliff.

**Andrew B. Lawson, Aberdeen,**  
*March 2000*

## **Part I**

# **The Nature of Spatial Epidemiology**



# 1

## Definitions, Terminology and Data Sets

Spatial epidemiology concerns the analysis of the spatial/geographical distribution of the incidence of disease. In its simplest form the subject concerns the use and interpretation of *maps* of the locations of disease cases, and the associated issues relating to map production and the statistical analysis of mapped data must apply within this subject. In addition, the nature of *disease* maps ensures that many epidemiological concepts also play an important role in the analysis. In essence, these two different aspects of the subject have their own impact on the methodology which has developed to deal with the many issues which arise in this area.

First, since mapped data are spatial in nature, the application of *spatial* statistical methods forms a core part of the subject area. The reason for this lies in the fact that the study of any data which are georeferenced (i.e. have a spatial/geographical location associated with them) may have properties which relate to the location of individual data items and also the surrounding data. For example, Figure 1.1 shows the total number of deaths from respiratory cancer found in 26 small areas (census tracts) in central Scotland over the period 1976–1983. This map displays a number of features which commonly arise when the geographical distribution of disease is examined. On this map the numbers (counts) of cases within each area are displayed. In some areas of the map the counts are similar to those found in the immediately surrounding areas (e.g. in the south and southeast of the map counts of 4 and 6 are recorded, while in the northwest of the map, lower counts are found in many areas). This similarity in the count data in groups of tracts is unlikely to have arisen from the allocation of a random sample of counts from a common statistical distribution. The counts may display some form of correlation in their levels based on their location, i.e. counts close to each other in space are similar. This form of correlation does not arise from the usual statistical models assumed to

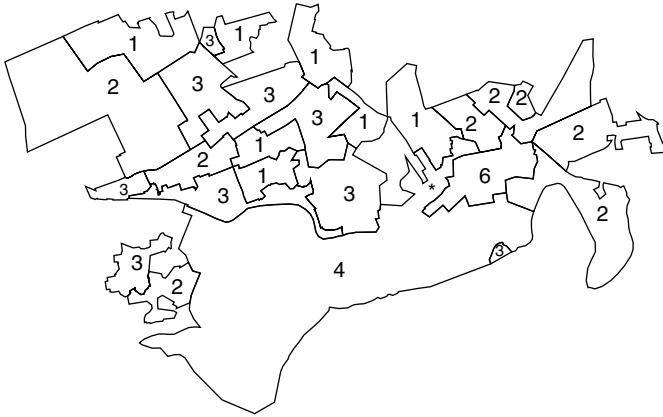


Figure 1.1 Falkirk: central Scotland respiratory cancer counts in 26 census enumeration districts over a fixed time period. \* Putative health hazard.

apply to independent observations found in, for example, clinical medical studies or other conventional statistical application areas. Hence, methods which apply to the analysis of these data must be able to address the possibility of such correlation existing in the mapped data under study. Another feature of this example, which commonly arises in the study of spatial epidemiology, is the irregular nature of the regions within which the counts are observed, i.e. the census tracts have irregular geographical boundaries. This may arise as a feature of the whole study region (*study window*) or may be found associated with tracts themselves. In some countries, notably in North America, small areas are often regular in shape and size and this feature simplifies the resulting analysis. However, in many other areas irregular region geometries are common. Finally, in some studies, the spatial distribution of cases or counts of disease are to be related to other locations on the map. For example, in Figure 1.1 the location of a potential (*putative*) environmental health hazard is also mapped (a metal-processing plant), and the focus of the study may be to assess the relationship of the disease incidence on the map to that location, perhaps to make inferences about the environmental risk in its vicinity.

The second feature which uniquely defines the study of spatial epidemiology is that the mapped data are often *discrete*. Unlike other areas of spatial statistical analysis, which are often focused on continuous data, e.g. geostatistical methods, the data found in spatial epidemiology often take the form of point locations (the address locations of cases of disease) or counts of disease within regions such as census tracts or, at larger scale, counties or municipalities. Hence, the mapped data often consist of cartesian coordinates in the form of a grid reference or longitude/latitude of an address of a case, or a count of cases within a region with the associated location of that region (either as a point location of a centroid or as a set of boundary line segments defining the region). Given this form of data format,



it is not surprising that models which have been developed for applications within this area are derived from stochastic point process theory (for case locations) and associated discrete probability distributions (for counts within arbitrary regions).

Finally, the epidemiological nature of these discrete spatial data leads to the derivation of models and methods which are related to conventional epidemiological studies. For example, the case-control study, where individual cases are matched to control individuals based on specific criteria, has parallels in spatial epidemiology where spatial control *distributions* are used to provide a locational control for cases. This is akin to the estimation of background hazard in survival studies. One fundamental epidemiological issue which arises in these studies is the incorporation of the local population which is at risk of contracting the disease in question. As we must control for the spatial variation in the underlying population, then we must be able to obtain good estimates of the population from which the cases or counts arise. This estimation often leads to the derivation of *expected* rates in the region count case and further to the estimation of the ratio of count to expected count/rate or the *relative risk*, in each area. Relative risk is a fundamental epidemiological concept (Clayton and Hills, 1993) in non-spatial epidemiological studies.

## 1.1 Map Hypotheses and Modelling Approaches

In any spatial epidemiological analysis, there will usually be a study focus which specifies the nature and style of the methods to be used. This focus will usually consist of a hypothesis or hypotheses about the nature of the spatial distribution of the disease which is to be examined, and it is convenient to categorise these hypotheses into three broad classes: *disease mapping*, *ecological analysis* and *disease clustering*. Usually, the distribution of cases of disease, whether in the form of counts or case address locations, can be thought to follow an underlying model, and the observed data may contain extra noise in the form of random variation around the model of interest. Often, the model will include aspects of the *null* (hypothesis) spatial distribution of the cases, which captures the ‘normal’ variation which is expected, and also aspects of the *alternative* spatial distribution. In much of spatial epidemiology, the focus of attention is on identifying features of the spatial distribution which are not captured by the null hypothesis distribution. This is mainly related to *excess spatial aggregation* of cases in areas of the map. That is, once the normal variation is allowed for, the residual spatial incidence *above* the normal incidence is the focus. Seldom is there any need to examine areas of lower aggregation than would be normally expected. Note that ‘normal’ variation is usually assumed to be defined by the underlying population distribution of the study region/window and cases are thought to arise in relation to the local variation in that distribution.

The first class, that of *disease mapping*, concerns the use of models to describe the overall disease distribution on the map. In disease mapping, often the object is simply to ‘clean’ the map of disease of the extra noise to uncover the underlying

structure. In that situation, the null hypothesis could be that the case distribution arises from an *unspecified* or partly specified null spatial distribution (which includes the population spatial distribution) and the object is to remove the extra noise/variation. In this sense disease mapping is close in spirit to image processing where *segmentation* usually describes the process of allocating pixels or groups of pixels to classes.

The second class, that of *ecological analysis*, concerns the analysis of the relation between the spatial distribution of disease incidence and measured explanatory factors. This is usually carried out at an aggregated spatial level, and usually concerns regional incidence compared to explanatory factors measured at regional or other levels of aggregation (Greenberg *et al.*, 1996). This contrasts with studies which use measurements made on individual subjects. However, many of the issues concerning interpretation of ecological studies are concerned with *change* in aggregation level and not aggregated data per se. For example, the *ecological fallacy* concerns making inference about individuals from analyses carried out at a higher scale, e.g. regional or country-wide level. Equally, the *atomistic fallacy* concerns making inferences about average characteristics from individual measurements. In what follows we assume a relatively wide definition of ecological, more in the sense of ecology itself, as any study which seeks to describe/explain the spatial distribution of disease based on the inclusion of explanatory variables. Two classic studies of this kind are presented by Cook and Pocock (1983), who examined the relation of cardiovascular incidence in the UK to a variety of variables (including water hardness, climate, location, socioeconomic and genetic factors and air pollution), and Donnelly (1995), who examined the respiratory health of school children and volatile organic compounds in the outdoor atmosphere. Note that this general definition can include the situation where case address locations are related to a pollution hazard via explanatory variables such as distance and direction from the hazard. In that case individual data are related to explanatory variables.

The final class, that of *disease clustering*, concerns the analysis of 'unusual' aggregations of disease, i.e. assessing whether there are any areas of elevated incidence of disease within a map. This type of analysis could take a variety of forms. First, the analysis could include the assessment of a complete map to ascertain whether the map is *clustered*. This is often termed *general clustering*. In this case, the null hypothesis would be that the disease map represents normal variation in incidence given the population distribution. The alternative hypothesis would include some specified clustering mechanism for the disease cases. This mechanism could be descriptive or include some notion of how the clusters form (e.g. clusters can form if infectious diseases are examined, and the contact rate of individuals can be modelled). General clustering is often treated as a form of *autocorrelation* and models for such effects are often employed. This form of clustering can be termed *non-specific* as it does not seek to determine where clusters are found but instead simply seeks to determine whether the pattern is clustered.

Second, *specific* cluster studies attempt to ascertain the locations of any clusters if they exist on the map. These clusters could have known (fixed) locations and the incidence of disease around these locations may be assessed for its relation to the location(s). Studies of putative pollution hazards fall within this category. This is often termed *focused* clustering. If the locations of clusters are unknown a priori, then the locations must also be estimated from the data; this is termed *non-focused* clustering. Often, ecological regression methods can be used in focused clustering studies, whereas, for non-focused studies, special methods must be constructed which allow the estimation of cluster locations and their form.

In all the above areas of study, fundamental to the methods employed is the inclusion of spatial location in the analysis and so spatial statistical methods are often employed to model the observed data; that epidemiological considerations should be employed in any study of the distribution of disease incidence, in that the concept of normal variation of disease (i.e. that generated from the population at risk from the disease) must be catered for in any model of incidence; and that methods used should be appropriate to the analysis of georeferenced discrete data.

## 1.2 Definitions and Data Examples

In this section, some basic definitions and concepts are introduced which are used throughout this book. In addition, a number of data examples make their first appearance and these will be referred to at various stages throughout the work.

In what follows we will mainly be concerned with data which are available within a single period of time. Hence, we do not provide notation for space-time problems here. Where such notation is appropriate, we provide it locally.

We define ‘epidemiology’ as the study of the occurrence of disease in relation to explanatory factors. A strict dictionary definition of the term implies the study of ‘epidemic diseases’. However, in this work we mainly restrict attention to *fixed* time period studies and do not directly examine the dynamic behaviour of disease incidence. This area has recently been reviewed in Mollison (1995), Daley and Gani (1999) and Andersson and Britton (2000). Some discussion of epidemic models appears in Chapter 10. Here the term ‘spatial epidemiology’ is defined to mean the study of the occurrence of disease in spatial locations and its explanatory factors. Usually, the disease to be examined occurs within a *map* and the data are expressed as a point location (case event) or are aggregated as a count of disease within a subregion of the map. Two examples of such data are provided in Figures 1.2 and 1.3. These two data types lead to different modelling approaches, and we make specific the following definitions as a basis for further discussion.

### 1.2.1 Case event data

We define the study window ( $W$ ), within which  $m$  disease case events occur at locations  $x_i$ ,  $i = 1, \dots, m$ . The area of  $W$  is denoted by  $|W|$ , Lebesgue measure on  $\mathbb{R}^2$ . Figure 1.4 displays these definitions.

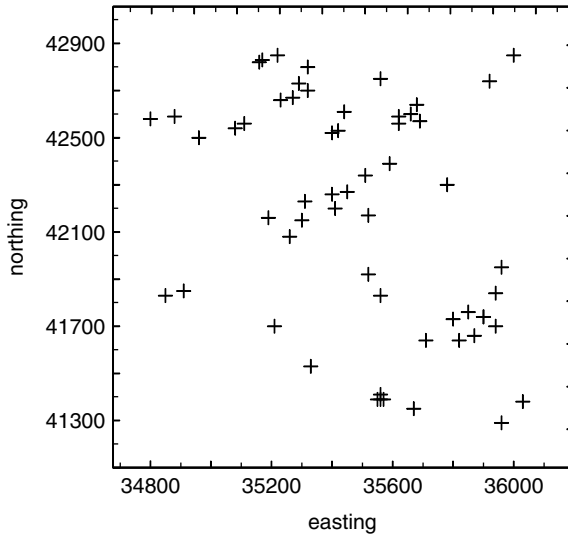


Figure 1.2 The locations of larynx cancer cases in an area of central Lancashire, UK, for the period 1974–1983.

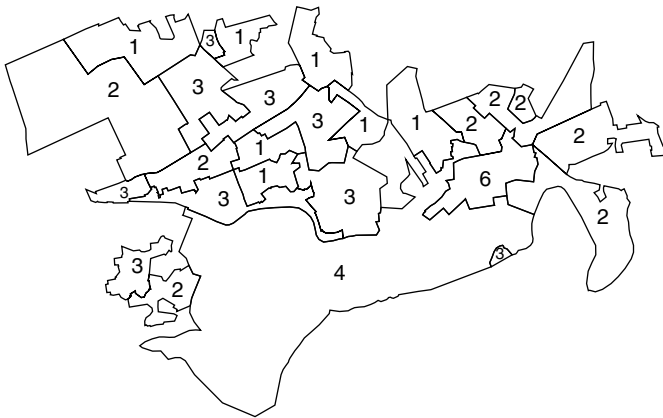


Figure 1.3 Respiratory cancer counts within census tracts (enumeration districts) of Falkirk, central Scotland, for the period 1978–1983.

## 1.2.2 Count data

We define the study window ( $W$ ) as above, within which  $m$  arbitrarily bounded subregions, wholly or in part, lie. The count in  $m$  subregion tracts is denoted  $n_i$ ,  $i = 1, \dots, m$ . In Figure 1.5, only regions 4, 5 and 6 are wholly within the

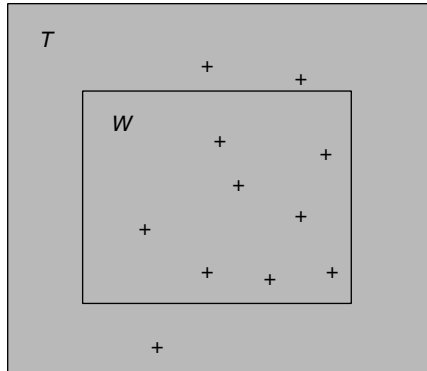


Figure 1.4 A notional study area ( $W$ ) and a guard area ( $T$ ).

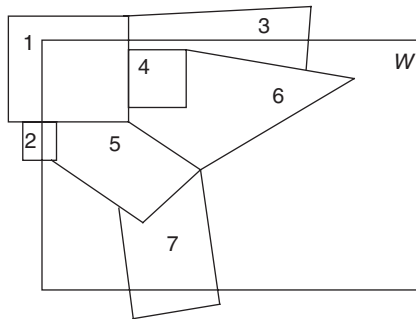


Figure 1.5 A study region within which counts are observed in subregions (tracts).

window. Regions 1, 2, 3 and 7 are cut by the window boundary. The effect of this region truncation will be discussed in detail later. However, it should be noted that, *usually*, the count available ( $n_i$ ) is from the complete region and *not* from the truncated region which appears in the study window.

Usually, the  $m$  subregions are politically defined administrative regions and are often tracts defined for the purposes of population censuses. We adopt the term ‘census tract’ to denote an arbitrarily defined region. In addition, the counts in census tracts are just an aggregation of case event data counted within the bounding tract boundaries. Hence, the data in Figure 1.5 could be derived from the data in Figure 1.4 by counting case events in census tract subregions of the window.

The object of analysis of case event or count data can define the type of summary measures used to describe the data. Usually, as a basic summary measure it is common to compute a local measure of relative risk, or to use a local measure of relative risk as the dependent variable in a more substantial analysis. Here,

relative risk is taken to mean the measure of excess risk found in relation to that supported purely by the local population, which is ‘at risk’. This population is sometimes called the ‘at-risk’ population or background. Relative risk is derived or computed from the relation of observed incidence to that which would be expected based on the ‘at-risk’ background. It is common practice within epidemiology to derive such risk estimates. In the case of spatial epidemiology it is common, when tract count data are available, to compute a standardised mortality (or morbidity) ratio (SMR), which is simply the ratio of the observed count within a tract to the expected count based on the ‘at-risk’ background. A ratio greater than 1.0 would suggest an excess of risk within the tract. These SMRs are often the basis for atlases of disease risk (see, for example, Pickle *et al.*, 1999).

### 1.3 Further Definitions

Some further definitions are required in relation to data which arise in such studies.

#### 1.3.1 Control events and processes

Often, an additional process or realisation of disease events is used to provide an estimate of the ‘background’ incidence of disease in an area. Define  $\mathbf{x}_{c_j}$ ,  $j = 1, \dots, m_c$ , to be these  $m_c$  control event locations. The use of such data will be detailed in a later section.

#### 1.3.2 Census tract information

The census tract count of a control disease is defined to be  $n_c$ .

Instead of using a control disease to represent ‘background’, the ‘expected’ incidence of disease can be used. This is usually based on known rates of disease in the population (Inskip *et al.*, 1983). Denote this expected incidence as  $e_i$ ,  $i = 1, \dots, m$ . The total population of a tract is  $p_i$ , while the extent of the tract is defined as  $a_i$ . The tract centroid, however defined, is denoted by  $\mathbf{x}_{n_i}$ .

For models involving explanatory variables measured at tract level, we define  $F$  as an  $m \times p$  matrix whose columns represent  $p$  explanatory variables, and  $\alpha$  as a  $p \times 1$  vector of parameters. (For case event models the row dimension of  $F$  will usually be  $m$  also.)

#### 1.3.3 Clustering definitions

In cases where clustering is studied, a number of additional definitions are required. First, cluster centre locations are defined as  $\mathbf{y}_j$ ,  $j = 1, \dots, k$ , where  $k$  is the number of centres in a suitably defined window. The term ‘parent’ is used here synonymously with cluster centre. This does not imply any genetic linkage with the observed data. The observed data belonging to a cluster are sometimes referred to as offspring. Again, there is no genetic linkage implied by this term. In addition,