Chong Gu

# Smoothing Spline ANOVA Models

*Second Edition*

Nonparametric function estimation with stochastic data, otherwise known as smoothing, has been studied by several generations of statisticians. Assisted by the ample computing power in today's servers, desktops, and laptops, smoothing methods have been finding their ways into everyday data analysis by practitioners. While scores of methods have proved successful for univariate smoothing, ones practical in multivariate settings number far less. Smoothing spline ANOVA models are a versatile family of smoothing methods derived through roughness penalties, that are suitable for both univariate and multivariate problems.

In this book, the author presents a treatise on penalty smoothing under a unified framework. Methods are developed for (i) regression with Gaussian and non-Gaussian responses as well as with censored life time data; (ii) density and conditional density estimation under a variety of sampling schemes; and (iii) hazard rate estimation with censored life time data and covariates. The unifying themes are the general penalized likelihood method and the construction of multivariate models with built-in ANOVA decompositions. Extensive discussions are devoted to model construction, smoothing parameter selection, computation, and asymptotic convergence.

Most of the computational and data analytical tools discussed in the book are implemented in R, an open-source platform for statistical computing and graphics. Suites of functions are embodied in the R package gss, and are illustrated throughout the book using simulated and real data examples.

This monograph will be useful as a reference work for researchers in theoretical and applied statistics as well as for those in other related disciplines. It can also be used as a text for graduate level courses on the subject. Most of the materials are accessible to a second year graduate student with a good training in calculus and linear algebra and working knowledge in basic statistical inferences such as linear models and maximum likelihood estimates.

**Chong Gu** received his Ph.D. from University of Wisconsin–Madison in 1989, and has been on the faculty in Department of Statistics, Purdue University since 1990. At various times during his career, he has held visiting appointments at University of British Columbia, University of Michigan, and National Institute of Statistical Sciences.

**Statistics**

▶ springer.com

# Contents