

Ezra Hauer

The Art of Regression Modeling in Road Safety



 Springer

The Art of Regression Modeling in Road Safety

Ezra Hauer

The Art of Regression Modeling in Road Safety

 Springer

Ezra Hauer
Professor Emeritus
University of Toronto
Toronto, ON, Canada

Additional material to this book can be downloaded from <http://extras.springer.com>

ISBN 978-3-319-12528-2 ISBN 978-3-319-12529-9 (eBook)
DOI 10.1007/978-3-319-12529-9
Springer Cham Heidelberg New York Dordrecht London

Library of Congress Control Number: 2014953294

© Springer International Publishing Switzerland 2015

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)

To Gordon Newell, mentor, in gratitude.

Preface

This book has two objectives. The first is to teach how to fit a multivariable statistical model to cross-sectional safety data using a simple spreadsheet. The second is to promote the understanding that is at the core of good modeling.

These twin objectives determine the flow and the structure of the book. After some preliminaries, a real data set is introduced. From here on and all the way to the last chapter, the same data are used to gradually build up a regression model. Along the way there are the “what and how” sections: what is an exploratory data analysis, how to use pivot tables, what is a curve-fitting spreadsheet and how to build one, how to use the Solver for parameter estimation, how to examine the quality of a fit by a CURE plot, what functions look like, etc. These are in support of the first objective, that of teaching how to fit a model to data. Interspersed between these are sections of a reflective nature. These support the “understanding” objective and speak about the “why, when, and whether” issues of modeling: why do we need to curve-fit, whether models be used in a cause–effect manner, when should a variable be added to the model equation, why it is important to know what function hides behind the data, etc.

Data about accidents, about the road, and about the traffic on it are routinely collected and maintained in databases. We know where reported accidents occurred and much about their circumstances. We also have information about many features of the road (grade, curve radius, lane width, speed limit, parking control, etc.) and of its traffic (daily volumes, percent of trucks, etc.). Inasmuch as these data pertain to what is observed to exist on a cross section of “units” (road segments, intersections, etc.) they are called “observational cross-sectional” data.

To be of use for evidence-based road safety management, the data need to be summarized and cast into the form of statistical models. Models serve three purposes:

1. To estimate how (un)safe are certain road segments, intersections, ramps, crossings, etc. thereby determining the size of the road safety problem that could perhaps be altered by interventions and design changes.

2. To estimate by how much safety has changed following an intervention or a design change.
3. To estimate by how much safety might be altered due to a change in some variable value.

The book focuses on the development of regression models for purposes 1 and 2. While the use of regression models for purpose 3 is commonplace, the trustworthiness of the result is in doubt. Still, it is possible that even this vexed purpose will be well served by the approach to modeling that is advocated here.

Who will use this book? Perhaps they will be graduate students with interest in road safety, perhaps professionals with responsibilities in data analysis, or perhaps others. I do not know how much math, probability, and statistics I can rely on. Some such knowledge is required by the very nature of the subject matter. There will be parts which some readers may find taxing. The hope is that judicious skipping will make the book accessible and useful for a variety of audiences. I tried to make the narration succinct; diversions, elaborations, and detail were relegated to footnotes. A glossary is provided to serve as a refresher of notation and acronyms; the index at the end is for finding the page on which a topic or concept is first introduced.

The book evolved from lecture notes used in a series of hands-on workshops and is richly laced with data-based illustrations, tables, and figures. The supporting materials are available for downloading. To access and download them, go to <http://extras.springer.com/> and enter the ISBN of this book. The ISBN (International Standard Book Number) is found just after the title page. The information is in five folders: Data, PowerPoint presentations, Problems, Solutions, and Spreadsheets. Together with the book these materials will be of interest to the reader, student, and instructor.

The book, of course, can be read as is. However, unlike textbooks in the past, it can also be used actively and creatively. The reader is invited to use the downloadable materials to see how results were obtained; to modify, expand, and enrich the analysis; and to use the same spreadsheets with other data.

Modeling in this book is built around the use of Excel spreadsheets and readers are assumed to have facility in their use. Information about less commonly known spreadsheet functionalities is provided where needed. Of course there are specialized and sophisticated statistical software packages which, once acquired and mastered, will do a good, perhaps a superior, job of model development. However, I find that the spreadsheet provides all the essentials; it makes for intimate contact with the data, it has adequate and flexible visualization, the “pivot tables” serve for exploration, an optimization tool does the curve-fitting, and it is a hospitable environment for writing custom pieces of code.

Model development is often presented as if it was a nearly algorithmic sequence of steps, an ordered progression of activities from “Start” to “End.” In my opinion such an approach tends to produce inferior results. It is better to think of model development as detective work with clues embedded in data. Like in a game of snakes and ladders, there are advances and setbacks whereby the modeler gradually moves towards a satisfactory outcome. Such work is well served by the atmosphere

of a spreadsheet. For all these reasons, the spreadsheet is my environment of choice for both instruction and creative modeling.

The modeling approach described in these pages may be thought old-fashioned, perhaps unsophisticated. Emphasis is on what is of essence. Papers describing novel modeling ideas and newfangled statistical techniques are being published daily and I make no attempt to capture the latest. In defense I only say that the quality of a meal depends more on the skill of the cook and the time spent on its preparation than on the modernity of the food processor. As has been said: "... second-rate minds grappling with first-rate problems can teach you more than first-rate minds lost in the shrubbery." (Lilla, 2013, p. xii).

Toronto, ON, Canada

Ezra Hauer

Reference

Lilla M in Foreword for Berlin I (2013) *Against the current*, 2nd ed. Princeton University Press, Princeton

Contents

1	What Is What	1
1.1	Units and Their Safety Property	1
1.2	Safety, Traits, and Populations	3
1.3	What $\hat{E}\{\mu\}$ and $\hat{\sigma}\{\mu\}$ Are Needed for	6
1.4	How $\hat{E}\{\mu\}$ and $\hat{\sigma}\{\mu\}$ Are Used: Numerical Examples	7
1.4.1	Data for Two Populations	8
1.4.2	Estimating $E\{\mu\}$ and $\sigma\{\mu\}$	8
1.4.3	How Many High- μ Units Are There?	11
1.4.4	The Performance of a Screen	13
1.4.5	Estimating the μ of a Unit	15
1.4.6	Is the Gamma Assumption Sensible?	16
1.5	The Chosen Perspective	18
1.6	Summary	19
	References	19
2	A Safety Performance Function for Real Populations	21
2.1	The Origin	21
2.2	The Estimate of $E\{\mu\}$	22
2.3	The Estimate of $\sigma\{\mu\}$	24
2.4	The Two σ 's; Homogeneity Versus Accuracy	25
2.5	Summary	28
	References	28
3	Exploratory Data Analysis	29
3.1	Introduction	29
3.2	The Data	31
3.3	The Pivot Table	33
3.4	Pausing for Reflection	38
3.5	Visualization	40
3.6	Terrain	43
3.7	Summary	44
	References	45

4	Curve-Fitting	47
4.1	Why Do We Need to Curve-Fit?	47
4.2	There is No Free Lunch	50
4.3	Kernel Regression	52
4.3.1	Bandwidth and Goodness of Fit	54
4.3.2	Adding a Variable	56
4.4	Summary	58
	References	59
5	Preparing for Parametric Curve-Fitting: The “Solver”	61
5.1	Optimization in Modeling	61
5.2	Using the Solver to Find Minima and Maxima	62
5.3	Solver for Curve-Fitting: An Example	66
5.4	Initial Guess and Parameter Scaling	69
5.5	Summary	70
	References	70
6	A First Parametric SPF	71
6.1	The Approach to Parametric SPF Modeling	71
6.2	A Simple Parametric SPF	75
6.3	Preparing and Using the First Curve-Fitting Spreadsheet	75
6.4	Modifying the Objective Function	77
6.5	Estimating $\sigma\{\mu\}$	79
6.6	The Accuracy of Parameter Estimates	81
6.6.1	The Statistical Inaccuracy of β_1	82
6.6.2	The Incompleteness of “Statistical Inaccuracy”	83
6.7	Regression, Design Choices, Interventions, and Safety Effect	84
6.7.1	A Road Design Example	85
6.7.2	A Speed-and-Safety Example	87
6.7.3	A Generalization	88
6.7.4	The Debate	89
6.8	Summary	95
	References	96
7	Which Fit Is Fitter	99
7.1	Goodness of Fit	99
7.2	The CURE Plot	101
7.3	The Bias-in-Fit	104
7.4	Leveling the Playing Field	106
7.5	When Is a CURE Plot Good Enough?	107
7.6	Comparing CURE Plots	109
7.7	Summary	110
	References	111

8	What to Optimize?	113
8.1	Introduction	113
8.2	Likelihood	115
8.2.1	The Parameter Behind Poisson Accident Counts	117
8.2.2	The Parameters Behind the NB Distribution	119
8.3	A Few Likelihood Functions	121
8.3.1	The Poisson Likelihood Function	121
8.3.2	The Negative Binomial Likelihood Function	124
8.3.3	The Negative Multinomial Likelihood Function	125
8.4	Alternative Objective Functions	126
8.5	Summary	131
	References	132
9	Adding Variables	135
9.1	When to Add a Variable	135
9.1.1	The Necessary Conditions	136
9.1.2	The Sufficient Condition	139
9.2	The Variable Introduction EDA: Is AADT Safety Related?	141
9.3	How to Add a Variable to the C-F Spreadsheet	144
9.4	The Omitted Variable Bias	145
9.5	A Few CURE Plots	148
9.6	Adding Variables: Terrain	150
9.7	Panel Data and the NM Likelihood	152
9.8	Panel Data and Alternative Objective Functions	155
9.9	Adding Another Variable: Year	158
9.10	Summary	159
	References	161
10	Choosing the Function Behind the Data	163
10.1	The Holy Grail	163
10.2	The Elusive $f()$: A Story with Morals	164
10.3	Enroute to the Multiplicative Model Equation	168
10.4	Trying for a Better Fit	170
10.4.1	Remedy I: A Bump Function for Segment Length	170
10.4.2	Remedy II: Alternative Functions	173
10.5	What Equations Look Like	174
10.6	Trying Various Functions	177
10.7	Parameter Proliferation	180
10.8	Options and Choices: Terrain Revisited	181
10.8.1	Fitting Separate SPFs	181
10.8.2	Making β_{Terrain} into a Function of Other Predictor Variables	185
10.9	Interaction	186
10.10	Summary	190
	References	191

11	Accuracies	193
11.1	Considerations	193
11.2	The Simulation Idea	195
11.3	The Idea Executed	196
11.3.1	Determining Standard Errors	197
11.3.2	How Accuracy Is Affected by the Addition of the Terrain Variable	198
11.3.3	How Accuracy Is Affected by the AADT “Error in Variables”	199
11.3.4	Study Design	200
11.4	Summary	201
	References	202
12	Closure	203
	Appendices	205
	Appendix A: Accident Counts on a Unit: The Poisson Assumption	205
	Appendix B: The Poisson Likelihood Function	207
	Appendix C: The Variance of μ 's and of Accident Counts in a Population of Units	207
	Appendix D: The Negative Binomial Distribution and the Gamma Assumption	209
	Appendix E: The Negative Binomial Likelihood Function	210
	Appendix F: The Conditional Expectation, $E\{\mu K = k\}$	212
	Appendix G: The Negative Multinomial Likelihood Function	212
	Appendix H: The Nadaraya Watson Kernel Regression	214
	Appendix I: The CURE Limits	215
	Appendix J: Towards Theory; First Steps	216
	Appendix K: The “Bump Function”	223
	Appendix L: Elasticity and CMFs for Multiplicative Single-Variable Functions	224
	Appendix M: Interaction Terms for Additive Linear and Multiplicative Power Models	224
	References	225
	Index	227

Glossary

Notational Conventions

$\sigma\{.\}$	Standard deviation of what the dot stands for
\wedge	A caret stands for “estimate of” the letter below
$E\{.\}$	Expected value of what the dot stands for
\ln	Natural logarithm
$P(K = k)$ or $P(k)$	Probability that the random variable K takes on the value k

Acronyms

AADT	Annual average daily traffic
AIC	Akaike information criterion
BIC	Bayesian information criterion
C-F	Curve-fitting
CMF	Crash modification factor or function
EDA	Exploratory data analysis
F&I	Fatal and injury
ML	Maximum likelihood
NB	Negative binomial
NM	Negative multinomial
N-W	Nadaraya-Watson
OVB	Omitted variable bias
pdf	Probability density function
PDF	Probability distribution function
PDO	Property damage only
RHR	Roadside hazard rating
SD	Squared differences (residuals)
SPF	Safety performance function

SSD	Sum of squared differences (residuals)
TAB	Total accumulated bias
VBA	Visual basic for applications
VIEDA	Variable introduction exploratory data analysis
vpd	Vehicles per day

Greek

α, β	Parameters
β_0	Scale parameter
β_i	Regression parameters $i = 1, 2, 3, \dots$
ε_{f, X_i}	Elasticity of function “ f ” with respect to change in variable X_i
θ	Gamma distributed random variable with mean = 1 and variance = $1/b$
λ	Mean number of reported accidents per unit of time
μ_i	Expected number of accidents of unit i
$\sigma\{\mu\}$	Standard deviation of the μ 's in a population of units
$\hat{\sigma}^2(i)$	Sum of sorted squared residuals from 1 to i
$\pm \hat{\sigma}'(i)$	Standard error for CURE plot at index = i

Latin

a, b	Parameters of the gamma and the negative binomial distributions
A, B, \dots	Traits
c	Vehicle “concentration”—the (average) number of vehicles/(lane-km)
$E\{\mu\}$	Expected value of the μ 's in a population of units
f	The expected number of reported accidents per lane per second
$f()$	The function linking predictor variables and parameters
h	Bandwidth in the N-W nonparametric regression
\bar{h}	Average headway
i	Counter (index) of units
j	Counter (index) of time periods
$K(.)$	Kernel function in the Nadaraya-Watson nonparametric regression
K, k	Count of accidents and a certain value of that count
$\mathcal{L}(.)$	Likelihood. The dot in the parenthesis is a placeholder for parameters.
L_i	Length of road segment i
$\ln(\mathcal{L}^*)$	Abridged log-likelihood
m	As subscript denotes “multivehicle”

N	Number of bins or levels of a trait
n	Number of units or, occasionally, of accident counts
NB	Negative binomial
NM	Negative multinomial
p	Probability of a vehicle to be in a crash in the next second
$P(\cdot)$	Probability of the event that the dot stands for
r	The probability of the crash to be reported
s	Standard error or, occasionally, the sum of accident counts
s	As subscript denotes “single-vehicle”
T	Duration of a time period
$V\{\cdot\}$	Variance of the random variable that the dot stands for
\bar{v}	Average speed
v_0	Average free-flow speed
v_b	Average speed at which a bottleneck begins to form
X_1, X_2, \dots	Variables in model equation
Δt	Small time interval

Abstract

Statistical models that use data to express the safety of populations of units (road segments, intersections, grade crossings, etc.) as a function their traits (traffic, geometry, operation) are nowadays called Safety Performance Functions; SPFs for short. To make this notion precise one has to say what is meant by “safety,” “unit,” “population,” and “trait.” Most importantly, one has to be clear about what exactly is the information that a SPF provides and what are its practical uses. These are then illustrated by a series of examples. The practical uses of SPFs call for an approach to modeling which is somewhat different from what is usually done.

1.1 Units and Their Safety Property

What should one call the safety of Bobcaygeon Road between the Scotch Line and Plantation Roads in Ontario or the safety of the intersection of Eglinton and Don Mills roads in Toronto? These questions refer to safety as a property of some elements of the real world to be called “units.” A road segment, an intersection, a vehicle, or a person is a “unit.” A key feature of a “unit” is that it may be involved in accidents¹ (crashes) or that crashes (accidents) may occur on it. While the count of

¹Those who prefer to use “crash” instead think that “accident” has connotations of being unavoidable, without cause, and thus unpreventable. This, they fear, might weaken the resolve to reduce crashes and their harm. Since the very purpose of studying road safety is to assist in the task of managing the frequency and severity of accidents, such an interpretation makes no sense in this book. In this book “crash” and “accident” describe the same event. One reason for not shunning the term “accident” is that it is the common currency in the community of transportation professionals. Another reason is that the word “accident” provides the proper associations for the randomness inherent in accident counts. The editor of the Canadian Oxford Dictionary (Barber 1999) says that: “No dictionary that I know of uses the word “unpreventable” in any of its definitions of the word “accident.” Were they to do so, the definition would be inaccurate and

accidents is an indicator of the safety of a unit, it is not identical to it. The safety property of a unit is defined to be *the number of accidents by type and severity, expected to occur on or to it in a specified period of time.*²

Two elements of this definition need clarification. First, the word “expected” usually corresponds to “average in the long run.” This works well for the idealized devices used in the instruction of probability theory: coin tosses, decks of unmarked cards, urns with balls, and fair dies. It works well because all these devices can be plausibly assumed not to change with repeated use. In contrast, the safety of real units changes with time.

Accordingly, one has to interpret term “expected” by a conditional and counterfactual statement as “what the limit of the long-term average would be if³ it was possible to freeze all the safety-related traits and circumstances of the unit.”

The second element in the definition of safety which requires clarification is the phrase “by type and severity.” This means that, generally, the safety property of a unit is not a single number. Thus, for example, the safety of the intersection of High and Main streets in the 2-year period 2012 and 2013 could be described by the array in Table 1.1.

From such an array one can obtain the row sums (5.00 rear-end, 2.40 angle, 0.42 single-vehicle, and 0.08 pedestrian accidents), column sums (4.80 PDO, 2.75 injury, and 0.35 fatal accident), and the total (7.90 accidents). Each accident may involve one or more drivers and their vehicles. Therefore for multi-vehicle categories such as rear-end and angle accidents the additional categorization by number of involved vehicles may be needed.⁴

not reflect the actual usage of the word. The defining terms that dominate are “unexpected,” “unforeseen,” “unintentional,” and “undesirable” Most people recognize that the things we refer to as “accidents” do indeed have causes, whether it be an unplanned pregnancy, slipping on a banana peel or the dog peeing on the rug and that accidents are preventable.” In some cases, such as when one speaks of Crash Modification Functions (CMFs) the use of “crash” is so well established that its use already seems natural. However, the exclusive use of “crash” is often a sign of advocacy, an impression I want to avoid. This is why both terms will be used interchangeably.

² For a more detailed discussion of how safety is defined see Chap. 3 in Hauer (1997).

³ As if the notion of “average in the long run” was not enough of an obstacle, the “would be if” phrase further distances the safety property from what is observable. Not only is safety not the same as the number of accidents, it is now something that can be imagined and perhaps estimated but can never be observed.

⁴ Alternatively, one may want to speak not of the number of accidents but of the number of accident-involved vehicles or drivers. Similarly, an injury accident is one in which one or more persons were injured. Thus one may wish to describe the safety of a unit not by the number of injury accidents but by the number of injured persons. To convert from “number of accidents” to “number of involved vehicles or drivers” one needs to know the ratio “involved vehicles/accident.” For the US from 1988 to 2005 the ratio remained steady at 1.724 ± 0.015 . To convert from “fatal accidents” to “persons killed” knowledge of the ratio “persons killed/fatal accident” is needed. For the US from 1988 to 2005 the ratio for fatalities/fatal crash was 1.113 ± 0.004 which is similar to that in Michigan (1.104 ± 0.013 and 1.114 ± 0.019 for those who had been drinking), in Ohio (1.088 ± 0.013), Wisconsin (1.130 ± 0.022). It is also similar to the ratio pedestrian fatalities/Pedestrian Fatal accidents which in Michigan was 1.080 ± 0.023 .