# Applied Multivariate Statistical Analysis

Wolfgang Karl Härdle • Léopold Simar

# Applied Multivariate Statistical Analysis

Fourth Edition

 Springer

Wolfgang Karl Härdle
C.A.S.E. Centre f. Appl. Stat. & Econ.
    School of Business and Economics
Humboldt-Universität zu Berlin
Berlin, Germany

Léopold Simar
Center of Operations Research &
    Econometrics (CORE)
Katholieke Univeristeit Leuven Inst.
    Statistics
Leuven, Belgium

The majority of chapters have quantlet codes in Matlab or R. These quantlets may be downloaded from http://extras.springer.com or via a link on http://springer.com/978-3-662-45170-0 and from www.quantlet.de

# Preface to the Fourth Edition

The fourth edition of this book on *Applied Multivariate Statistical Analysis* offers a new sub-chapter on Variable Selection by using least absolute shrinkage and selection operator (LASSO) and its general form the so-called Elastic Net.

All pictures and numerical examples have been now calculated in the (almost) standard language R & MATLAB. The code for each picture is indicated with a small **Q** sign near the picture, e.g. **Q** MVAdenbank denotes the corresponding quantlet for reproduction of Fig. 1.9, where we display the densities of the diagonal of genuine and counterfeit bank notes. We believe that these publicly available quantlets (see also http://sfb649.wiwi.hu-berlin.de/quantnet/) create a valuable contribution to distribution of knowledge in the statistical science. The symbols and notations have also been standardised. In the preparation of the fourth edition, we received valuable input from Dedy Dwi Prastyo, Petra Burdejova, Sergey Nasekin and Awdesch Melzer. We would like to thank them.

Berlin, Germany                                                    Wolfgang Karl Härdle
Louvain la Neuve, Belgium                                          Léopold Simar
January 2014

# Preface to the Third Edition

The third edition of this book on *Applied Multivariate Statistical Analysis* offers the following new features.

1. A new Chap. 8 on Regression Models has been added.
2. Almost all numerical examples have been reproduced in MATLAB or R.

The chapter on regression models focuses on a core business of multivariate statistical analysis. This contribution has not been subject of a prominent discussion in earlier editions of this book. We now take the opportunity to cover classical themes of ANOVA and ANCOVA analysis. Categorical responses are presented in Sect. 8.2. The spectrum of log linear models for contingency tables is presented in Sect. 8.2.2, and applications to count data, e.g. in the economic and medical science are presented there. Logit models are discussed in great detail, and the numerical implementation in terms of matrix manipulations is presented.

The majority of pictures and numerical examples has been now calculated in the (almost) standard language R & MATLAB. The code for each picture is indicated with a small  sign near the picture, e.g.  MVAdenbank denotes the corresponding quantlet for reproduction of Fig. 1.9, where we display the densities of the diagonal of genuine and counterfeit bank notes. We believe that these publicly available quantlets (see also www.quantlet.com) create a valuable contribution to distribution of knowledge in the statistical science. The symbols and notations have also been standardised. In the preparation of the third edition, we received valuable input from Song Song, Weining Wang and Mengmeng Guo. We would like to thank them.

Berlin, Germany                                    Wolfgang Karl Härdle
Louvain la Neuve, Belgium                                   Léopold Simar
June 2011

# Contents

**Part IV Appendix**