ROBERT NISBET

JOHN ELDER

GARY MINER

# HANDBOOK OF

# Statistical Analysis
# & Data Mining
# Applications

# HANDBOOK OF STATISTICAL ANALYSIS AND DATA MINING APPLICATIONS

ROBERT NISBET
*Pacific Capital Bankcorp N.A.*
*Santa Barbara, CA*

JOHN ELDER
*Elder Research, Inc., Charlottesville, VA*

GARY MINER
*StatSoft, Inc., Tulsa, Oklahoma*

ELSEVIER

AMSTERDAM • BOSTON • HEIDELBERG • LONDON
NEW YORK • OXFORD • PARIS • SAN DIEGO
SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO
Academic Press is an imprint of Elsevier

# Table of Contents

## II

# THE ALGORITHMS IN DATA MINING AND TEXT MINING, THE ORGANIZATION OF THE THREE MOST COMMON DATA MINING TOOLS, AND SELECTED SPECIALIZED AREAS USING DATA MINING

# III

# TUTORIALS—STEP-BY-STEP CASE STUDIES AS A STARTING POINT TO LEARN HOW TO DO DATA MINING ANALYSES

## Guest Authors of the Tutorials

### A. How to Use Data Miner Recipe

### B. Data Mining for Aviation Safety

### C. Predicting Movie Box-Office Receipts

### D. Detecting Unsatisfied Customers: A Case Study

### E. Credit Scoring

## IV

## MEASURING TRUE COMPLEXITY, THE "RIGHT MODEL FOR THE RIGHT USE," TOP MISTAKES, AND THE FUTURE OF ANALYTICS

# Foreword 1

This book will help the novice user become familiar with data mining. Basically, data mining is doing data analysis (or statistics) on data sets (often large) that have been obtained from potentially many sources. As such, the miner may not have control of the input data, but must rely on sources that have gathered the data. As such, there are problems that every data miner must be aware of as he or she begins (or completes) a mining operation. I strongly resonated to the material on "The Top 10 Data Mining Mistakes," which give a worthwhile checklist:

- Ensure you have a response variable and predictor variables—and that they are correctly measured.
- Beware of overfitting. With scads of variables, it is easy with most statistical programs to fit incredibly complex models, but they cannot be reproduced. It is good to save part of the sample to use to test the model. Various methods are offered in this book.
- Don't use only one method. Using only linear regression can be a problem. Try dichotomizing the response or categorizing it to remove nonlinearities in the response variable. Often, there are clusters of values at zero, which messes up any normality assumption. This, of course, loses information, so you may want to categorize a continuous response variable and use an alternative to regression. Similarly, predictor variables may need to be treated as factors rather than linear predictors. A classic example is using marital status or race as a linear predictor when there is no order.
- Asking the wrong question—when looking for a rare phenomenon, it may be helpful to identify the most common pattern. These may lead to complex analyses, as in item 3, but they may also be conceptually simple. Again, you may need to take care that you don't overfit the data.
- Don't become enamored with the data. There may be a substantial history from earlier data or from domain experts that can help with the modeling.
- Be wary of using an outcome variable (or one highly correlated with the outcome variable) and becoming excited about the result. The predictors should be "proper" predictors in the sense that (a) they are measured prior to the outcome and (b) are not a function of the outcome.
- Do not discard outliers without solid justification. Just because an observation is out of line with others is insufficient reason to ignore it. You must check the circumstances that led to the value. In any event, it is useful to conduct the analysis with the observation(s) included and excluded to determine the sensitivity of the results to the outlier.

- Extrapolating is a fine way to go broke—the best example is the stock market. Stick within your data, and if you must go outside, put plenty of caveats. Better still, restrain the impulse to extrapolate. Beware that pictures are often far too simple and we can be misled. Political campaigns oversimplify complex problems ("My opponent wants to raise taxes"; "My opponent will take us to war") when the realities may imply we have some infrastructure needs that can be handled only with new funding, or we have been attacked by some bad guys.

Be wary of your data sources. If you are combining several sets of data, they need to meet a few standards:

- The definitions of variables that are being merged should be identical. Often they are close but not exact (especially in meta-analysis where clinical studies may have somewhat different definitions due to different medical institutions or laboratories).
- Be careful about missing values. Often when multiple data sets are merged, missing values can be induced: one variable isn't present in another data set, what you thought was a unique variable name was slightly different in the two sets, so you end up with two variables that both have a lot of missing values.
- How you handle missing values can be crucial. In one example, I used complete cases and lost half of my sample—all variables had at least 85% completeness, but when put together the sample lost half of the data. The residual sum of squares from a stepwise regression was about 8. When I included more variables using mean replacement, almost the same set of predictor variables surfaced, but the residual sum of squares was 20. I then used multiple imputation and found approximately the same set of predictors but had a residual sum of squares (median of 20 imputations) of 25. I find that mean replacement is rather optimistic but surely better than relying on only complete cases. If using stepwise regression, I find it useful to replicate it with a bootstrap or with multiple imputation. However, with large data sets, this approach may be expensive computationally.

To conclude, there is a wealth of material in this handbook that will repay study.

*Peter A. Lachenbruch, Ph.D.,*
*Oregon State University*
*Past President, 2008, American Statistical Society*
*Professor, Oregon State University*
*Formerly: FDA and professor at Johns Hopkins University;*
*UCLA, and University of Iowa, and*
*University of North Carolina Chapel Hill*

# Foreword 2

A November 2008 search on Amazon.com for "data mining" books yielded over 15,000 hits—including 72 to be published in 2009. Most of these books either describe data mining in very technical and mathematical terms, beyond the reach of most individuals, or approach data mining at an introductory level without sufficient detail to be useful to the practitioner. The *Handbook of Statistical Analysis and Data Mining Applications* is the book that strikes the right balance between these two treatments of data mining.

This volume is not a theoretical treatment of the subject—the authors themselves recommend other books for this—but rather contains a description of data mining principles and techniques in a series of "knowledge-transfer" sessions, where examples from real data mining projects illustrate the main ideas. This aspect of the book makes it most valuable for practitioners, whether novice or more experienced.

While it would be easier for everyone if data mining were merely a matter of finding and applying the correct mathematical equation or approach for any given problem, the reality is that both "art" and "science" are necessary. The "art" in data mining requires experience: when one has seen and overcome the difficulties in finding solutions from among the many possible approaches, one can apply newfound wisdom to the next project. However, this process takes considerable time and, particularly for data mining novices, the iterative process inevitable in data mining can lead to discouragement when a "textbook" approach doesn't yield a good solution.

This book is different; it is organized with the practitioner in mind. The volume is divided into four parts. Part I provides an overview of analytics from a historical perspective and frameworks from which to approach data mining, including CRISP-DM and SEMMA. These chapters will provide a novice analyst an excellent overview by defining terms and methods to use, and will provide program managers a framework from which to approach a wide variety of data mining problems. Part II describes algorithms, though without extensive mathematics. These will appeal to practitioners who are or will be involved with day-to-day analytics and need to understand the qualitative aspects of the algorithms. The inclusion of a chapter on text mining is particularly timely, as text mining has shown tremendous growth in recent years.

Part III provides a series of tutorials that are both domain-specific and software-specific. Any instructor knows that examples make the abstract concept more concrete, and these tutorials accomplish exactly that. In addition, each tutorial shows how the solutions were developed using popular data mining software tools, such as Clementine, Enterprise Miner, Weka, and *STATISTICA*. The step-by-step specifics will assist practitioners in learning not only how to approach a wide variety of problems, but also how to use these software

products effectively. Part IV presents a look at the future of data mining, including a treatment of model ensembles and "The Top 10 Data Mining Mistakes," from the popular presentation by Dr. Elder.

However, the book is best read a few chapters at a time while actively doing the data mining rather than read cover-to-cover (a daunting task for a book this size). Practitioners will appreciate tutorials that match their business objectives and choose to ignore other tutorials. They may choose to read sections on a particular algorithm to increase insight into that algorithm and then decide to add a second algorithm after the first is mastered. For those new to a particular software tool highlighted in the tutorials section, the step-by-step approach will operate much like a user's manual. Many chapters stand well on their own, such as the excellent "History of Statistics and Data Mining" and "The Top 10 Data Mining Mistakes" chapters. These are broadly applicable and should be read by even the most experienced data miners.

The *Handbook of Statistical Analysis and Data Mining Applications* is an exceptional book that should be on every data miner's bookshelf or, better yet, found lying open next to their computer.

*Dean Abbott*
*President*
*Abbott Analytics*
*San Diego, California*

# Preface

Data mining scientists in research and academia may look askance at this book because it does not present algorithm theory in the commonly accepted mathematical form. Most articles and books on data mining and knowledge discovery are packed with equations and mathematical symbols that only experts can follow. Granted, there is a good reason for insistence on this formalism. The underlying complexity of nature and human response requires teachers and researchers to be extremely clear and unambiguous in their terminology and definitions. Otherwise, ambiguities will be communicated to students and readers, and their understanding will not penetrate to the essential elements of any topic. Academic areas of study are not called *disciplines* without reason.

This rigorous approach to data mining and knowledge discovery builds a fine foundation for academic studies and research by experts. Excellent examples of such books are

- *The Handbook of Data Mining*, 2003, by Nong Ye (Ed.). Mahwah, New Jersey: Lawrence Erlbaum Associates.
- *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd edition, 2009, by T. Hastie, R. Tibshirani, & J. Friedman. New York: Springer-Verlag.

Books like these were especially necessary in the early days of data mining, when analytical tools were relatively crude and required much manual configuration to make them work right. Early users had to understand the tools in depth to be able to use them productively. These books are still necessary for the college classroom and research centers. Students must understand the theory behind these tools in the same way that the developers understood it so that they will be able to build new and improved versions.

Modern data mining tools, like the ones featured in this book, permit ordinary business analysts to follow a path through the data mining process to create models that are "good enough." These less-than-optimal models are far better in their ability to leverage faint patterns in databases to solve problems than the ways it used to be done. These tools provide default configurations and automatic operations, which shield the user from the technical complexity underneath. They provide one part in the crude analogy to the automobile interface. You don't have to be a chemical engineer or physicist who understands moments of force to be able to operate a car. All you have to do is learn to turn the key in the ignition, step on the gas and the brake at the right times, turn the wheel to change direction in a safe manner, and *voila*, you are an expert user of the very complex technology