

Statistics for Biology and Health

John O'Quigley

Proportional Hazards Regression

 Springer

John O'Quigley

Proportional Hazards Regression

Statistics for Biology and Health

Series Editors:

M. Gail

K. Krickeberg

J. Samet

A. Tsiatis

W. Wong

Statistics for Biology and Health

- Bacchieri/Cioppa*: Fundamentals of Clinical Research
- Borchers/Buckland/Zucchini*: Estimating Animal Abundance: Closed Populations
- Burzykowski/Molenberghs/Buyse*: The Evaluation of Surrogate Endpoints
- Duchateau/Janssen*: The Frailty Model
- Everitt/Rabe-Hesketh*: Analyzing Medical Data Using S-PLUS
- Ewens/Grant*: Statistical Methods in Bioinformatics: An Introduction, 2nd ed.
- Gentleman/Carey/Huber/Irizarry/Dudoit*: Bioinformatics and Computational Biology Solutions Using R and Bioconductor
- Hougaard*: Analysis of Multivariate Survival Data
- Keyfitz/Caswell*: Applied Mathematical Demography, 3rd ed.
- Klein/Moeschberger*: Survival Analysis: Techniques for Censored and Truncated Data, 2nd ed.
- Kleinbaum/Klein*: Survival Analysis: A Self-Learning Text, 2nd ed.
- Kleinbaum/Klein*: Logistic Regression: A Self-Learning Text, 2nd ed.
- Lange*: Mathematical and Statistical Methods for Genetic Analysis, 2nd ed.
- Manton/Singer/Suzman*: Forecasting the Health of Elderly Populations
- Martynussen/Scheike*: Dynamic Regression Models for Survival Data
- Moyé*: Multiple Analyses in Clinical Trials: Fundamentals for Investigators
- Nielsen*: Statistical Methods in Molecular Evolution
- O'Quigley*: Proportional Hazards Regression
- Parmigiani/Garrett/Irizarry/Zeger*: The Analysis of Gene Expression Data: Methods and Software
- Proschan/Lan Wittes*: Statistical Monitoring of Clinical Trials: A Unified Approach
- Siegmund/Yakir*: The Statistics of Gene Mapping
- Simon/Korn/McShane/Radmacher/Wright/Zhao*: Design and Analysis of DNA Microarray Investigations
- Sorensen/Gianola*: Likelihood, Bayesian, and MCMC Methods in Quantitative Genetics
- Stallard/Manton/Cohen*: Forecasting Product Liability Claims: Epidemiology and Modeling in the Manville Asbestos Case
- Sun*: The Statistical Analysis of Interval-censored Failure Time Data
- Therneau/Grambsch*: Modeling Survival Data: Extending the Cox Model
- Ting*: Dose Finding in Drug Development
- Vittinghoff/Glidden/Shiboski/McCulloch*: Regression Methods in Biostatistics: Linear, Logistic, Survival, and Repeated Measures Models
- Wu/Ma/Casella*: Statistical Genetics of Quantitative Traits: Linkage, Maps, and QTL
- Zhang/Singer*: Recursive Partitioning in the Health Sciences
- Zuur/Ieno/Smith*: Analysing Ecological Data

John O'Quigley

Proportional Hazards Regression

 Springer

John O'Quigley
Institut Curie
26 rue d'Ulm
75005 Paris, France

Series Editors

M. Gail
National Cancer Institute
Rockville, MD 20892
USA

K. Krickeberg
Le Chatelet
F-63270 Manglieu
France

J. Sarnet
Department of Epidemiology
School of Public Health
Johns Hopkins University
615 Wolfe Street
Baltimore, MD 21205-2103
USA

A. Tsiatis
Department of Statistics
North Carolina State University
Raleigh, NC 27695
USA

W. Wong
Department of Statistics
Stanford University
Stanford, CA 94305-4065
USA

ISBN 978-0-387-25148-6 e-ISBN 978-0-387-68639-4
DOI: 10.1007/978-0-387-68639-4

Library of Congress Control Number: 2008920048

© 2008 Springer Science+Business Media, LLC

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed on acid-free paper.

9 8 7 6 5 4 3 2 1

springer.com

To Máire Proinsias

Preface

Proportional hazards models and their extensions (models with time-dependent covariates, models with time dependent regression coefficients, models with random coefficients and any mixture of these) can be used to characterize just about any applied problem to which the techniques of survival analysis are appropriate. This simple observation enables us to find an elegant statistical expression for all plausible practical situations arising in the analysis of survival data. We have a single unifying framework. In consequence, a solid understanding of the framework itself offers the statistician the ability to tackle the thorniest of questions which may arise when dealing with survival data.

The main goal of this text is not to present or review the very substantial amount of research that has been carried out on proportional hazards and related models. Rather, the goal is to consider the many questions which are of interest in a regression analysis of survival data (prediction, goodness of fit, model construction, inference and interpretation in the presence of misspecified models) from the standpoint of the proportional hazards and the non-proportional hazards models.

This standpoint is essentially mathematical in that the aim is to put all of the inferential questions on a firm conceptual footing. However, unlike the current widely accepted approach based almost entirely on counting processes, stochastic integrals and the martingale central limit theorem for multivariate counting processes, we mostly work with much more classic and better known central limit theorems. In particular we appeal to theorems dealing with sums of independent but not

necessarily identically distributed univariate random variables and, of particular interest, the functional central limit theorem establishing the Brownian motion limit for standardized univariate sums. Delicate measure theoretic arguments borrowed from mathematical analysis can then be wholly avoided. Admittedly, some tricky situations can still be more readily resolved by making an appeal to the martingale central limit theorem and it may also be true that the use of martingale techniques for multivariate counting processes affords greater generality. Nonetheless, in the author's view, a very high percentage of practical problems can be tackled by making an appeal to the arguably less general but more standard and well known central limit theorems for univariate sums.

Mathematicians always strive for the greatest generality and, in the author's view - at least as far as survival analysis is concerned - this undertaking has not been without some unfortunate drawbacks. The measure theoretic underpinning of the counting processes and martingale approach is quite opaque. The subject is very difficult and while outstanding efforts have been made across the globe in leading statistics and biostatistics departments to explain the essential ideas behind the material, few would claim that students, other than the small minority already well steered in mathematical analysis, ever really fully grasp just what is going on. This is a situation that we need be concerned about. The author, having taught such courses in a number of institutions, speculates that even for the most successful students entering research careers and publishing articles in our leading journals it is not easy for them to do other than reiterate well rehearsed near inscrutable arguments. Their work, reviewed by their peers - alumni survivors of similar courses - may clear the publishing hurdle and achieve technical excellence but somehow, along the way, creativity is stifled.

In brief, there is a real danger of the subject of survival analysis sustaining itself from within and reluctant to absorb input from without. The pressure to focus so much attention on the resolution of mathematical subtleties by non-mathematicians has led us away from those areas where we have traditionally done well ... abstract modeling of medical, biological, physical and social phenomena. The technical demands of those in the area of survival analysis are such that it is becoming difficult for them to construct or challenge models via considerations other than those pertaining to the correct application of

an abstruse theory. Those not in the area, a large and diverse pool of potential contributors, will typically throw up their hands and say that they are not sufficiently comfortable with survival analysis to make criticism of substance. A somewhat ambitious goal, or hope, for this text is to help change this state of affairs.

Work on this book began during the author's thesis. I would like to acknowledge the input of Dr Salah Rashid, a visiting surgeon to the University of Leeds Medical School, for asking a lot of awkward questions to a, then, very inexperienced statistician. Among these questions were 'how much of the variation is explained by the predictors', 'why would you assume that the strength of effect remains the same through time' and 'what is the relative importance of biological measurements to clinical measurements.' I believe that I can now attempt an answer to some of these questions although I fear, taking rather longer than expected to answer the clinician's concern - in this case some twenty odd years - that my good friend Dr Rashid may have moved on to other questions. Much of my thesis was based on collaborative work with Dr Rashid and his comments and questions then, and for decades to follow, provided an invaluable source of food for thought. I share the debt we all owe to Professor Sir David Cox for his great vision and scientific imagination, making all of this work possible, but also a personal debt to Sir David for having so very kindly agreed to be the external examiner on my own Ph.d thesis and for having patiently explained issues which, alone, I was unable to resolve.

My career at the Institut National de la Santé et de la Recherche Médicale in France was made possible thanks to the unfailing support of Professor Daniel Schwartz, one of the founders of the modern theory of clinical trials and I thank him warmly for that as well as for numerous discussions on parametric survival models, especially as they relate to problems in human fertility. A number of Professor Schwartz's colleagues, Joseph Lellouch, Denis Hémon, Alfred Spira in particular, were of great assistance to me in gaining understanding of the role played by survival analysis in quantitative epidemiology. However, my good fortune did not end there and I would like to offer my warm appreciation for the support, help and advice offered by Ross Prentice in inviting me to work at the Fred Hutchinson Cancer Research Center, Seattle during the late eighties, an opportunity which brought me into contact with a remarkable number of major contributors to the area of survival analysis. Among these I would like to express my gratitude to

Ross himself alongside Norman Breslow, John Crowley, Tom Fleming, Suresh Moolgavkar, Margaret Pepe and Steve Self, all of whom showed great generosity and forbearance in discussing general concepts along with their own ideas on different aspects of survival analysis.

Competing with Seattle as the world's leader in survival analysis is the Department of Biostatistics in the Harvard School of Public Health. I was given the opportunity of spending several months there in 1999 and would like to thank Nan Laird, the then department chair, for that. Many visits and collaborations followed and, although these were not in the area of survival analysis, I took advantage of the proximity to talk to those who have left quite a mark in the field. In particular I would like to offer my thanks to Victor DeGruttola, Dave Harrington, Michael Hughes, Steve Lagakos, David Schoenfeld, L.J. Wei, Marvin Zelen, all of whom, from very demanding schedules, gave time to the exchange of ideas.

I have been very fortunate in coming into contact with quite a number of the most creative researchers in this area, so many of whom have shown such scholarly patience in explaining their own views to a keen, always enthusiastic but often slow listener. I hesitate to list them since I will surely miss out some names and then, of course, if I were to credit all of the writings which have greatly helped me, this preface would take a good third of the whole book. But let me include a special mention for Janez Stare and Ronghui Xu who will recognize within these pages much which stems from our extensive collaborations. And, as an extension to this special mention, let me thank those colleagues whose kindness, as well as exceptional talent, helped provide part of a hard-to-define support structure without which this enduring task would most likely have been abandoned many years ago. I have in mind Jacques Bénichou, Claude Chastang (whose colorful view of statistics as well as life in general is so sorely missed), Michel Chavance, Philippe Flandre, Catherine Hill, Joe Ibrahim, Richard Kay, John Kent, Susanne May, Thierry Moreau, Loki Natarajan, Fabienne Pessione, Maja Pohar, Catherine Quantin, Peter Sassieni, Michael Schemper, Martin Schumacher, Lesley Struthers and Joe Whittaker. The many students who followed my course on Survival Analysis in the Department of Mathematics at the University of California at San Diego, the course upon which the skeleton of this book ended up being based, are sincerely thanked for their enthusiasm and obstinate questioning.

Finally, although the very word statistics, let alone proportional hazards regression, would leave them quite at a loss, this work owes its greatest debt to those closest to me - my nearest and dearest - for a contribution which involved untold patience and tender indulgence. My warmest gratitude goes to them.

Paris, February 2007.

Contents

1	Introduction	1
1.1	Summary	1
1.2	Motivation	1
1.3	Objectives	6
1.4	Controversies	7
1.5	Data sets	10
1.6	Use as a graduate text	10
1.7	Exercises and class projects	11
2	Background: Probability	13
2.1	Summary	13
2.2	Motivation	14
2.3	Integration and measure	14
2.4	Random variables and probability measure	17
2.5	Distributions and densities	19
2.6	Expectation	24
2.7	Order statistics and their expectations	26
2.8	Entropy and variance	32
2.9	Approximations	36
2.10	Stochastic processes	41
2.11	Brownian motion	42
2.12	Counting processes and martingales	51
2.13	Exercises and class projects	59
3	Background: General inference	63
3.1	Summary	63
3.2	Motivation	64
3.3	Limit theorems for sums of random variables	64
3.4	Functional Central Limit Theorem	68

3.5	Empirical distribution function	71
3.6	Inference for martingales and stochastic integrals . . .	74
3.7	Estimating equations	80
3.8	Inference using resampling techniques	88
3.9	Explained variation	91
3.10	Exercises and class projects	99
4	Background: Survival analysis	103
4.1	Summary	103
4.2	Motivation	103
4.3	Basic tools	104
4.4	Some potential models	112
4.5	Censoring	120
4.6	Competing risks as a particular type of censoring . . .	124
4.7	Exercises and class projects	125
5	Marginal survival	129
5.1	Summary	129
5.2	Motivation	129
5.3	Maximum likelihood estimation	130
5.4	Empirical estimate (no censoring)	136
5.5	Empirical estimate (with censoring)	138
5.6	Exercises and class projects	148
6	Regression models and subject heterogeneity	151
6.1	Summary	151
6.2	Motivation	152
6.3	General or nonproportional hazards model	153
6.4	Proportional hazards model	154
6.5	The Cox regression model	155
6.6	Modeling multivariate problems	165
6.7	Partially proportional hazards models	172
6.8	Non proportional hazards model with intercept	183
6.9	Time-dependent covariates	186
6.10	Time-dependent covariates and non proportional hazards models	188
6.11	Proportional hazards models in epidemiology	189
6.12	Exercises and class projects	199

7	Inference: Estimating equations	203
7.1	Summary	203
7.2	Motivation	204
7.3	The observations	205
7.4	Main theorem	207
7.5	The estimating equations	219
7.6	Consistency and asymptotic normality of $\tilde{\beta}$	225
7.7	Interpretation for β^* as average effect	226
7.8	Exercises and class projects	228
8	Inference: Functions of Brownian motion	231
8.1	Summary	231
8.2	Motivation	232
8.3	Brownian motion approximations	233
8.4	Non and partially proportional hazards models	238
8.5	Tests based on functions of Brownian motion	239
8.6	Multivariate model	249
8.7	Graphical representation of regression effects	254
8.8	Operating characteristics of tests	258
8.9	Goodness-of-fit tests	260
8.10	Exercises and class projects	263
9	Inference: Likelihood	267
9.1	Summary	267
9.2	Motivation	267
9.3	Likelihood solution for parametric models	268
9.4	Likelihood solution for exponential models	270
9.5	Semi-parametric likelihood solution	275
9.6	Other likelihood expressions	280
9.7	Goodness-of-fit of likelihood estimates	287
9.8	Exercises and class projects	292
10	Inference: Stochastic integrals	295
10.1	Summary	295
10.2	Motivation	295
10.3	Counting process framework to the model	296
10.4	Some nonparametric statistics	298
10.5	Stochastic integral representation of score statistic	299
10.6	Exercises and class projects	308