Michel Denuit
Donatien Hainaut
Julien Trufin

# Effective Statistical Learning Methods for Actuaries I

## GLMs and Extensions

Springer

# Springer Actuarial

## Springer Actuarial Lecture Notes

This subseries of Springer Actuarial includes books with the character of lecture notes. Typically these are research monographs on new, cutting-edge developments in actuarial science; sometimes they may be a glimpse of a new field of research activity, or presentations of a new angle in a more classical field.

In the established tradition of Lecture Notes, the timeliness of a manuscript can be more important than its form, which may be informal, preliminary or tentative.

More information about this subseries at http://www.springer.com/series/15682

Michel Denuit · Donatien Hainaut ·
Julien Trufin

# Effective Statistical Learning Methods for Actuaries I

GLMs and Extensions

Springer

Michel Denuit
Institut de Statistique, Biostatistique et
Sciences Actuarielles
Université Catholique de Louvain
Louvain-la-Neuve, Belgium

Donatien Hainaut
Institut de Statistique, Biostatistique et
Sciences Actuarielles
Université Catholique de Louvain
Louvain-la-Neuve, Belgium

Julien Trufin
Département de Mathématiques
Université Libre de Bruxelles
Brussels, Belgium

# Preface

The present material is written for students enrolled in actuarial master programs and practicing actuaries, who would like to gain a better understanding of insurance data analytics. It is built in three volumes, starting from the celebrated Generalized Linear Models, or GLMs and continuing with tree-based methods and neural networks.

After an introductory chapter, this first volume starts with a recap' of the basic statistical aspects of insurance data analytics and summarizes the state of the art using GLMs and their various extensions: GAMs, mixed models and credibility, and some nonlinear versions, or GNMs. Analytical tools from Extreme Value Theory are also presented to deal with tail events that arise in liability insurance or survival analysis. This book also goes beyond mean modeling, considering volatility modeling (double GLMs) and the general modeling of location, scale and shape parameters (GAMLSS).

Throughout this book, we alternate between methodological aspects and numerical illustrations or case studies to demonstrate practical applications of the proposed techniques. The numerous examples cover all areas of insurance, not only property and casualty but also life and health, being based on real data sets from the industry or collected by regulators.

The R statistical software has been found convenient to perform the analyses throughout this book. It is a free language and environment for statistical computing and graphics. In addition to our own R code, we have benefited from many R packages contributed by the members of the very active community of R-users. We provide the readers with information about the resources available in R throughout the text as well as in the closing section to each chapter. The open-source statistical software R is freely available from https://www.r-project.org/.

The technical requirements to understand the material are kept at a reasonable level so that this text is meant for a broad readership. We refrain from proving all results but rather favor an intuitive approach with supportive numerical illustrations, providing the reader with relevant references where all justifications can be found, as well as more advanced material. These references are gathered in a dedicated section at the end of each chapter.

The three authors are professors of actuarial mathematics at the universities of Brussels and Louvain-la-Neuve, Belgium. Together, they accumulate decades of teaching experience related to the topics treated in the three books, in Belgium and throughout Europe and Canada. They are also scientific directors at Detralytics, a consulting office based in Brussels.

Within Detralytics as well as on behalf of actuarial associations, the authors have had the opportunity to teach the material contained in the three volumes of "Effective Statistical Learning Methods for Actuaries" to various audiences of practitioners. The feedback received from the participants to these short courses greatly helped to improve the exposition of the topic. Throughout their contacts with the industry, the authors also implemented these techniques in a variety of consulting and R&D projects. This makes the three volumes of "Effective Statistical Learning Methods for Actuaries" the ideal support for teaching students and CPD events for professionals.

Louvain-la-Neuve, Belgium                                                    Michel Denuit
Louvain-la-Neuve, Belgium                                                Donatien Hainaut
Brussels, Belgium                                                              Julien Trufin
June 2019

# Notation

Here are a few words on the notation and terminology used throughout the book. For the most part, the notation conforms to what is usual in mathematical statistics as well as insurance mathematics. Unless stated otherwise,

| | |
|---|---|
| $n =$ | Number of observations (or data points). |
| $y_i =$ | The $i$th observed value of the response (or outcome, dependent variable), observed response for individual $i$, realization of the random variable $Y_i, i = 1, \ldots, n$. |
| $\widehat{y}_i =$ | Fitted value, predicted response for the $i$th individual. |
| $\bar{y} =$ | Average, or sample mean of the $n$ observed responses $y_1, \ldots, y_n$. |
| $\boldsymbol{y} =$ | Column vector of all $n$ response values. All vectors are column ones, by convention. |
| $p =$ | Number of features (sometimes also called covariates, independent variables, regressors, as well as explanatory variables or predictors when they influence the response). |
| $x_{ij} =$ | Value of the $j$th feature for the $i$th data point, $i = 1, \ldots, n$, $j = 1, \ldots, p$, realization of the random variable $X_{ij}$. |
| $\boldsymbol{x}_i =$ | Column vector of the $p$ features for the $i$th data point. |
| $\boldsymbol{X} =$ | Matrix of the $p$ features for all data points, known as the design matrix with $n$ rows and $p$ columns ($\boldsymbol{x}_i$ is the $i$th row of $\boldsymbol{X}$) |
| | or random vector with components $X_j, j = 1, \ldots, p$ (clear from the context). |
| $\boldsymbol{X}^\top =$ | The transpose of the design matrix $\boldsymbol{X}$, with $p$ rows and $n$ columns. |
| $\mathrm{I}[\cdot] =$ | Indicator function of an event (equal to 1 if the event appearing within the brackets is realized and to 0 otherwise). |
| $\mathrm{P}[\cdot] =$ | Probability of an event. |
| $\mathrm{E}[\cdot] =$ | Expectation of a random variable. |
| $\mathrm{Var}[\cdot] =$ | Variance of a random variable. |
| $\mathrm{Cov}[\cdot, \cdot] =$ | Covariance of a pair of random variables. |

$\widehat{\theta} =$                  Estimator or estimate (clear from the context) of the unknown
parameter $\theta$ (parameters are denoted by Greek letters).

    The real line is denoted as $(-\infty, \infty)$. We use the symbol $\sim$ to mean "is distributed as", $\approx$ to mean "is approximately equal to" or "is approximately distributed as" (clear from the context), and $\propto$ to mean "is proportional to".

# Contents