

# Geospatial Analysis

A Comprehensive Guide to Principles,  
Techniques and Software Tools

- Third Edition -

Michael J de Smith

Michael F Goodchild

Paul A Longley

# Copyright and licensing information

Copyright © 2007, 2008, 2009 All Rights reserved. Third Edition. Issue version: 3.12

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 UK Copyright Designs and Patents Act 1998 or with the written permission of the authors. The moral right of the authors has been asserted. Copies of this document are available in printed, electronic book, and web-accessible formats.

## Disclaimer

This publication is designed to offer accurate and authoritative information in regard to the subject matter. It is provided on the understanding that it is not supplied as a form of professional or advisory service. References to software products, datasets or publications are purely made for information purposes and the inclusion or exclusion of any such item does not imply recommendation or otherwise of the product or material in question.

## Licensing and ordering

For ordering (printed and special PDF versions), licensing and contact details please refer to the Guide's website: [www.spatialanalysisonline.com](http://www.spatialanalysisonline.com)

Published by Matador (an imprint of Troubador Publishing Ltd) on behalf of The Winchelsea Press

Postal address: Matador, 9 De Montfort Mews, Leicester, LE1 7FW, UK. Tel: +44 (0)116 255 9311

Matador email address: [books@troubador.co.uk](mailto:books@troubador.co.uk); Web: [www.troubador.co.uk](http://www.troubador.co.uk)

Chinese edition: Publishing House of Electronics Industry, Beijing, PRC, [www.phei.com.cn](http://www.phei.com.cn)

## British Library Cataloguing in Publishing Data

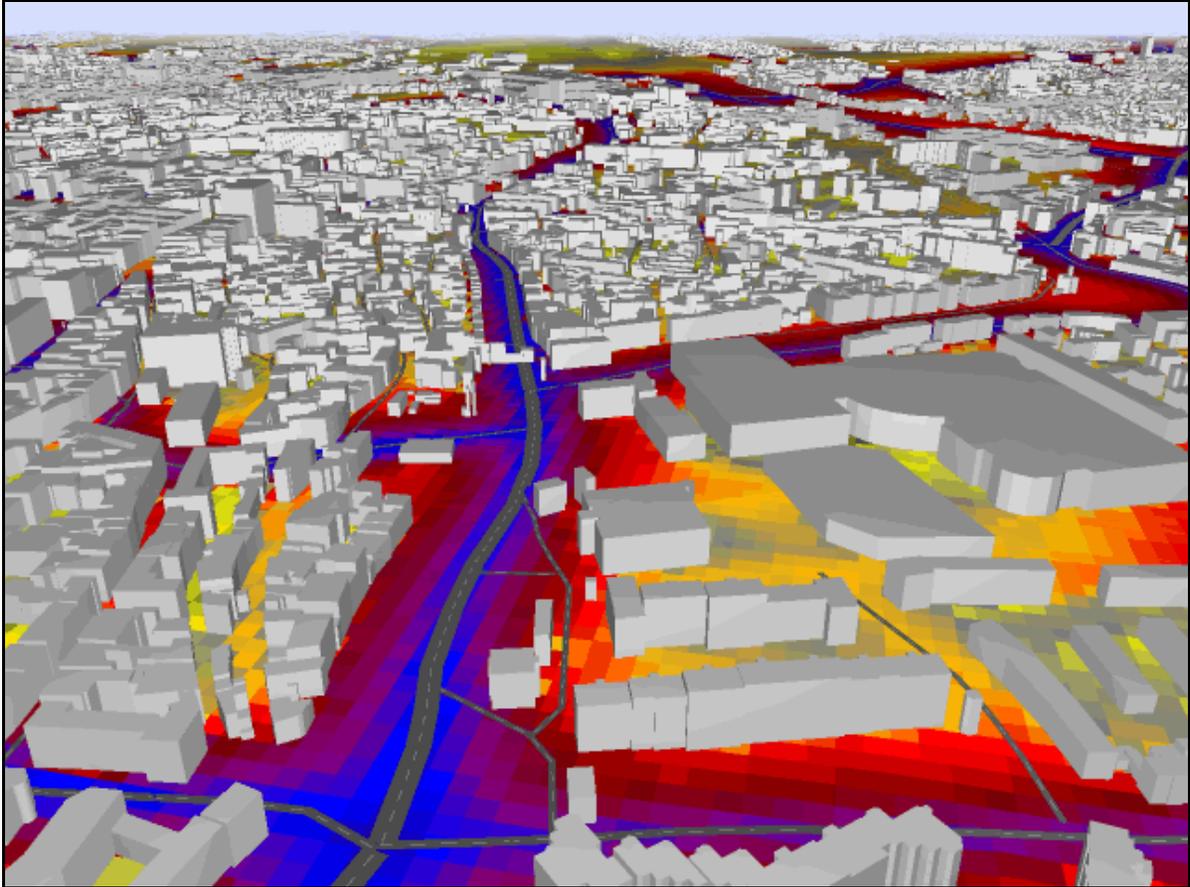
ISBN 13: 9781848761582 (soft cover version)

ISBN 13: 9781848761575 (hard cover version)

Production of this edition has been generously supported by a Fellowship award to Dr de Smith from the UK HEFCE SPLINT initiative:



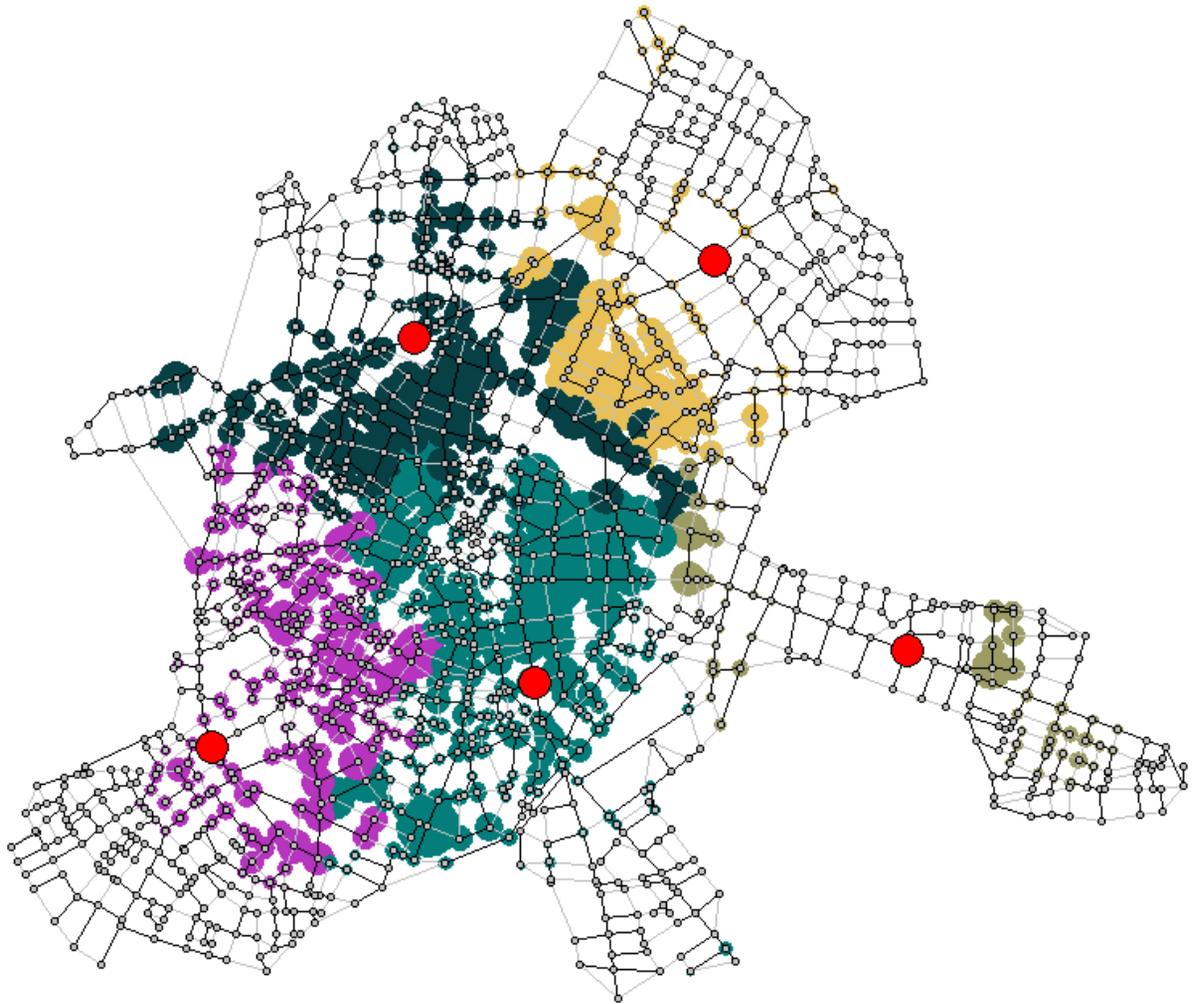
3D visualization of modeled road-related noise levels  
in an urbanized area



*Visualization using CadnaA software, courtesy of Accon GmbH & DataKustik GmbH.  
See Section 3.3 for more details*

## Optimized service center location and allocated demand

### Tripolis, in Arcadia, Greece



*Coverage or  $p$ -center location optimization problem. See Section 7.4.2 for more details. Map produced using *S-Distance* software, courtesy of S A Sirigos, Greece*

# Contents, figures and tables

Copyright and licensing information .....	ii
Contents, figures and tables .....	v
Foreword to the Third Edition .....	xix
Acknowledgements .....	xxi
<b>1 Introduction and terminology .....</b>	<b>23</b>
1.1 Motivation and Media .....	23
1.1.1 Guide overview .....	23
1.1.2 Spatial analysis, GIS and software tools .....	25
1.1.3 Intended audience and scope .....	29
1.2 Software tools and Companion Materials .....	31
1.2.1 GIS and related software tools .....	31
1.2.1.1 Sample software products .....	32
1.2.1.2 Software performance .....	32
1.2.2 Suggested reading .....	33
1.3 Terminology and Abbreviations .....	36
1.3.1 Definitions .....	36
1.4 Common Measures and Notation .....	43
1.4.1 Notation .....	43
1.4.2 Statistical measures and related formulas .....	45
1.4.2.1 Counts and specific values .....	45
1.4.2.2 Measures of centrality .....	46
1.4.2.3 Measures of spread .....	47
1.4.2.4 Measures of distribution shape .....	50
1.4.2.5 Measures of complexity and dimensionality .....	50
1.4.2.6 Common distributions .....	51
1.4.2.7 Data transforms and back transforms .....	52
1.4.2.8 Selected functions .....	53
1.4.2.9 Matrix expressions .....	54
<b>2 Conceptual Frameworks for Spatial Analysis .....</b>	<b>57</b>
2.1 The geospatial perspective .....	57
2.2 Basic Primitives .....	58
2.2.1 Place .....	58
2.2.2 Attributes .....	59
2.2.3 Objects .....	61
2.2.4 Maps .....	62
2.2.5 Multiple properties of places .....	63
2.2.6 Fields .....	64
2.2.7 Networks .....	65
2.2.8 Density estimation .....	65
2.2.9 Detail, resolution, and scale .....	65
2.2.10 Topology .....	66
2.3 Spatial Relationships .....	68
2.3.1 Co-location .....	68
2.3.2 Distance, direction and spatial weights matrices .....	68

2.3.3	Multidimensional scaling .....	69
2.3.4	Spatial context .....	70
2.3.5	Neighborhood .....	70
2.3.6	Spatial heterogeneity .....	71
2.3.7	Spatial dependence .....	71
2.3.8	Spatial sampling .....	72
2.3.9	Spatial interpolation .....	72
2.3.10	Smoothing and sharpening .....	73
2.3.11	First- and second-order processes .....	73
<b>2.4</b>	<b>Spatial Statistics .....</b>	<b>75</b>
2.4.1	Spatial probability .....	75
2.4.2	Probability density .....	75
2.4.3	Uncertainty .....	75
2.4.4	Statistical inference.....	76
<b>2.5</b>	<b>Spatial Data Infrastructure .....</b>	<b>78</b>
2.5.1	Geoportals .....	78
2.5.2	Metadata .....	79
2.5.3	Interoperability .....	79
2.5.4	Conclusion .....	79
<b>3</b>	<b>Methodological Context .....</b>	<b>81</b>
3.1	Spatial analysis as a process .....	81
3.2	Analytical methodologies .....	83
3.3	Spatial analysis and the PPDAC model .....	87
3.3.1	Problem: Framing the question .....	88
3.3.2	Plan: Formulating the approach.....	90
3.3.3	Data: Data acquisition .....	91
3.3.4	Analysis: Analytical methods and tools .....	93
3.3.5	Conclusions: Delivering the results .....	95
3.4	Geospatial analysis and model building .....	96
3.5	The changing context of GIScience .....	102
<b>4</b>	<b>Building Blocks of Spatial Analysis .....</b>	<b>105</b>
4.1	Spatial Data Models and Methods .....	105
4.2	Geometric and Related Operations .....	107
4.2.1	Length and area for vector data .....	107
4.2.2	Length and area for raster datasets .....	109
4.2.3	Surface area .....	111
4.2.3.1	Projected surfaces .....	111
4.2.3.2	Terrestrial (unprojected) surface area .....	113
4.2.4	Line Smoothing and point-weeding .....	114
4.2.5	Centroids and centers.....	116
4.2.5.1	Polygon centroids and centers .....	116
4.2.5.2	Point sets .....	119
4.2.5.3	Lines.....	121
4.2.6	Point (object) in polygon (PIP) .....	121
4.2.7	Polygon decomposition .....	123
4.2.8	Shape .....	124
4.2.9	Overlay and combination operations .....	125
4.2.10	Areal interpolation .....	128
4.2.11	Districting and re-districting.....	130
4.2.12	Classification and clustering.....	135
4.2.12.1	Univariate classification schemes .....	135

4.2.12.2	Multivariate classification and clustering .....	138
4.2.12.3	Multi-band image classification .....	140
4.2.12.4	Uncertainty and image processing.....	145
4.2.12.5	Hyperspectral image classification .....	146
4.2.13	Boundaries and zone membership .....	149
4.2.13.1	Convex hulls.....	149
4.2.13.2	Non-convex hulls .....	150
4.2.13.3	Minimum Bounding Rectangles (MBRs) .....	152
4.2.13.4	Fuzzy boundaries .....	153
4.2.13.5	Breaklines and natural boundaries .....	156
4.2.14	Tessellations and triangulations.....	157
4.2.14.1	Delaunay Triangulation.....	157
4.2.14.2	TINs – Triangulated irregular networks .....	158
4.2.14.3	Voronoi/Thiessen polygons .....	159
<b>4.3</b>	<b>Queries, Computations and Density .....</b>	<b>163</b>
4.3.1	Spatial selection and spatial queries .....	163
4.3.2	Simple calculations.....	163
4.3.3	Ratios, indices, normalization, standardization and rate smoothing.....	166
4.3.4	Density, kernels and occupancy.....	171
4.3.4.1	Point density .....	171
4.3.4.2	Kernel density for networks .....	178
4.3.4.3	Line and intersection densities .....	179
4.3.4.4	Cartograms .....	179
<b>4.4</b>	<b>Distance Operations .....</b>	<b>184</b>
4.4.1	Metrics.....	186
4.4.1.1	Introduction .....	186
4.4.1.2	Terrestrial distances.....	187
4.4.1.3	Extended Euclidean and $L_p$ -metric distances .....	188
4.4.2	Cost distance .....	190
4.4.2.1	Accumulated cost surfaces and least cost paths .....	191
4.4.2.2	Distance transforms.....	196
4.4.3	Network distance .....	202
4.4.4	Buffering .....	203
4.4.4.1	Vector buffering.....	203
4.4.4.2	Raster buffering .....	205
4.4.4.3	Hybrid buffering.....	205
4.4.4.4	Network buffering.....	205
4.4.5	Distance decay models .....	205
<b>4.5</b>	<b>Directional Operations .....</b>	<b>210</b>
4.5.1	Directional analysis – overview.....	210
4.5.2	Directional analysis of linear datasets .....	210
4.5.3	Directional analysis of point datasets .....	215
4.5.4	Directional analysis of surfaces .....	217
<b>4.6</b>	<b>Grid Operations and Map Algebra .....</b>	<b>219</b>
4.6.1	Operations on single and multiple grids .....	219
4.6.2	Linear spatial filtering .....	220
4.6.3	Non-linear spatial filtering.....	223
4.6.4	Erosion and dilation .....	224
<b>5</b>	<b>Data Exploration and Spatial Statistics .....</b>	<b>225</b>
<b>5.1</b>	<b>Statistical Methods and Spatial Data .....</b>	<b>225</b>
5.1.1	Descriptive statistics .....	227

5.1.2	Spatial sampling .....	228
5.1.2.1	Sampling frameworks .....	230
5.1.2.2	Declustering .....	234
<b>5.2</b>	<b>Exploratory Spatial Data Analysis .....</b>	<b>236</b>
5.2.1	EDA, ESDA and ESTDA .....	236
5.2.2	Outlier detection .....	238
5.2.2.1	Mapped histograms.....	238
5.2.2.2	Box plots.....	239
5.2.3	Cross tabulations and conditional choropleth plots.....	241
5.2.4	ESDA and mapped point data .....	243
5.2.5	Trend analysis of continuous data .....	244
5.2.6	Cluster hunting and scan statistics .....	244
<b>5.3</b>	<b>Grid-based Statistics .....</b>	<b>247</b>
5.3.1	Overview of grid-based statistics .....	247
5.3.2	Crosstabulated grid data, the Kappa Index and Cramer's V statistic .....	248
5.3.3	Quadrat analysis of grid datasets .....	250
5.3.4	Landscape Metrics .....	252
5.3.4.1	Non-spatial landscape metrics .....	254
5.3.4.2	Spatial landscape metrics.....	255
<b>5.4</b>	<b>Point Sets and Distance Statistics .....</b>	<b>259</b>
5.4.1	Basic distance-derived statistics .....	259
5.4.2	Nearest neighbor methods .....	260
5.4.3	Pairwise distances .....	263
5.4.4	Hot spot and cluster analysis .....	267
5.4.4.1	Hierarchical nearest neighbor clustering.....	268
5.4.4.2	K-means clustering.....	269
5.4.4.3	Kernel density clustering.....	269
5.4.4.4	Spatio-temporal clustering .....	270
5.4.5	Proximity matrix comparisons .....	272
<b>5.5</b>	<b>Spatial Autocorrelation .....</b>	<b>274</b>
5.5.1	Autocorrelation, time series and spatial analysis .....	274
5.5.2	Global spatial autocorrelation.....	276
5.5.2.1	Join counts and the analysis of nominal-valued spatial data .....	276
5.5.2.2	Moran <i>I</i> and Geary <i>C</i> .....	282
5.5.2.3	Weighting models and lags.....	289
5.5.3	Local indicators of spatial association (LISA) .....	290
5.5.4	Significance tests for autocorrelation indices.....	292
<b>5.6</b>	<b>Spatial Regression .....</b>	<b>294</b>
5.6.1	Regression overview.....	294
5.6.2	Simple regression and trend surface modeling .....	299
5.6.3	Geographically Weighted Regression (GWR) .....	301
5.6.4	Spatial autoregressive and Bayesian modeling.....	305
5.6.4.1	Spatial autoregressive modeling .....	305
5.6.4.2	Conditional autoregressive and Bayesian modeling.....	308
5.6.5	Spatial filtering models.....	311
<b>6</b>	<b>Surface and Field Analysis .....</b>	<b>313</b>
<b>6.1</b>	<b>Modeling Surfaces .....</b>	<b>313</b>
6.1.1	Test datasets .....	313
6.1.2	Surfaces and fields .....	314
6.1.3	Raster models .....	315
6.1.4	Vector models.....	318

6.1.5	Mathematical models .....	319
6.1.6	Statistical and fractal models .....	320
<b>6.2</b>	<b>Surface Geometry .....</b>	<b>323</b>
6.2.1	Gradient, slope and aspect .....	323
6.2.1.1	Slope .....	323
6.2.1.2	Aspect .....	325
6.2.2	Profiles and curvature .....	328
6.2.2.1	Profiles and cross-sections .....	328
6.2.2.2	Curvature and morphometric analysis .....	328
6.2.2.3	Profile curvature .....	331
6.2.2.4	Plan curvature .....	331
6.2.2.5	Tangential curvature .....	332
6.2.2.6	Longitudinal and cross-sectional curvature .....	332
6.2.2.7	Mean, maximum and minimum curvature .....	332
6.2.3	Directional derivatives .....	332
6.2.4	Paths on surfaces .....	333
6.2.5	Surface smoothing .....	334
6.2.6	Pit filling .....	335
6.2.7	Volumetric analysis .....	336
<b>6.3</b>	<b>Visibility .....</b>	<b>337</b>
6.3.1	Viewsheds and RF propagation .....	337
6.3.2	Line of sight .....	340
6.3.3	Isovist analysis and space syntax .....	341
6.3.3.1	Isovists .....	341
6.3.3.2	Space syntax .....	343
<b>6.4</b>	<b>Watersheds and Drainage .....</b>	<b>345</b>
6.4.1	Overview of watersheds and drainage .....	345
6.4.2	Drainage modeling .....	345
6.4.3	D-infinity model .....	346
6.4.4	Drainage modeling case study .....	347
6.4.4.1	Flow accumulation .....	347
6.4.4.2	Stream network construction .....	347
6.4.4.3	Stream basin construction .....	348
<b>6.5</b>	<b>Gridding, Interpolation and Contouring .....</b>	<b>349</b>
6.5.1	Overview of gridding and interpolation .....	349
6.5.2	Gridding and interpolation methods .....	350
6.5.3	Contouring .....	356
<b>6.6</b>	<b>Deterministic Interpolation Methods .....</b>	<b>358</b>
6.6.1	Inverse distance weighting (IDW) .....	359
6.6.2	Natural neighbor .....	361
6.6.3	Nearest-neighbor .....	363
6.6.4	Radial basis and spline functions .....	363
6.6.5	Modified Shepard .....	364
6.6.6	Triangulation with linear interpolation .....	365
6.6.7	Triangulation with spline-like interpolation .....	365
6.6.8	Rectangular or bi-linear interpolation .....	366
6.6.9	Profiling .....	366
6.6.10	Polynomial regression .....	366
6.6.11	Minimum curvature .....	367
6.6.12	Moving average .....	367
6.6.13	Local polynomial .....	367
6.6.14	Topogrid/Topo to raster .....	368

<b>6.7</b>	<b>Geostatistical Interpolation Methods</b>	<b>369</b>
6.7.1	Core concepts	369
6.7.1.1	Geostatistics	369
6.7.1.2	Geostatistical references	370
6.7.1.3	Semivariance	370
6.7.1.4	Sample size	371
6.7.1.5	Support	371
6.7.1.6	Declustering	371
6.7.1.7	Variogram	372
6.7.1.8	Stationarity	372
6.7.1.9	Sill, range and nugget	372
6.7.1.10	Transformation	373
6.7.1.11	Anisotropy	374
6.7.1.12	Indicator semivariance	375
6.7.1.13	Cross-semivariance	375
6.7.1.14	Comments on geostatistical software packages	376
6.7.1.15	Semivariance modeling	377
6.7.1.16	Fractal analysis	381
6.7.1.17	Madograms and Rodograms	381
6.7.1.18	Periodograms and Fourier analysis	381
6.7.2	Kriging interpolation	382
6.7.2.1	Core process	382
6.7.2.2	Goodness of fit	384
6.7.2.3	Simple Kriging	384
6.7.2.4	Ordinary Kriging	385
6.7.2.5	Universal Kriging	386
6.7.2.6	Median-Polishing and Kriging	386
6.7.2.7	Indicator Kriging	386
6.7.2.8	Probability Kriging	386
6.7.2.9	Disjunctive Kriging	386
6.7.2.10	Non-stationary Modeling and Stratified Kriging	387
6.7.2.11	Co-Kriging	387
6.7.2.12	Factorial Kriging	387
6.7.2.13	Conditional simulation	387
<b>7</b>	<b>Network and Location Analysis</b>	<b>391</b>
7.1	Introduction to Network and Location Analysis	391
7.1.1	Overview of network and location analysis	391
7.1.2	Terminology	391
7.1.3	Source data	393
7.1.4	Algorithms and computational complexity theory	395
7.2	Key Problems in Network and Location Analysis	397
7.2.1	Overview – network analysis	397
7.2.1.1	Key problems in network analysis	398
7.2.1.2	Network analysis software	401
7.2.1.3	Key problems in location analysis	404
7.2.2	Heuristic and meta-heuristic algorithms	406
7.2.2.1	Greedy heuristics and local search	407
7.2.2.2	Interchange heuristics	408
7.2.2.3	Metaheuristics	409
7.2.2.4	Tabu search	409
7.2.2.5	Cross-entropy (CE) methods	410

7.2.2.6	Simulated annealing .....	410
7.2.2.7	Lagrangian multipliers and Lagrangian relaxation .....	411
7.2.2.8	Ant systems and ant colony optimization (ACO) .....	414
<b>7.3</b>	<b>Network Construction, Optimal Routes and Optimal Tours .....</b>	<b>416</b>
7.3.1	Minimum spanning tree .....	416
7.3.2	Gabriel network.....	417
7.3.3	Steiner trees.....	419
7.3.4	Shortest (network) path problems .....	420
7.3.4.1	Overview of shortest path problems .....	420
7.3.4.2	Dantzig algorithm .....	421
7.3.4.3	Dijkstra algorithm .....	421
7.3.4.4	A* algorithm .....	422
7.3.4.5	GIS implementations of SPAs .....	422
7.3.4.6	Further SPAs applications.....	424
7.3.5	Tours, travelling salesman problems and vehicle routing.....	425
7.3.5.1	Capacitated vehicle routing .....	428
<b>7.4</b>	<b>Location and Service Area Problems .....</b>	<b>430</b>
7.4.1	Location problems .....	430
7.4.2	Larger p-median and p-center problems.....	433
7.4.2.1	Simple heuristics .....	433
7.4.2.2	Lagrangian relaxation.....	433
7.4.2.3	Comparison of alternative p-median heuristics .....	436
7.4.3	Service areas .....	438
7.4.3.1	Travel time zones .....	439
<b>7.5</b>	<b>Arc Routing .....</b>	<b>441</b>
7.5.1	Network traversal problems .....	441
<b>8</b>	<b>Geocomputational methods and modeling .....</b>	<b>445</b>
<b>8.1</b>	<b>Introduction to Geocomputation .....</b>	<b>445</b>
8.1.1	Geocomputational methods.....	445
8.1.2	Modeling dynamic processes within GIS.....	446
8.1.2.1	Representing time and change within GIS.....	447
8.1.2.2	Linkage/coupling versus integration/embedding .....	449
<b>8.2</b>	<b>Geosimulation .....</b>	<b>452</b>
8.2.1	Introduction to geosimulation .....	452
8.2.2	Cellular automata (CA) .....	452
8.2.3	Agents and agent-based models.....	456
8.2.3.1	Agent-based models .....	456
8.2.3.2	Agents .....	457
8.2.4	Applications of agent-based models.....	459
8.2.5	Advantages of agent-based models .....	462
8.2.6	Limitations of agent-based models .....	463
8.2.7	Explanation or prediction? .....	464
8.2.8	Developing an agent-based model .....	466
8.2.9	Types of simulation/modeling (s/m) systems for agent-based modeling.....	467
8.2.10	Guidelines for choosing a simulation/modeling (s/m) system.....	469
8.2.11	Simulation/modeling (s/m) systems for agent-based modeling.....	470
8.2.12	Verification and calibration of agent-based models .....	482
8.2.13	Validation and analysis of agent-based model outputs.....	484
<b>8.3</b>	<b>Artificial Neural Networks (ANN).....</b>	<b>486</b>
8.3.1	Introduction to artificial neural networks .....	486
8.3.1.1	Multi-level perceptrons (MLP).....	487

8.3.1.2	Learning and back-propagation for MLPs .....	489
8.3.1.3	MLP Example 1: Function approximation .....	492
8.3.1.4	MLP Example 2: Landcover change modeling (LCM) .....	494
8.3.1.5	MLP Example 3: Spatial interaction modeling .....	497
8.3.2	Radial basis function networks .....	499
8.3.3	Self organizing networks .....	501
8.3.3.1	Self Organizing Maps (SOMs) .....	501
8.3.3.2	SOM unsupervised classification of hyper-spectral image data .....	502
8.3.3.3	TSP optimization using SOM concepts.....	507
<b>8.4</b>	<b>Genetic Algorithms and Evolutionary Computing .....</b>	<b>509</b>
8.4.1	Genetic algorithms – introduction .....	509
8.4.2	Genetic algorithm components.....	510
8.4.2.1	Encoding or representation.....	510
8.4.2.2	Fitness function.....	511
8.4.2.3	Population initialization.....	512
8.4.2.4	Selection.....	512
8.4.2.5	Reproduction.....	513
8.4.2.6	Crossover .....	513
8.4.2.7	Mutation .....	514
8.4.2.8	Local search.....	514
8.4.2.9	Termination .....	514
8.4.3	Example GA applications .....	514
8.4.3.1	GA Example 1: TSP .....	514
8.4.3.2	GA Example 2: Clustering.....	514
8.4.3.3	GA Example 3: Map labeling.....	515
8.4.3.4	GA Example 4: Optimum location .....	517
8.4.4	Evolutionary computing and genetic programming .....	518
<b>Afterword.....</b>		<b>519</b>
<b>References .....</b>		<b>521</b>
<b>CATMOG Guides .....</b>		<b>538</b>
<b>R-Project spatial statistics software packages .....</b>		<b>540</b>
<b>Fragstats landscape metrics .....</b>		<b>543</b>
<b>Web links .....</b>		<b>545</b>
Associations and academic bodies .....		545
Online technical dictionaries/definitions .....		546
Spatial data, test data and spatial information sources .....		546
Selected national and international data and information sources .....		546
Transport planning and analysis.....		547
Network analysis test datasets.....		547
Statistics and Spatial Statistics links .....		547
Other GIS web sites .....		547
<b>Index.....</b>		<b>549</b>

# List of Figures

Figure 1-1 3D Physical GIS models.....	28
Figure 2-1 Attribute tables - spatial datasets.....	59
Figure 2-2 Cyclic attribute data – Wind direction, single location .....	61
Figure 2-3 An example map showing points, lines, and areas appropriately symbolized.....	62
Figure 2-4 Layers and overlay .....	63
Figure 2-5 Noise level raster .....	64
Figure 2-6 Filled contour view of field data .....	64
Figure 2-7 Topological relationships .....	67
Figure 2-8 Spatial weights computation .....	69
Figure 2-9 Three alternative ways of defining neighborhood, using simple GIS functions .....	70
Figure 2-10 Simple interpolation modeling .....	73
Figure 2-11 Four distinct patterns of twelve points in a study area .....	74
Figure 2-12 The process of statistical inference .....	76
Figure 2-13 Pupil performance and school catchment area in the East Riding of Yorkshire, UK....	79
Figure 3-1 Analytical process – Mitchell .....	83
Figure 3-2 Analytical process – Draper .....	84
Figure 3-3 PPDAC as an iterative process .....	86
Figure 3-4 Noise map, Augsburg .....	87
Figure 3-5 Simple GIS graphical model (ESRI ArcGIS).....	96
Figure 3-6 Dynamic residential growth model (Idrisi) .....	98
Figure 3-7 Modeling wildfire risks, Arizona, USA .....	99
Figure 4-1 Area calculation using Simpson’s rule .....	107
Figure 4-2 3x3 grid neighborhood.....	110
Figure 4-3 5x5 grid neighborhood.....	110
Figure 4-4 Planimetric and surface area of a 3D triangle.....	111
Figure 4-5 DEM surface area.....	112
Figure 4-6 Surface model of DEM .....	113
Figure 4-7 Smoothing techniques .....	115
Figure 4-8 Triangle centroid.....	117
Figure 4-9 Polygon centroid (M2) and alternative polygon centers .....	117
Figure 4-10 Center and centroid positioning.....	118
Figure 4-11 Polygon center selection .....	119
Figure 4-12 Point set centers.....	120
Figure 4-13 Point in polygon – tests and special cases .....	122
Figure 4-14 Skeletonised convex polygon .....	123
Figure 4-15 GRASS overlay operations, v.overlay .....	126
Figure 4-16 Areal interpolation from census areas to a single grid cell .....	129
Figure 4-17 Proportionally assigned population values.....	129
Figure 4-18 Grouping data – Zone arrangement effects on voting results.....	132
Figure 4-19 Creating postcode polygons.....	134
Figure 4-20 Automated Zone Procedure (AZP) .....	134
Figure 4-21 AZP applied to part of Manchester, UK .....	135
Figure 4-22 Jenks Natural Breaks algorithm .....	138
Figure 4-23 SPOT Band 1 image histogram – distinct peaks highlighted (CLUSTER) .....	145
Figure 4-24 2D map of Cuprite mining district, Western Nevada, USA .....	147
Figure 4-25 3D hypercube visualization of Cuprite mining district, Western Nevada, USA.....	147
Figure 4-26 Single class assignment from spectral angle analysis .....	148
Figure 4-27 Convex hull of sample point set.....	149
Figure 4-28 Alpha hulls .....	151
Figure 4-29 Interpolation within “centroid” MBR .....	153

Figure 4-30 Point locations inside and outside bounding polygon.....	153
Figure 4-31 Sigmoidal fuzzy membership functions .....	155
Figure 4-32 Delaunay triangulation of spot height locations.....	158
Figure 4-33 Voronoi regions generated in ArcGIS and MATLAB.....	160
Figure 4-34 Voronoi cells for a homogeneous grid using a 3x3 distance transform.....	161
Figure 4-35 Network-based Voronoi regions – Shibuya district, Tokyo.....	162
Figure 4-36 Cell-by-cell or Local operations .....	165
Figure 4-37 Map algebra: Index creation .....	166
Figure 4-38 Normalization within ArcGIS .....	169
Figure 4-39 Quantile map of normalized SIDS data .....	169
Figure 4-40 Excess risk rate map for SIDS data .....	170
Figure 4-41 Point data .....	172
Figure 4-42 Simple linear (box or uniform) kernel smoothing.....	172
Figure 4-43 Univariate Normal kernel smoothing and cumulative densities.....	173
Figure 4-44 Alternative univariate kernel density functions .....	173
Figure 4-45 2D Normal kernel .....	174
Figure 4-46 Kernel density map, Lung Case data, 3D visualization .....	175
Figure 4-47 Univariate kernel density functions, unit bandwidth .....	177
Figure 4-48 Cartogram creation using basic Dorling algorithm .....	180
Figure 4-49 Cartogram creation using Dougenik, Chrisman and Niemeyer algorithm .....	181
Figure 4-50 World Population as a Cartogram .....	181
Figure 4-51 Cartograms of births data, 1974 .....	182
Figure 4-52 Hexagonal cartogram showing UK mortality data, age group 20-24.....	183
Figure 4-53 Alternative measures of terrain distance.....	184
Figure 4-54 Glasgow's Clockwork Orange Underground .....	186
Figure 4-55 Great circle and constant bearing paths, Boston to Bristol.....	187
Figure 4-56 p-metric circles .....	189
Figure 4-57 Cost distance model.....	191
Figure 4-58 Cost surface as grid .....	191
Figure 4-59 Grid resolution and cost distance .....	192
Figure 4-60 Accumulated cost surface and least cost paths.....	193
Figure 4-61 Alternative route selection by ACS .....	194
Figure 4-62 Steepest path vs tracked path.....	195
Figure 4-63 3x3 Distance transformation - scan elements.....	196
Figure 4-64 5x5 DT mask .....	198
Figure 4-65 Distance transform, single point .....	198
Figure 4-66 Urban traffic modeling.....	199
Figure 4-67 Notting Hill Carnival routes .....	199
Figure 4-68 Alternative routes selected by gradient constrained DT .....	200
Figure 4-69 Hellisheiði power plant pipeline route selection .....	201
Figure 4-70 Shortest and least time paths .....	202
Figure 4-71 Simple buffering .....	203
Figure 4-72 Manifold: Buffer operations.....	204
Figure 4-73 Manifold: Buffering options .....	204
Figure 4-74 Inverse distance decay models, $\alpha/d^{\beta}$ .....	208
Figure 4-75 Exponential distance decay models, $\alpha e^{-\beta d}$ .....	208
Figure 4-76 Directional analysis of streams .....	212
Figure 4-77 Two-variable wind rose .....	214
Figure 4-78 Standard distance circle and ellipses.....	215
Figure 4-79 Correlated Random Walk simulation .....	216
Figure 4-80 Slope and aspect plot, Mt St Helens data, USA .....	217
Figure 4-81 Wind flow grid simulation using WindNinja.....	218

Figure 4-82 Dilation and erosion operations .....	224
Figure 5-1 Point-based sampling schemes.....	229
Figure 5-2 Grid generation examples .....	231
Figure 5-3 Grid sampling examples within hexagonal grid, 1 hectare area.....	231
Figure 5-4 Random point generation examples – ArcGIS .....	233
Figure 5-5 Random point samples on a network.....	233
Figure 5-6 Brushing and linking, GeoDa.....	236
Figure 5-7 Parallel coordinate plot.....	237
Figure 5-8 Star plot .....	237
Figure 5-9 Histogram linkage .....	239
Figure 5-10 Simple box plot .....	239
Figure 5-11 Mapped box plot, GeoDa .....	240
Figure 5-12 Conditional Choropleth mapping.....	242
Figure 5-13 Exploratory analysis of radioactivity data .....	243
Figure 5-14 Trend analysis of radioactivity dataset .....	244
Figure 5-15 Quadrat counts.....	250
Figure 5-16 Texture analysis – variability.....	253
Figure 5-17 Nearest Neighbor distribution .....	260
Figure 5-18 Ripley’s K function computation .....	263
Figure 5-19 Ripley K function, shown as transformed L function plot.....	265
Figure 5-20 Thomas Poisson Cluster Process (20 clusters, SD=0.03, mean=5).....	265
Figure 5-21 Lung cancer incidence data.....	267
Figure 5-22 Lung cancer NNh clusters.....	269
Figure 5-23 KDE cancer incidence mapping.....	270
Figure 5-24 Time series of stock price and volume data.....	274
Figure 5-25 Join count patterns.....	277
Figure 5-26 Join count computation .....	278
Figure 5-27 Homogeneous and non-homogenous probability images .....	279
Figure 5-28 Grouping and size effects .....	281
Figure 5-29 Irregular lattice dataset .....	282
Figure 5-30 Adjacency matrix, W.....	283
Figure 5-31 Moran’s I computation.....	285
Figure 5-32 Revised source data .....	285
Figure 5-33 Sample dataset and Moran I analysis .....	287
Figure 5-34 Moran I (co)variance cloud, lag 1 .....	288
Figure 5-35 Local Moran I computation .....	291
Figure 5-36 LISA map, Moran I.....	291
Figure 5-37 Significance tests for revised sample dataset.....	293
Figure 5-38 Georgia educational attainment: GWR residuals map, Gaussian adaptive kernel.....	303
Figure 6-1 East Sussex test surface, OS TQ81NE .....	313
Figure 6-2 Pentland Hills test surface.....	314
Figure 6-3 Linear regression surface fit to NT04 spot heights .....	315
Figure 6-4 Raster file neighborhoods.....	316
Figure 6-5 Vector models of TQ81NE.....	319
Figure 6-6 First, second and third order mathematical surfaces .....	320
Figure 6-7 Random and fractal grids .....	321
Figure 6-8 Pseudo-random surfaces .....	322
Figure 6-9 8-triangle slope computation .....	324
Figure 6-10 Gradient and sampling resolution.....	325
Figure 6-11 Slope computation output.....	325
Figure 6-12 Frequency distribution of aspect values .....	326
Figure 6-13 Aspect computation output .....	327

Figure 6-14 Profile of NS transect, TQ81NE.....	328
Figure 6-15 Multiple profile computation .....	328
Figure 6-16 Surface morphology .....	329
Figure 6-17 Path smoothing – vertical profile.....	334
Figure 6-18 Grid smoothing .....	335
Figure 6-19 Viewshed computation.....	338
Figure 6-20 3D Urban radio wave propagation modeling using Cellular Expert and ArcGIS .....	339
Figure 6-21 Radio frequency viewshed .....	340
Figure 6-22 Line of sight analysis .....	341
Figure 6-23 Viewsheds and lines of sight on a synthetic (Gaussian) surface.....	341
Figure 6-24 Isovist analysis, Street network, central London .....	342
Figure 6-25 Axial lines and connectivity.....	343
Figure 6-26 Depthmap – Gallery space visibility map.....	344
Figure 6-27 D-Infinity flow assignment .....	346
Figure 6-28 Flow direction and accumulation .....	347
Figure 6-29 Stream identification .....	348
Figure 6-30 Watersheds and basins.....	348
Figure 6-31 Contour plots for alternative interpolation methods – generated with Surfer 8 .....	353
Figure 6-32 Linear interpolation of contours .....	356
Figure 6-33 Contour computation output.....	357
Figure 6-34 IDW as surface plot.....	359
Figure 6-35 Contour plots for alternative IDW methods, OS NT04 .....	360
Figure 6-36 Natural Neighbor interpolation – computation of weights .....	362
Figure 6-37 Clough-Tocher TIN interpolation .....	366
Figure 6-38 Regression fitting to test dataset OS NT04 .....	367
Figure 6-39 Sample variogram.....	371
Figure 6-40 Sill, range and nugget.....	373
Figure 6-41 Data transformation for Normality.....	374
Figure 6-42 Anisotropy 2D map, zinc data .....	375
Figure 6-43 Indicator variograms .....	375
Figure 6-44 Variogram models – graphs.....	380
Figure 6-45 Fractal analysis of TQ81NE.....	381
Figure 6-46 Ordinary Kriging of zinc dataset.....	385
Figure 6-47 Conditional simulation of untransformed zinc test dataset.....	389
Figure 7-1 Network topologies .....	393
Figure 7-2 Visualization of lane/movement simulation (Dynameq) .....	403
Figure 7-3 LP Solution graphs .....	413
Figure 7-4 Minimum Spanning Tree.....	416
Figure 7-5 Gabriel network construction .....	417
Figure 7-6 Relative neighborhood network and related constructions .....	418
Figure 7-7 Steiner MST construction .....	419
Figure 7-8 Dantzig shortest path algorithm .....	421
Figure 7-9 Salt Lake City – Sample networking problems and solutions .....	423
Figure 7-10 Shortest obstacle-avoiding path.....	424
Figure 7-11 MST, TSP and related problems .....	427
Figure 7-12 Heuristic solution and dual circuit TSP examples .....	428
Figure 7-13 Tanker delivery tours .....	429
Figure 7-14 Optimum facility location on a network – LOLA solution.....	432
Figure 7-15 Comparison of heuristic p-median solutions, Tripolis, Greece.....	437
Figure 7-16 Facility location in Tripolis, Greece, planar model .....	438
Figure 7-17 Service area definition.....	439
Figure 7-18 Travel-time or drive-time zones .....	440

Figure 7-19 Routing directions .....	441
Figure 7-20 Arc routing.....	442
Figure 8-1 Game of Life Model .....	453
Figure 8-2 Heatbugs Model.....	453
Figure 8-3 Moore and von Neumann neighborhoods .....	455
Figure 8-4 Schelling segregation model.....	460
Figure 8-5 Pedestrian movement simulation – Subway hall model .....	461
Figure 8-6 Model development balance.....	470
Figure 8-7 Geometric and Locational Features of the Notting Hill Carnival Swarm model.....	472
Figure 8-8 RepastS point and click modeling and runtime environments .....	474
Figure 8-9 RepastCity – importing GIS network data into Repast Symphony .....	475
Figure 8-10 Repast Symphony – agent-based model visualized using NASA’s Worldwind.....	475
Figure 8-11 StarLogo TNG drag and drop programming interface and 3D view of a simulation ...	478
Figure 8-12 Outputs from OBEUS: Schelling residential dynamics model .....	479
Figure 8-13 AgentSheets: The Boulder Mountain Biking Advisor .....	481
Figure 8-14 An urban and transport dynamics model developed in AnyLogic (2006).....	482
Figure 8-15 Simple 3-5-2 feedforward artificial neural network .....	486
Figure 8-16 MLP 3-5-2 with bias nodes.....	487
Figure 8-17 ANN hidden node structure .....	487
Figure 8-18 Sample activation functions .....	488
Figure 8-19 MLP: Test data and fitted model.....	493
Figure 8-20 MLP: RMSE curves.....	493
Figure 8-21 Land cover, 1986, Chiquitania.....	495
Figure 8-22 Distance raster (meters), anthropogenic disturbance, Chiquitania .....	495
Figure 8-23 MLP Classifier – Idrisi.....	496
Figure 8-24 Transition potential map .....	497
Figure 8-25 MLP trip distribution model 1.....	498
Figure 8-26 MLP trip distribution model 2.....	499
Figure 8-27 Radial basis function NN model .....	500
Figure 8-28 SOM grids.....	502
Figure 8-29 SOM classification of remotely-sensed hyperspectral data .....	504
Figure 8-30 Self Organizing Map (SOM) classification – Idrisi .....	506
Figure 8-31 SOM classified 3-band image.....	507
Figure 8-32 Rank score transform .....	512

# List of Tables

Table 1-1 Selected terminology.....	36
Table 1-2 Notation and symbology .....	43
Table 1-3 Common formulas and statistical measures .....	45
Table 4-1 Geographic data models .....	105
Table 4-2 OGC OpenGIS Simple Features Specification – Principal Methods.....	106
Table 4-3 Spatial overlay methods, Manifold GIS .....	127
Table 4-4 Regional employment data – grouping affects .....	131
Table 4-5 Selected univariate classification schemes.....	136
Table 4-6 Image classification facilities – Selected classifiers .....	142
Table 4-7 Selected MATLab/GRASS planar geometric analysis functions.....	159
Table 4-8 Widely used univariate kernel density functions .....	176
Table 4-9 Interpretation of p-values .....	190
Table 4-10 3x3 Chamfer metrics.....	197
Table 4-11 Linear spatial filters .....	222
Table 5-1 Implications of Data Models.....	226
Table 5-2 Description of methods for analysis of spatial data in ecology.....	227
Table 5-3 Voronoi-based ESDA .....	244
Table 5-4 Sample statistical tools for grid data – Idrisi .....	247
Table 5-5 Simple 2-way contingency table.....	248
Table 5-6 Simple Chi-square frequency table computation .....	251
Table 5-7 NN Statistics and study area size .....	262
Table 5-8 Join count analysis results.....	278
Table 5-9 Join count mean and variance formulas .....	280
Table 5-10 Tabulated lattice data.....	283
Table 5-11 Selected regression analysis terminology.....	298
Table 5-12 Georgia dataset – global regression estimates and diagnostics .....	303
Table 5-13 Georgia dataset – comparative regression estimates and diagnostics .....	307
Table 6-1 Morphometric features – a simplified classification.....	331
Table 6-2 Gridding and interpolation methods .....	351
Table 6-3 Variogram models (univariate, isotropic) .....	379
Table 7-1 Network analysis terminology.....	392
Table 7-2 Some key optimization problems in network analysis .....	399
Table 7-3 Sample network analysis problem parameters.....	400
Table 7-4 Routing functionality in selected logistics software packages.....	402
Table 7-5 Taxonomy of location analysis problems.....	404
Table 8-1 Agent-based modeling and GIS coupling .....	450
Table 8-2 Agents and environments.....	467
Table 8-3 Comparison of open source simulation/modeling toolkits .....	471
Table 8-4 Comparison of shareware/freeware simulation/modeling systems.....	476
Table 8-5 Comparison of proprietary simulation/modeling systems.....	480
Table 8-6 W weights matrix, Chiquitania MLP model.....	497
Table 8-7 SOM neighborhood and learning rate functions .....	505

# Foreword to the Third Edition

---

## Changes from the 2<sup>nd</sup> edition:

This 3<sup>rd</sup> edition includes the following principal changes from earlier editions: the print version is now provided in full color; embedded links (applicable to the PDF and Web versions) are highlighted in the body of the printed text in dark blue type; there has been extensive revision and updating of Chapters 3 and 8; removal of tables of software products – these are now all referenced via the accompanying website page which is regularly updated: [www.spatialanalysisonline.com/software.html](http://www.spatialanalysisonline.com/software.html); addition of new sections, including geospatial modeling, Knox and Mantel tests, cartogram production, and space syntax; extended discussion of geovisualization techniques in many sections; expansion and revision of the sections covering distance transforms (DTs) and join count statistics (JCS); inclusion of many new diagrams and examples (e.g. fire, pollution, wind and noise modeling) and re-implementation of many existing diagrams; addition of links to downloadable PDF files of the out-of-print **CATMOG** publication series – although rather ‘historic’ those referenced include useful discussions of many of the topics mentioned in this Guide; provision of Google mapping links for many placenames cited; and updates to many sections to reflect recent technological advances, particularly in software tools. These include the addition of information on analytical tools not previously covered such as **PySal**, **PASSaGE**, **SAM**, and **R Spatial** (see further the selected summary “R-Project spatial statistics software packages” at the end of this Guide).

**Geospatial Analysis: A Comprehensive Guide to Principles, Techniques and Software Tools** originated as a document to accompany the spatial analysis module of the MSc in Geographic Information Science at University College London delivered by the principal author, Dr Mike de Smith. As is often the case, from its conception through to completion of the first draft it developed a life of its own, growing into a substantial Guide designed for use by a wide audience. Once several of the chapters had been written – notably those covering the building blocks of spatial analysis and on surface analysis – the project was discussed with Professors Longley and Goodchild. They kindly agreed to contribute to the contents of the Guide itself. As such, this Guide may be seen as a companion to the pioneering book on *Geographic Information Systems and Science* by Longley, Goodchild, Maguire and Rhind, particularly the chapters of that work which deal with spatial analysis and modeling. Their participation has also facilitated links with broader “spatial literacy” and spatial analysis programmes. Notable amongst these are the *GIS&T Body of Knowledge* materials provided by the Association of American Geographers at [www.aag.org/bok/](http://www.aag.org/bok/) together with the spatial educational programmes provided at [www.spatial-literacy.org](http://www.spatial-literacy.org), [www.spatial.ucsb.edu](http://www.spatial.ucsb.edu), [www.ncgia.ucsb.edu](http://www.ncgia.ucsb.edu) and [www.csiss.org](http://www.csiss.org).

The three formats in which this Guide has been published: Printed, Web and E-book (PDF) versions have proved to be extremely popular, encouraging us to seek to improve and extend the material and associated resources further. Many academics and industry professionals have provided helpful comments on previous editions, and universities in several parts of the world have now developed courses which make use of the Guide and the accompanying resources. Workshops based on these materials have been run in Ireland, the USA, East Africa, Italy and Japan, and a Chinese version of the Guide has been published by the Publishing House of Electronics Industry, Beijing, PRC, [www.phei.com.cn](http://www.phei.com.cn) in 2009.

A unique, ongoing, feature of this Guide is its independent evaluation of software, in particular the set of readily available tools and packages for conducting various forms of geospatial analysis. To our knowledge, there is no similarly extensive resource that is available in printed or electronic form. We remain convinced that there is a need for guidance on where to find and how to apply

selected tools. Inevitably, some topics have been omitted, primarily where there is little or no readily available commercial or open source software to support particular analytical operations. Other topics, whilst included, have been covered relatively briefly and/or with limited examples, reflecting the inevitable constraints of time and the authors' limited access to some of the available software resources.

Every effort has been made to ensure the information provided is up-to-date, accurate, compact, comprehensive and representative – we do not claim it to be exhaustive. However, with fast-moving changes in the software industry and in the development of new techniques it would be impractical and uneconomic to publish the material in a conventional manner. Accordingly the Guide has been prepared without intermediary typesetting. This has enabled the time between producing the text and delivery in electronic (web, e-book) and printed formats to be greatly reduced, thereby ensuring that the work is as current as possible. It also enables the work to be updated on a regular basis, with embedded hyperlinks to external resources and suppliers (highlighted and activated in the Web and PDF versions), thus making the Guide a more dynamic and extensive resource than would otherwise be possible. This approach does come with some minor disadvantages. These include: the need to provide rather more subsections to chapters and keywording of terms than would normally be the case in order to support topic selection within the web-based version; and the need for careful use of symbology and embedded graphic symbols at various points within the text to ensure that the web-based output correctly displays Greek letters and other symbols across a range of web browsers.

As with the previous editions, comments and suggestions regarding the scope, detailed content and associated materials (e.g. case studies) are welcome and amendments will be made available via the Guide web site, [www.spatialanalysisonline.com](http://www.spatialanalysisonline.com). We would like to thank all those users of the web site, electronic version of the Guide and of the printed book, for their comments and suggestions which have assisted us in producing this third edition.

Mike de Smith, Edinburgh ♦ Mike Goodchild, Santa Barbara ♦ Paul Longley, London

June 2009 (3<sup>rd</sup> Edition)

## Acknowledgements

The authors would like to express their particular thanks to the following individuals and organizations: [Accon GmbH](#), Greifenberg, Germany for permission to use the noise mapping images on the inside cover of this Guide and in Figure 3-4; Prof D Martin for permission to use Figure 4-19 and Figure 4-20; Prof D Dorling and colleagues for permission to use Figure 4-50 and Figure 4-52; Dr K McGarigal for permission to use the [Fragstats](#) summary in Section 5.3.4; Dr H Kristinsson, Faculty of Engineering, University of Iceland for permission to use Figure 4-69; Dr S Rana, formerly of the Center for Transport Studies, University College London for permission to use Figure 6-24; Prof B Jiang, Department of Technology and Built Environment of [University of Gävle](#), Sweden for permission to use the Axwoman software and sample data in Section 6.3.3.2; Dr G Dubois, European Commission (EC), Joint Research Center Directorate (DG JRC) for comments on parts of Chapter 6 and permission to use material from the original [AI-Geostats](#) website; [Geovariances \(France\)](#) for provision of an evaluation copy of their [Isatis](#) geostatistical software; F O'Sullivan for use of Figure 6-41; Profs A Okabe, K Okunuki and S Shiode ([Center for Spatial Information Science](#), Tokyo University, Japan) for use of their [SANET](#) software and sample data; and S A Sirigos, University of Thessaly, Greece for permission to use his Tripolis dataset in the Figure at the front of this Guide, the provision of his [S-Distance](#) software, and comments on part of Chapter 7. Sections 8.1 and 8.2 of Chapter 8 are substantially derived from material researched and written by Christian Castle and Andrew Crooks (and updated for the latest edition by Andrew) with the financial support of the Economic and Social Research Council (ESRC), Camden Primary Care Trust (PCT), and the Greater London Authority (GLA) Economics Unit. The front cover has been designed by Dr Alex Singleton.

A number of the maps displayed in this Guide, notably those in Chapter 6, have been created using GB Ordnance Survey data provided via the [EDINA](#) Digimap/JISC service. These datasets and other GB OS data illustrated is © Crown Copyright. Every effort has been made to acknowledge and establish copyright of materials used in this publication. Anyone with a query regarding any such item should contact the authors via the Guide's website, [www.spatialanalysisonline.com](http://www.spatialanalysisonline.com)



# 1 Introduction and terminology

## 1.1 Motivation and Media

Our objective in producing this Guide is to be comprehensive in terms of concepts and techniques (but not necessarily exhaustive), representative and independent in terms of software tools, and above all practical in terms of application and implementation. However, we believe that it is no longer appropriate to think of a standard, discipline-specific textbook as capable of satisfying every kind of new user need. Accordingly, an innovative feature of our approach here is the range of formats and channels through which we disseminate the material.

### 1.1.1 Guide overview

In this Guide we address the full spectrum of spatial analysis and associated modeling techniques that are provided within currently available and widely used geographic information systems (GIS) and associated software. Collectively such techniques and tools are often now described as *geospatial analysis*, although we use the more common form, *spatial analysis*, in most of our discussions.

The term ‘GIS’ is widely attributed to Roger Tomlinson and colleagues, who used it in 1963 to describe their activities in building a digital natural resource inventory system for Canada (Tomlinson 1967, 1970). The history of the field has been charted in an edited volume by Foresman (1998) containing contributions by many of its early protagonists. A timeline of many of the formative influences upon the field up to the year 2000 is available via: [www.casa.ucl.ac.uk/gistimeline/](http://www.casa.ucl.ac.uk/gistimeline/); a printed summary is provided by Longley *et al.* (2005, and Chapter 1 of the forthcoming 3<sup>rd</sup> edition), and useful background information may be found at the GIS History Project web site based at the NCGIA (Buffalo): <http://www.ncgia.buffalo.edu/gishist/>. Each of these sources makes the unassailable point

that the success of GIS as an area of activity has fundamentally been driven by the success of its applications in solving real world problems. Many applications are illustrated in Longley *et al.* (Chapter 2 of the forthcoming 3<sup>rd</sup> edition, “A gallery of applications”). In a similar vein the web site for this Guide provides companion material focusing on applications. Amongst these are a series of sector-specific *case studies* drawing on recent work in and around London (UK), together with a number of international case studies.

In order to cover such a wide range of topics, this Guide has been divided into a number of main sections or chapters. These are then further subdivided, in part to identify distinct topics as closely as possible, facilitating the creation of a web site from the text of the Guide. Hyperlinks embedded within the document enable users of the web and PDF versions of this document to navigate around the Guide and to external sources of information, data, software, maps, and reading materials.

Chapter 2 provides an introduction to spatial thinking, recently described by some as “spatial literacy”, and addresses the central issues and problems associated with spatial data that need to be considered in any analytical exercise. In practice, real-world applications are likely to be governed by the organizational *practices* and *procedures* that prevail with respect to particular *places*. Not only are there wide differences in the volume and remit of data that the public sector collects about population characteristics in different parts of the world, but there are differences in the ways in which data are collected, assembled and disseminated (e.g. general purpose censuses versus statistical modeling of social surveys, property registers and tax payments). There are also differences in the ways in which different data holdings can legally be merged and the purposes for which data may be used – particularly with regard to health and law enforcement data. Finally, there are geographical differences in the cost of geographically referenced data. Some organizations, such as the US Geological Survey, are bound by statute to limit charges

for data to sundry costs such as media used for delivering data while others, such as most national mapping organizations in Europe, are required to exact much heavier charges in order to recoup much or all of the cost of data creation. Analysts may already be aware of these contextual considerations through local knowledge, and other considerations may become apparent through browsing metadata catalogues. GIS applications must by definition be sensitive to context, since they represent unique locations on the Earth's surface.

This initial discussion is followed in Chapter 3 by an examination of the methodological background to GIS analysis. Initially we examine a number of formal methodologies and then apply ideas drawn from these to the specific case of spatial analysis. A process known by its initials, PPDAC (Problem, Plan, Data, Analysis, Conclusions) is described as a methodological framework that may be applied to a very wide range of spatial analysis problems and projects. We conclude Chapter 3 with a discussion on model-building, with particular reference to the various types of model that can be constructed to address geospatial problems.

Subsequent Chapters present the various analytical methods supported within widely available software tools. The majority of the methods described in Chapter 4 and many of those in Chapter 6 are implemented as standard facilities in modern commercial GIS packages such as [ArcGIS](#), [MapInfo](#), [Manifold](#), [TNTMips](#) and [Geomedia](#). Many are also provided in more specialized GIS products such as [Idrisi](#), [GRASS](#), [Terraseer](#) and [ENVI](#). In addition we discuss a number of more specialized tools, designed to address the needs of specific sectors or technical problems that are otherwise not well-supported within the core GIS packages at present. Chapter 5, which focuses on statistical methods, and Chapters 7 and 8 which address Network and Location Analysis, and Geocomputation, are much less commonly supported in GIS packages, but may provide loose- or close-coupling with such systems, depending upon the application area. In all instances we

provide detailed examples and commentary on software tools that are readily available.

As noted above, throughout this Guide examples are drawn from and refer to specific products – these have been selected purely as examples and are not intended as recommendations. Extensive use has also been made of tabulated information, providing abbreviated summaries of techniques and formulas for reasons of both compactness and coverage. These tables are designed to provide a quick reference to the various topics covered and are, therefore, not intended as a substitute for fuller details on the various items covered. We provide limited discussion of novel 2D and 3D mapping facilities, and the support for digital globe formats (e.g. [KML](#) and [KMZ](#)), which is increasingly being embedded into general-purpose and specialized data analysis toolsets. These developments confirm the trend towards integration of geospatial data and presentation layers into mainstream software systems and services, both terrestrial and planetary (see, for example, the [KML](#) images of Mars DEMs at the end of this Guide).

Just as all datasets and software packages contain errors, known and unknown, so too do all books and websites, and the authors of this Guide expect that there will be errors despite our best efforts to remove these! Some may be genuine errors or misprints, whilst others may reflect our use of specific versions of software packages and their documentation. Inevitably with respect to the latter, new versions of the packages that we have used to illustrate this Guide will have appeared even before publication, so specific examples, illustrations and comments on scope or restrictions may have been superseded. In all cases the user should review the documentation provided with the software version they plan to use, check release notes for changes and known bugs, and look at any relevant online services (e.g. user/developer forums and blogs on the web) for additional materials and insights.

The interactive web version of this Guide may be accessed via the associated Internet site: [www.spatialanalysisonline.com](http://www.spatialanalysisonline.com). The contents and sample sections of the PDF version may

also be accessed from this site. In both cases the information is regularly updated. The Internet is now well established as society's principal mode of information exchange, and most aspiring GIS users are accustomed to searching for material that can easily be customized to specific needs. Our objective for such users is to provide an independent, reliable and authoritative first port of call for conceptual, technical, software and applications material that addresses the panoply of new user requirements.

Readers wishing to obtain a more in-depth understanding of the background to many of the topics covered in this Guide should review the Suggested Reading topic (Section 1.2.2). Those seeking examples of software tools that might be used for geospatial analysis should refer to Section 1.1.2.

### 1.1.2 Spatial analysis, GIS and software tools

Given the vast range of spatial analysis techniques that have been developed over the past half century many topics can only be covered to a limited depth, whilst others have been omitted because they are not implemented in current mainstream GIS products. This is a rapidly changing field and increasingly GIS packages are including analytical tools as standard built-in facilities or as optional *toolsets*, *add-ins* or *analysts*. In many instances such facilities are provided by the original software suppliers (commercial vendors or collaborative non-commercial development teams) whilst in other cases facilities have been developed and are provided by third parties. Many products offer software development kits (SDKs), programming languages and language support, scripting facilities and/or special interfaces for developing one's own analytical tools or variants.

In addition, a wide variety of web-based or web-deployed tools have become available, enabling datasets to be analyzed and mapped, including dynamic interaction and drill-down capabilities, without the need for local GIS software installation. These tools include the widespread use of [Java](#) applets, [Flash](#)-based

mapping, [AJAX](#) and [Web 2.0](#) applications, and interactive [Virtual Globe](#) explorers, some of which are described in this Guide. They provide an illustration of the direction that many toolset and service providers are taking.

Throughout this Guide there are numerous examples of the use of software tools that facilitate geospatial analysis. In addition, some subsections of the Guide and the software section of the accompanying website, provide summary information about such tools and links to their suppliers. Commercial software products rarely provide access to source code or full details of the algorithms employed. Typically they provide references to books and articles on which procedures are based, coupled with online help and "white papers" describing their parameters and applications. This means that results produced using one package on a given dataset can rarely be exactly matched to those produced using any other package or through hand-crafted coding. There are many reasons for these inconsistencies including: differences in the software architectures of the various packages and the algorithms used to implement individual methods; errors in the source materials or their interpretation; coding errors; inconsistencies arising out of the ways in which different GIS packages model, store and manipulate information; and differing treatments of special cases (e.g. missing values, boundaries, adjacency, obstacles, distance computations etc.).

Non-commercial packages sometimes provide source code and test data for some or all of the analytical functions provided, although it is important to understand that "non-commercial" often does not mean that users can download the full source code. Source code greatly aids understanding, reproducibility and further development. Such software will often also provide details of known bugs and restrictions associated with functions – although this information may also be provided with commercial products it is generally less transparent. In this respect non-commercial software may meet the requirements of scientific rigor more fully than many commercial offerings, but is often

provided with limited documentation, training tools, cross-platform testing and/or technical support, and thus is generally more demanding on the users and system administrators. In many instances open source and similar not-for-profit GIS software may also be less generic, focusing on a particular form of spatial representation (e.g. a grid or raster spatial model). Like some commercial software, it may also be designed with particular application areas in mind, such as addressing problems in hydrology or epidemiology.

The process of selecting software tools encourages us to ask: (i) “what is meant by geospatial analysis techniques?” and (ii) “what should we consider to be GIS software?” To some extent the answer to the second question is the simpler, if we are prepared to be guided by self-selection. For our purposes we focus principally on products that claim to provide geographic information systems capabilities, supporting at least 2D mapping (display and output) of raster (grid based) and/or vector (point/line/polygon based) data, with a minimum of basic map manipulation facilities. We concentrate our review on a number of the products most widely used or with the most readily accessible analytical facilities. This leads us beyond the realm of pure GIS. For example: we use examples drawn from packages that do not directly provide mapping facilities (e.g. [Crimestat](#)) but which provide input and/or output in widely used GIS map-able formats; products that include some mapping facilities but whose primary purpose is spatial or spatio-temporal data exploration and analysis (e.g. [GS+](#), [STIS](#), [GeoDa](#), [PySal](#)); and products that are general- or special-purpose analytical engines incorporating mapping capabilities (e.g. [MATLab](#) with the Mapping Toolbox, [WinBUGS](#) with [GeoBUGS](#)) – for more details on these and other example software tools, please see the website page:

[www.spatialanalysisonline.com/software.html](http://www.spatialanalysisonline.com/software.html)

The more difficult of the two questions above is the first – what should be considered as “geospatial analysis”? In conceptual terms, the phrase identifies the subset of techniques that

are applicable when, as a minimum, data can be referenced on a two-dimensional frame and relate to terrestrial activities. The results of geospatial analysis will change if the location or extent of the frame changes, or if objects are repositioned within it: if they do not, then “everywhere is nowhere”, location is unimportant, and it is simpler and more appropriate to use conventional, *aspatial*, techniques.

Many GIS products apply the term (geo)spatial analysis in a very narrow context. In the case of vector-based GIS this typically means operations such as: map overlay (combining two or more maps or map layers according to predefined rules); simple buffering (identifying regions of a map within a specified distance of one or more features, such as towns, roads or rivers); and similar basic operations. This reflects (and is reflected in) the use of the term *spatial analysis* within the Open Geospatial Consortium (OGC) “simple feature specifications” (see further Table 4-2). For raster-based GIS, widely used in the environmental sciences and [remote sensing](#), this typically means a range of actions applied to the grid cells of one or more maps (or images) often involving filtering and/or algebraic operations (*map algebra*). These techniques involve processing one or more raster layers according to simple rules resulting in a new map layer, for example replacing each cell value with some combination of its neighbors’ values, or computing the sum or difference of specific attribute values for each grid cell in two matching raster datasets. Descriptive statistics, such as cell counts, means, variances, maxima, minima, cumulative values, frequencies and a number of other measures and distance computations are also often included in this generic term “spatial analysis”.

However, at this point only the most basic of facilities have been included, albeit those that may be the most frequently used by the greatest number of GIS professionals. To this initial set must be added a large variety of statistical techniques (descriptive, exploratory, explanatory and predictive) that

have been designed specifically for spatial and spatio-temporal data. Today such techniques are of great importance in social and political sciences, despite the fact that their origins may often be traced back to problems in the environmental and life sciences, in particular ecology, geology and epidemiology. It is also to be noted that spatial statistics is largely an observational science (like astronomy) rather than an experimental science (like agronomy or pharmaceutical research). This aspect of geospatial science has important implications for analysis, particularly the application of a range of statistical methods to spatial problems.

Limiting the definition of geospatial analysis to 2D mapping operations and spatial statistics remains too restrictive for our purposes. There are other very important areas to be considered. These include: surface analysis – in particular analyzing the properties of physical surfaces, such as gradient, aspect and visibility, and analyzing surface-like data “fields”; network analysis – examining the properties of natural and man-made networks in order to understand the behavior of flows within and around such networks; and locational analysis. GIS-based network analysis may be used to address a wide range of practical problems such as route selection and facility location, and problems involving flows such as those found in hydrology. In many instances location problems relate to networks and as such are often best addressed with tools designed for this purpose, but in others existing networks may have little or no relevance or may be impractical to incorporate within the modeling process. Problems that are not specifically network constrained, such as new road or pipeline routing, regional warehouse location, mobile phone mast positioning, pedestrian movement or the selection of rural community health care sites, may be effectively analyzed (at least initially) without reference to existing physical networks. Locational analysis “in the plane” is also applicable where suitable network datasets are not available, or are too large or expensive to be utilized, or where the location algorithm is very complex or involves the

examination or simulation of a very large number of alternative configurations.

A further important aspect of geospatial analysis is visualization (or *geovisualization*) – the use, creation and manipulation of images, maps, diagrams, charts, 3D static and dynamic views, high resolution satellite imagery and digital globes, and their associated tabular datasets (see further, [Slocum et al., 2008](#), [Dodge et al., 2008](#), [Longley et al. \(Chapter 13 in the forthcoming 3<sup>rd</sup> edition\)](#) and the work of the [GeoVista](#) project team). For further insights into how some of these developments may be applied, see [Andrew Hudson-Smith \(2008\) “Digital Geography: Geographic visualization for urban environments”](#) and [Martin Dodge and Rob Kitchin’s earlier “Atlas of Cyberspace”](#) which is now available as a free downloadable document.

GIS packages and web-based services increasingly incorporate a range of such tools, providing static or rotating views, draping images over 2.5D surface representations, providing animations and fly-throughs, dynamic linking and brushing and spatio-temporal visualizations. This latter class of tools has been, until recently, the least developed, reflecting in part the limited range of suitable compatible datasets and the limited set of analytical methods available, although this picture is changing rapidly. One recent example is the availability of image time series from NASA’s Earth Observation Satellites, yielding vast quantities of data on a daily basis (e.g. [Aqua mission](#), commenced 2002; [Terra mission](#), commenced 1999).

Geovisualization is the subject of ongoing research by the International Cartographic Association (ICA), [Commission on Geovisualization](#), who have organized a series of workshops and publications addressing developments in geovisualization, notably with a cartographic focus.

As datasets, software tools and processing capabilities develop, 3D geometric and photo-realistic visualization are becoming a *sine qua non* of modern geospatial systems and services – see [Andy Hudson-Smith’s “Digital Urban”](#)

blog for a regularly updated commentary on this field. We expect to see an explosion of tools and services and datasets in this area over the coming years – many examples are included as illustrations in this Guide. Other examples readers may wish to explore include: the static and dynamic visualizations at [3DNature](#) and similar sites; the 2D and 3D [Atlas of Switzerland](#); Urban 3D modeling programmes such as [LandExplorer](#) and [CityGML](#); and the integration of GIS technologies and data with digital globe software, e.g. data from [Digital Globe](#) and [GeoEye](#), and Earth-based frameworks such as [Google Earth](#), [Microsoft Virtual Earth](#), [NASA Worldwind](#) and [Edushi](#) (Chinese). There are also automated translators between GIS packages such as ArcGIS and digital Earth models (see for example [Arc2Earth](#)).

These novel visualization tools and facilities augment the core tools utilized in spatial analysis throughout many parts of the analytical process: exploration of data; identification of patterns and relationships; construction of models; dynamic interaction with models; and communication of results – see, for example, the recent work of the city of [Portland, Oregon](#), who have used 3D visualization to communicate the results of zoning, crime analysis and other key local variables to the public. Another example is the 3D visualizations provided as part of the web-accessible [London Air Quality](#) network (see example at the front of this Guide). These are designed to enable:

- users to visualize air pollution in the areas that they work, live or walk
- transport planners to identify the most polluted parts of London.
- urban planners to see how building density affects pollution concentrations in the City and other high density areas, and
- students to understand pollution sources and dispersion characteristics

Physical 3D models and hybrid physical-digital models are also being developed and applied to practical analysis problems. For example: 3D physical models constructed from plaster,

wood, paper and plastics have been used for many years in architectural and engineering planning projects; hybrid sandtables are being used to help firefighters in California visualize the progress of wildfires (see Figure 1-1A); very large sculptured solid terrain models (e.g. see [STM](#)) are being used for educational purposes, to assist land use modeling programmes, and to facilitate participatory 3D modeling in less-developed communities ([P3DM](#)); and 3D digital printing technology is being used to rapidly generate 3D landscapes and cityscapes from GIS, CAD and/or VRML files with planning, security, architectural, archaeological and geological applications (see Figure 1-1B and the websites of [Z corporation](#), [Landprint](#) and [Dimension Printing](#) for more details. To create large landscape models multiple individual prints, which are typically only around 20cm x 20cm x 5cm, are made, in much the same manner as raster file mosaics.

**Figure 1-1 3D Physical GIS models**

A. Sand-in-a-box model, [Albuquerque, USA](#)



B. 3D GIS printing



GIS software, notably in the commercial sphere, is driven primarily by demand and applicability, as manifest in willingness to pay.

Hence, to an extent, the facilities available often reflect commercial and resourcing realities (including the development of improvements in processing and display hardware, and the ready availability of high quality datasets) rather than the status of development in geospatial science. Indeed, there may be many capabilities available in software packages that are provided simply because it is extremely easy for the designers and programmers to implement them, especially those employing object-oriented programming and data models. For example, a given operation may be provided for polygonal features in response to a well-understood application requirement, which is then easily enabled for other features (e.g. point sets, polylines) despite the fact that there may be no known or likely requirement for the facility.

Despite this cautionary note, for specific well-defined or *core* problems, software developers will frequently utilize the most up-to-date research on algorithms in order to improve the quality (accuracy, optimality) and efficiency (speed, memory usage) of their products. For further information on algorithms and data structures, see the online [NIST Dictionary of algorithms and data structures](#).

Furthermore, the quality, variety and efficiency of spatial analysis facilities provide an important discriminator between commercial offerings in an increasingly competitive and open market for software. However, the ready availability of analysis tools does not imply that one product is necessarily better or more complete than another – it is the selection and application of *appropriate* tools in a manner that is *fit for purpose* that is important. Guidance documents exist in some disciplines that assist users in this process, e.g. [Perry et al. \(2002\)](#) dealing with ecological data analysis, and to a significant degree we hope that this Guide will assist users from many disciplines in the selection process.

### 1.1.3 Intended audience and scope

This Guide has been designed to be accessible to a wide range of readers – from undergraduates and postgraduates studying GIS

and spatial analysis, to GIS practitioners and professional analysts. It is intended to be much more than a cookbook of formulas, algorithms and techniques – its aim is to provide an explanation of the key techniques of spatial analysis using examples from widely available software packages. It stops short, however, of attempting a systematic evaluation of competing software products. A substantial range of application examples are provided, but any specific selection inevitably illustrates only a small subset of the huge range of facilities available. Wherever possible, examples have been drawn from non-academic sources, highlighting the growing understanding and acceptance of GIS technology in the commercial and government sectors.

The scope of this Guide incorporates the various spatial analysis topics included within the [NCGIA Core Curriculum](#) (Goodchild and Kemp, 1990) and as such may provide a useful accompaniment to GIS Analysis courses based closely or loosely on this programme. More recently the Education Committee of the University Consortium for Geographic Information Science ([UCGIS](#)) in conjunction with the Association of American Geographers ([AAG](#)) has produced a comprehensive “Body of Knowledge” (BoK) document, which is available from the [AAG website](#) (<http://www.aag.org/bok/>). This Guide covers materials that primarily relate to the BoK sections **CF**: Conceptual Foundations; **AM**: Analytical Methods and **GC**: Geocomputation. In the general introduction to the **AM** knowledge area the authors of the BoK summarize this component as follows:

“This knowledge area encompasses a wide variety of operations whose objective is to derive analytical results from geospatial data. Data analysis seeks to understand both first-order (environmental) effects and second-order (interaction) effects. Approaches that are both data-driven (exploration of geospatial data) and model-driven (testing hypotheses and creating models) are included. Data-driven techniques derive summary descriptions of data, evoke insights about characteristics of data, contribute to the development of research hypotheses,

and lead to the derivation of analytical results. The goal of model-driven analysis is to create and test geospatial process models. In general, model-driven analysis is an advanced knowledge area where previous experience with exploratory spatial data analysis would constitute a desired prerequisite.” (BoK, p83 of the e-book version).

## 1.2 Software tools and Companion Materials

### 1.2.1 GIS and related software tools

The GIS software and analysis tools that an individual, group or corporate body chooses to use will depend very much on the purposes to which they will be put. There is an enormous difference between the requirements of academic researchers and educators, and those with responsibility for planning and delivery of emergency control systems or large scale physical infrastructure projects. The spectrum of products that may be described as a GIS includes (amongst others):

- highly specialized, sector specific packages: for example civil engineering design and costing systems; satellite image processing systems; and utility infrastructure management systems
- transportation and logistics management systems
- civil and military control room systems
- systems for visualizing the built environment for architectural purposes, for public consultation or as part of simulated environments for interactive gaming
- land registration systems
- census data management systems
- commercial location services and Digital Earth models

The list of software functions and applications is long and in some instances suppliers would not describe their offerings as a GIS. In many cases such systems fulfill specific operational needs, solving a well-defined subset of spatial problems and providing mapped output as an incidental but essential part of their operation. Many of the capabilities may be found in generic GIS products. In other instances a specialized package may utilize a GIS engine for the display and in some cases processing of spatial data (directly, or indirectly through interfacing or file

input/output mechanisms). For this reason, and in order to draw a boundary around the present work, reference to application-specific GIS will be limited.

A number of GIS packages and related toolsets have particularly strong facilities for processing and analyzing binary, grayscale and color images. They may have been designed originally for the processing of remote sensed data from satellite and aerial surveys, but many have developed into much more sophisticated and complete GIS tools, e.g. Clark Lab's [Idrisi](#) software; MicroImage's [TNTMips](#) product set; the [ERDAS](#) suite of products; and [ENVI](#) with associated packages such as [RiverTools](#). Alternatively, image handling may have been deliberately included within the original design parameters for a generic GIS package (e.g. [Manifold](#)), or simply be toolsets for image processing that may be combined with mapping tools (e.g. the [MATLab](#) Image Processing Toolbox). Whatever their origins, a central purpose of such tools has been the capture, manipulation and interpretation of image data, rather than spatial analysis per se, although the latter inevitably follows from the former.

In this Guide we do not provide a separate chapter on image processing, despite its considerable importance in GIS, focusing instead on those areas where image processing tools and concepts are applied for spatial analysis (e.g. surface analysis). We have adopted a similar position with respect to other forms of data capture, such as field and geodetic survey systems and data cleansing software – although these incorporate analytical tools, their primary function remains the recording and georeferencing of datasets, rather than the analysis of such datasets once stored.

For most GIS professionals, spatial analysis and associated modeling is an infrequent activity. Even for those whose job focuses on analysis the range of techniques employed tends to be quite narrow and application focused. GIS consultants, researchers and academics on the other hand are continually exploring and developing analytical techniques. For the first

group and for consultants, especially in commercial environments, the imperatives of financial considerations, timeliness and corporate policy loom large, directing attention to: delivery of solutions within well-defined time and cost parameters; working within commercial constraints on the cost and availability of software, datasets and staffing; ensuring that solutions are fit for purpose/meet client and end-user expectations and agreed standards; and in some cases, meeting “political” expectations.

For the second group of users it is common to make use of a variety of tools, data and programming facilities developed in the academic sphere. Increasingly these make use of non-commercial wide-ranging spatial analysis software libraries, such as the **R-Spatial** project (in “R”); **PySal** (in “Python”); and **Splancs** (in “S”).

#### 1.2.1.1 Sample software products

The principal products we have included in this latest edition of the Guide are included on the accompanying website’s [software page](#). Many of these products are free whilst others are available (at least in some form) for a small fee for all or selected groups of users. Others are licensed at varying per user prices, from a few hundred to over a thousand US dollars per user. Our tests and examples have largely been carried out using desktop/Windows versions of these software products. Different versions that support Unix-based operating systems and more sophisticated back-end database engines have not been utilized. In the context of this Guide we do not believe these selections affect our discussions in any substantial manner, although such issues may have performance and systems architecture implications that are extremely important for many users. OGC compliant software products are listed on the OGC resources web page: <http://www.opengeospatial.org/resource/products/compliant>. To quote from the OGC: “The OGC Compliance Testing Program provides a formal process for testing compliance of products that implement OpenGIS® Standards. Compliance Testing determines that a specific product implementation of a particular OpenGIS®

Standard complies with all mandatory elements as specified in the standard and that these elements operate as described in the standard.”

#### 1.2.1.2 Software performance

Suppliers should be able to provide advice on performance issues (e.g. see the **ESRI** web site, “Services” area for relevant documents relating to their products) and in some cases such information is provided within product Help files (e.g. see the Performance Tips section within the **Manifold GIS** help file). Some analytical tasks are very processor- and memory-hungry, particularly as the number of elements involved increases. For example, vector overlay and buffering is relatively fast with a few objects and layers, but slows appreciably as the number of elements involved increases. This increase is generally at least linear with the number of layers and features, but for some problems grows in a highly non-linear (i.e. geometric) manner. Many optimization tasks, such as optimal routing through networks or trip distribution modeling, are known to be extremely hard or impossible to solve optimally and methods to achieve a best solution with a large dataset can take a considerable time to run (see Section 7.1.4 for a fuller discussion of this topic). Similar problems exist with the processing and display of raster files, especially large images or sets of images. Geocomputational methods, some of which are beginning to appear within GIS packages and related toolsets, are almost by definition computationally intensive. This certainly applies to large-scale (Monte Carlo) simulation models, cellular automata and agent-based models and some raster-based optimization techniques, especially where modeling extends into the time domain.

A frequent criticism of GIS software is that it is over-complicated, resource-hungry and requires specialist expertise to understand and use. Such criticisms are often valid and for many problems it may prove simpler, faster and more transparent to utilize specialized tools for the analytical work and draw on the strengths of GIS in data management and mapping to provide input/output and

visualization functionality. Example approaches include: (i) using high-level programming facilities within a GIS (e.g. macros, scripts, VBA, [Python](#)) - many add-ins are developed in this way; (ii) using wide-ranging programmable spatial analysis software libraries and toolsets that incorporate GIS file reading, writing and display, such as the [R-Spatial](#) and [PySal](#) projects noted earlier; (iii) using general purpose data processing toolsets (e.g. [MATLab](#), Excel, Python's [Matplotlib](#), Numeric Python ([Numpy](#)) and other libraries from [Enthought](#); or (iv) directly utilizing mainstream programming languages (e.g. Java, C++). The advantage of these approaches is control and transparency, the disadvantages are that software development is never trivial, is often subject to frustrating and unforeseen delays and errors, and generally requires ongoing maintenance. In some instances analytical applications may be well-suited to parallel or grid-enabled processing - as for example is the case with [GWR](#) (see [Harris et al., 2006](#)).

At present there are no standardized tests for the quality, speed and accuracy of GIS procedures. It remains the buyer's and user's responsibility and duty to evaluate the software they wish to use for the specific task at hand, and by systematic controlled tests or by other means establish that the product and facility within that product they choose to use is truly fit for purpose – *caveat emptor!* Details of how to obtain these products are provided on the [software page](#) of the website that accompanies this book. The list maintained on [Wikipedia](#) is also a useful source of information and links, although is far from being complete or independent. A number of trade magazines and websites (such as [Geoplace](#) and [Geocommunity](#)) provide ad hoc reviews of GIS software offerings, especially new releases, although coverage of analytical functionality may be limited.

### 1.2.2 Suggested reading

There are numerous excellent modern books on GIS and spatial analysis, although few address software facilities and developments. Hypertext links are provided here, and throughout the text where they are cited, to

the more recent publications and web resources listed.

As a background to this Guide any readers unfamiliar with GIS are encouraged to first tackle “Geographic Information Systems and Science” (GISSc) by [Longley et al. \(2005\)](#) - this work is soon to be updated in a 3<sup>rd</sup> edition. GISSc seeks to provide a comprehensive and highly accessible introduction to the subject as a whole. The GB Ordnance Survey's “[GIS Files](#)” document, downloadable from their website also provides an excellent brief introduction to GIS and its application.

Some of the basic mathematics and statistics of relevance to GIS analysis is covered in [Dale \(2005\)](#) and [Allan \(2004\)](#). For detailed information on datums and map projections, see [Iliffe and Lott \(2008\)](#). A useful online resource for those involved in data analysis, particularly with a statistical content, is the [e-Handbook of Statistical Methods](#) produced by the US National Institute on Standards and Technology, NIST). The more informally produced set of articles on statistical topics provided under the [Wikipedia](#) umbrella are also an extremely useful resource. These works, and the mathematics reference site, [Mathworld](#), are referred to (with hypertext links) at various points throughout this document. For more specific sources on geostatistics and associated software packages, the European Commission's AI-GEOSTATS website ([www.ai-geostats.org](#)) is highly recommended, as is the web site of the [Center for Computational Geostatistics \(CCG\)](#) at the University of Alberta. For those who find mathematics and statistics something of a mystery, [de Smith \(2006\)](#) and [Bluman \(2003\)](#) provide useful starting points. For guidance on how to avoid the many pitfalls of statistical data analysis readers are recommended the material in the classic work by [Huff \(1993\)](#) “How to lie with statistics”, and the 2008 book by [Blastland and Dilnot](#) “The tiger that isn't”.

A relatively new development has been the increasing availability of out-of-print published books, articles and guides as free downloads in PDF format. These include: the series of 59 short guides published under the CATMOG

umbrella (Concepts and Methods in Modern Geography), published between 1975 and 1995, most of which are now available at the [QMRG website](#) (a full list of all the guides is provided at the end of this book); the [AutoCarto archives](#) (1972-1997); the [Atlas of Cyberspace](#) by Dodge and Kitchin; and [Fractal Cities](#), by Batty and Longley.

Undergraduates and MSc programme students will find [Burrough and McDonnell \(1998\)](#) provides excellent coverage of many aspects of geospatial analysis, especially from an environmental sciences perspective. Valuable guidance on the relationship between spatial process and spatial modeling may be found in [Cliff and Ord \(1981\)](#) and [Bailey and Gatrell \(1995\)](#). The latter provides an excellent introduction to the application of statistical methods to spatial data analysis. [O'Sullivan and Unwin \(2003\)](#) is a more broad-ranging book covering the topic the authors describe as "Geographic Information Analysis". This work is best suited to advanced undergraduates and first year postgraduate students. In many respects a deeper and more challenging work is [Haining's \(2003\)](#) "Spatial Data Analysis – Theory and Practice". This book is strongly recommended as a companion to the present Guide for postgraduate researchers and professional analysts involved in using GIS in conjunction with statistical analysis.

However, these authors do not address the broader spectrum of geospatial analysis and associated modeling as we have defined it. For example, problems relating to networks and location are often not covered and the literature relating to this area is scattered across many disciplines, being founded upon the mathematics of graph theory, with applications ranging from electronic circuit design to computer networking and from transport planning to the design of complex molecular structures. Useful books addressing this field include [Miller and Shaw \(2001\)](#) "Geographic Information Systems for Transportation" (especially Chapters 3, 5 and 6), and [Rodrigue et al. \(2006\)](#) "The geography of transport systems" (see further: <http://people.hofstra.edu/geotrans/>).

As companion reading on these topics for the present Guide we suggest the two volumes from the Handbooks in Operations Research and Management Science series by [Ball et al. \(1995\)](#): "Network Models", and "Network Routing". These rather expensive volumes provide collections of reviews covering many classes of network problems, from the core optimization problems of shortest paths and arc routing (e.g. street cleaning), to the complex problems of dynamic routing in variable networks, and a great deal more besides. This is challenging material and many readers may prefer to seek out more approachable material, available in a number of other books and articles, e.g. [Ahuja et al. \(1993\)](#), [Mark Daskin's](#) excellent book "Network and Discrete Location" (1995) and the earlier seminal works by [Haggett and Chorley \(1969\)](#), and [Scott \(1971\)](#), together with the widely available online materials accessible via the Internet. Final recommendations here are [Stephen Wise's](#) excellent [GIS Basics \(2002\)](#) and [Worboys and Duckham \(2004\)](#) which address GIS from a computing perspective. Both these volumes covers many topics, including the central issues of data modeling and data structures, key algorithms, system architectures and interfaces.

Many recent books described as covering (geo)spatial analysis are essentially edited collections of papers or brief articles. As such most do not seek to provide comprehensive coverage of the field, but tend to cover information on recent developments, often with a specific application focus (e.g. health, transport, archaeology). The latter is particularly common where these works are selections from sector- or discipline-specific conference proceedings, whilst in other cases they are carefully chosen or specially written papers. Classic amongst these is [Berry and Marble \(1968\)](#) "Spatial Analysis: A reader in statistical geography". More recent examples include "GIS, Spatial Analysis and Modeling" edited by [Maguire, Batty and Goodchild \(2005\)](#), and the excellent (but costly) compendium work "[The SAGE handbook of Spatial Analysis](#)" edited by [Fotheringham and Rogerson \(2008\)](#).

A second category of companion materials to the present work is the extensive product-specific documentation available from software suppliers. Some of the online help files and product manuals are excellent, as are associated example data files, tutorials, worked examples and white papers (see for example, ESRI's GIS site: <http://www.gis.com/> which provides a wide-ranging guide to GIS. In many instances we utilize these to illustrate the capabilities of specific pieces of software and to enable readers to replicate our results using readily available materials. In addition some suppliers, notably ESRI, have a substantial publishing operation, including more general (i.e. not product specific) books of relevance to the present work. Amongst their publications we strongly recommend the "ESRI Guide to GIS Analysis Volume 1: Geographic patterns and relationships" (1999) by Andy Mitchell, which is full of valuable tips and examples. This is a basic introduction to GIS Analysis, which he defines in this context as "a process for looking at geographic patterns and relationships between features". Mitchell's [Volume 2 \(July 2005\)](#) covers more advanced techniques of data analysis, notably some of the more accessible and widely supported methods of spatial statistics, and is equally highly recommended. A number of the topics covered in his Volume 2 also appear in this Guide. Those considering using Open Source software should investigate the recent books by [Neteler and Mitasova \(2008\)](#), [Tyler Mitchell \(2005\)](#) and [Sherman \(2008\)](#).

In parallel with the increasing range and sophistication of spatial analysis facilities to be found within GIS packages, there has been a major change in spatial analytical techniques. In large measure this has come about as a result of technological developments and the related availability of software tools and detailed publicly available datasets. One aspect of this has been noted already – the move towards network-based location modeling where in the past this would have been unfeasible. More general shifts can be seen in the move towards local rather than simply global analysis, for example in the field of exploratory data analysis; in the increasing use of advanced forms of visualization as an

aid to analysis and communication; and in the development of a wide range of computationally intensive and simulation methods that address problems through micro-scale processes (geocomputational methods). These trends are addressed at many points throughout this Guide.

## 1.3 Terminology and Abbreviations

GIS, like all disciplines, utilizes a wide range of terms and abbreviations, many of which have well-understood and recognized meanings. For a large number of commonly used terms online dictionaries have been developed, for example: those created by the Association for Geographic Information (AGI); the Open Geospatial Consortium (OGC); and by various software suppliers. The latter includes many terms and definitions that are particular to specific products, but remain a valuable resource. The University of California maintains an online dictionary of abbreviations and acronyms used in GIS, cartography and

remote sensing. Web site details for each of these are provided at the end of this Guide.

### 1.3.1 Definitions

Geospatial analysis utilizes many of these terms, but many others are drawn from disciplines such as mathematics and statistics. The result that the same terms may mean entirely different things depending on their context and in many cases, on the software provider utilizing them. In most instances terms used in this Guide are defined on the first occasion they are used, but a number warrant defining at this stage. Table 1-1 provides a selection of such terms, utilizing definitions from widely recognized sources where available and appropriate.

Table 1-1 Selected terminology

Term	Definition
Adjacency	The sharing of a common side or boundary by two or more polygons (AGI). Note that adjacency may also apply to features that lie either side of a common boundary where these features are not necessarily polygons
Arc	Commonly used to refer to a straight line segment connecting two nodes or vertices of a polyline or polygon. Arcs may include segments or circles, spline functions or other forms of smooth curve. In connection with graphs and networks, arcs may be directed or undirected, and may have other attributes (e.g. cost, capacity etc.)
Artifact	A result (observation or set of observations) that appears to show something unusual (e.g. a spike in the surface of a 3D plot) but which is of no significance. Artifacts may be generated by the way in which data have been collected, defined or re-computed (e.g. resolution changing), or as a result of a computational operation (e.g. rounding error or substantive software error). Linear artifacts are sometimes referred to as “ghost lines”
Aspect	The direction in which slope is maximized for a selected point on a surface (see also, Gradient and Slope)
Attribute	A data item associated with an individual object (record) in a spatial database. Attributes may be explicit, in which case they are typically stored as one or more fields in tables linked to a set of objects, or they may be implicit (sometimes referred to as <i>intrinsic</i> ), being either stored but hidden or computed as and when required (e.g. polyline length, polygon centroid). Raster/grid datasets typically have a single explicit attribute (a value) associated with each cell, rather than an attribute table containing as many records as there are cells in the grid
Azimuth	The horizontal direction of a vector, measured clockwise in degrees of rotation from the positive Y-axis, for example, degrees on a compass (AGI)
Azimuthal Projection	A type of map projection constructed as if a plane were to be placed at a tangent to the Earth's surface and the area to be mapped were projected onto the plane.

Term	Definition
	All points on this projection keep their true compass bearing (AGI)
(Spatial) Autocorrelation	The degree of relationship that exists between two or more (spatial) variables, such that when one changes, the other(s) also change. This change can either be in the same direction, which is a positive autocorrelation, or in the opposite direction, which is a negative autocorrelation (AGI). The term autocorrelation is usually applied to ordered datasets, such as those relating to time series or spatial data ordered by distance band. The existence of such a relationship suggests but does not definitely establish causality
Cartogram	A cartogram is a form of map in which some variable such as Population Size or Gross National Product typically is substituted for land area. The geometry or space of the map is distorted in order to convey the information of this alternate variable. Cartograms use a variety of approaches to map distortion, including the use of continuous and discrete regions. The term cartogram (or <i>linear cartogram</i> ) is also used on occasion to refer to maps that distort distance for particular display purposes, such as the London Underground map
Choropleth	A thematic map [i.e. a map showing a theme, such as soil types or rainfall levels] portraying properties of a surface using area symbols such as shading [or color]. Area symbols on a choropleth map usually represent categorized classes of the mapped phenomenon (AGI)
Conflation	A term used to describe the process of combining (merging) information from two data sources into a single source, reconciling disparities where possible (e.g. by rubber-sheeting – see below). The term is distinct from <i>concatenation</i> which refers to combinations of data sources (e.g. by overlaying one upon another) but retaining access to their distinct components
Contiguity	The topological identification of adjacent polygons by recording the left and right polygons of each arc. Contiguity is not concerned with the exact locations of polygons, only their relative positions. Contiguity data can be stored in a table, matrix or simply as [i.e. in] a list, that can be cross-referenced to the relevant co-ordinate data if required (AGI).
Curve	A one-dimensional geometric object stored as a sequence of points, with the subtype of curve specifying the form of interpolation between points. A curve is simple if it does not pass through the same point twice (OGC). A LineString (or polyline – see below) is a subtype of a curve
Datum	Strictly speaking, the singular of <i>data</i> . In GIS the word datum usually relates to a reference level (surface) applying on a nationally or internationally defined basis from which elevation is to be calculated. In the context of terrestrial geodesy datum is usually defined by a model of the Earth or section of the Earth, such as WGS84 (see below). The term is also used for horizontal referencing of measurements; see <a href="#">Iliffe and Lott (2008)</a> for full details
DEM	Digital elevation model (a DEM is a particular kind of DTM, see below)
DTM	Digital terrain model
EDM	Electronic distance measurement
EDA, ESDA	Exploratory data analysis/Exploratory spatial data analysis
Ellipsoid/Spheroid	An ellipse rotated about its minor axis determines a spheroid (sphere-like object),

## TERMINOLOGY AND ABBREVIATIONS

Term	Definition
	also known as an ellipsoid of revolution (see also, WGS84)
Feature	Frequently used within GIS referring to point, line (including polyline and mathematical functions defining arcs), polygon and sometimes text (annotation) objects (see also, vector)
Geoid	An imaginary shape for the Earth defined by mean sea level and its imagined continuation under the continents at the same level of gravitational potential (AGI)
Geodemographics	The analysis of people by where they live, in particular by type of neighborhood. Such localized classifications have been shown to be powerful discriminators of consumer behavior and related social and behavioral patterns
Geospatial	Referring to location relative to the Earth's surface. "Geospatial" is more precise in many GI contexts than "geographic," because geospatial information is often used in ways that do not involve a graphic representation, or map, of the information. OGC
Geostatistics	Statistical methods developed for and applied to geographic data. These statistical methods are required because geographic data do not usually conform to the requirements of standard statistical procedures, due to spatial autocorrelation and other problems associated with spatial data (AGI). The term is widely used to refer to a family of tools used in connection with spatial interpolation (prediction) of (piecewise) continuous datasets and is widely applied in the environmental sciences. Spatial statistics is a term more commonly applied to the analysis of discrete objects (e.g. points, areas) and is particularly associated with the social and health sciences
Geovisualization	A family of techniques that provide visualizations of spatial and spatio-temporal datasets, extending from static, 2D maps and cartograms, to representations of 3D using perspective and shading, solid terrain modeling and increasingly extending into dynamic visualization interfaces such as linked windows, digital globes, fly-throughs, animations, virtual reality and immersive systems. Geovisualization is the subject of ongoing research by the International Cartographic Association (ICA), <a href="#">Commission on Geovisualization</a>
GIS-T	GIS applied to transportation problems
GPS/ DGPS	Global positioning system; Differential global positioning system – DGPS provides improved accuracy over standard GPS by the use of one or more fixed reference stations that provide corrections to GPS data
Gradient	Used in spatial analysis with reference to surfaces (scalar fields). Gradient is a vector field comprised of the <i>aspect</i> (direction of maximum slope) and <i>slope</i> computed in this direction (magnitude of rise over run) at each point of the surface. The magnitude of the gradient (the slope or inclination) is sometimes itself referred to as the gradient (see also, Slope and Aspect)
Graph	A collection of vertices and edges (links between vertices) constitutes a graph. The mathematical study of the properties of graphs and paths through graphs is known as graph theory
Heuristic	A term derived from the same Greek root as Eureka, heuristic refers to procedures for finding solutions to problems that may be difficult or impossible to solve by direct means. In the context of optimization heuristic algorithms are systematic procedures that seek a good or near optimal solution to a well-defined

Term	Definition
	problem, but not one that is necessarily optimal. They are often based on some form of intelligent trial and error or search procedure
<i>iid</i>	An abbreviation for “independently and identically distributed”. Used in statistical analysis in connection with the distribution of errors or residuals
Invariance	In the context of GIS invariance refers to properties of features that remain unchanged under one or more (spatial) transformations
Kernel	Literally, the core or central part of an item. Often used in computer science to refer to the central part of an operating system, the term kernel in geospatial analysis refers to methods (e.g. density modeling, local grid analysis) that involve calculations using a well-defined local neighborhood (block of cells, radially symmetric function)
Layer	A collection of geographic entities of the same type (e.g. points, lines or polygons). Grouped layers may combine layers of different geometric types
Map algebra	A range of actions applied to the grid cells of one or more maps (or images) often involving filtering and/or algebraic operations. These techniques involve processing one or more raster layers according to simple rules resulting in a new map layer, for example replacing each cell value with some combination of its neighbors’ values, or computing the sum or difference of specific attribute values for each grid cell in two matching raster datasets
Mashup	A recently coined term used to describe websites whose content is composed from multiple (often distinct) data sources, such as a mapping service and property price information, constructed using programmable interfaces to these sources (as opposed to simple compositing or embedding)
MBR/ MER	Minimum bounding rectangle/Minimum enclosing (or envelope) rectangle (of a feature set)
Planar/non-planar/planar enforced	Literally, lying entirely within a plane surface. A polygon set is said to be planar enforced if every point in the set lies in exactly one polygon, or on the boundary between two or more polygons. See also, planar graph. A graph or network with edges crossing (e.g. bridges/underpasses) is non-planar
Planar graph	If a graph can be drawn in the plane (embedded) in such a way as to ensure edges only intersect at points that are vertices then the graph is described as planar
Pixel/image	Picture element – a single defined point of an image. Pixels have a “color” attribute whose value will depend on the encoding method used. They are typically either binary (0/1 values), grayscale (effectively a color mapping with values, typically in the integer range [0,255]), or color with values from 0 upwards depending on the number of colors supported. Image files can be regarded as a particular form of raster or grid file
Polygon	A closed figure in the plane, typically comprised of an ordered set of connected vertices, $v_1, v_2, \dots, v_{n-1}, v_n = v_1$ where the connections (edges) are provided by straight line segments. If the sequence of edges is not self-crossing it is called a simple polygon. A point is inside a simple polygon if traversing the boundary in a clockwise direction the point is always on the right of the observer. If every pair of points inside a polygon can be joined by a straight line that also lies inside the polygon then the polygon is described as being convex (i.e. the interior is a connected point set). The OGC definition of a polygon is “a planar surface defined by 1 exterior boundary and 0 or more interior boundaries. Each interior boundary

Term	Definition
	defines a hole in the polygon”
Polyhedral surface	A Polyhedral surface is a contiguous collection of polygons, which share common boundary segments (OGC). See also, Tesseral/Tessellation
Polyline	An ordered set of connected vertices, $v_1, v_2, \dots, v_{n-1}, v_n \neq v_1$ where the connections (edges) are provided by straight line segments. The vertex $v_1$ is referred to as the start of the polyline and $v_n$ as the end of the polyline. The OGC specification uses the term LineString which it defines as: a curve with linear interpolation between points. Each consecutive pair of points defines a line segment
Raster/grid	A data model in which geographic features are represented using discrete cells, generally squares, arranged as a (contiguous) rectangular grid. A single grid is essentially the same as a two-dimensional matrix, but is typically referenced from the lower left corner rather than the norm for matrices, which are referenced from the upper left. Raster files may have one or more values (attributes or bands) associated with each cell position or pixel
Resampling	<ol style="list-style-type: none"> <li>1. Procedures for (automatically) adjusting one or more raster datasets to ensure that the grid resolutions of all sets match when carrying out combination operations. Resampling is often performed to match the coarsest resolution of a set of input rasters. Increasing resolution rather than decreasing requires an interpolation procedure such as bicubic spline.</li> <li>2. The process of reducing image dataset size by representing a group of pixels with a single pixel. Thus, pixel count is lowered, individual pixel size is increased, and overall image geographic extent is retained. Resampled images are “coarse” and have less information than the images from which they are taken. Conversely, this process can also be executed in the reverse (AGI)</li> <li>3. In a statistical context the term resampling (or re-sampling) is sometimes used to describe the process of selecting a subset of the original data, such that the samples can reasonably be expected to be independent</li> </ol>
Rubber sheeting	A procedure to adjust the co-ordinates all of the data points in a dataset to allow a more accurate match between known locations and a few data points within the dataset. Rubber sheeting ... preserves the interconnectivity or topology, between points and objects through stretching, shrinking or re-orienting their interconnecting lines (AGI). Rubber-sheeting techniques are widely used in the production of Cartograms ( <i>op. cit.</i> )
Slope	The amount of <i>rise</i> of a surface (change in elevation) divided by the distance over which this rise is computed (the <i>run</i> ), along a straight line transect in a specified direction. The run is usually defined as the <i>planar</i> distance, in which case the slope is the $\tan()$ function. Unless the surface is flat the slope at a given point on a surface will (typically) have a maximum value in a particular direction (depending on the surface and the way in which the calculations are carried out). This direction is known as the <i>aspect</i> . The <i>vector</i> consisting of the slope and aspect is the <i>gradient</i> of the surface at that point (see also, Gradient and Aspect)
Spatial econometrics	A subset of econometric methods that is concerned with spatial aspects present in cross-sectional and space-time observations. These methods focus in particular on two forms of so-called spatial effects in econometric models, referred to as spatial dependence and spatial heterogeneity (Anselin, 1988, 2006)
Spheroid	A flattened (oblate) form of a sphere, or ellipse of revolution. The most widely used model of the Earth is that of a spheroid, although the detailed form is

Term	Definition
	slightly different from a true spheroid
SQL/Structured Query Language	Within GIS software SQL extensions known as spatial queries are frequently implemented. These support queries that are based on spatial relationships rather than simply attribute values
Surface	A 2D geometric object. A simple surface consists of a single 'patch' that is associated with one exterior boundary and 0 or more interior boundaries. Simple surfaces in 3D are isomorphic to planar surfaces. Polyhedral surfaces are formed by 'stitching' together simple surfaces along their boundaries (OGC). Surfaces may be regarded as scalar fields, i.e. fields with a single value, e.g. elevation or temperature, at every point
Tesseral/Tessellation	A gridded representation of a plane surface into disjoint polygons. These polygons are normally either square (raster), triangular (TIN – see below), or hexagonal. These models can be built into hierarchical structures, and have a range of algorithms available to navigate through them. A (regular or irregular) 2D tessellation involves the subdivision of a 2-dimensional plane into polygonal tiles (polyhedral blocks) that completely cover a plane (AGI). The term lattice is sometimes used to describe the complete division of the plane into regular or irregular disjoint polygons. More generally the subdivision of the plane may be achieved using arcs that are not necessarily straight lines
TIN	Triangulated irregular network. A form of the tesseral model based on triangles. The vertices of the triangles form irregularly spaced nodes. Unlike the grid, the TIN allows dense information in complex areas, and sparse information in simpler or more homogeneous areas. The TIN dataset includes topological relationships between points and their neighboring triangles. Each sample point has an X,Y coordinate and a surface, or Z-Value. These points are connected by edges to form a set of non-overlapping triangles used to represent the surface. TINs are also called irregular triangular mesh or irregular triangular surface model (AGI)
Topology	The relative location of geographic phenomena independent of their exact position. In digital data, topological relationships such as connectivity, adjacency and relative position are usually expressed as relationships between nodes, links and polygons. For example, the topology of a line includes its from- and to-nodes, and its left and right polygons (AGI). In mathematics, a property is said to be <i>topological</i> if it survives stretching and distorting of space
Transformation 1. Map	Map transformation: A computational process of converting an image or map from one coordinate system to another. Transformation ... typically involves rotation and scaling of grid cells, and thus requires resampling of values (AGI)
Transformation 2. Affine	Affine transformation: When a map is digitized, the X and Y coordinates are initially held in digitizer measurements. To make these X,Y pairs useful they must be converted to a real world coordinate system. The affine transformation is a combination of linear transformations that converts digitizer coordinates into Cartesian coordinates. The basic property of an affine transformation is that parallel lines remain parallel (AGI, with modifications). The principal affine transformations are contraction, expansion, dilation, reflection, rotation, shear and translation

## TERMINOLOGY AND ABBREVIATIONS

Term	Definition
Transformation 3. Data	Data transformation (see also, subsection 6.7.1.10): A mathematical procedure (usually a one-to-one mapping or function) applied to an initial dataset to produce a result dataset. An example might be the transformation of a set of sampled values $\{x_i\}$ using the $\log()$ function, to create the set $\{\log(x_i)\}$ . Affine and map transformations are examples of mathematical transformations applied to coordinate datasets. Note that operations on transformed data, e.g. checking whether a value is within 10% of a target value, is not equivalent to the same operation on untransformed data, even after back transformation
Transformation 4. Back	Back transformation: If a set of sampled values $\{x_i\}$ has been transformed by a one-to-one mapping function $f()$ into the set $\{f(x_i)\}$ , and $f()$ has a one-to-one inverse mapping function $f^{-1}()$ , then the process of computing $f^{-1}\{f(x_i)\}=\{x_i\}$ is known as back transformation. Example $f()=\ln()$ and $f^{-1}=\exp()$
Vector	<p>1. Within GIS the term vector refers to data that are comprised of lines or arcs, defined by beginning and end points, which meet at nodes. The locations of these nodes and the topological structure are usually stored explicitly. Features are defined by their boundaries only and curved lines are represented as a series of connecting arcs. Vector storage involves the storage of explicit topology, which raises overheads, however it only stores those points which define a feature and all space outside these features is “non-existent” (AGI)</p> <p>2. In mathematics the term refers to a directed line, i.e. a line with a defined origin, direction and orientation. The same term is used to refer to a single column or row of a matrix, in which case it is denoted by a bold letter, usually in lower case</p>
Viewshed	Regions of visibility observable from one or more observation points. Typically a viewshed will be defined by the numerical or color coding of a raster image, indicating whether the (target) cell can be seen from (or probably seen from) the (source) observation points. By definition a cell that can be viewed from a specific observation point is inter-visible with that point (each location can see the other). Viewsheds are usually determined for optically defined visibility within a maximum range
WGS84	World Geodetic System, 1984 version. This models the Earth as a spheroid with major axis 6378.137 kms and flattening factor of 1:298.257, i.e. roughly 0.3% flatter at the poles than a perfect sphere. One of a number of such global models

*Note: Where cited, references are drawn from the Association for Geographic Information (AGI), and the Open Geospatial Consortium (OGC). Square bracketed text denotes insertion by the present authors into these definitions. For OGC definitions see: Open Geospatial Consortium Inc (2006) in References section*

## 1.4 Common Measures and Notation

### 1.4.1 Notation

Throughout this Guide a number of terms and associated formulas are used that are common to many analytical procedures. In this section we provide a brief summary of those that fall into this category. Others, that are more specific to a particular field of analysis, are

treated within the section to which they primarily apply. Many of the measures we list will be familiar to readers, since they originate from standard single variable (univariate) statistics. For brevity we provide details of these in tabular form. In order to clarify the expressions used here and elsewhere in the text, we use the notation shown in Table 1-2. *Italics* are used within the text and formulas to denote variables and parameters, as well as selected terms.

**Table 1-2 Notation and symbology**

$[a,b]$	A closed interval of the Real line, for example $[0,1]$ means the set of all values between 0 and 1, including 0 and 1
$(a,b)$	An open interval of the Real line, for example $(0,1)$ means the set of all values between 0 and 1, NOT including 0 and 1. This should not be confused with the notation for coordinate pairs, $(x,y)$ , or its use within bivariate functions such as $f(x,y)$ , or in connection with graph edges (see below) – the meaning should be clear from the context
$(i,j)$	In the context of graph theory, which forms the basis for network analysis, this pairwise notation is often used to define an edge connecting the two vertices $i$ and $j$
$(x,y)$	A (spatial) data pair, usually representing a pair of coordinates in two dimensions. Terrestrial coordinates are typically Cartesian (i.e. in the plane, or <i>planar</i> ) based on a pre-specified projection of the sphere, or Spherical (latitude, longitude). Spherical coordinates are often quoted in positive or negative degrees from the Equator and the Greenwich meridian, so may have the ranges $[-90,+90]$ for latitude (north-south measurement) and $[-180,180]$ for longitude (east-west measurement)
$(x,y,z)$	A (spatial) data triple, usually representing a pair of coordinates in two dimensions, plus a third coordinate (usually height or depth) or an attribute value, such as soil type or household income
$\{x_i\}$	A set of $n$ values $x_1, x_2, x_3, \dots, x_n$ , typically continuous ratio-scaled variables in the range $(-\infty, \infty)$ or $[0, \infty)$ . The values may represent measurements or attributes of distinct objects, or values that represent a collection of objects (for example the population of a census tract)
$\{X_i\}$	An ordered set of $n$ values $X_1, X_2, X_3, \dots, X_n$ , such that $X_i \leq X_{i+1}$ for all $i$
$\mathbf{X}, \mathbf{x}$	The use of bold symbols in expressions indicates matrices (upper case) and vectors (lower case)
$\{f_i\}$	A set of $k$ frequencies ( $k \leq n$ ), derived from a dataset $\{x_i\}$ . If $\{x_i\}$ contains discrete values, some of which occur multiple times, then $\{f_i\}$ represents the number of occurrences or the <i>count</i> of each distinct value. $\{f_i\}$ may also represent the number of occurrences of values that lie in a range or set of ranges, $\{r_i\}$ . If a dataset contains $n$ values, then the sum $\sum f_i = n$ . The set $\{f_i\}$ can also be written $f(x_i)$ . If $\{f_i\}$ is regarded as a set of weights (for example attribute values) associated with the $\{x_i\}$ , it may be written as the set $\{w_i\}$ or $w(x_i)$
$\{p_i\}$	A set of $k$ probabilities ( $k \leq n$ ), estimated from a dataset or theoretically derived. With a finite set of values $\{x_i\}$ , $p_i = f_i/n$ . If $\{x_i\}$ represents a set of $k$ classes or ranges then $p_i$ is the probability of

	finding an occurrence in the $i^{\text{th}}$ class or range, i.e. the proportion of events or values occurring in that class or range. The sum $\sum p_i = 1$ . If a set of frequencies, $\{f_i\}$ , have been standardized by dividing each value $f_i$ by their sum, $\sum f_i$ , then $\{p_i\}$ is equivalent to $\{f_i\}$
$\Sigma$	Summation symbol, e.g. $x_1 + x_2 + x_3 + \dots + x_n$ . If no limits are shown the sum is assumed to apply to all subsequent elements, otherwise upper and/or lower limits for summation are provided
$\Pi$	Product symbol, e.g. $x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n$ . If no limits are shown the product is assumed to apply to all subsequent elements, otherwise upper and/or lower limits for multiplication are provided
$\hat{\phantom{x}}$	Used here in conjunction with Greek symbols (directly above) to indicate a value is an estimate of the true population value. Sometimes referred to as “hat”
$\sim$	Is distributed as, for example $y \sim N(0, 1)$ means the variable $y$ has a distribution that is Normal with a mean of 0 and standard deviation of 1
!	Factorial symbol. $z = x!$ means $z = x(x-1)(x-2)\dots 1$ . $x \geq 0$ . Usually applied to integer values of $x$ . May be defined for fractional values of $x$ using the Gamma function (Table 1-3, Section 1.4.2.8)
$\equiv$	‘Equivalent to’ symbol
$\approx$	‘Approximately equal to’ symbol
$\in$	‘Belongs to’ symbol, e.g. $x \in [0, 2]$ means that $x$ belongs to/is drawn from the set of all values in the closed interval $[0, 2]$ ; $x \in \{0, 1\}$ means that $x$ can take the values 0 and 1
$\leq$	Less than or equal to, represented in the text where necessary by $\leq$ (provided in this form to support display by some web browsers)
$\geq$	Greater than or equal to, represented in the text where necessary by $\geq$ (provided in this form to support display by some web browsers)

## 1.4.2 Statistical measures and related formulas

Table 1-3 provides a list of common measures (univariate statistics) applied to datasets, and associated formulas for calculating the measure from a sample dataset in summation form (rather than integral form) where necessary. In some instances these formulas are adjusted to provide estimates of the population values rather than those obtained from the sample of data one is working on.

Many of the measures can be extended to two-dimensional forms in a very straightforward manner, and thus they provide the basis for numerous standard formulas in spatial statistics. For a number of univariate statistics (variance, skewness, kurtosis) we refer to the notion of (estimated) *moments* about the mean. These are computations of the form

$$\sum (x_i - \bar{x})^r, r = 1, 2, 3, \dots$$

When  $r=1$  this summation will be 0, since this is just the difference of all values from the mean. For values of  $r>1$  the expression provides measures that are useful for describing the shape (spread, skewness, peakedness) of a distribution, and simple variations on the formula are used to define

the correlation between two or more datasets (the *product moment* correlation). The term *moment* in this context comes from physics, i.e. like ‘momentum’ and ‘moment of inertia’, and in a spatial (2D) context provides the basis for the definition of a centroid – the center of mass or center of gravity of an object, such as a polygon (see further, Section 4.2.5).

### Table 1-3 Common formulas and statistical measures

This table of measures has been divided into 9 subsections for ease of use. Each is provided with its own subheading:

- Counts and specific values
- Measures of centrality
- Measures of spread
- Measures of distribution shape
- Measures of complexity and dimensionality
- Common distributions
- Data transforms and back transforms
- Selected functions
- Matrix expressions

### 1.4.2.1 Counts and specific values

Measure	Definition	Expression(s)
Count	The number of data values in a set	$Count(\{x_i\})=n$
Top $m$ , Bottom $m$	The set of the largest (smallest) $m$ values from a set. May be generated via an SQL command	$Top_m\{x_i\}=\{X_{n-m+1}, \dots, X_{n-1}, X_n\};$ $Bot_m\{x_i\}=\{X_1, X_2, \dots, X_m\};$
Variety	The number of distinct i.e. different data values in a set. Some packages refer to the variety as diversity, which should not be confused with information theoretic and other diversity measures	
Majority	The most common i.e. most frequent data values in a set. Similar to mode (see below), but often applied to raster datasets at the neighborhood or zonal level. For general datasets the term should only be	

Measure	Definition	Expression(s)
	applied to cases where a given class is 50%+ of the total	
Minority	The least common i.e. least frequently occurring data values in a set. Often applied to raster datasets at the neighborhood or zonal level	
Maximum, <i>Max</i>	The maximum value of a set of values. May not be unique	$Max\{x_i\}=X_n$
Minimum, <i>Min</i>	The minimum value of a set of values. May not be unique	$Min\{x_i\}=X_1$
Sum	The sum of a set of data values	$\sum_{i=1}^n x_i$

**1.4.2.2 Measures of centrality**

Measure	Definition	Expression(s)
Mean (arithmetic)	The arithmetic average of a set of data values (also known as the <i>sample mean</i> where the data are a sample from a larger population). Note that if the set $\{f_j\}$ are regarded as weights rather than frequencies the result is known as the <i>weighted mean</i> . Other mean values include the geometric and harmonic mean. The population mean is often denoted by the symbol $\mu$ . In many instances the sample mean is the best (unbiased) estimate of the population mean and is sometimes denoted by $\mu$ with a ^ symbol above it) or as a variable such as $x$ with a bar above it.	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ $\bar{x} = \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i}$ $\bar{x} = \sum_{i=1}^n p_i x_i$
Mean (harmonic)	The harmonic mean, $H$ , is the mean of the reciprocals of the data values, which is then adjusted by taking the reciprocal of the result. The harmonic mean is less than or equal to the geometric mean, which is less than or equal to the arithmetic mean	$H = \left( \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i} \right)^{-1}$
Mean (geometric)	The geometric mean, $G$ , is the mean defined by taking the products of the data values and then adjusting the value by taking the $n^{\text{th}}$ root of the result. The geometric mean is greater than or equal to the harmonic mean and is less than or equal to the arithmetic mean	$G = \left( \prod_{i=1}^n x_i \right)^{1/n}$ <p style="text-align: right;">, hence</p> $\log(G) = \frac{1}{n} \sum_{i=1}^n \log(x_i)$
Mean (power)	The general (limit) expression for mean values. Values for $p$ give the following means: $p=1$ arithmetic; $p=2$ root mean	$M = \left( \frac{1}{n} \sum_{i=1}^n x_i^p \right)^{1/p}$

Measure	Definition	Expression(s)
	square; $p=-1$ harmonic. Limit values for $p$ (i.e. as $p$ tends to these values) give the following means: $p=0$ geometric; $p=-\infty$ minimum; $p=\infty$ maximum	
Trim-mean, $TM$ , $t$ , Olympic mean	The mean value computed with a specified percentage (proportion), $t/2$ , of values removed from each tail to eliminate the highest and lowest outliers and extreme values. For small samples a specific number of observations (e.g. 1) rather than a percentage, may be ignored. In general an equal number, $k$ , of high and low values should be removed and the number of observations summed should equal $n(1-t)$ expressed as an integer. This variant is sometimes described as the Olympic mean, as is used in scoring Olympic gymnastics for example	$TM = \frac{1}{n(1-t)} \sum_{i=nt/2}^{n(1-t/2)} X_i, t \in [0,1]$
Mode	The most common or frequently occurring value in a set. Where a set has one dominant value or range of values it is said to be unimodal; if there are several commonly occurring values or ranges it is described as multi-modal. Note that arithmetic mean-mode $\approx$ 3 (arithmetic mean-median) for many unimodal distributions	
Median, $Med$	The middle value in an ordered set of data if the set contains an odd number of values, or the average of the two middle values if the set contains an even number of values. For a continuous distribution the median is the 50% point (0.5) obtained from the cumulative distribution of the values or function	$Med\{x_i\}=X_{(n+1)/2}; n \text{ odd}$ $Med\{x_i\}=(X_{n/2}+X_{n/2+1})/2; n \text{ even}$
Mid-range, $MR$	The middle value of the Range	$MR\{x_i\}=\text{Range}/2$
Root mean square (RMS)	The root of the mean of squared data values. Squaring removes negative values	$\sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}$

#### 1.4.2.3 Measures of spread

Measure	Definition	Expression(s)
Range	The difference between the maximum and minimum values of a set	$Range\{x_i\}=X_n-X_1$
Lower quartile (25%), $LQ$	In an ordered set, 25% of data items are less than or equal to the upper bound of this range. For a continuous distribution	$LQ=\{X_1, \dots, X_{(n+1)/4}\}$

Measure	Definition	Expression(s)
	the LQ is the set of values from 0% to 25% (0.25) obtained from the cumulative distribution of the values or function. Treatment of cases where $n$ is even and $n$ is odd, and when $i$ runs from 1 to $n$ or 0 to $n$ vary	
Upper quartile (75%), $UQ$	In an ordered set 75% of data items are less than or equal to the upper bound of this range. For a continuous distribution the UQ is the set of values from 75% (0.75) to 100% obtained from the cumulative distribution of the values or function. Treatment of cases where $n$ is even and $n$ is odd, and when $i$ runs from 1 to $n$ or 0 to $n$ vary	$UQ=\{X_{3(n+1)/4}, \dots X_n\}$
Inter-quartile range, $IQR$	The difference between the lower and upper quartile values, hence covering the middle 50% of the distribution. The inter-quartile range can be obtained by taking the median of the dataset, then finding the median of the upper and lower halves of the set. The IQR is then the difference between these two secondary medians	$IQR=UQ-LQ$
Trim-range, $TR, t$	The range computed with a specified percentage (proportion), $t/2$ , of the highest and lowest values removed to eliminate outliers and extreme values. For small samples a specific number of observations (e.g. 1) rather than a percentage, may be ignored. In general an equal number, $k$ , of high and low values are removed (if possible)	$TR_t=X_{n(1-t/2)}-X_{nt/2}, t \in [0,1]$ $TR_{50\%}=IQR$
Variance, $Var, \sigma^2, s^2, \mu_2$	The average squared difference of values in a dataset from their population mean, $\mu$ , or from the sample mean (also known as the sample variance where the data are a sample from a larger population). Differences are squared to remove the effect of negative values (the summation would otherwise be 0). The third formula is the frequency form, where frequencies have been standardized, i.e. $\sum f_i=1$ . $Var$ is a function of the 2 <sup>nd</sup> moment about the mean. The population variance is often denoted by the symbol $\mu_2$ or $\sigma^2$ . The estimated population variance is often denoted by $s^2$ or by $\hat{\sigma}^2$ with a ^ symbol above it	$Var = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$ $Var = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ $Var = \sum_{i=1}^n f_i (x_i - \bar{x})^2$ $Var = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})$ $s^2 = \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

Measure	Definition	Expression(s)
Standard deviation, $SD$ , $s$ or $RMSD$	The square root of the variance, hence it is the Root Mean Squared Deviation (RMSD). The population standard deviation is often denoted by the symbol $\sigma$ . $SD^*$ shows the estimated population standard deviation (sometimes denoted by $\hat{\sigma}$ with a ^ symbol above it or by $s$ )	$SD = \sqrt{Var} = \sigma$ $SD = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$ $SD^* = \hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$
Standard error of the mean, $SE$	The estimated standard deviation of the mean values of $n$ samples from the same population. It is simply the sample standard deviation reduced by a factor equal to the square root of the number of samples, $n > 1$	$SE = \frac{SD}{\sqrt{n}}$
Root mean squared error, $RMSE$	The standard deviation of samples from a known set of true values, $x_i^*$ . If $x_i^*$ are estimated by the mean of sampled values $RMSE$ is equivalent to $RMSD$	$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - x_i^*)^2}$
Mean deviation/error, $MD$ or $ME$	The mean deviation of samples from the known set of true values, $x_i^*$	$MD = \frac{1}{n} \sum_{i=1}^n (x_i - x_i^*)$
Mean absolute deviation/error, $MAD$ or $MAE$	The mean absolute deviation of samples from the known set of true values, $x_i^*$	$MAE = \frac{1}{n} \sum_{i=1}^n  x_i - x_i^* $
Covariance, $Cov$	Literally the pattern of common (or co-) variation observed in a collection of two (or more) datasets, or partitions of a single dataset. Note that if the two sets are the same the covariance is the same as the variance	$Cov(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ $Cov(x, x) = Var(x)$
Correlation/product moment or Pearson's correlation coefficient, $r$	A measure of the similarity between two (or more) paired datasets. The correlation coefficient is the ratio of the covariance to the product of the standard deviations. If the two datasets are the same or perfectly matched this will give a result=1	$r = Cov(x, y) / SD_x SD_y$ $r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$
Coefficient of variation, $CV$	The ratio of the standard deviation to the mean, sometime computed as a percentage. If this ratio is close to 1, and the distribution is strongly left skewed, it may suggest the underlying distribution is Exponential. Note, mean values close to 0 may produce unstable results	$CV = SD / \bar{x}$

Measure	Definition	Expression(s)
Variance mean ratio, VMR	The ratio of the variance to the mean, sometime computed as a percentage. If this ratio is close to 1, and the distribution is unimodal and relates to count data, it may suggest the underlying distribution is Poisson. Note, mean values close to 0 may produce unstable results	$VMR = Var / \bar{x}$

**1.4.2.4 Measures of distribution shape**

Measure	Definition	Expression(s)
Skewness, $\alpha_3$	If a frequency distribution is unimodal and symmetric about the mean it has a skewness of 0. Values greater than 0 suggest skewness of a unimodal distribution to the right, whilst values less than 0 indicate skewness to the left. A function of the 3 <sup>rd</sup> moment about the mean (denoted by $\alpha_3$ with a ^ symbol above it for the sample skewness)	$\alpha_3 = \frac{1}{n\sigma^3} \sum_{i=1}^n (x_i - \mu)^3$ $\alpha_3 = \frac{1}{n\hat{\sigma}^3} \sum_{i=1}^n (x_i - \bar{x})^3$ $\hat{\alpha}_3 = \frac{n}{(n-1)(n-2)\hat{\sigma}^3} \sum_{i=1}^n (x_i - \bar{x})^3$
Kurtosis, $\alpha_4$	A measure of the peakedness of a frequency distribution. More pointy distributions tend to have high kurtosis values. A function of the 4 <sup>th</sup> moment about the mean. It is customary to subtract 3 from the raw kurtosis value (which is the kurtosis of the Normal distribution) to give a figure relative to the Normal (denoted by $\alpha_4$ with a ^ symbol above it for the sample kurtosis)	$\alpha_4 = \frac{1}{n\sigma^4} \sum_{i=1}^n (x_i - \mu)^4$ $\alpha_4 = \frac{1}{n\hat{\sigma}^4} \sum_{i=1}^n (x_i - \bar{x})^4$ $\hat{\alpha}_4 = \frac{a}{\hat{\sigma}^4} \sum_{i=1}^n (x_i - \bar{x})^4 - b \quad \text{where}$ $a = \frac{n(n+1)}{(n-1)(n-2)(n-3)}, \quad b = \frac{3(n-1)^2}{(n-2)(n-3)}$

**1.4.2.5 Measures of complexity and dimensionality**

Measure	Definition	Expression(s)
Information statistic (Entropy), $I$ (Shannon's)	A measure of the amount of pattern, disorder or <i>information</i> , in a set $\{x_i\}$ where $p_i$ is the proportion of events or values occurring in the $i^{th}$ class or range. Note that if $p_i=0$ then $p_i \log_2(p_i)$ is 0. $I$ takes values in the range $[0, \log_2(k)]$ . The lower value means all data falls into 1 category, whilst the upper means all data are evenly spread	$I = - \sum_{i=1}^k p_i \log_2(p_i)$

Measure	Definition	Expression(s)
Information statistic (Diversity), $Div$	Shannon's entropy statistic (see above) standardized by the number of classes, $k$ , to give a range of values from 0 to 1	$Div = \frac{-\sum_{i=1}^k p_i \log_2(p_i)}{\log_2(k)}$
Dimension (topological), $D_T$	Broadly, the number of (intrinsic) coordinates needed to refer to a single point anywhere on the object. The dimension of a point=0, a rectifiable line=1, a surface=2 and a solid=3. See text for fuller explanation. The value 2.5 (often denoted 2.5D) is used in GIS to denote a planar region over which a single-valued attribute has been defined at each point (e.g. height). In mathematics topological dimension is now equated to a definition similar to cover dimension (see below)	$D_T=0,1,2,3,\dots$
Dimension (capacity, cover or fractal), $D_C$	Let $N(h)$ represent the number of small elements of edge length $h$ required to cover an object. For a line, length 1, each element has length $1/h$ . For a plane surface each element (small square of side length $1/h$ ) has area $1/h^2$ , and for a volume, each element is a cube with volume $1/h^3$ . More generally $N(h)=1/h^D$ , where $D$ is the topological dimension, so $N(h)=h^{-D}$ and thus $\log(N(h))=-D\log(h)$ and so $D_C=-\log(N(h))/\log(h)$ . $D_C$ may be fractional, in which case the term <i>fractal</i> is used	$D_C = -\lim_{h \rightarrow 0^+} \frac{\ln N(h)}{\ln(h)}$ $D_C \geq 0$

#### 1.4.2.6 Common distributions

Measure	Definition	Expression(s)
Uniform (continuous)	All values in the range are equally likely. Mean= $a/2$ , variance= $a^2/12$ . Here we use $f(x)$ to denote the probability distribution associated with continuous valued variables $x$ , also described as a <i>probability density function</i>	$f(x) = \frac{1}{a}; x \in [0, a]$
Binomial (discrete)	The terms of the Binomial give the probability of $x$ successes out of $n$ trials, for example 3 heads in 10 tosses of a coin, where $p$ =probability of success and $q=1-p$ =probability of failure. Mean, $m=np$ , variance= $npq$ . Here we use $p(x)$ to denote the probability distribution associated with discrete valued variables $x$	$p(x) = \frac{n!}{(n-x)!x!} p^x q^{1-x}; x=1,2,\dots,n$
Poisson (discrete)	An approximation to the Binomial when $p$ is very small and $n$ is large ( $>100$ ), but the mean $m=np$ is fixed and finite (usually not	$p(x) = \frac{m^x}{x!} e^{-m}; x=1,2,\dots,n$

## INTRODUCTION AND TERMINOLOGY

Measure	Definition	Expression(s)
	large). Mean=variance= $m$	
Normal (continuous)	The distribution of a measurement, $x$ , that is subject to a large number of independent, random, additive errors. The Normal distribution may also be derived as an approximation to the Binomial when $p$ is not small (e.g. $p \approx 1/2$ ) and $n$ is large. If $\mu$ =mean and $\sigma$ =standard deviation, we write $N(\mu, \sigma)$ as the Normal distribution with these parameters. The Normal- or z-transform $z=(x-\mu)/\sigma$ changes (normalizes) the distribution so that it has a zero mean and unit variance, $N(0,1)$ . The distribution of $n$ mean values of independent random variables drawn from <i>any</i> underlying distribution is also Normal ( <i>Central Limit Theorem</i> )	$f(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}; z \in [-\infty, \infty]$

### 1.4.2.7 Data transforms and back transforms

Measure	Definition	Expression(s)
Log	If the frequency distribution for a dataset is broadly unimodal and left-skewed, the natural log transform (logarithms base e) will adjust the pattern to make it more symmetric/similar to a Normal distribution. For variates whose values may range from 0 upwards a value of 1 is often added to the transform. Back transform with the exp() function	$z = \ln(x)$ or $z = \ln(x+1)$ <i>n.b.</i> $\ln(x) = \log_e(x) = \log_{10}(x) * \log_{10}(e)$ $x = \exp(z)$ or $x = \exp(z) - 1$
Square root (Freeman-Tukey)	A transform that may adjust the dataset to make it more similar to a Normal distribution. For variates whose values may range from 0 upwards a value of 1 is often added to the transform. For $0 \leq x \leq 1$ (e.g. rate data) the combined form of the transform is often used, and is known as the Freeman-Tukey (FT) transform	$z = \sqrt{x}$ , or $z = \sqrt{x+1}$ , or $z = \sqrt{x} + \sqrt{x+1}$ (FT) $x = z^2$ , or $x = z^2 - 1$
Logit	Often used to transform binary response data, such as survival/non-survival or present/absent, to provide a continuous value in the range $(-\infty, \infty)$ , where $p$ is the proportion of the sample that is 1 (or 0). The inverse or back-transform is shown as $p$ in terms of $z$ . This transform avoids concentration of values at the ends of the range. For samples where proportions $p$ may take the values 0 or 1 a modified form of the transform may be used. This is typically achieved by adding $1/2n$ to the	$z = \ln\left(\frac{p}{1-p}\right), p \in [0,1]$ $p = \frac{e^z}{1 + e^z}$

Measure	Definition	Expression(s)
	numerator and denominator, where $n$ is the sample size. Often used to correct S-shaped ( <i>logistic</i> ) relationships between response and explanatory variables	
Normal, z-transform	This transform normalizes or standardizes the distribution so that it has a zero mean and unit variance. If $\{x_i\}$ is a set of $n$ sample <i>mean</i> values from <i>any</i> probability distribution with mean $\mu$ and variance $\sigma^2$ then the z-transform shown here as $z_2$ will be distributed $N(0,1)$ for large $n$ (Central Limit Theorem). The divisor in this instance is the standard error. In both instances the standard deviation must be non-zero	$z_1 = \frac{(x - \mu)}{\sigma}$ $z_2 = \frac{(x - \mu)}{\sigma/\sqrt{n}}$
Box-Cox, power transforms	A family of transforms defined for positive data values only, that often can make datasets more Normal; $k$ is a parameter. The inverse or back-transform is also shown as $x$ in terms of $z$	$z = \frac{(x^k - 1)}{k}, k > 0, x > 0$ $x = (kz + 1)^{1/k}, k > 0$
Angular transforms (Freeman-Tukey)	A transform for proportions, $p$ , designed to spread the set of values near the end of the range. $k$ is typically 0.5. Often used to correct S-shaped relationships between response and explanatory variables. If $p=x/n$ then the Freeman-Tukey (FT) version of this transform is the averaged version shown. This is a variance-stabilizing transform	$z = \sin^{-1}(p^k),$ $p = \sin(z)^{1/k}$ $z = \sin^{-1}\left(\sqrt{\frac{x}{n+1}}\right) +$ $\sin^{-1}\left(\sqrt{\frac{x+1}{n+1}}\right) \text{ (FT)}$

#### 1.4.2.8 Selected functions

Measure	Definition	Expression(s)
Bessel function of the first kind	Bessel functions occur as the solution to specific differential equations. They are described with reference to a parameter known as the order, shown as a subscript. For integer orders Bessel functions can be represented as an infinite series. Order 0 and Order 1 expansions are shown here. The graph of a Bessel function is similar to a dampening sine wave. Usage in spatial analysis arises in connection with directional statistics and spline curve fitting. See the <a href="#">Mathworld</a> website entry for more details	$I_0(\kappa) = \sum_{i=0}^{\infty} \frac{(-1)^i (\kappa/2)^{2i}}{(i!)^2}, \text{ and}$ $I_1(\kappa) = \frac{\kappa}{2} \sum_{i=0}^{\infty} \frac{(-1)^i (\kappa/2)^{2i+1}}{i!(i+1)!}$
Exponential integral function, $E_1(x)$	A definite integral function. Used in association with spline curve fitting. See the <a href="#">Mathworld</a> website entry for more	$E_1(x) = \int_1^{\infty} \frac{e^{-tx}}{t} dt$

Measure	Definition	Expression(s)
	details	
Gamma function, $\Gamma$	A widely used definite integral function. For integer values of $x$ : $\Gamma(x)=(x-1)!$ and $\Gamma(x/2)=(x/2-1)!$ so $\Gamma(3/2)=(1/2)!/2=(\sqrt{\pi})/2$ See the <a href="#">Mathworld</a> website entry for more details	$\Gamma(x) = \int_0^{\infty} x^{1/2} e^{-x} dx$ $\Gamma(1/2) = \sqrt{\pi}$

**1.4.2.9 Matrix expressions**

Measure	Definition	Expression(s)
Identity	A matrix with diagonal elements 1 and off-diagonal elements 0	$I = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 1 \end{bmatrix}$
Determinant	Determinants are only defined for square matrices. Let $A$ be an $n$ by $n$ matrix with elements $\{a_{ij}\}$ . The matrix $M_{ij}$ here is a subset of $A$ known as the <i>minor</i> , formed by eliminating row $i$ and column $j$ from $A$ . An $n$ by $n$ matrix, $A$ , with $\text{Det}=0$ is described as <i>singular</i> , and such a matrix has no inverse. If $\text{Det}(A)$ is very close to 0 it is described as <i>ill-conditioned</i>	$ A , \text{Det}(A)$ $ A  = \sum_i^n a_{ij} a^{ij}$ , where $a^{ij} = (-1)^{i+j} M_{ij}$
Inverse	The matrix equivalent of division in conventional algebra. For a matrix, $A$ , to be invertible its determinant must be non-zero, and ideally not very close to zero. A matrix that has an inverse is by definition non-singular. A symmetric real-valued matrix is <i>positive definite</i> if all its eigenvalues are positive, whereas a <i>positive semi-definite</i> matrix allows for some eigenvalues to be 0. A matrix, $A$ , that is invertible satisfies the relation $AA^{-1}=I$	$A^{-1}$
Transpose	A matrix operation in which the rows and columns are transposed, i.e. in which elements $a_{ij}$ are swapped with $a_{ji}$ for all $i, j$ . The inverse of a transposed matrix is the same as the transpose of the matrix inverse	$A^T$ or $A'$ $(A^T)^{-1} = (A^{-1})^T$
Symmetric	A matrix in which element $a_{ij}=a_{ji}$ for all $i, j$	$A=A^T$
Trace	The sum of the diagonal elements of a matrix, $a_{ii}$ – the sum of the eigenvalues of a matrix equals its trace	$\text{Tr}(A)$

Measure	Definition	Expression(s)
Eigenvalue, Eigenvector	If $A$ is a real-valued $k$ by $k$ square matrix and $x$ is a non-zero real-valued vector, then a scalar $\lambda$ that satisfies the equation shown in the adjacent column is known as an eigenvalue of $A$ and $x$ is an eigenvector of $A$ . There are $k$ eigenvalues of $A$ , each with a corresponding eigenvector. The matrix $A$ can be decomposed into three parts, as shown, where $E$ is a matrix of its eigenvectors and $D$ is a diagonal matrix of its eigenvalues	$(A-\lambda I)x=0$ $A=EDE^{-1}$ (diagonalization)