# Regression Modeling with Actuarial and Financial Applications

Edward W. Frees

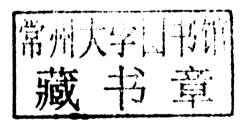
CAMBRIDGE



# Regression Modeling with Actuarial and Financial Applications

### EDWARD W. FREES

University of Wisconsin, Madison





### CAMBRIDGE UNIVERSITY PRESS

Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore, São Paulo, Delhi, Dubai, Tokyo

Cambridge University Press
32 Avenue of the Americas, New York, NY 10013-2473, USA

www.cambridge.org
Information on this title: www.cambridge.org/9780521135962

© Edward W. Frees 2010

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 2010

Printed in the United States of America

A catalog record for this publication is available from the British Library.

Library of Congress Cataloging in Publication data

Frees, Edward W.

Regression modeling with actuarial and financial applications / Edward W. (Jed) Frees.

p. cm.

Includes index.

ISBN 978-0-521-76011-9 (hardback : alk. paper) 1. Insurance – Statistical methods. 2. Finance – Statistical methods. 3. Regression analysis. I. Title.

HG8781.F67 2010 519.5'36–dc22 2009032791

ISBN 978-0-521-76011-9 Hardback ISBN 978-0-521-13596-2 Paperback

Additional resources for this publication at http://research.bus.wisc.edu/RegActuaries

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party Internet Web sites referred to in this publication and does not guarantee that any content on such Web sites is, or will remain, accurate or appropriate.

# Regression Modeling with Actuarial and Financial Applications

Statistical techniques can be used to address new situations. This is important in a rapidly evolving risk management and financial world. Analysts with a strong statistical background understand that a large data set can represent a treasure trove of information to be mined and can yield a strong competitive advantage.

This book provides budding actuaries and financial analysts with a foundation in multiple regression and time series. Readers will learn about these statistical techniques using data on the demand for insurance, lottery sales, foreign exchange rates, and other applications. Although no specific knowledge of risk management or finance is presumed, the approach introduces applications in which statistical techniques can be used to analyze real data of interest. In addition to the fundamentals, this book describes several advanced statistical topics that are particularly relevant to actuarial and financial practice, including the analysis of longitudinal, two-part (frequency/severity), and fat-tailed data.

Datasets with detailed descriptions, sample statistical software scripts in R and SAS, and tips on writing a statistical report, including sample projects, can be found on the book's Web site: http://research.bus.wisc.edu/RegActuaries.

### Christopher Daykin, Independent Consultant and Actuary Angus Macdonald, Heriot-Watt University

The International Series on Actuarial Science, published by Cambridge University Press in conjunction with the Institute of Actuaries and the Faculty of Actuaries, will contain textbooks for students taking courses in or related to actuarial science, as well as more advanced works designed for continuing professional development or for describing and synthesizing research. The series will be a vehicle for publishing books that reflect changes and developments in the curriculum, that encourage the introduction of courses on actuarial science in universities, and that show how actuarial science can be used in all areas in which there is long-term financial risk.

There is an old saying, attributed to Sir Issac Newton:

"If I have seen far, it is by standing on the shoulders of giants."

I dedicate this book to the memory of two giants who helped me, and everyone who knew them, see farther and live better lives:

James C. Hickman and Joseph P. Sullivan.

### **Preface**

Actuaries and other financial analysts quantify situations using data — we are "numbers" people. Many of our approaches and models are stylized, based on years of experience and investigations performed by legions of analysts. However, the financial and risk management world evolves rapidly. Many analysts are confronted with new situations in which tried-and-true methods simply do not work. This is where a toolkit like regression analysis comes in.

Regression is the study of relationships among variables. It is a generic statistics discipline that is not restricted to the financial world – it has applications in the fields of social, biological, and physical sciences. You can use regression techniques to investigate large and complex data sets. To familiarize you with regression, this book explores many examples and data sets based on actuarial and financial applications. This is not to say that you will not encounter applications outside of the financial world (e.g., an actuary may need to understand the latest scientific evidence on genetic testing for underwriting purposes). However, as you become acquainted with this toolkit, you will see how regression can be applied in many (and sometimes new) situations.

### Who Is This Book For?

This book is written for financial analysts who face uncertain events and wish to quantify the events using empirical information. No industry knowledge is assumed, although readers will find the reading much easier if they have an interest in the applications discussed here! This book is designed for students who are just being introduced to the field as well as industry analysts who would like to brush up on old techniques and (for the later chapters) get an introduction to new developments.

To read this book, I assume knowledge comparable to a one-semester introduction to probability and statistics – Appendix A1 provides a brief review to brush up if you are rusty. Actuarial students in North America will have a one-year introduction to probability and statistics – this type of introduction will help readers grasp concepts more quickly than a one-semester background. Finally, readers will find matrix, or linear, algebra helpful though not a prerequisite for reading this text.

Different readers are interested in understanding statistics at different levels. This book is written to accommodate the "armchair actuary," that is, one who passively reads and does not get involved by attempting the exercises in the text. Consider an analogy to football or any other game. Just like the armchair quarterback of football,

xiv Preface

there is a great deal that you can learn about the game just by watching. However, if you want to sharpen your skills, you have to go out and play the game. If you do the exercises or reproduce the statistical analyses in the text, you will become a better player. Still, this text interweaves examples with the basic principles. Thus, even the armchair actuary can obtain a solid understanding of regression techniques through this text.

### What Is This Book About?

The table of contents provides an overview of the topics covered, which are organized into four parts. The first part introduces linear regression. This is the core material of the book, with refreshers on mathematical statistics, distributions, and matrix algebra woven in as needed.

The second part is devoted to topics in time series. Why integrate time series topics into a regression book? The reasons are simple, yet compelling: most accounting, financial, and economic data become available over time. Although cross-sectional inferences are useful, business decisions need to be made in real time with currently available data. Chapters 7–10 introduce time series techniques that can be readily accomplished using regression tools (and there are many).

Nonlinear regression is the subject of the third part. Many modern-day predictive modeling tools are based on nonlinear regression – these are the workhorses of statistical shops in the financial and risk management industry.

The fourth part concerns actuarial applications, topics that I have found relevant in my research and consulting work in financial risk management. The first four chapters of this part consist of variations of regression models that are particularly useful in risk management. The last two chapters focus on communications, specifically report writing and designing graphs. Communicating information is an important aspect of every technical discipline, and statistics is certainly no exception.

### **How Does This Book Deliver Its Message?**

### Chapter Development

Each chapter has several examples interwoven with theory. In chapters in which a model is introduced, I begin with an example and discuss the data analysis without regard to the theory. This analysis is presented at an intuitive level, without reference to a specific model. This is straightforward, because it amounts to little more than curve fitting. The goal is to have students summarize data sensibly without having the notion of a model obscure good data analysis. Then, an introduction to the theory is provided in the context of the introductory example. One or more additional examples follow that reinforce the theory already introduced and provide a context for explaining additional theory. In Chapters 5 and 6, which do introduce not models but rather techniques for analysis, I begin with an introduction of the technique. This introduction is then followed by an example that reinforces

Preface xv

the explanation. In this way, the data analysis can be easily omitted without loss of continuity, if time is a concern.

### Real Data

Many of the exercises ask the reader to work with real data. The need for working with real data is well documented; for example, see Hogg (1972) or Singer and Willett (1990). Some criteria of Singer and Willett for judging a good data set include authenticity, availability of background information, interest and relevance to substantive learning, and availability of elements with which readers can identify. Of course, there are some important disadvantages to working with real data. Data sets can quickly become outdated. Further, the ideal data set for illustrating a specific statistical issue is difficult to find. This is because with real data, almost by definition, several issues occur simultaneously. This makes it difficult to isolate a specific aspect. I particularly enjoy working with large data sets. The larger the data set, the greater the need for statistics to summarize the information content.

The larger the data set, the greater the need for statistics to summarize the information content.

### Statistical Software and Data

My goal in writing this text is to reach a broad group of students and industry analysts. Thus, to avoid excluding large segments, I chose not to integrate any specific statistical software package into the text. Nonetheless, because of the applications orientation, it is critical that the methodology presented be easily accomplished using readily available packages. For the course taught at the University of Wisconsin, I use the statistical packages SAS and R. On the book's Web site, at

http://research.bus.wisc.edu/RegActuaries,

users will find scripts written in SAS and R for the analyses presented in the text. The data are available in text format, allowing readers to employ any statistical packages that they wish. When you see a display such as this in the margin, you will also be able to find this data set (*TermLife*) on the book's Web site.

® EMPIRICAL Filename is "TermLife"

### **Technical Supplements**

The technical supplements reinforce and extend the results in the main body of the text by giving a more formal, mathematical treatment of the material. This treatment is, in fact, a supplement because the applications and examples are described in the main body of the text. For readers with sufficient mathematical background, the supplements provide additional material that is useful in communicating to technical audiences. The technical supplements provide a deeper, and broader, coverage of applied regression analysis.

I believe that analysts should have an idea of "what is going on under the hood," or "how the engine works." Most of these topics will be omitted from the first reading of the material. However, as you work with regression, you will be confronted with questions such as, "Why?" and you will need to get into the details

xvi Preface

to see exactly how a certain technique works. Further, the technical supplements provide a menu of optional items that an instructor may wish to cover.

### Suggested Courses

There is a wide variety of topics that can go into a regression course. Here are some suggested courses. The course that I teach at the University of Wisconsin is the first on the list in the following table.

Audience	Nature of Course	Suggested Chapters
One-year background in probability and statistics	Survey of regression and time series models	Chapters 1–8, 11–13, 20–21, main body of text only
One-year background in probability and statistics	Regression and time series models	Chapters 1–8, 20–21, selected portions of technical supplements
One-year background in probability and statistics	Regression modeling	Chapters 1–6, 11–13, 20–21, selected portions of technica supplements
Background in statistics and linear regression	Actuarial regression models	Chapters 10–21, selected portions of technical supplements

In addition to the previously suggested courses, this book is designed as supplemental reading for a time series course as well as a reference book for industry analysts. My hope is that college students who use the beginning parts of the book in their university courses will find the later chapters helpful in their industry positions. In this way I hope to promote lifelong learning!

### Acknowledgments

It is appropriate to begin the acknowledgment section by thanking the students in the actuarial program here at the University of Wisconsin; students are important partners in the knowledge creation and dissemination business at universities. Through their questions and feedback, I have learned a tremendous amount over the years. I have also benefited from excellent assistance from those who have helped me pull together all the pieces for this book, specifically Missy Pinney, Peng Shi, Yunjie (Winnie) Sun, and Ziyan Xie.

I have enjoyed working with several former students and colleagues on regression problems in recent years, including Katrien Antonio, Jie Gao, Paul Johnson, Margie Rosenberg, Jiafeng Sun, Emil Valdez, and Ping Wang. Their contributions are reflected indirectly throughout the text. Because of my long association with the University of Wisconsin–Madison, I am reluctant to go back further in time and provide a longer list, for fear of missing important individuals. I have also been fortunate to have a more recent association with the Insurance Services Office (ISO). Colleagues at ISO have provided me with important insights into applications.

*Preface* xvii

Through this text that features applications of regression into actuarial and financial industry problems, I hope to encourage the fostering of additional partnerships between academia and industry.

I am pleased to acknowledge detailed reviews that I received from my colleagues Tim Welnetz and Margie Rosenberg. I also wish to thank Bob Miller for permission to include our joint work on designing effective graphs in Chapter 21. Bob has taught me a lot about regression over the years.

Moreover, I am happy to acknowledge financial support through the Assurant Health Professorship in Actuarial Science at the University of Wisconsin–Madison.

Saving the most important for last, I thank my family for their support. Ten thousand thanks to my mother, Mary; my brothers Randy, Guy, and Joe; my wife, Deirdre; and our sons, Nathan and Adam.

# Contents

Pr	eface		page xiii
1	Regi	ression and the Normal Distribution	1
	1.1	What Is Regression Analysis?	1
	1.2	Fitting Data to a Normal Distribution	3
	1.3	Power Transforms	7
	1.4	Sampling and the Role of Normality	8
	1.5	Regression and Sampling Designs	10
	1.6	Actuarial Applications of Regression	12
	1.7	Further Reading and References	13
	1.8	Exercises	14
	1.9	Technical Supplement – Central Limit Theorem	18
Pa	rt I	Linear Regression	
2	Basi	c Linear Regression	23
	2.1	Correlations and Least Squares	23
	2.2	Basic Linear Regression Model	29
	2.3	Is the Model Useful? Some Basic Summary Measures	32
	2.4	Properties of Regression Coefficient Estimators	35
	2.5	Statistical Inference	37
	2.6	Building a Better Model: Residual Analysis	41
	2.7	Application: Capital Asset Pricing Model	46
	2.8	Illustrative Regression Computer Output	51
	2.9	Further Reading and References	54
		Exercises	54
	2.11	Technical Supplement – Elements of Matrix Algebra	62
3		tiple Linear Regression – I	70
	3.1	Method of Least Squares	70
	3.2	Linear Regression Model and Properties of Estimators	76
	3.3	Estimation and Goodness of Fit	81
	3.4	Statistical Inference for a Single Coefficient	85
	3.5	Some Special Explanatory Variables	92
	3.6	Further Reading and References	100
	3.7	Exercises	101

viii Contents

4	Mul	tiple Linear Regression – II	107
	4.1	The Role of Binary Variables	107
	4.2	Statistical Inference for Several Coefficients	113
	4.3	One Factor ANOVA Model	120
	4.4	Combining Categorical and Continuous Explanatory Variables	126
	4.5	Further Reading and References	133
	4.6	Exercises	133
	4.7	Technical Supplement – Matrix Expressions	138
5	Vari	able Selection	148
	5.1	An Iterative Approach to Data Analysis and Modeling	148
	5.2	Automatic Variable Selection Procedures	149
	5.3	Residual Analysis	153
	5.4	Influential Points	160
	5.5	Collinearity	165
	5.6	Selection Criteria	171
	5.7	Heteroscedasticity	175
	5.8	Further Reading and References	179
	5.9	Exercises	180
	5.10	Technical Supplements for Chapter 5	182
6	Inte	rpreting Regression Results	189
	6.1	What the Modeling Process Tells Us	190
	6.2	The Importance of Variable Selection	196
	6.3	The Importance of Data Collection	198
	6.4	Missing Data Models	205
	6.5	Application: Risk Managers' Cost-Effectiveness	209
	6.6	Further Reading and References	218
	6.7	Exercises	219
	6.8	Technical Supplements for Chapter 6	222
Pa	rt II	Topics in Time Series	
7	Mod	leling Trends	227
	7.1	Introduction	227
	7.2	Fitting Trends in Time	229
	7.3	Stationarity and Random Walk Models	236
	7.4	Inference Using Random Walk Models	238
	7.5	Filtering to Achieve Stationarity	243
	7.6	Forecast Evaluation	245
	7.7	Further Reading and References	248
	7.8	Exercises	249
8	Auto	ocorrelations and Autoregressive Models	251
	8.1	Autocorrelations	251
	8.2	Autoregressive Models of Order One	254

Contents ix

	8.3	Estimation and Diagnostic Checking	256
	8.4	Smoothing and Prediction	258
	8.5	Box-Jenkins Modeling and Forecasting	260
	8.6	Application: Hong Kong Exchange Rates	265
	8.7	Further Reading and References	269
	8.8	Exercises	270
9	Fore	casting and Time Series Models	273
	9.1	Smoothing with Moving Averages	273
	9.2	Exponential Smoothing	275
	9.3	Seasonal Time Series Models	278
	9.4	Unit Root Tests	284
	9.5	ARCH/GARCH Models	285
	9.6	Further Reading and References	288
10	Long	gitudinal and Panel Data Models	289
		What Are Longitudinal and Panel Data?	289
		Visualizing Longitudinal and Panel Data	291
		Basic Fixed Effects Models	293
		Extended Fixed Effects Models	296
		Random Effects Models	299
	10.6	Further Reading and References	301
Par	t III	Topics in Nonlinear Regression	
11	Cate	gorical Dependent Variables	305
	11.1	Binary Dependent Variables	305
		Logistic and Probit Regression Models	307
		Inference for Logistic and Probit Regression Models	312
		Application: Medical Expenditures	315
		Nominal Dependent Variables	318
		Ordinal Dependent Variables	325
		Further Reading and References	328
		Exercises	329
	11.9	Technical Supplements – Likelihood-Based Inference	337
12		nt Dependent Variables	343
		Poisson Regression	343
		Application: Singapore Automobile Insurance	348
		Overdispersion and Negative Binomial Models	352
		Other Count Models	354
		Further Reading and References	359
	12.6	Exercises	360
13		eralized Linear Models	362
		Introduction	362
	13.2	GLM Model	364

X	Contents

	13.3 Estimation	367
	13.4 Application: Medical Expenditures	371
	13.5 Residuals	374
	13.6 Tweedie Distribution	375
	13.7 Further Reading and References	376
	13.8 Exercises	377
	13.9 Technical Supplements – Exponential Family	378
14	Survival Models	383
	14.1 Introduction	383
	14.2 Censoring and Truncation	385
	14.3 Accelerated Failure Time Model	390
	14.4 Proportional Hazards Model	392
	14.5 Recurrent Events	395
	14.6 Further Reading and References	397
15	Miscellaneous Regression Topics	399
	15.1 Mixed Linear Models	399
	15.2 Bayesian Regression	403
	15.3 Density Estimation and Scatterplot Smoothing	406
	15.4 Generalized Additive Models	409
	15.5 Bootstrapping	410
		412
	15.6 Further Reading and References	412
Par	t IV Actuarial Applications	412
Par 16		417
	t IV Actuarial Applications	
	et IV Actuarial Applications  Frequency-Severity Models	417
	rt IV Actuarial Applications Frequency-Severity Models 16.1 Introduction	<b>417</b> 417
	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model	<b>417</b> 417 418
	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model	417 417 418 421 424 427
	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References	417 417 418 421 424 427 429
	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model	417 417 418 421 424 427
	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References	417 417 418 421 424 427 429
16	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References 16.7 Exercises	417 418 421 424 427 429 432
16	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References 16.7 Exercises  Fat-Tailed Regression Models 17.1 Introduction 17.2 Transformations	417 418 421 424 427 429 432 433 433
16	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References 16.7 Exercises  Fat-Tailed Regression Models 17.1 Introduction 17.2 Transformations 17.3 Generalized Linear Models	417 418 421 424 427 429 432 433 433 434 437
16	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References 16.7 Exercises  Fat-Tailed Regression Models 17.1 Introduction 17.2 Transformations 17.3 Generalized Linear Models 17.4 Generalized Distributions	417 418 421 424 427 429 432 433 433 434 437 442
16	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References 16.7 Exercises  Fat-Tailed Regression Models 17.1 Introduction 17.2 Transformations 17.3 Generalized Linear Models 17.4 Generalized Distributions 17.5 Quantile Regression	417 418 421 424 427 429 432 433 433 434 437 442
16	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References 16.7 Exercises  Fat-Tailed Regression Models 17.1 Introduction 17.2 Transformations 17.3 Generalized Linear Models 17.4 Generalized Distributions 17.5 Quantile Regression 17.6 Extreme Value Models	417 418 421 424 427 429 432 433 433 434 437 442 446 448
16	Frequency-Severity Models 16.1 Introduction 16.2 Tobit Model 16.3 Application: Medical Expenditures 16.4 Two-Part Model 16.5 Aggregate Loss Model 16.6 Further Reading and References 16.7 Exercises  Fat-Tailed Regression Models 17.1 Introduction 17.2 Transformations 17.3 Generalized Linear Models 17.4 Generalized Distributions 17.5 Quantile Regression	417 418 421 424 427 429 432 433 433 434 437 442

*Contents* xi

18	Credibility and Bonus-Malus 18.1 Risk Classification and Experience Rating 18.2 Credibility 18.3 Credibility and Regression 18.4 Bonus-Malus 18.5 Further Reading and References	452 452 453 458 464 465
19	Claims Triangles 19.1 Introduction 19.2 Regression Using Functions of Time as Explanatory Variables 19.3 Using Past Developments 19.4 Further Reading and References 19.5 Exercises	467 467 471 475 477 478
20	Report Writing: Communicating Data Analysis Results 20.1 Overview 20.2 Methods for Communicating Data 20.3 How to Organize 20.4 Further Suggestions for Report Writing 20.5 Case Study: Swedish Automobile Claims 20.6 Further Reading and References 20.7 Exercises	481 481 482 486 490 491 503 504
21	Designing Effective Graphs 21.1 Introduction 21.2 Graphic Design Choices Make a Difference 21.3 Design Guidelines 21.4 Empirical Foundations for Guidelines 21.5 Concluding Remarks 21.6 Further Reading and References	505 506 508 513 520 526 526
Brie	ef Answers to Selected Exercises	529
App	endix 1: Basic Statistical Inference	547
App	endix 2: Matrix Algebra	551
App	endix 3: Probability Tables	554
Inde	ex	559

## Regression and the Normal Distribution

Chapter Preview. Regression analysis is a statistical method that is widely used in many fields of study, with actuarial science being no exception. This chapter provides an introduction to the role of the normal distribution in regression, the use of logarithmic transformations in specifying regression relationships, and the sampling basis that is critical for inferring regression results to broad populations of interest.

### 1.1 What Is Regression Analysis?

Statistics is about data. As a discipline, it is about the collection, summarization, and analysis of data to make statements about the real world. When analysts collect data, they are really collecting information that is quantified, that is, transformed to a numerical scale. There are easy, well-understood rules for reducing the data, through either numerical or graphical summary measures. These summary measures can then be linked to a theoretical representation, or model, of the data. With a model that is calibrated by data, statements about the world can be made.

Statistical methods have had a major impact on several fields of study:

- In the area of data collection, the careful design of sample surveys is crucial to market research groups and to the auditing procedures of accounting firms.
- Experimental design is a subdiscipline devoted to data collection. The focus of experimental design is on constructing methods of data collection that will extract information in the most efficient way possible. This is especially important in fields such as agriculture and engineering where each observation is expensive, possibly costing millions of dollars.
- Other applied statistical methods focus on managing and predicting data.
   Process control deals with monitoring a process over time and deciding when intervention is most fruitful. Process control helps manage the quality of goods produced by manufacturers.
- Forecasting is about extrapolating a process into the future, whether it be sales of a product or movements of an interest rate.

Regression analysis is a statistical method used to analyze data. As we will see, the distinguishing feature of this method is the ability to make statements about

Statistics is about the collection, summarization, and analysis of data to make statements about the real world.