# Data Mining with SPSS Modeler

Tilo Wendler • Sören Gröttrup

# Data Mining with SPSS Modeler

Theory, Exercises and Solutions

Springer

Tilo Wendler
HTW Berlin
University of Applied Sciences
Berlin
Germany

Sören Gröttrup
Stuttgart
Germany

# Preface

Data Analytics, Data Mining and Big Data are terms often used in everyday business. Companies collect more and more data and store it in databases, with the hope of finding helpful patterns that can improve business. Shortly after deciding to more use of such data, managers often confess that analysing these datasets is resource-consuming and anything but easy. Involving the firm's IT-experts leads to a discussion regarding which tools to use. Very few applications are available in the marketplace that are appropriate for handling large datasets in a professional way. Two commercial products worth mentioning are 'Enterprise Miner' by SAS and 'SPSS Modeler' by IBM.

At first glance, these applications are easy to use. After a while, however, many users realize that more difficult questions require deeper statistical knowledge. Many people are interested in gaining such statistical skills and applying them, using one of the data mining tools offered by the industry.

This book will help users to become familiar with a wide range of statistical concepts or algorithms and apply them to concrete datasets. After a short statistical overview of how the procedures work and what assumptions to keep in mind, step-by-step procedures show how to find the solutions with the SPSS Modeler.

Features of This Book
– Easy to read
– Standardised chapter structure, including exercises and solutions
– All datasets are provided as downloads and explained in detail
– Template streams help the reader focus on the interesting parts of the stream and leave out recurring tasks
– Complete solution streams are ready to use
– Each example includes step-by-step explanations
– Short explanations of the most important statistical assumptions used when applying the algorithms are included
– Hundreds of screenshots are included, to ensure successful application of the algorithms to the datasets
– Exercises teach how to secure and systematise this knowledge
– Explanations and solutions are provided for all exercises
– Skills acquired through solving the exercises allow the user to create his/her own streams

Berlin, Germany                                                         Tilo Wendler
Stuttgart, Germany                                                   Sören Gröttrup

# Contents