

Springer Series in Statistics

# Smoothing Spline ANOVA Models

*Second Edition*

 Springer

# Springer Series in Statistics 297

*Advisors:*

P. Bickel, P. Diggle, S. Feinberg, U. Gather,  
I. Olkin, S. Zeger

For further volumes:

<http://www.springer.com/series/692>



Chong Gu

# Smoothing Spline ANOVA Models

Second Edition

 Springer

Chong Gu  
Department of Statistics  
Purdue University  
West Lafayette, IN 47907  
USA

ISSN 0172-7397  
ISBN 978-1-4614-5368-0 ISBN 978-1-4614-5369-7 (eBook)  
DOI 10.1007/978-1-4614-5369-7  
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2012950795

© Springer Science+Business Media New York 2013

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media ([www.springer.com](http://www.springer.com))

*To my father  
For the books and the bookcases*



# Preface to the First Edition

Thirty years have passed since the pioneering work of [Kimeldorf and Wahba \(1970a, 1970b, 1971\)](#) and [Good and Gaskins \(1971\)](#), and during this time, a rich body of literature has been developed on smoothing methods with roughness penalties. There have been two books solely devoted to the subject prior to this one, of which [Wahba \(1990\)](#) compiled an excellent synthesis for work up to that date, and [Green and Silverman \(1994\)](#) provided a mathematically gentler introduction to the field through regression models that are largely univariate.

Much has happened in the past decade, and more has been done with the penalty method than just regression. In this book, I have tried to assemble a comprehensive treatment of penalty smoothing under a unified framework. Treated are (i) regression with Gaussian and non-Gaussian responses as well as with censored lifetime data, (ii) density and conditional density estimation under a variety of sampling schemes, and (iii) hazard rate estimation with censored lifetime data and covariates. The unifying themes are the general penalized likelihood method and the construction of multivariate models with certain ANOVA decompositions built in. Extensive discussions are devoted to model (penalty) construction, smoothing parameter selection, computation, and asymptotic convergence. There are, however, many omissions, and the selection and treatment of topics solely reflect my personal preferences and views. Most of the materials have appeared in the literature, but a few items are new, as noted in the bibliographic notes at the end of the chapters.



An adequate treatment of model construction in the context requires some elementary knowledge of reproducing kernel Hilbert spaces, of which a self-contained introduction is included early in the book; the materials should be accessible to a second-year graduate student with a good training in calculus and linear algebra. Also assumed is a working knowledge of basic statistical inference such as linear models, maximum likelihood estimates, etc. To better understand materials on hazard estimation, prior knowledge of basic survival analysis would also help.

Most of the computational and data analytical tools discussed in the book are implemented in R, an open-source clone of the popular S/Splus language. Code for regression is reasonably polished and user-friendly and has been distributed in the R package `gss` available through CRAN, the Comprehensive R Archive Network, with the master site at

<http://cran.r-project.org>

The use of `gss` facilities is illustrated in the book through simulated and real-data examples.

Remaining on my wish list are (i) polished, user-friendly software tools for density estimation and hazard estimation, (ii) fast computation via approximate solutions of penalized likelihood problems, and (iii) handling of parametric random effects such as those appearing in longitudinal models and hazard models with frailty. All of the above are under active development and could be addressed in a later edition of the book or, sooner than that, in later releases of `gss`.

The book was conceived in Spring 1996 when I was on leave at the Department of Statistics, University of Michigan, which offered me the opportunity to teach a course on the subject. Work on the book has been on and off since then, with much of the progress being made in the 1997–1998 academic year during my visit at the National Institute of Statistical Sciences, and in Fall 2000 when I was teaching a course on the subject at Purdue.

I am indebted to Grace Wahba, who taught me smoothing splines, and to Doug Bates, who taught me statistical computing. Bill Studden carefully read various drafts of Chaps. 1, 2, and 4; his questions alerted me to numerous accounts of mathematical sloppiness in the text and his suggestions led to much improved presentations. Detailed comments and suggestions by Nancy Heckman on a late draft helped me to fix numerous problems throughout the first five chapters and to shape the final organization of the book (e.g., the inclusion of §1.4). For various ways in which they helped, I would also like to thank Mary Ellen Bock, Jerry Davis, Nels Grevstad, Wensheng Guo, Alan Karr, Youngju Kim, Ping Ma, Jerry Sacks, Jingyuan Wang, Yuedong Wang, Jeff Wu, Dong Xiang, Liqing Yan, and the classes at Michigan and Purdue. Last but not least, I would like to thank the R Core Team, for creating a most enjoyable platform for statistical computing.

# Preface

When the first edition was published a decade ago, I wrote in the Preface:

Remaining on my wish list are (i) polished, user-friendly software tools for density estimation and hazard estimation, (ii) fast computation via approximate solutions of penalized likelihood problems, and (iii) handling of parametric random effects such as those appearing in longitudinal models and hazard models with frailty.

I am happy to report that the wishes have been fulfilled, plus some more, and it is time to present an updated treatise on smoothing methods with roughness penalties.

The developments of software tools embodied in an R package `gss` have gone a long way in the past decade, with the user-interface polished, functionality expanded, and/or numerical efficiency improved from release to release. The primary objective of this new edition is to introduce extensive software illustrations to complement the theoretical and methodological discussions, so the reader not only can read about the methods but also can use them in everyday data analysis.

Newly developed theoretical, methodological, and computational techniques are integrated in a few new chapters and new sections, along with some previously omitted entries; due modifications are made in related chapters and sections to maintain coherence. Empirical studies are expanded, reorganized, and mostly rerun using the latest software.

Two appendices are also added. One appendix outlines the overall design of the R package `gss`. The other presents some conceptual critiques on a few issues concerning smoothing methods at large, which are potentially controversial.

Much of the new materials that went into this edition were taken from or inspired by collaborations or communications with Pang Du, Anouschka Foltz, Chun Han, Young-Ju Kim, Yi Lin, Ping Ma, Christophe Pouzat, Jingyuan Wang, and Tonglin Zhang, to whom I owe thanks. I can not thank enough the R Core Team, for creating and maintaining a most enjoyable platform for statistical computing.

West Lafayette, Indiana  
August 2011

Chong Gu

# Contents

<b>Preface to the First Edition</b>	<b>vii</b>
<b>Preface</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Estimation Problem and Method . . . . .	2
1.1.1 Cubic Smoothing Spline . . . . .	2
1.1.2 Penalized Likelihood Method . . . . .	4
1.2 Notation . . . . .	5
1.3 Decomposition of Multivariate Functions . . . . .	6
1.3.1 ANOVA Decomposition and Averaging Operator . . . . .	6
1.3.2 Multiway ANOVA Decomposition . . . . .	7
1.3.3 Multivariate Statistical Models . . . . .	10
1.4 Case Studies . . . . .	12
1.4.1 Water Acidity in Lakes . . . . .	12
1.4.2 AIDS Incubation . . . . .	14
1.4.3 Survival After Heart Transplant . . . . .	15
1.5 Scope . . . . .	17
1.6 Bibliographic Notes . . . . .	19
1.7 Problems . . . . .	20

<b>2</b>	<b>Model Construction</b>	<b>23</b>
2.1	Reproducing Kernel Hilbert Spaces . . . . .	24
2.1.1	Hilbert Spaces and Linear Subspaces . . . . .	24
2.1.2	Riesz Representation Theorem . . . . .	29
2.1.3	Reproducing Kernel and Non-Negative Definite Function . . . . .	29
2.2	Smoothing Splines on $\{1, \dots, K\}$ . . . . .	32
2.3	Polynomial Smoothing Splines on $[0, 1]$ . . . . .	34
2.3.1	A Reproducing Kernel in $\mathcal{C}^{(m)}[0, 1]$ . . . . .	34
2.3.2	Computation of Polynomial Smoothing Splines . . . . .	36
2.3.3	Another Reproducing Kernel in $\mathcal{C}^{(m)}[0, 1]$ . . . . .	37
2.4	Smoothing Splines on Product Domains . . . . .	40
2.4.1	Tensor Product Reproducing Kernel Hilbert Spaces . . . . .	40
2.4.2	Reproducing Kernel Hilbert Spaces on $\{1, \dots, K\}^2$ . . . . .	41
2.4.3	Reproducing Kernel Hilbert Spaces on $[0, 1]^2$ . . . . .	42
2.4.4	Reproducing Kernel Hilbert Spaces on $\{1, \dots, K\} \times [0, 1]$ . . . . .	44
2.4.5	Multiple-Term Reproducing Kernel Hilbert Spaces: General Form . . . . .	45
2.5	Bayes Model . . . . .	48
2.5.1	Shrinkage Estimates as Bayes Estimates . . . . .	48
2.5.2	Polynomial Splines as Bayes Estimates . . . . .	49
2.5.3	Smoothing Splines as Bayes Estimates . . . . .	51
2.6	Minimization of Penalized Functional . . . . .	51
2.6.1	Existence of Minimizer . . . . .	52
2.6.2	Penalized and Constrained Optimization . . . . .	53
2.7	Bibliographic Notes . . . . .	54
2.8	Problems . . . . .	56
<b>3</b>	<b>Regression with Gaussian-Type Responses</b>	<b>61</b>
3.1	Preliminaries . . . . .	62
3.2	Smoothing Parameter Selection . . . . .	64
3.2.1	Unbiased Estimate of Relative Loss . . . . .	65
3.2.2	Cross-Validation and Generalized Cross-Validation . . . . .	67
3.2.3	Restricted Maximum Likelihood Under Bayes Model . . . . .	70
3.2.4	Weighted and Replicated Data . . . . .	72
3.2.5	Empirical Performance . . . . .	74
3.3	Bayesian Confidence Intervals . . . . .	75
3.3.1	Posterior Distribution . . . . .	76
3.3.2	Confidence Intervals on Sampling Points . . . . .	78
3.3.3	Across-the-Function Coverage . . . . .	78

3.4	Computation: Generic Algorithms . . . . .	79
3.4.1	Algorithm for Fixed Smoothing Parameters . . . . .	80
3.4.2	Algorithm for Single Smoothing Parameter . . . . .	80
3.4.3	Algorithm for Multiple Smoothing Parameters . . . . .	82
3.4.4	Calculation of Posterior Variances . . . . .	84
3.5	Efficient Approximation . . . . .	85
3.5.1	Preliminaries . . . . .	85
3.5.2	Bayes Model . . . . .	86
3.5.3	Computation . . . . .	88
3.5.4	Empirical Choice of $q$ . . . . .	90
3.5.5	Numerical Accuracy . . . . .	92
3.6	Software . . . . .	93
3.6.1	RKPACK . . . . .	93
3.6.2	R Package <code>gss</code> : <code>ssanova</code> and <code>ssanova0 Suites</code> . . . . .	94
3.7	Model Checking Tools . . . . .	98
3.7.1	Cosine Diagnostics . . . . .	98
3.7.2	Examples . . . . .	99
3.7.3	Concepts and Heuristics . . . . .	103
3.8	Square Error Projection . . . . .	104
3.9	Case Studies . . . . .	106
3.9.1	Nitrogen Oxides in Engine Exhaust . . . . .	106
3.9.2	Ozone Concentration in Los Angeles Basin . . . . .	107
3.10	Computation: Special Algorithms . . . . .	111
3.10.1	Fast Algorithm for Polynomial Splines . . . . .	112
3.10.2	Iterative Algorithms and Monte Carlo Cross-Validation . . . . .	114
3.11	Bibliographic Notes . . . . .	115
3.12	Problems . . . . .	118
<b>4</b>	<b>More Splines</b> . . . . .	<b>125</b>
4.1	Partial Splines . . . . .	126
4.2	Splines on the Circle . . . . .	127
4.2.1	Periodic Polynomial Splines . . . . .	127
4.2.2	Splines as Low-Pass Filters . . . . .	128
4.2.3	More on Asymptotics of §3.2 . . . . .	130
4.3	Thin-Plate Splines . . . . .	134
4.3.1	Semi-Kernels for Thin-Plate Splines . . . . .	135
4.3.2	Reproducing Kernels for Thin-Plate Splines . . . . .	136
4.3.3	Tensor Product Splines with Thin-Plate Marginals . . . . .	139
4.3.4	Case Study: Water Acidity in Lakes . . . . .	140
4.4	Splines on the Sphere . . . . .	143
4.4.1	Spherical Harmonics . . . . .	143
4.4.2	Laplacian on the Sphere and Spherical Splines . . . . .	144

4.4.3	Reproducing Kernels in Closed Forms . . . . .	146
4.4.4	Case Study: Global Temperature Map . . . . .	147
4.5	L-Splines . . . . .	149
4.5.1	Trigonometric Splines . . . . .	150
4.5.2	Chebyshev Splines . . . . .	153
4.5.3	General Construction . . . . .	157
4.5.4	Case Study: Weight Loss of Obese Patient . . . . .	161
4.5.5	Fast Algorithm . . . . .	165
4.6	Bibliographic Notes . . . . .	166
4.7	Problems . . . . .	167
<b>5</b>	<b>Regression with Responses from Exponential Families</b>	<b>175</b>
5.1	Preliminaries . . . . .	176
5.2	Smoothing Parameter Selection . . . . .	177
5.2.1	Performance-Oriented Iteration . . . . .	178
5.2.2	Direct Cross-Validation . . . . .	181
5.3	Inferential Tools . . . . .	184
5.3.1	Approximate Bayesian Confidence Intervals . . . . .	185
5.3.2	Kullback-Leibler Projection . . . . .	186
5.4	Software, Customization, and Empirical Performance . . . . .	187
5.4.1	R Package <code>gss</code> : <code>gssanova</code> , <code>gssanova0</code> , and <code>gssanova1</code> Suites . . . . .	187
5.4.2	Binomial Family . . . . .	188
5.4.3	Poisson Family . . . . .	191
5.4.4	Gamma Family . . . . .	193
5.4.5	Inverse Gaussian Family . . . . .	196
5.4.6	Negative Binomial Family . . . . .	199
5.5	Case Studies . . . . .	202
5.5.1	Eruption Time of Old Faithful . . . . .	202
5.5.2	Spectrum of Yearly Sunspots . . . . .	203
5.5.3	Progression of Diabetic Retinopathy . . . . .	205
5.5.4	Colorectal Cancer Mortality Rate . . . . .	208
5.6	Bibliographic Notes . . . . .	210
5.7	Problems . . . . .	212
<b>6</b>	<b>Regression with Correlated Responses</b>	<b>215</b>
6.1	Models for Correlated Data . . . . .	216
6.1.1	Random Effects . . . . .	216
6.1.2	Stationary Time Series . . . . .	216
6.2	Mixed-Effect Models and Penalized Joint Likelihood . . . . .	217
6.2.1	Smoothing Matrices . . . . .	218
6.2.2	Bayes Model . . . . .	219
6.2.3	Optimality of Generalized Cross-Validation . . . . .	219
6.2.4	Empirical Performance . . . . .	221

6.2.5	Non-Gaussian Regression . . . . .	222
6.2.6	R Package <code>gss</code> : Optional Argument <code>random</code> . . . . .	222
6.3	Penalized Likelihood with Correlated Data . . . . .	223
6.3.1	Bayes Model . . . . .	223
6.3.2	Extension of Cross-Validation . . . . .	225
6.3.3	Optimality of Cross-Validation . . . . .	226
6.3.4	Empirical Performance . . . . .	228
6.3.5	R Package <code>gss</code> : <code>ssanova9</code> Suite . . . . .	230
6.4	Case Studies . . . . .	231
6.4.1	Treatment of Bacteriuria . . . . .	231
6.4.2	Ozone Concentration in Los Angeles Basin . . . . .	232
6.5	Bibliographic Notes . . . . .	233
6.6	Problems . . . . .	235
<b>7</b>	<b>Probability Density Estimation</b> . . . . .	<b>237</b>
7.1	Preliminaries . . . . .	238
7.2	Poisson Intensity . . . . .	242
7.3	Smoothing Parameter Selection . . . . .	243
7.3.1	Kullback-Leibler Loss . . . . .	243
7.3.2	Cross-Validation . . . . .	244
7.3.3	Empirical Performance . . . . .	246
7.4	Computation, Inference, and Software . . . . .	247
7.4.1	Newton Iteration . . . . .	247
7.4.2	Numerical Integration . . . . .	248
7.4.3	Kullback-Leibler Projection . . . . .	250
7.4.4	R Package <code>gss</code> : <code>ssden</code> Suite . . . . .	250
7.5	Case Studies . . . . .	253
7.5.1	Buffalo Snowfall . . . . .	253
7.5.2	Eruption Time of Old Faithful . . . . .	254
7.5.3	AIDS Incubation . . . . .	255
7.6	Biased Sampling and Random Truncation . . . . .	257
7.6.1	Biased and Truncated Samples . . . . .	257
7.6.2	Penalized Likelihood Estimation . . . . .	258
7.6.3	Empirical Performance . . . . .	260
7.6.4	R Package <code>gss</code> : <code>ssden</code> Suite . . . . .	260
7.6.5	Case Study: AIDS Incubation . . . . .	262
7.7	Conditional Densities . . . . .	263
7.7.1	Penalized Likelihood Estimation . . . . .	263
7.7.2	Empirical Performance of Cross-validation . . . . .	265
7.7.3	Kullback-Leibler Projection . . . . .	266
7.7.4	R Package <code>gss</code> : <code>sscdn</code> Suite . . . . .	266
7.7.5	Case Study: Penny Thickness . . . . .	268
7.8	Regression with Cross-Classified Responses . . . . .	269
7.8.1	Logistic Regression . . . . .	269
7.8.2	Log-Linear Regression Models . . . . .	271



7.8.3	Bayesian Confidence Intervals for $y$ -Contrasts . . . . .	271
7.8.4	Mixed-Effect Models for Correlated Data . . . . .	272
7.8.5	Empirical Performance of Cross-Validation . . . . .	273
7.8.6	R Package <code>gss: sllrm</code> Suite . . . . .	274
7.8.7	Case Study: Eyetracking Experiments . . . . .	275
7.9	Response-Based Sampling . . . . .	278
7.9.1	Response-Based Samples . . . . .	278
7.9.2	Penalized Likelihood Estimation . . . . .	279
7.10	Bibliographic Notes . . . . .	280
7.11	Problems . . . . .	282
<b>8</b>	<b>Hazard Rate Estimation</b> . . . . .	<b>285</b>
8.1	Preliminaries . . . . .	286
8.2	Smoothing Parameter Selection . . . . .	288
8.2.1	Kullback-Leibler Loss and Cross-Validation . . . . .	289
8.2.2	Empirical Performance . . . . .	291
8.3	Inference and Software . . . . .	292
8.3.1	Bayesian Confidence Intervals . . . . .	292
8.3.2	Kullback-Leibler Projection . . . . .	293
8.3.3	Frailty Models for Correlated Data . . . . .	293
8.3.4	R Package <code>gss: sshzd</code> Suite . . . . .	293
8.4	Case Studies . . . . .	295
8.4.1	Treatments of Gastric Cancer . . . . .	295
8.4.2	Survival After Heart Transplant . . . . .	297
8.5	Penalized Partial Likelihood . . . . .	299
8.5.1	Partial Likelihood and Biased Sampling . . . . .	299
8.5.2	Inference . . . . .	300
8.5.3	R Package <code>gss: sscox</code> Suite . . . . .	300
8.5.4	Case Study: Survival After Heart Transplant . . . . .	302
8.6	Models Parametric in Time . . . . .	303
8.6.1	Location-Scale Families and Accelerated Life Models . . . . .	303
8.6.2	Kullback-Leibler and Cross-Validation . . . . .	305
8.6.3	Weibull Family . . . . .	305
8.6.4	Log Normal Family . . . . .	309
8.6.5	Log Logistic Family . . . . .	311
8.6.6	Case Study: Survival After Heart Transplant . . . . .	314
8.7	Bibliographic Notes . . . . .	316
8.8	Problems . . . . .	317
<b>9</b>	<b>Asymptotic Convergence</b> . . . . .	<b>319</b>
9.1	Preliminaries . . . . .	319
9.2	Rates for Density Estimates . . . . .	322
9.2.1	Linear Approximation . . . . .	323
9.2.2	Approximation Error and Main Results . . . . .	325

9.2.3	Efficient Approximation . . . . .	327
9.2.4	Convergence Under Incorrect Model . . . . .	330
9.2.5	Estimation Under Biased Sampling . . . . .	331
9.2.6	Estimation of Conditional Density . . . . .	332
9.2.7	Estimation Under Response-Based Sampling . . . . .	332
9.3	Rates for Hazard Estimates . . . . .	333
9.3.1	Martingale Structure . . . . .	333
9.3.2	Linear Approximation . . . . .	334
9.3.3	Approximation Error and Main Results . . . . .	335
9.3.4	Efficient Approximation . . . . .	338
9.3.5	Convergence Under Incorrect Model . . . . .	341
9.4	Rates for Regression Estimates . . . . .	341
9.4.1	General Formulation . . . . .	341
9.4.2	Linear Approximation . . . . .	342
9.4.3	Approximation Error and Main Result . . . . .	343
9.4.4	Efficient Approximation . . . . .	345
9.4.5	Convergence Under Incorrect Model . . . . .	347
9.5	Bibliographic Notes . . . . .	348
9.6	Problems . . . . .	349
<b>10</b>	<b>Penalized Pseudo Likelihood</b>	<b>351</b>
10.1	Density Estimation on Product Domains . . . . .	352
10.1.1	Pseudo and Genuine Likelihoods . . . . .	352
10.1.2	Preliminaries . . . . .	353
10.1.3	Smoothing Parameter Selection . . . . .	354
10.1.4	Square Error Projection . . . . .	357
10.1.5	R Package <code>gss: ssden1</code> Suite . . . . .	358
10.1.6	Case Study: Transcription Factor Association . . . . .	359
10.2	Density Estimation: Asymptotic Convergence . . . . .	360
10.2.1	Linear Approximation . . . . .	361
10.2.2	Approximation Error and Main Results . . . . .	361
10.2.3	Efficient Approximation . . . . .	363
10.3	Conditional Density Estimation . . . . .	364
10.3.1	Preliminaries . . . . .	365
10.3.2	Smoothing Parameter Selection . . . . .	366
10.3.3	Square Error Projection . . . . .	369
10.3.4	R Package <code>gss: ssden1</code> Suite . . . . .	369
10.3.5	Case Study: Penny Thickness . . . . .	371
10.3.6	Asymptotic Convergence . . . . .	372
10.4	Hazard Estimation . . . . .	372
10.4.1	Preliminaries . . . . .	373
10.4.2	Smoothing Parameter Selection . . . . .	374
10.4.3	Inference . . . . .	375
10.4.4	R Package <code>gss: sshzd1</code> Suite . . . . .	376
10.4.5	Case Study: Survival After Heart Transplant . . . . .	378

10.5 Hazard Estimation: Asymptotic Convergence . . . . .	378
10.5.1 Linear Approximation . . . . .	379
10.5.2 Approximation Error and Main Results . . . . .	380
10.5.3 Efficient Approximation . . . . .	381
10.6 Bibliographic Notes . . . . .	383
10.7 Problems . . . . .	384
<b>A R Package gss</b>	<b>387</b>
A.1 Model Construction . . . . .	387
A.1.1 Marginal Configurations . . . . .	388
A.1.2 Construction of Interaction Terms . . . . .	389
A.1.3 Custom Types . . . . .	390
A.2 Modeling and Data Analytical Tools . . . . .	391
A.3 Numerical Engines . . . . .	393
<b>B Conceptual Critiques</b>	<b>395</b>
B.1 Model Indexing . . . . .	395
B.2 Optimal and Cross-Validation Indices . . . . .	397
B.3 Loss, Risk, and Smoothing Parameter Selection . . . . .	398
B.4 Degrees of Freedom . . . . .	400
<b>References</b>	<b>403</b>
<b>Author Index</b>	<b>417</b>
<b>Subject Index</b>	<b>421</b>

# 1

## Introduction

Data and models are two sources of information in a statistical analysis. Data carry noise but are “unbiased,” whereas models, effectively a set of constraints, help to reduce noise but are responsible for “biases.” Representing the two extremes on the spectrum of “bias-variance” trade-off are standard parametric models and constraint-free nonparametric “models” such as the empirical distribution for a probability density. In between the two extremes, there exist scores of nonparametric or semiparametric models, of which most are also known as smoothing methods. A family of such nonparametric models in a variety of stochastic settings can be derived through the penalized likelihood method, forming the subject of this book.

The general penalized likelihood method can be readily abstracted from the cubic smoothing spline as the solution to a minimization problem, and its applications in regression, density estimation, and hazard estimation set out the subject of study (§1.1). Some general notation is set in §1.2. Multivariate statistical models can often be characterized through function decompositions similar to the classical analysis of variance (ANOVA) decomposition, which we discuss in §1.3. To illustrate the potential applications of the methodology, previews of selected case studies are presented in §1.4. Brief summaries of the chapters to follow are given in §1.5.

## 1.1 Estimation Problem and Method

The problem to be addressed in this book is flexible function estimation based on stochastic data. To allow for flexibility in the estimation of  $\eta$ , say, soft constraints of the form  $J(\eta) \leq \rho$  are used in lieu of the rigid constraints of parametric models, where  $J(\eta)$  quantifies the roughness of  $\eta$  and  $\rho$  sets the allowance; an example of  $J(\eta)$  for  $\eta$  on  $[0, 1]$  is  $\int_0^1 (d^2\eta/dx^2)^2 dx$ . Solving the constrained maximum likelihood problem by the Lagrange method, one is led to the penalized likelihood method.

In what follows, a brief discussion of the cubic smoothing spline helps to motivate the idea, and a simple simulation illustrates the role of  $\rho$  through the Lagrange multiplier, better known as the smoothing parameter in the context. Following a straightforward abstraction, the penalized likelihood method is exemplified in regression, density estimation, and hazard estimation.

### 1.1.1 Cubic Smoothing Spline

Consider a regression problem  $Y_i = \eta(x_i) + \epsilon_i$ ,  $i = 1, \dots, n$ , where  $x_i \in [0, 1]$  and  $\epsilon_i \sim N(0, \sigma^2)$ . In a classical parametric regression analysis,  $\eta$  is assumed to be of form  $\eta(x, \beta)$ , known up to the parameters  $\beta$ , which are to be estimated from the data. When  $\eta(x, \beta)$  is linear in  $\beta$ , one has a standard linear model. A parametric model characterizes a set of rigid constraints on  $\eta$ . The dimension of the model space (i.e., the number of unknown parameters) is typically much smaller than the sample size  $n$ .

To avoid possible model misspecification in a parametric analysis, otherwise known as bias, an alternative approach to estimation is to allow  $\eta$  to vary in a high-dimensional (possibly infinite) function space, leading to various nonparametric or semiparametric estimation methods. A popular approach to the nonparametric estimation of  $\eta$  is via the minimization of a penalized least squares score,

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \eta(x_i))^2 + \lambda \int_0^1 \ddot{\eta}^2 dx, \quad (1.1)$$

with  $\ddot{\eta} = d^2\eta/dx^2$ , where the first term discourages the lack of fit of  $\eta$  to the data, the second term penalizes the roughness of  $\eta$ , and the smoothing parameter  $\lambda$  controls the trade-off between the two conflicting goals. The minimization of (1.1) is implicitly over functions with square integrable second derivatives. The minimizer  $\eta_\lambda$  of (1.1) is called a cubic smoothing spline. As  $\lambda \rightarrow 0$ ,  $\eta_\lambda$  approaches the minimum curvature interpolant. As  $\lambda \rightarrow \infty$ ,  $\eta_\lambda$  approaches the simple linear regression line. Note that the linear polynomials  $\{f : f = \beta_0 + \beta_1 x\}$  form the so-called null space of the roughness penalty  $\int_0^1 \ddot{f}^2 dx$ ,  $\{f : \int_0^1 \ddot{f}^2 dx = 0\}$ .

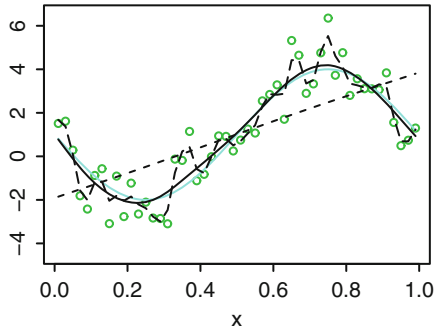


FIGURE 1.1. Cubic smoothing splines. The test function is in the *faded line* and the estimates are in the *solid*, *dashed*, and *long-dashed lines*. The data are superimposed as *circles*.

To illustrate, consider a simple simulation with  $x_i = (i - 0.5)/50$ ,  $i = 1, \dots, 50$ ,  $\eta(x) = 1 + 3 \sin(2\pi x - \pi)$ , and  $\sigma^2 = 1$ . The estimate  $\eta_\lambda$  was calculated at  $\log_{10} n\lambda = 0, -3, -6$ . Plotted in Fig. 1.1 are the test function (faded line), the estimates (solid, dashed, and long-dashed lines), and the data (circles). The rough fit corresponds to  $\log_{10} n\lambda = -6$ , the near straight line to  $\log_{10} n\lambda = 0$ , and the close fit to  $\log_{10} n\lambda = -3$ .

An alternative derivation of the cubic smoothing spline is through a constrained least squares problem, which solves

$$\min \frac{1}{n} \sum_{i=1}^n (Y_i - \eta(x_i))^2, \quad \text{subject to} \quad \int_0^1 \ddot{\eta}^2 dx \leq \rho, \quad (1.2)$$

for some  $\rho \geq 0$ . The solution to (1.2) usually falls on the boundary of the permissible region,  $\int_0^1 \ddot{\eta}^2 dx = \rho$ , and by the Lagrange method, it can be calculated as the minimizer of (1.1) with an appropriate Lagrange multiplier  $\lambda$ . Thus, up to the choices of  $\lambda$  and  $\rho$ , a penalized least squares problem with a penalty proportional to  $\int_0^1 \ddot{\eta}^2 dx$  is equivalent to a constrained least squares problem subject to a soft constraint of the form  $\int_0^1 \ddot{\eta}^2 dx \leq \rho$ ; see, e.g., Schoenberg (1964). See also §2.6.2.

Defined as the solution to a penalized optimization problem, a smoothing spline is also known as a natural spline in the numerical analysis literature. The minimizer  $\eta_\lambda$  of (1.1) is called a cubic spline because it is a piecewise cubic polynomial. It is three times differentiable, with the third derivative jumping at the knots  $\xi_1 < \xi_2 < \dots < \xi_q$ , the ordered distinctive sampling points  $x_i$ , and it is linear beyond the first knot  $\xi_1$  and the last knot  $\xi_q$ . See Schumaker (1981, Chap. 8) for a comprehensive treatment of smoothing splines from a numerical analytical perspective. See also de Boor (1978).

### 1.1.2 Penalized Likelihood Method

The cubic smoothing spline of (1.1) is a specialization of the general penalized likelihood method in univariate Gaussian regression. To estimate a function of interest  $\eta$  on a generic domain  $\mathcal{X}$  using stochastic data, one may use the minimizer of

$$L(\eta|\text{data}) + \frac{\lambda}{2}J(\eta), \quad (1.3)$$

where  $L(\eta|\text{data})$  is usually taken as the minus log likelihood of the data and  $J(f)$  is a quadratic roughness functional with a null space  $\mathcal{N}_J = \{f : J(f) = 0\}$  of low dimension; see §2.1.1 for the definition of quadratic functional. The solution of (1.3) is the maximum likelihood estimate in a model space  $\mathcal{M}_\rho = \{f : J(f) \leq \rho\}$  for some  $\rho \geq 0$ , and the smoothing parameter  $\lambda$  in (1.3) is the Lagrange multiplier. See §2.6.2 for a detailed discussion of the role of  $\lambda$  as a Lagrange multiplier.

A few examples of penalized likelihood estimation follow.

**Example 1.1 (Response data regression)** Assume

$$Y|x \sim \exp \left\{ (y\eta(x) - b(\eta(x))) / a(\phi) + c(y, \phi) \right\},$$

an exponential family density with a modeling parameter  $\eta$  and a possibly unknown nuisance parameter  $\phi$ . Observing independent data  $(x_i, Y_i)$ ,  $i = 1, \dots, n$ , the method estimates  $\eta$  via the minimization of

$$-\frac{1}{n} \sum_{i=1}^n \{Y_i \eta(x_i) - b(\eta(x_i))\} + \frac{\lambda}{2} J(\eta). \quad (1.4)$$

When the density is Gaussian, (1.4) reduces to a penalized least squares problem; see Problem 1.1. Penalized least squares regression for Gaussian-type responses is the subject of Chap. 3. Penalized likelihood regression for non-Gaussian responses will be studied in Chap. 5.  $\square$

**Example 1.2 (Density estimation)** Observing independent and identically distributed samples  $X_i$ ,  $i = 1, \dots, n$  from a probability density  $f(x)$  supported on a bounded domain  $\mathcal{X}$ , the method estimates  $f$  by  $e^{\eta} / \int_{\mathcal{X}} e^{\eta} dx$ , where  $\eta$  minimizes

$$-\frac{1}{n} \sum_{i=1}^n \left\{ \eta(X_i) - \log \int_{\mathcal{X}} e^{\eta(x)} dx \right\} + \frac{\lambda}{2} J(\eta). \quad (1.5)$$

A side condition, say  $\int_{\mathcal{X}} \eta dx = 0$ , shall be imposed on  $\eta$  for a one-to-one transform  $f \leftrightarrow e^{\eta} / \int_{\mathcal{X}} e^{\eta} dx$ . Penalized likelihood density estimation is the subject of Chap. 7.  $\square$

**Example 1.3 (Hazard estimation)** Let  $T$  be the lifetime of an item with survival function  $S(t|u) = P(T > t|u)$ , possibly dependent on a covariate  $U$ . The hazard function is defined as  $e^{\eta(t,u)} = -\partial \log S(t|u)/\partial t$ . Let  $Z$  be the left-truncation time and  $C$  be the right-censoring time, independent of  $T$  and of each other. Observing  $(U_i, Z_i, X_i, \delta_i)$ ,  $i = 1, \dots, n$ , where  $X = \min(T, C)$ ,  $\delta = I_{[T \leq C]}$ , and  $Z < X$ , the method estimates the log hazard  $\eta$  via the minimization of

$$-\frac{1}{n} \sum_{i=1}^n \left\{ \delta_i \eta(X_i, U_i) - \int_{Z_i}^{X_i} e^{\eta(t, U_i)} dt \right\} + \frac{\lambda}{2} J(\eta); \quad (1.6)$$

see Problem 1.2 for the derivation of the likelihood. Penalized likelihood hazard estimation will be studied in Chap. 8.  $\square$

The two basic components of a statistical model, the deterministic part and the stochastic part, are well separated in (1.3). The structure of the deterministic part is determined by the construction of  $J(\eta)$  for  $\eta$  on a domain  $\mathcal{X}$ , of which a comprehensive treatment is presented in Chap. 2. The stochastic part is reflected in the likelihood  $L(\eta|\text{data})$  and determines, among other things, the natural measures with which the performance of the estimate is to be assessed. The minimizer of (1.3) with a varying  $\lambda$  defines a family of estimates, and from the cubic spline simulation shown in Fig. 1.1, we have seen how differently the family members may behave. Data-driven procedures for the proper selection of the smoothing parameter are crucial to the practicability of penalized likelihood estimation, to which extensive discussion will be devoted in the settings of regression, density estimation, and hazard estimation in their respective chapters.

## 1.2 Notation

Listed below is some general notation used in this book. Context-specific or subject-specific notation may differ from that listed here, in which case every effort will be made to avoid possible confusion.

Domains are usually denoted by  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ , etc., or subscripted as  $\mathcal{X}_1, \mathcal{X}_2$ , etc. Points on domains are usually denoted by  $x \in \mathcal{X}, y \in \mathcal{Y}$ , or  $x_1, x_2, y \in \mathcal{X}$ . Points on product domains are denoted by  $x_1, x_2, y \in \mathcal{X} \times \mathcal{X}_2$ , with  $x_{1(1)}, x_{2(1)}, y_{(1)} \in \mathcal{X}_1$  and  $x_{1(2)}, x_{2(2)}, y_{(2)} \in \mathcal{X}_2$ , or by  $z = (x, y) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ , with  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Ordinary subscripts are used to denote multiple points on a domain, but *not* coordinates of a point on a product domain.

Function spaces are usually denoted by  $\mathcal{H}, \mathcal{G}$ , etc. Functions in function spaces are usually denoted by  $f, g, h \in \mathcal{H}, \eta, \phi, \xi \in \mathcal{H}$ , etc. Derivatives of a univariate function  $f(x)$  are denoted by  $\dot{f} = df/dx$ ,  $\ddot{f} = d^2f/dx^2$ ,



or by the general notation  $f^{(m)} = d^m f/dx^m$ . Derivatives of multivariate functions  $f(x_{(1)}, x_{(2)})$  on  $\mathcal{X}_1 \times \mathcal{X}_2$  or  $g(x, y)$  on  $\mathcal{X} \times \mathcal{Y}$  are denoted by  $f_{(112)}^{(3)} = \partial^3 f/\partial x_{(1)}^2 \partial x_{(2)}$ ,  $\ddot{g}_{(xy)} = \partial^2 g/\partial x \partial y$ , etc.

Matrices are denoted by the standard notation of uppercase letters. Vectors, however, are often *not* denoted by boldface letters in this book. For a point on a product domain  $\mathcal{X} = \prod_{\gamma=1}^{\Gamma} \mathcal{X}_{\gamma}$ , we write  $x = (x_{(1)}, \dots, x_{(\Gamma)})$ . For a function on domain  $\mathcal{X} = \{1, \dots, K\}$ , we write  $f = (f(1), \dots, f(K))^T$ , which may be used as a vector in standard matrix arithmetic. Boldface vectors are used where confusion may result otherwise. For example,  $\mathbf{1} = (1, \dots, 1)^T$  is used to denote a vector of all one's, and  $\mathbf{c} = (c_1, \dots, c_n)^T$  is used to encapsulate subscripted coefficients. In formulas concerning matrix computation, vectors are always set in boldface.

The standard  $O_p$ ,  $o_p$  notation is used in the asymptotic analyses of §§3.2, 4.2.3, 5.2, 6.2, 6.3, Chap. 9, §§10.2, and 10.5. If  $P(|X| > KY) \rightarrow 0$  for some constant  $K < \infty$ , we write  $X = O_p(Y)$ , and when  $P(|X| > \epsilon Y) \rightarrow 0$ ,  $\forall \epsilon > 0$ , we denote  $X = o_p(Y)$ .

### 1.3 Decomposition of Multivariate Functions

An important aspect of statistical modeling, which distinguishes it from mere function approximation, is the interpretability of the results. Of great utility are decomposition of multivariate functions similar to the classical analysis of variance (ANOVA) decomposition and the associated notions of main effect and interaction. Higher-order interactions are often excluded in practical estimation to control model complexity; the exclusion of all interactions yields the popular additive models. Selective exclusion of certain interactions also characterizes many interesting statistical models in a variety of stochastic settings.

Casting the classical one-way ANOVA decomposition as the decomposition of functions on a discrete domain, a simple averaging operator is introduced to facilitate the generalization of the notion to arbitrary domains. Multiway ANOVA decomposition is then defined, with the identifiability of the terms assured by side conditions specified through the averaging operators. Examples are given and a proposition is proved concerning certain intrinsic structures that are independent of the side conditions. The utility and implication of selective term trimming in an ANOVA decomposition are then briefly discussed in the context of regression, density estimation, and hazard estimation.

#### 1.3.1 ANOVA Decomposition and Averaging Operator

Consider a standard one-way ANOVA model,  $Y_{ij} = \mu_i + \epsilon_{ij}$ , where  $\mu_i$  are the treatment means at treatment levels  $i = 1, \dots, K$  and  $\epsilon_{ij}$  are

independent normal errors. Writing  $\mu_i = \mu + \alpha_i$ , one has the “overall mean”  $\mu$  and the treatment effect  $\alpha_i$ . The identifiability of  $\mu$  and  $\alpha_i$  are assured through a side condition, of which common choices include  $\alpha_1 = 0$  with level 1 treated as the control and  $\sum_{i=1}^K \alpha_i = 0$  with all levels treated symmetrically.

The one-way ANOVA model can be recast as  $Y_j = f(x_j) + \epsilon_j$ , where  $f(x)$  is defined on the discrete domain  $\mathcal{X} = \{1, \dots, K\}$ ; the treatment levels are now coded by  $x$  and the subscript  $j$  labels the observations. The ANOVA decomposition  $\mu_i = \mu + \alpha_i$  in the standard ANOVA model notation can be written as

$$f(x) = Af + (I - A)f = f_\emptyset + f_x,$$

where  $A$  is an averaging operator that “averages out” the argument  $x$  to return a constant function and  $I$  is the identity operator. For example, with  $Af = f(1)$ , one has  $f(x) = f(1) + \{f(x) - f(1)\}$ , corresponding to  $\alpha_1 = 0$ . With  $Af = \sum_{x=1}^K f(x)/K = \bar{f}$ , one has  $f(x) = \bar{f} + (f(x) - \bar{f})$ , corresponding to  $\sum_{i=1}^K \alpha_i = 0$ . Note that applying  $A$  to a constant function returns that constant, hence the name “averaging.” It follows that  $A(Af) = Af$ ,  $\forall f$ , or, simply,  $A^2 = A$ . The constant term  $f_\emptyset = Af$  is the “overall mean” and the term  $f_x = (I - A)f$  is the treatment effect, or “contrast,” that satisfies the side condition  $Af_x = 0$ .

On a continuous domain, say  $\mathcal{X} = [a, b]$ , one may similarly define an ANOVA decomposition  $f(x) = Af + (I - A)f = f_\emptyset + f_x$  through an appropriately defined averaging operator  $A$ , where  $f_x$  satisfies the side condition  $Af_x = 0$ . For example, with  $Af = f(a)$ , one has  $f(x) = f(a) + \{f(x) - f(a)\}$ . Similarly, with  $Af = \int_a^b f dx / (b - a)$ , one has  $f(x) = \int_a^b f dx / (b - a) + \{f(x) - \int_a^b f dx / (b - a)\}$ .

### 1.3.2 Multiway ANOVA Decomposition

Now consider a function  $f(x) = f(x_{(1)}, \dots, x_{(\Gamma)})$  on a product domain  $\mathcal{X} = \prod_{\gamma=1}^{\Gamma} \mathcal{X}_\gamma$ , where  $x_{(\gamma)} \in \mathcal{X}_\gamma$  denotes the  $\gamma$ th coordinate of  $x \in \mathcal{X}$ . Let  $A_\gamma$  be an averaging operator on  $\mathcal{X}_\gamma$  that averages out  $x_{(\gamma)}$  from the active argument list and satisfies  $A_\gamma^2 = A_\gamma$ ;  $A_\gamma f$  is constant on the  $\mathcal{X}_\gamma$  axis but not necessarily an overall constant function. An ANOVA decomposition of  $f$  can be defined as

$$f = \left\{ \prod_{\gamma=1}^{\Gamma} (I - A_\gamma + A_\gamma) \right\} f = \sum_{\mathcal{S}} \left\{ \prod_{\gamma \in \mathcal{S}} (I - A_\gamma) \prod_{\gamma \notin \mathcal{S}} A_\gamma \right\} f = \sum_{\mathcal{S}} f_{\mathcal{S}}, \quad (1.7)$$

where  $\mathcal{S} \subseteq \{1, \dots, \Gamma\}$  enlists the active arguments in  $f_{\mathcal{S}}$  and the summation is over all of the  $2^\Gamma$  subsets of  $\{1, \dots, \Gamma\}$ . The term  $f_\emptyset = \prod A_\gamma f$  is a constant, the term  $f_\gamma = f_{\{\gamma\}} = (I - A_\gamma) \prod_{\alpha \neq \gamma} A_\alpha f$  is the  $x_{(\gamma)}$  main effect,

the term  $f_{\gamma,\delta} = f_{\{\gamma,\delta\}} = (I - A_\gamma)(I - A_\delta) \prod_{\alpha \neq \gamma,\delta} A_\alpha f$  is the  $x_{(\gamma)}-x_{(\delta)}$  interaction, and so forth. The terms of such a decomposition satisfy the side conditions  $A_\gamma f_S = 0, \forall S \ni \gamma$ . The choices of  $A_\gamma$ , or the side conditions on each axes, are open to specification.

The domains  $\mathcal{X}_\gamma$  are generic in the above discussion; in particular, they can be product domains themselves. As a matter of fact, the ANOVA decomposition of (1.7) can also be defined recursively through a series of nested constructions with  $\Gamma = 2$ ; see, e.g., Problem 1.3.

The ANOVA decomposition can be built into penalized likelihood estimation through the proper construction of the roughness functional  $J(f)$ ; details are to be found in §2.4.

**Example 1.4** When  $\Gamma = 2$ ,  $\mathcal{X}_1 = \{1, \dots, K_1\}$ , and  $\mathcal{X}_2 = \{1, \dots, K_2\}$ , the decomposition reduces to a standard two-way ANOVA decomposition. With averaging operators  $A_1 f = f(1, x_{(2)})$  and  $A_2 f = f(x_{(1)}, 1)$ , one has

$$\begin{aligned} f_\emptyset &= A_1 A_2 f = f(1, 1), \\ f_1 &= (I - A_1) A_2 f = f(x_{(1)}, 1) - f(1, 1), \\ f_2 &= A_1 (I - A_2) f = f(1, x_{(2)}) - f(1, 1), \\ f_{1,2} &= (I - A_1)(I - A_2) f \\ &= f(x_{(1)}, x_{(2)}) - f(x_{(1)}, 1) - f(1, x_{(2)}) + f(1, 1). \end{aligned}$$

With  $A_\gamma f = \sum_{x_{(\gamma)}=1}^{K_\gamma} f(x_{(1)}, x_{(2)}) / K_\gamma$ ,  $\gamma = 1, 2$ , one similarly has

$$\begin{aligned} f_\emptyset &= A_1 A_2 f = f_{..}, \\ f_1 &= (I - A_1) A_2 f = f_{x_{(1)}\cdot} - f_{..}, \\ f_2 &= A_1 (I - A_2) f = f_{\cdot x_{(2)}} - f_{..}, \\ f_{1,2} &= (I - A_1)(I - A_2) f \\ &= f(x_{(1)}, x_{(2)}) - f_{x_{(1)}\cdot} - f_{\cdot x_{(2)}} + f_{..}, \end{aligned}$$

where  $f_{..} = \sum_{x_{(1)}, x_{(2)}} f(x_{(1)}, x_{(2)}) / K_1 K_2$ ,  $f_{x_{(1)}\cdot} = \sum_{x_{(2)}} f(x_{(1)}, x_{(2)}) / K_2$ , and  $f_{\cdot x_{(2)}} = \sum_{x_{(1)}} f(x_{(1)}, x_{(2)}) / K_1$ . One may also use different averaging operators on different axes; see Problem 1.4.  $\square$

**Example 1.5** Consider  $\Gamma = 2$  and  $\mathcal{X}_1 = \mathcal{X}_2 = [0, 1]$ . With  $A_1 f = f(0, x_{(2)})$  and  $A_2 f = f(x_{(1)}, 0)$ , one has

$$\begin{aligned} f_\emptyset &= A_1 A_2 f = f(0, 0), \\ f_1 &= (I - A_1) A_2 f = f(x_{(1)}, 0) - f(0, 0), \\ f_2 &= A_1 (I - A_2) f = f(0, x_{(2)}) - f(0, 0), \\ f_{1,2} &= (I - A_1)(I - A_2) f \\ &= f(x_{(1)}, x_{(2)}) - f(x_{(1)}, 0) - f(0, x_{(2)}) + f(0, 0). \end{aligned}$$

With  $A_\gamma f = \int_0^1 f dx_{(\gamma)}$ ,  $\gamma = 1, 2$ , one has

$$\begin{aligned} f_\emptyset &= A_1 A_2 f = \int_0^1 \int_0^1 f dx_{(1)} dx_{(2)}, \\ f_1 &= (I - A_1) A_2 f = \int_0^1 (f - \int_0^1 f dx_{(1)}) dx_{(2)}, \\ f_2 &= A_1 (I - A_2) f = \int_0^1 (f - \int_0^1 f dx_{(2)}) dx_{(1)}, \\ f_{1,2} &= (I - A_1)(I - A_2) f \\ &= f - \int_0^1 f dx_{(2)} - \int_0^1 f dx_{(1)} + \int_0^1 \int_0^1 f dx_{(1)} dx_{(2)}. \end{aligned}$$

Similar results with different averaging operators on different axes are also straightforward; see Problem 1.5.  $\square$

In standard ANOVA models, higher-order terms are frequently eliminated, whereas main effects and lower-order interactions are estimated from the data. One learns not to drop the  $x_{(1)}$  and  $x_{(2)}$  main effects if the  $x_{(1)}-x_{(2)}$  interaction is considered, however, and not to drop the  $x_{(1)}-x_{(2)}$  interaction when the  $x_{(1)}-x_{(2)}-x_{(3)}$  interaction is included. Although the ANOVA decomposition as defined in (1.7) obviously depends on the averaging operators  $A_\gamma$ , certain structures are independent of the particular choices of  $A_\gamma$ . Specifically, for any index set  $\mathcal{I}$ , if  $f_S = 0$ ,  $\forall S \supseteq \mathcal{I}$  with a particular set of  $A_\gamma$ , then the structure also holds for any other choices of  $A_\gamma$ , as the following proposition asserts.

**Proposition 1.1** *For any two sets of averaging operators  $A_\gamma$  and  $\tilde{A}_\gamma$  satisfying  $A_\gamma^2 = A_\gamma$  and  $\tilde{A}_\gamma^2 = \tilde{A}_\gamma$ ,  $\prod_{\gamma \in \mathcal{I}} (I - A_\gamma) f = 0$  if and only if  $\prod_{\gamma \in \mathcal{I}} (I - \tilde{A}_\gamma) f = 0$ , where  $\mathcal{I}$  is any index set.*

Note that the condition  $\prod_{\gamma \in \mathcal{I}} (I - A_\gamma) f = 0$  means that  $f_S = 0$ ,  $\forall S \supseteq \mathcal{I}$ . For example,  $(I - A_1) f = 0$  implies that all terms involving  $x_{(1)}$  vanish, and  $(I - A_1)(I - A_2) f = 0$  means that all terms involving both  $x_{(1)}$  and  $x_{(2)}$  disappear. Model structures that can be characterized through constraints of the form  $\prod_{\gamma \in \mathcal{I}} (I - A_\gamma) f = 0$  permit a term  $f_S$  only when all of its “subset terms,”  $f_{S'}$  for  $S' \subset S$ , are permitted. A simple corollary of the proposition is the obvious fact that an additive model remains an additive model regardless of the side conditions.

*Proof of Proposition 1.1:* It is easy to see that  $(I - \tilde{A}_\gamma) A_\gamma = 0$ . Suppose  $\prod_{\gamma \in \mathcal{I}} (I - A_\gamma) f = 0$  and define the ANOVA decomposition in (1.7) using  $A_\gamma$ . Now, for any nonzero term  $f_S$  in (1.7), one has  $S \not\supseteq \mathcal{I}$ , so there exists  $\gamma \in \mathcal{I}$  but  $\gamma \notin S$ , hence  $f_S = [\cdots A_\gamma \cdots] f$ . The corresponding  $(I - \tilde{A}_\gamma)$  in  $\prod_{\gamma \in \mathcal{I}} (I - \tilde{A}_\gamma)$  then annihilates the term. It follows that all nonzero ANOVA terms in (1.7) are annihilated by  $\prod_{\gamma \in \mathcal{I}} (I - \tilde{A}_\gamma)$ , so  $\prod_{\gamma \in \mathcal{I}} (I - \tilde{A}_\gamma) f = 0$ . The converse is true by symmetry.  $\square$

### 1.3.3 Multivariate Statistical Models

Many multivariate statistical models can be characterized by selective term elimination in an ANOVA decomposition. Some of such models are discussed below.

#### *Curse of Dimensionality and Additive Models*

Recall the classical ANOVA models with  $\mathcal{X}_\gamma$  discrete. In practical data analysis, one usually includes only the main effects, with the possible addition of a few lower-order interactions. Higher-order interactions are less interpretable yet more difficult to estimate, as they usually consume many more degrees of freedom than the lower-order terms. Models with only main effects included are called additive models.

The difficulty associated with function estimation in high-dimensional spaces may be perceived through the sparsity of the space. Take  $\mathcal{X}_\gamma = [0, 1]$ , for example, a  $k$ -dimensional cube with each side of length 0.5 has volume  $0.5^k$ . Assume a uniform distribution of the data and consider a piecewise constant function with jumps only possible at  $x_{(\gamma)} = 0.5$ . To estimate such a function in 1 dimension with two pieces, one has information from 50% of the data per piece, in 2 dimensions with four pieces, 25% per piece, in 3 dimensions with eight pieces, 12.5% per piece, etc. The lack of data due to the sparsity of high-dimensional spaces is often referred to as the curse of dimensionality. Alternatively, the curse of dimensionality may also be characterized by the explosive increase in the number of parameters, or the degrees of freedom, that one would need to approximate a function well in a high-dimensional space. To achieve the flexibility of a five-piece piecewise polynomial in 1 dimension, for example, one would end up with 125 pieces in 3 dimensions by taking products of the pieces in 1 dimension.

To combat the curse of dimensionality in multivariate function estimation, one needs to eliminate higher-order interactions to control model complexity. As with classical ANOVA models, additive models with the possible addition of second-order interactions are among the most popular models used in practice.

#### *Conditional Independence and Graphical Models*

To simplify notation, the marginal domains will be denoted by  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $\mathcal{Z}$ , etc., in the rest of the section instead of the subscripted  $\mathcal{X}$  used in (1.7).

Consider a probability density  $f(x)$  of a random variable  $X$  on a domain  $\mathcal{X}$ . Writing

$$f(x) = \frac{e^{\eta(x)}}{\int_{\mathcal{X}} e^{\eta(x)} dx}, \quad (1.8)$$

known as the logistic density transform, the log density  $\eta(x)$  is free of the positivity and unity constraints,  $f(x) > 0$  and  $\int_{\mathcal{X}} f(x) dx = 1$ , that  $f(x)$

must satisfy. The transform is not one-to-one, though, as  $e^{\eta(x)}/\int_{\mathcal{X}} e^{\eta(x)} dx = e^{C+\eta(x)}/\int_{\mathcal{X}} e^{C+\eta(x)} dx$  for any constant  $C$ . The transform can be made one-to-one, however, by imposing a side condition  $A_x \eta = 0$  for some averaging operator  $A_x$  on  $\mathcal{X}$ ; this can be achieved by eliminating the constant term in a one-way ANOVA decomposition  $\eta = A_x \eta + (I - A_x) \eta = \eta_0 + \eta_x$ .

For a joint density  $f(x, y)$  of random variables  $(X, Y)$  on a product domain  $\mathcal{X} \times \mathcal{Y}$ , one may write

$$f(x, y) = \frac{e^{\eta(x, y)}}{\int_{\mathcal{X}} dx \int_{\mathcal{Y}} e^{\eta(x, y)} dy} = \frac{e^{\eta_x + \eta_y + \eta_{x, y}}}{\int_{\mathcal{X}} dx \int_{\mathcal{Y}} e^{\eta_x + \eta_y + \eta_{x, y}} dy},$$

where  $\eta_x$ ,  $\eta_y$ , and  $\eta_{x, y}$  are the main effects and interaction of  $\eta(x, y)$  in an ANOVA decomposition; the constant is eliminated in the rightmost expression for a one-to-one transform. The conditional distribution of  $Y$  given  $X$  has a density

$$f(y|x) = \frac{e^{\eta(x, y)}}{\int_{\mathcal{Y}} e^{\eta(x, y)} dy} = \frac{e^{\eta_y + \eta_{x, y}}}{\int_{\mathcal{Y}} e^{\eta_y + \eta_{x, y}} dy}, \quad (1.9)$$

where the logistic conditional density transform is one-to-one only for the rightmost expression with the side conditions  $A_y(\eta_y + \eta_{x, y}) = 0$ ,  $\forall x \in \mathcal{X}$ , where  $A_y$  is the averaging operator on  $\mathcal{Y}$  that help to define the ANOVA decomposition. The independence of  $X$  and  $Y$ , denoted by  $X \perp Y$ , is characterized by  $\eta_{x, y} = 0$ , or  $(I - A_x)(I - A_y)\eta = 0$ .

The domains  $\mathcal{X}$  and  $\mathcal{Y}$  are generic in (1.9); in particular, they can be product domains themselves. Substituting  $(y, z)$  for  $y$  in (1.9), one has

$$f(y, z|x) = \frac{e^{\eta_y + \eta_z + \eta_{y, z} + \eta_{x, y} + \eta_{x, z} + \eta_{x, y, z}}}{\int_{\mathcal{Y}} dy \int_{\mathcal{Z}} e^{\eta_y + \eta_z + \eta_{y, z} + \eta_{x, y} + \eta_{x, z} + \eta_{x, y, z}} dz},$$

where  $\eta_{(y, z)}$  is expanded out as  $\eta_y + \eta_z + \eta_{y, z}$  and  $\eta_{x, (y, z)}$  is expanded out as  $\eta_{x, y} + \eta_{x, z} + \eta_{x, y, z}$ ; see Problem 1.3. The conditional independence of  $Y$  and  $Z$  given  $X$ , denoted by  $(Y \perp Z)|X$ , is characterized by  $\eta_{y, z} + \eta_{x, y, z} = 0$ , or  $(I - A_y)(I - A_z)\eta = 0$ .

Now, consider the joint density of four random variables  $(U, V, Y, Z)$ , with  $(U \perp V)|(Y, Z)$  and  $(Y \perp Z)|(U, V)$ . It can be shown that such a structure is characterized by  $\eta_{u, v} + \eta_{y, z} + \eta_{u, v, y} + \eta_{u, v, z} + \eta_{u, y, z} + \eta_{v, y, z} + \eta_{u, v, y, z} = 0$  in an ANOVA decomposition, or  $(I - A_u)(I - A_v)\eta = (I - A_y)(I - A_z)\eta = 0$ ; see Problem 1.7.

As noted above, the ANOVA decompositions in the log density  $\eta$  that characterize conditional independence structures are all of the type covered in Proposition 1.1. The elimination of lower-order terms in (1.8) and (1.9) for one-to-one transforms only serve to remove technical redundancies introduced by the ‘‘overparameterization’’ of  $f(x)$  or  $f(y|x)$  by the corresponding unrestricted  $\eta$ .