

## Semiparametric Regression

Semiparametric regression is concerned with the flexible incorporation of nonlinear functional relationships in regression analyses. Any application area that uses regression analysis can benefit from semiparametric regression. Assuming only a basic familiarity with ordinary parametric regression, this user-friendly book explains the techniques and benefits of semiparametric regression in a concise and modular fashion. The authors make liberal use of graphics and examples plus case studies taken from environmental, financial, and other applications. They include practical advice on implementation and pointers to relevant software.

This book is suitable as a textbook for students with little background in regression as well as a reference book for statistically oriented scientists – such as biostatisticians, econometricians, quantitative social scientists, and epidemiologists – with a good working knowledge of regression and the desire to begin using more flexible semiparametric models. Even experts on semiparametric regression should find something new here.

David Ruppert is the Andrew Schultz, Jr., Professor of Engineering (School of Operations Research and Industrial Engineering) and Professor of Statistical Science at Cornell University. He has served as editor for a number of prestigious series and journals and has published some 80 articles of his own as well as co-authoring two popular books, *Transformation and Weighting in Regression* and *Measurement Error in Nonlinear Models*. He is also winner of the Wilcoxon Prize for best practical applications paper in technometrics and an elected Fellow of the American Statistical Association and the Institute of Mathematical Statistics.

M. P. Wand is Professor of Statistics at the University of New South Wales in Sydney, Australia. He has held faculty appointments at Harvard University, Rice University, and Texas A&M University. Dr. Wand is a Fellow of the American Statistical Association and has served as an associate editor for the *Journal of the American Statistical Association* and *Biometrika*. He is winner of the P. A. P. Moran Medal for statistical research.

R. J. Carroll is Distinguished Professor of Statistics, Nutrition and Toxicology at Texas A&M University. Among his many honors are the COPSS Presidents' Award, the Fisher Lecture, the Snedecor Award, and the Wilcoxon Prize. He is an elected Fellow of the American Statistical Association and the Institute of Mathematical Statistics as well as an elected member of the International Statistical Institute.

---

CAMBRIDGE SERIES IN STATISTICAL AND PROBABILISTIC MATHEMATICS

---

*Editorial Board*

- R. Gill (Department of Mathematics, Utrecht University)  
B. D. Ripley (Department of Statistics, University of Oxford)  
S. Ross (Department of Industrial Engineering, University of California, Berkeley)  
M. Stein (Department of Statistics, University of Chicago)  
D. Williams (School of Mathematical Sciences, University of Bath)

This series of high-quality upper-division textbooks and expository monographs covers all aspects of stochastic applicable mathematics. The topics range from pure and applied statistics to probability theory, operations research, optimization, and mathematical programming. The books contain clear presentations of new developments in the field and also of the state of the art in classical methods. While emphasizing rigorous treatment of theoretical methods, the books also contain applications and discussions of new techniques made possible by advances in computational practice.

*Already published*

1. *Bootstrap Methods and Their Application*, by A. C. Davison and D. V. Hinkley
2. *Markov Chains*, by J. Norris
3. *Asymptotic Statistics*, by A. W. van der Vaart
4. *Wavelet Methods for Time Series Analysis*, by Donald B. Percival and Andrew T. Walden
5. *Bayesian Methods*, by Thomas Leonard and John S. J. Hsu
6. *Empirical Processes in M-Estimation*, by Sara van de Geer
7. *Numerical Methods of Statistics*, by John F. Monahan
8. *A User's Guide to Measure Theoretic Probability*, by David Pollard
9. *The Estimation and Tracking of Frequency*, by B. G. Quinn and E. J. Hannan
10. *Data Analysis and Graphics using R*, by John Maindonald and John Braun
11. *Statistical Models*, by A. C. Davison

Cambridge University Press  
978-0-521-78050-6 - Semiparametric Regression  
David Ruppert, M. P. Wand and R. J. Carroll  
Frontmatter  
[More information](#)

---

# Semiparametric Regression

---

DAVID RUPPERT

*Cornell University*

M. P. WAND

*Harvard University*

R. J. CARROLL

*Texas A&M University*



CAMBRIDGE  
UNIVERSITY PRESS

Cambridge University Press  
 978-0-521-78050-6 - Semiparametric Regression  
 David Ruppert, M. P. Wand and R. J. Carroll  
 Frontmatter  
[More information](#)

CAMBRIDGE UNIVERSITY PRESS  
 Cambridge, New York, Melbourne, Madrid, Cape Town,  
 Singapore, São Paulo, Delhi, Tokyo, Mexico City

Cambridge University Press  
 The Edinburgh Building, Cambridge CB2 8RU, UK

Published in the United States of America by  
 Cambridge University Press, New York

[www.cambridge.org](http://www.cambridge.org)  
 Information on this title: [www.cambridge.org/9780521780506](http://www.cambridge.org/9780521780506)

© David Ruppert, M. P. Wand, R. J. Carroll 2003

This publication is in copyright. Subject to statutory exception  
 and to the provisions of relevant collective licensing agreements,  
 no reproduction of any part may take place without the written  
 permission of Cambridge University Press.

First published 2003  
 Reprinted 2005, 2006, 2008, 2009

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

*Library of Congress Cataloguing in Publication data*

Ruppert, David, 1948–  
 Semiparametric regression / David Ruppert, M. P. Wand, R. J. Carroll.  
 p.cm.

Includes bibliographical references and index.

ISBN 0-521-78050-0 – ISBN 0-521-78516-2 (pb.)

I. Regression analysis. 2. Nonparametric statistics. I. Wand, M. P. (Matthew P).  
 II. Carroll, Raymond J. III. Title.

QA278.2.R87 2003  
 519.5'36 – dc21 2002041460

ISBN 978-0-521-78050-6 Hardback

ISBN 978-0-521-78516-7 Hardback

Cambridge University Press has no responsibility for the persistence or  
 accuracy of URLs for external or third-party internet websites referred to in  
 this publication, and does not guarantee that any content on such websites is,  
 or will remain, accurate or appropriate. Information regarding prices, travel  
 timetables, and other factual information given in this work is correct at  
 the time of first printing but Cambridge University Press does not guarantee  
 the accuracy of such information thereafter.

Cambridge University Press  
978-0-521-78050-6 - Semiparametric Regression  
David Ruppert, M. P. Wand and R. J. Carroll  
Frontmatter  
[More information](#)

---

*To Anne, with love*  
— David

*To my wife's parents, Ayhan and Recep*  
— Matt

*To Brett and Jeb*  
— Raymond

Cambridge University Press  
978-0-521-78050-6 - Semiparametric Regression  
David Ruppert, M. P. Wand and R. J. Carroll  
Frontmatter  
[More information](#)

---

# Contents

---

<i>Preface</i>	<i>page</i> xiii
<i>Guide to Notation</i>	xv
<b>1 Introduction</b>	<b>1</b>
1.1 Assessing the Carcinogenicity of Phenolphthalein	3
1.2 Salinity and Fishing in North Carolina	4
1.3 Management of a Retirement Fund	5
1.4 Biomonitoring of Airborne Mercury	7
1.5 Term Structure of Interest Rates	7
1.6 Air Pollution and Mortality in Milan: The Harvesting Effect	11
<b>2 Parametric Regression</b>	<b>15</b>
2.1 Introduction	15
2.2 Linear Regression Models	15
2.3 Regression Diagnostics	20
2.4 Inference	28
2.5 Parametric Additive Models	36
2.6 Model Selection	44
2.7 Polynomial Regression Models	46
2.8 Nonlinear Regression	48
2.9 Transformations in Regression	51
2.10 Bibliographic Notes	55
2.11 Summary of Formulas	55
<b>3 Scatterplot Smoothing</b>	<b>57</b>
3.1 Introduction	57
3.2 Preliminary Ideas	58
3.3 Practical Implementation	62
3.4 Automatic Knot Selection	64
3.5 Penalized Spline Regression	65
3.6 Quadratic Spline Bases	67
3.7 Other Spline Models and Bases	69
3.8 Other Penalties	74
3.9 General Definition of a Penalized Spline	75
3.10 Linear Smoothers	76
3.11 Error of a Smoother	76
	vii

3.12 Rank of a Smoother	78
3.13 Degrees of Freedom of a Smoother	80
3.14 Residual Degrees of Freedom	82
3.15 Other Approaches to Scatterplot Smoothing	84
3.16 Choosing a Scatterplot Smoother	87
3.17 Bibliographical Notes	88
3.18 Summary of Formulas	89
<b>4 Mixed Models</b>	<b>91</b>
4.1 Introduction	91
4.2 Mixed Models	91
4.3 Prediction	95
4.4 The Linear Mixed Model (LMM)	98
4.5 Estimation and Prediction in LMM	98
4.6 Estimated BLUP (EBLUP)	101
4.7 Standard Error Estimation	102
4.8 Hypothesis Testing	104
4.9 Penalized Splines as BLUPs	108
4.10 Bibliographical Notes	110
4.11 Summary of Formulas	110
<b>5 Automatic Scatterplot Smoothing</b>	<b>112</b>
5.1 Introduction	112
5.2 The Likelihood Approach	113
5.3 The Model Selection Approach	114
5.4 Caveats of Automatic Parameter Selection	120
5.5 Choosing the Knots and Basis Functions	123
5.6 Automatic Selection of the Number of Knots	127
5.7 Bibliographical Notes	131
5.8 Summary of Formulas	131
<b>6 Inference</b>	<b>133</b>
6.1 Introduction	133
6.2 Variability Bands	133
6.3 Confidence and Prediction Intervals	135
6.4 Inference for Penalized Splines	137
6.5 Simultaneous Confidence Bands	142
6.6 Testing the Adequacy of Parametric Models	145
6.7 Testing for No Effect	149
6.8 Inference Using First Derivatives	151
6.9 Testing for Existence of a Feature	156
6.10 Bibliographical Notes	158
6.11 Summary of Formulas	159
<b>7 Simple Semiparametric Models</b>	<b>161</b>
7.1 Introduction	161
7.2 Beyond Scatterplot Smoothing	161



<i>Contents</i>	ix
7.3 Semiparametric Binary Offset Model	162
7.4 Additivity and Interactions	164
7.5 General Parametric Component	164
7.6 Inference	167
7.7 Bibliographical Notes	168
<b>8 Additive Models</b>	<b>170</b>
8.1 Introduction	170
8.2 Fitting an Additive Model	171
8.3 Degrees of Freedom	174
8.4 Smoothing Parameter Selection	176
8.5 Hypothesis Testing	181
8.6 Model Selection	183
8.7 Bibliographical Notes	185
<b>9 Semiparametric Mixed Models</b>	<b>186</b>
9.1 Introduction	186
9.2 Additive Mixed Models	186
9.3 Subject-Specific Curves	191
9.4 Bibliographical Notes	192
<b>10 Generalized Parametric Regression</b>	<b>194</b>
10.1 Introduction	194
10.2 Binary Response Data	194
10.3 Logistic Regression	195
10.4 Other Generalized Linear Models	197
10.5 Iteratively Reweighted Least Squares	200
10.6 Hat Matrix, Degrees of Freedom, and Standard Errors	201
10.7 Overdispersion and Variance Functions: Pseudolikelihood	201
10.8 Generalized Linear Mixed Models	203
10.9 Deviance	209
10.10 Technical Details	210
10.11 Bibliographical Notes	213
<b>11 Generalized Additive Models</b>	<b>214</b>
11.1 Introduction	214
11.2 Generalized Scatterplot Smoothing	215
11.3 Generalized Additive Mixed Models	217
11.4 Degrees-of-Freedom Approximations	219
11.5 Automatic Smoothing Parameter Selection	220
11.6 Hypothesis Testing	220
11.7 Model Selection	221
11.8 Density Estimation	221
11.9 Bibliographical Notes	222
<b>12 Interaction Models</b>	<b>223</b>
12.1 Introduction	223
12.2 Binary-by-Continuous Interaction Models	224

x	<i>Contents</i>
12.3	Factor-by-Curve Interactions in Additive Models 226
12.4	Varying Coefficient Models 234
12.5	Continuous-by-Continuous Interactions 235
12.6	Bibliographical Notes 237
<b>13</b>	<b>Bivariate Smoothing 238</b>
13.1	Introduction 238
13.2	Choice of Bivariate Basis Functions 240
13.3	Kriging 242
13.4	General Radial Smoothing 248
13.5	Default Automatic Bivariate Smoother 256
13.6	Geoadditive Models 258
13.7	Additive Plus Interaction Models 259
13.8	Generalized Bivariate Smoothing 259
13.9	Appendix: Equivalence of BLUP using $\mathbf{Z}_R$ and $\mathbf{Z}_P$ 259
13.10	Bibliographical Notes 260
<b>14</b>	<b>Variance Function Estimation 261</b>
14.1	Introduction 261
14.2	Formulation 263
14.3	Application to the LIDAR Data 264
14.4	Quasilikelihood and Variance Functions 266
14.5	Bibliographical Notes 267
<b>15</b>	<b>Measurement Error 268</b>
15.1	Introduction 268
15.2	Formulation 269
15.3	The Expectation Maximization (EM) Algorithm 270
15.4	Simulated Example Revisited 273
15.5	Sensitivity Analysis Example 273
15.6	Bibliographical Notes 275
<b>16</b>	<b>Bayesian Semiparametric Regression 276</b>
16.1	Introduction 276
16.2	General Framework 277
16.3	Scatterplot Smoothing 279
16.4	Linear Mixed Models 285
16.5	Generalized Linear Mixed Models 288
16.6	Rao–Blackwellization 291
16.7	Bibliographical Notes 292
<b>17</b>	<b>Spatially Adaptive Smoothing 293</b>
17.1	Introduction 293
17.2	A Local Penalty Method 294
17.3	Completely Automatic Algorithm 295
17.4	Bayesian Inference 296

<i>Contents</i>	xi
17.5 Simulations	298
17.6 LIDAR Example	304
17.7 Additive Models	305
17.8 Bibliographical Notes	307
<b>18 Analyses</b>	<b>308</b>
18.1 Cancer Rates on Cape Cod	308
18.2 Assessing the Carcinogenicity of Phenolphthalein	308
18.3 Salinity and Fishing in North Carolina	308
18.4 Management of a Retirement Fund	313
18.5 Biomonitoring of Airborne Mercury	314
18.6 Term Structure of Interest Rates	315
18.7 Air Pollution and Mortality in Milan: The Harvesting Effect	319
<b>19 Epilogue</b>	<b>320</b>
19.1 Introduction	320
19.2 Minimalist Statistics	320
19.3 Some Omitted Topics	321
19.4 Future Research	325
<b>A Technical Complements</b>	<b>326</b>
A.1 Introduction	326
A.2 Matrix Definitions and Results	326
A.3 Linear Algebra	331
A.4 Probability Definitions and Results	333
A.5 Maximum Likelihood Estimation	335
A.6 Bibliographical Notes	335
<b>B Computational Issues</b>	<b>336</b>
B.1 Fast Computation of Penalized Spline Smooths	336
B.2 Computation of Covariance Matrix Estimators	351
B.3 Software	353
<i>Bibliography</i>	361
<i>Author Index</i>	375
<i>Notation Index</i>	380
<i>Example Index</i>	381
<i>Subject Index</i>	382

Cambridge University Press  
978-0-521-78050-6 - Semiparametric Regression  
David Ruppert, M. P. Wand and R. J. Carroll  
Frontmatter  
[More information](#)

---

## Preface

---

The primary aim of this book is to guide researchers needing to flexibly incorporate nonlinear relationships into their regression analyses. Flexible nonlinear regression is traditionally known as *nonparametric regression*; it differs from parametric regression in that the shape of the functional relationships are not predetermined but can adjust to capture unusual or unexpected features of the data.

Almost all existing regression texts treat either parametric or nonparametric regression exclusively. The level of exposition between books of either type differs quite alarmingly. In this book we argue that nonparametric regression can be viewed as a relatively simple extension of parametric regression and treat the two together. We refer to this combination as *semiparametric regression*. Our approach to semiparametric regression is based on penalized regression splines and mixed models. Indeed, every model in this book is a special case of the linear mixed model or its generalized counterpart. This makes the methodology modular and is in keeping with our general philosophy of *minimalist statistics* (see Section 19.2), where the amount of methodology, terminology, and so on is kept to a minimum. This is the first smoothing book that makes use of the mixed model representation of smoothers.

Unlike many other texts on nonparametric regression, this book is very much problem-driven. Examples from our collaborative research (and elsewhere) have driven the selection of material and emphases and are used throughout the book.

The book is suitable for several audiences. One audience consists of students or working scientists with only a moderate background in regression, though familiarity with matrix and linear algebra is assumed. Marginal notes and the appendices are intended for beginners, especially those from interface disciplines. We make liberal use of graphics because visualization is a particularly effective tool for acquiring intuition in a new subject.

Another audience that we are aiming at consists of statistically oriented scientists (e.g., biostatisticians, econometricians, quantitative social scientists, and epidemiologists) who have a good working knowledge of linear models and the desire to begin using more flexible semiparametric models. There are many connections between linear and nonparametric regression. Our goal is to exploit them and the reader's knowledge of linear models to provide a foundation for understanding nonparametric modeling.

There is enough new material to be of interest even to experts on smoothing, and they are a third possible audience.

There are several competing approaches to nonparametric modeling: smoothing splines (e.g., Eubank 1988, 1999; Wahba 1990; Green and Silverman 1994); series-based smoothers, including wavelets (Tarter and Lock 1993; Ogden 1996); kernel methods, including local regression (Wand and Jones 1995; Fan and Gijbels 1996); and regression splines (Friedman 1991; Stone et al. 1997; Hansen and Kooperberg 2002). All four approaches can be used effectively and have their devotees. We believe that the nature of the data should play a role in the choice among them. For example, wavelets are more suited to highly oscillatory functions. Apart from this, the choice of a nonparametric regression method is a matter somewhat of individual taste and background. Based on our motivating applications and personal tastes, the approach to nonparametric regression used throughout this book is what we call *penalized splines*, although they are also labeled as *P-splines*, *pseudosplines*, and *low-rank spline smoothers* in the literature. Penalized splines are quite similar to smoothing splines; in fact, they are a generalization of smoothing splines that allow more flexible choices of the spline model, the basis functions for that model, and the penalty.

Penalized splines have close ties with ridge regression, mixed models, and Bayesian statistics, ties that were discovered by researchers working on smoothing splines. These ties allow techniques from mixed models – for example, (restricted) maximum likelihood estimation and likelihood ratio tests – to be added to penalized spline methodology. Similarly, Bayesian techniques based on Markov chain Monte Carlo provide what we believe to be the most satisfactory approach to fitting complex semiparametric models as well as the direction that semiparametric regression is most likely to take in the future. This book includes introductions to mixed models and to Bayesian modeling.

### *Acknowledgments*

We are especially grateful to Ciprian Crainiceanu and Bhaswati Ganguli for their assistance in the preparation of this book. Ciprian wrote the WinBugs program in Appendix B and wrote the programs used for simulations-based  $p$ -values for likelihood ratio tests. Several other of our colleagues and collaborators have contributed to the book in various ways. We would like to thank Marc Aerts, Babette Brumback, Tianxi Cai, Gerda Claeskens, Brent Coull, Maria Durban, Garrett Fitzmaurice, Jonathan French, Robert Gentleman, Bob Gray, Nick Horton, Joe Ibrahim, Erin Kammann, Göran Kauermann, Robert Kohn, Nan Laird, Nick Lange, Mary Lindstrom, Long Ngo, Doug Nychka, Michael O’Connell, Helen Parise, José Pinheiro, Louise Ryan, Misha Salganik, Joel Schwartz, John Staudenmayer, Sally Thurston, Carrie Wager, Naisyin Wang, Jim Ware, Antonella Zanobetti, and Yihua Zhao for their collaboration, interest, and comments.

We thank Lauren Cowles for being a very supportive and patient editor. We would like to express our gratitude to Tom Ryan and Misha Salganik for sending us errata.

The second author lovingly acknowledges the support of his wife, Handan, and children, Declan and Jaida, throughout this project. Support of the Department of Biostatistics, Harvard University, is also gratefully acknowledged.

## Guide to Notation

---

This chapter gives a brief overview of notational conventions used in the book. Please see the Notation Index for more specialized notation.

The symbol “ $\equiv$ ” means “equal by definition”.

We use both lower- and uppercase letters (e.g.,  $x$ ,  $X$ , and  $\lambda$ ) to denote scalar quantities, either fixed or random. Lowercase bold letters (e.g.,  $\mathbf{x}$  and  $\boldsymbol{\lambda}$ ) will be used for vectors. Uppercase bold fonts (e.g.,  $\mathbf{X}$  and  $\mathbf{\Lambda}$ ) will denote matrices. The entries of a vector or matrix use the same letter and case as the vector or matrix itself but are not bold. Thus,

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

and

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}.$$

If a matrix is partitioned then the submatrices are in bold; for example,

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}.$$

We will indicate the row index of a matrix to the right and the column index below, as in:

$$\mathbf{C} = \left[ c_{ik} \right]_{\substack{1 \leq k \leq K \\ 1 \leq i \leq n}}.$$

The transpose of  $\mathbf{A}$  is denoted by  $\mathbf{A}^T$ . If  $\mathbf{A}$  is an invertible square matrix, then  $\mathbf{A}^{-1}$  denotes its inverse. Any vector is assumed to be a column, so its transpose is a row.

The *norm* of a vector  $\mathbf{x}$  is denoted by  $\|\mathbf{x}\|$ ; that is,

$$\|\mathbf{x}\| \equiv \sqrt{\mathbf{x}^T \mathbf{x}}.$$

The real line will be denoted by  $\mathbb{R}$ , and  $d$ -dimensional space will be denoted by  $\mathbb{R}^d$ .

For a function  $f(x)$  of a scalar  $x$ ,

$$f^{(r)}(x) \equiv (d^r/dx^r) f(x),$$

the  $r$ th derivative of  $f(x)$ .

If  $f(\mathbf{x})$  is a function from  $\mathbb{R}^d$  to  $\mathbb{R}$  then the *derivative vector* is a  $1 \times d$  row vector with  $j$ th entry equal to  $(\partial/\partial x_j)f(\mathbf{x})$ , the partial derivative of  $f(\mathbf{x})$  with respect to  $x_j$ , and is denoted by

$$Df(\mathbf{x}).$$

The *Hessian matrix* is a  $d \times d$  matrix whose  $(i, j)$  entry is equal to

$$\frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x});$$

it is denoted by

$$Hf(\mathbf{x}).$$

If  $x$  and  $y$  are random variables, then  $E(x)$ ,  $\text{var}(x)$ , and  $\text{st.dev.}(x)$  are the mean, variance, and standard deviation of  $x$ , and  $\text{cov}(x, y)$  is the covariance between  $x$  and  $y$ .  $\text{Cov}(\mathbf{x})$  is the covariance matrix of a random vector  $\mathbf{x}$ ; see Appendix A for its definition.