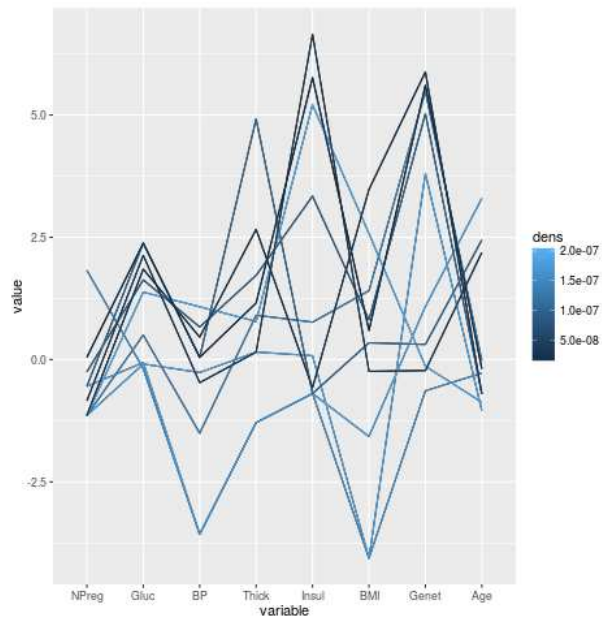# Statistical Regression and Classification

## From Linear Models to Machine Learning

Norman Matloff

University of California, Davis

Outlier Hunt

2

# Contents

i

# Preface

Why write yet another regression book? There is a plethora of books out there already, written by authors whom I greatly admire, and whose work I myself have found useful. I might cite the books by Harrell [60] and Fox [49], among many, many excellent examples. Note that I am indeed referring to general books on regression analysis, as opposed to more specialized work such as [65] and [75], which belong to a different genre. My book here is intended for a traditional (though modernized) regression course, rather than one on statistical learning.

Yet, I felt there is an urgent need for a different kind of book. So, why is this regression book different from all other regression books? First, **it modernizes the standard treatment of regression methods**. In particular:

- The book supplements classical regression models with introductory material on machine learning methods.

- Recognizing that these days, classification is the focus of many applications, the book covers this topic in detail, especially the multiclass case.

- In view of the voluminous nature of many modern datasets, there is a chapter on Big Data.

- There is much more hands-on involvement of computer usage.

Other major senses in which this book differs from others are:

- Though presenting the material in a mathematically precise manner, the book aims to provide much needed practical insight for the practicing analyst, remedying the "too many equations, too few explanations" problem.

For instance, the book not only shows how the math works for trans-
formations of variables, but also raises points on why one might refrain
from applying transformations.

- The book features a recurring interplay between parametric and non-
parametric methods. For instance, in an example involving currency
data, the book finds that the fitted linear model predicts substantially
more poorly than a k-nearest neighbor fit, suggesting deficiencies in
the linear model. Nonparametric analysis is then used to further in-
vestigate, providing parametric model assessment in a manner that is
arguably more insightful than classical residual plots.

- For those interested in computing issues, many of the book's chap-
ters include optional sections titled Computational Complements, on
topics such as data wrangling, views of package source code, parallel
computing and so on.

  Also, many of the exercises are code-oriented. In particular, in such
  exercises the reader is asked to write "mini-CRAN" functions,[1] short
  but useful library functions that can be applied to practical regression
  analysis. Here is an example exercise of this kind:

  > Write a function stepAR2() that works similarly to stepAIC(),
  > except that this new function uses adjusted $R^2$ as its crite-
  > rion for adding or deleting a predictor. The call form will
  > be

  ```
  stepAR2(lmobj,direction='fwd',
      nsteps=ncol(lmobj$model)-1)
  ```

  > where the arguments are...Predictors will be added/deleted
  > one at a time, according to which one maximizes adjusted
  > $R^2$. The return value will be an S3 object of type 'stepr2',
  > with sole component a data frame of $\widehat{\beta}_i$ values (0s meaning
  > the predictor is not currently in the prediction equation),
  > one row per model. There will also be an $R^2$ column. Write
  > a summary() function for this class that shows the actions
  > taken at each step of the process.

- For those who wish to go into more depth on mathematical topics,
there are Mathematical Complements sections at the end of most
chapters, and math-oriented exercises. The material ranges from
straightforward computation of mean squared error to esoteric top-
ics such as a proof of the Tower Property, $E\left[E(V|U_1, U_2) \mid U_1\right] = E(V \mid U_1)$, a result that is used in the text.

---

[1]CRAN is the online repository of user-contributed R code.

As mentioned, **this is still a book on traditional regression analysis.** In contrast to [65], this book is aimed at a traditional regression course. Except for Chapters 10 and 11, the primary methodology used is linear and generalized linear parametric models, covering both the Description and Prediction goals of regression methods. We are just as interested in Description applications of regression, such as measuring the gender wage gap in Silicon Valley, as we are in forecasting tomorrow's demand for bike rentals. An entire chapter is devoted to measuring such effects, including discussion of Simpson's Paradox, multiple inference, and causation issues. The book's examples are split approximately equally in terms of Description and Prediction goals. Issues of model fit play a major role.

The book includes more than 75 full examples, using real data. But concerning the above comment regarding "too many equations, too few explanations,", merely including examples with real data is not enough to truly tell the story in a way that will be useful in practice. Rather few books go much beyond presenting the formulas and techniques, thus leaving the hapless practitioner to his own devices. Too little is said in terms of what the equations really mean in a practical sense, what can be done with regard to the inevitable imperfections of our models, which techniques are too much the subject of "hype," and so on.

As a nonstatistician, baseball great Yogi Berra, put it in his inimitable style, "In theory there is no difference between theory and practice. In practice there is." This book aims to remedy this gaping deficit. It develops the material in a manner that is *mathematically precise* yet always maintains as its top priority — borrowing from a book title of the late Leo Breiman — *a view toward applications.*

In other words:

> The philosophy of this book is to not only prepare the analyst to know *how* to do something, but also to understand *what* she is doing. For successful application of data science techniques, the latter is just as important as the former.

Some further examples of how this book differs from the other regression books:

**Intended audience and chapter coverage:**

This book is aimed at both practicing professionals and use in the classroom. It aims to be both accessible and valuable to this diversity of readership.

In terms of classroom use, with proper choice of chapters and appendices, the book could be used as a text tailored to various discipline-specific audiences and various levels, undergraduate or graduate. I would recommend that the core of any course consist of most sections of Chapters 1-4 (excluding the Math and Computational Complements sections), with coverage of at least introductory sections of Chapters 5, 6, 7, 8 and 9 for all audiences. Beyond that, different types of disciplines might warrant different choices of further material. For example:

- **Statistics students:** Depending on level, at least some of the Mathematical Complements and math-oriented exercises should be involved. There might be more emphasis on Chapters 6, 7 and 9.

- **Computer science students:** Here one would cover more of classification, machine learning and Big Data material, Chapters 5, 8, 10, 11 and 12. Also, one should cover the Computational Complements sections and associated "mini-CRAN" code exercises.

- **Economics/social science students:** Here there would be heavy emphasis on the Description side, Chapters 6 and 7, with special emphasis on topics such as Instrumental Variables and Propensity Matching in Chapter 7. Material on generalized linear models and logistic regression, in Chapter 4 and parts of Chapter 5, might also be given emphasis.

- **Student class level:** The core of the book could easily be used in an undergraduate regression course, but aimed at students with background in calculus and matrix algebra, such as majors in statistics, math or computer science. A graduate course would cover more of the chapters on advanced topics, and would likely cover more of the Mathematical Complements sections.

- **Level of mathematical sophistication:** In the main body of the text, i.e., excluding the Mathematical Complements sections, basic matrix algebra is used throughout, but use of calculus is minimal. As noted, for those instructors who want the mathematical content, it is there in the Mathematical Complements sections, but the main body of the text requires only the matrix algebra and a little calculus.

The reader must of course be familiar with terms like *confidence interval*, *significance test* and *normal distribution*. Many readers will have had at least some prior exposure to regression analysis, but this is not assumed, and the subject is developed from the beginning.

The reader is assumed to have some prior experience with R, but at a minimal level: familiarity with function arguments, loops, if-else and vector/matrix operations and so on. For those without such background, there are many gentle tutorials on the Web, as well as a leisurely introduction in a statistical context in [21]. Those with programming experience can also read the quick introduction in the appendix of [102]. My book [95] gives a detailed treatment of R as a programming language, but that level of sophistication is certainly not needed for the present book.

**A comment on the field of machine learning:**

Mention should be made of the fact that this book's title includes both the word *regression* and the phrase *machine learning*. The latter phrase is included to reflect that the book includes some introductory material on machine learning, in a regression context.

Much has been written on a perceived gap between the statistics and machine learning communities [24]. This gap is indeed real, but work has been done to reconcile them [16], and in any event, the gap is actually not as wide as people think.

My own view is that machine learning (ML) consists of the development of regression models with the Prediction goal. Typically nonparametric (or what I call semi-parameteric) methods are used. Classification models are more common than those for predicting continuous variables, and it is common that more than two classes are involved, sometimes a great many classes. All in all, though, it's still regression analysis, involving the conditional mean of $Y$ given $X$ (reducing to $P(Y = 1|X)$ in the classification context).

One often-claimed distinction between statistics and ML is that the former is based on the notion of a sample from a population whereas the latter is concerned only with the content of the data itself. But this difference is more perceived than real. The idea of cross-validation is central to ML methods, and since that approach is intended to measure how well one's model generalizes beyond our own data, it is clear that ML people do think in terms of samples after all. Similar comments apply to ML's citing the variance-vs.-bias tradeoff, overfitting and so on

So, at the end of the day, we all are doing regression analysis, and this book takes this viewpoint.

**Code and software:**

The book also makes use of some of my research results and associated software. The latter is in my package **regtools**, available from CRAN [98].

A number of other packages from CRAN are used. Note that typically we use only the default values for the myriad arguments available in many functions; otherwise we could fill an entire book devoted to each package! Cross-validation is suggested for selection of tuning parameters, but with a warning that it too can be problematic.

In some cases, the **regtools** source code is also displayed within the text, so as to make clear exactly what the algorithms are doing. Similarly, data wrangling/data cleaning code is shown, not only for the purpose of "hands-on" learning, but also to highlight the importance of those topics.

**Thanks:**

Conversations with a number of people have directly or indirectly enhanced the quality of this book, among them Charles Abromaitis, Stuart Ambler, Doug Bates, Oleksiy Budilovsky, Yongtao Cao, Tony Corke, Tal Galili, Frank Harrell, Harlan Harris, Benjamin Hofner, Jiming Jiang, Hyunseung Kang, Martin Mächler, Erin McGinnis, John Mount, Richard Olshen, Pooja Rajkumar, Ariel Shin, Chuck Stone, Jessica Tsoi, Yu Wu, Yihui Xie, Yingkang Xie, Achim Zeileis and Jiaping Zhang.

A seminar presentation by Art Owen introduced me to the application of random effects models in recommender systems, a provocative blend of old and new. This led to the MovieLens examples and other similar examples in the book, as well as a vigorous new research interest for me. Art also led me to two Stanford statistics PhD students, Alex Chin and Jing Miao, who each read two of the chapters in great detail. Special thanks also go to Nello Cristianini, Hui Lin, Ira Sharenow and my old friend Gail Gong for their detailed feedback.

Thanks go to my editor, John Kimmel, for his encouragement and much-appreciated patience, and to the internal reviewers, David Giles and Robert Gramacy. Of course, I cannot put into words how much I owe to my wonderful wife Gamis and our daughter Laura, both of whom inspire all that I do, including this book project.

**Website:**

Code, errata, extra examples and so on are available at

*http://heather.cs.ucdavis.edu/regclass.html.*

**A final comment:**

My career has evolved quite a bit over the years. I wrote my dissertation in abstract probability theory [104], but turned my attention to applied statistics soon afterward. I was one of the founders of the Department of

Statistics at UC Davis. Though a few years later I transferred into the new Computer Science Department, I am still a statistician, and much of my CS research has been statistical, e.g., [100]. Most important, my interest in regression has remained strong throughout those decades.

I published my first research papers on regression methodology way back in the 1980s, and the subject has captivated me ever since. My long-held wish has been to write a regression book, and thus one can say this work is 30 years in the making. I hope you find its goals both worthy and attained. Above all, I simply hope you find it an interesting read.

# List of Symbols and Abbreviations

$Y$: the response variable
$X$: vector of predictor variables
$\widetilde{X}$: $X$ with a 1 prepended
$X^{(j)}$: the $j^{th}$ predictor variable
$n$: number of observations
$p$: number of predictors
$Y_i$: value of the response variable in observation $i$
$X_i$: vector of predictors in observation $i$
$X_i^{(j)}$: value of the $j^{th}$ predictor variable in observation $i$
$A$: $n \times (p+1)$ matrix of the predictor data in a linear model
$D$: length-$n$ vector of the response data in a linear model
$H$: the *hat matrix*, $A(A'A)^{-1}A'$
$\mu(t)$: the regression function $E(Y|X=t)$
$\sigma^2(t)$: $Var(Y|X=t)$
$\widehat{\mu}(t)$: estimated value of $\mu(t)$
$\beta$: vector of coefficients in a linear/generalized linear model
$\widehat{\beta}$: estimated value of $\beta$
$'$: matrix transpose
$I$: multiplicative identity matrix
k-NN: k-Nearest Neighbor method
MSE: Mean Squared (Estimation) Error
MSPE: Mean Squared Prediction Error
CART: Classification and Regression Trees
SVM: Support Vector Machine
NN: neural network
PCA: Principal Components Analysis
NMF: Nonnegative Matrix Factorization
OVA: One vs. All classification

AVA: All vs. All classification
LDA: Linear Discriminant Analysis