





















## CONTENTS

---

<b>Preface</b>	<b>ix</b>
<b>Acknowledgments</b>	<b>xi</b>
<b>1. Introduction</b>	<b>1</b>
Reference	6
<b>2. Processing the Information and Getting to Know Your Data</b>	<b>7</b>
2.1 Example 1: 2006 Birth Data	7
2.2 Example 2: Alumni Donations	17
2.3 Example 3: Orange Juice	31
References	39
<b>3. Standard Linear Regression</b>	<b>40</b>
3.1 Estimation in R	43
3.2 Example 1: Fuel Efficiency of Automobiles	43
3.3 Example 2: Toyota Used-Car Prices	47
Appendix 3.A The Effects of Model Overfitting on the Average Mean Square Error of the Regression Prediction	53
References	54
<b>4. Local Polynomial Regression: a Nonparametric Regression     Approach</b>	<b>55</b>
4.1 Model Selection	56
4.2 Application to Density Estimation and the Smoothing of Histograms	58
4.3 Extension to the Multiple Regression Model	58
4.4 Examples and Software	58
References	65
<b>5. Importance of Parsimony in Statistical Modeling</b>	<b>67</b>
5.1 How Do We Guard Against False Discovery	67
References	70

<b>6. Penalty-Based Variable Selection in Regression Models with Many Parameters (LASSO)</b>	<b>71</b>
6.1 Example 1: Prostate Cancer	74
6.2 Example 2: Orange Juice	78
References	82
<b>7. Logistic Regression</b>	<b>83</b>
7.1 Building a Linear Model for Binary Response Data	83
7.2 Interpretation of the Regression Coefficients in a Logistic Regression Model	85
7.3 Statistical Inference	85
7.4 Classification of New Cases	86
7.5 Estimation in R	87
7.6 Example 1: Death Penalty Data	87
7.7 Example 2: Delayed Airplanes	92
7.8 Example 3: Loan Acceptance	100
7.9 Example 4: German Credit Data	103
References	107
<b>8. Binary Classification, Probabilities, and Evaluating Classification Performance</b>	<b>108</b>
8.1 Binary Classification	108
8.2 Using Probabilities to Make Decisions	108
8.3 Sensitivity and Specificity	109
8.4 Example: German Credit Data	109
<b>9. Classification Using a Nearest Neighbor Analysis</b>	<b>115</b>
9.1 The $k$ -Nearest Neighbor Algorithm	116
9.2 Example 1: Forensic Glass	117
9.3 Example 2: German Credit Data	122
Reference	125
<b>10. The Naïve Bayesian Analysis: a Model for Predicting a Categorical Response from Mostly Categorical Predictor Variables</b>	<b>126</b>
10.1 Example: Delayed Airplanes	127
Reference	131
<b>11. Multinomial Logistic Regression</b>	<b>132</b>
11.1 Computer Software	134
11.2 Example 1: Forensic Glass	134

11.3	Example 2: Forensic Glass Revisited	141
	Appendix 11.A Specification of a Simple Triplet Matrix	147
	References	149
<b>12.</b>	<b>More on Classification and a Discussion on Discriminant Analysis</b>	<b>150</b>
12.1	Fisher's Linear Discriminant Function	153
12.2	Example 1: German Credit Data	154
12.3	Example 2: Fisher Iris Data	156
12.4	Example 3: Forensic Glass Data	157
12.5	Example 4: MBA Admission Data	159
	Reference	160
<b>13.</b>	<b>Decision Trees</b>	<b>161</b>
13.1	Example 1: Prostate Cancer	167
13.2	Example 2: Motorcycle Acceleration	179
13.3	Example 3: Fisher Iris Data Revisited	182
<b>14.</b>	<b>Further Discussion on Regression and Classification Trees, Computer Software, and Other Useful Classification Methods</b>	<b>185</b>
14.1	R Packages for Tree Construction	185
14.2	Chi-Square Automatic Interaction Detection (CHAID)	186
14.3	Ensemble Methods: Bagging, Boosting, and Random Forests	188
14.4	Support Vector Machines (SVM)	192
14.5	Neural Networks	192
14.6	The R Package Rattle: A Useful Graphical User Interface for Data Mining	193
	References	195
<b>15.</b>	<b>Clustering</b>	<b>196</b>
15.1	<i>k</i> -Means Clustering	196
15.2	Another Way to Look at Clustering: Applying the Expectation-Maximization (EM) Algorithm to Mixtures of Normal Distributions	204
15.3	Hierarchical Clustering Procedures	212
	References	219
<b>16.</b>	<b>Market Basket Analysis: Association Rules and Lift</b>	<b>220</b>
16.1	Example 1: Online Radio	222
16.2	Example 2: Predicting Income	227
	References	234

<b>17. Dimension Reduction: Factor Models and Principal Components</b>	<b>235</b>
17.1 Example 1: European Protein Consumption	238
17.2 Example 2: Monthly US Unemployment Rates	243
<b>18. Reducing the Dimension in Regressions with Multicollinear Inputs: Principal Components Regression and Partial Least Squares</b>	<b>247</b>
18.1 Three Examples	249
References	257
<b>19. Text as Data: Text Mining and Sentiment Analysis</b>	<b>258</b>
19.1 Inverse Multinomial Logistic Regression	259
19.2 Example 1: Restaurant Reviews	261
19.3 Example 2: Political Sentiment	266
Appendix 19.A Relationship Between the Gentzkow Shapiro Estimate of “Slant” and Partial Least Squares	268
References	271
<b>20. Network Data</b>	<b>272</b>
20.1 Example 1: Marriage and Power in Fifteenth Century Florence	274
20.2 Example 2: Connections in a Friendship Network	278
References	292
<b>Appendix A: Exercises</b>	<b>293</b>
Exercise 1	294
Exercise 2	294
Exercise 3	296
Exercise 4	298
Exercise 5	299
Exercise 6	300
Exercise 7	301
<b>Appendix B: References</b>	<b>338</b>
<b>Index</b>	<b>341</b>

# Introduction

Today's statistics applications involve enormous data sets: many cases (rows of a data spreadsheet, with a row representing the information on a studied case) and many variables (columns of the spreadsheet, with a column representing the outcomes on a certain characteristic across the studied cases). A case may be a certain item such as a purchase transaction, or a subject such as a customer or a country, or an object such as a car or a manufactured product. The information that we collect varies across the cases, and the explanation of this variability is central to the tools that we study in this book. Many variables are typically collected on each case, but usually only a few of them turn out to be useful. The majority of the collected variables may be irrelevant and represent just noise. It is important to find those variables that matter and those that do not.

Here are a few types of data sets that one encounters in data mining. In marketing applications, we observe the purchase decisions, made over many time periods, of thousands of individuals who select among several products under a variety of price and advertising conditions. Social network data contains information on the presence of links among thousands or millions of subjects; in addition, such data includes demographic characteristics of the subjects (such as gender, age, income, race, and education) that may have an effect on whether subjects are "linked" or not. Google has extensive information on 100 million users, and Facebook has data on even more. The recommender systems developed by firms such as Netflix and Amazon use available demographic information and the detailed purchase/rental histories from millions of customers. Medical data sets contain the outcomes of thousands of performed procedures, and include information on their characteristics such as the type of procedure and its outcome, and the location where and the time when the procedure has been performed.

While traditional statistics applications focus on relatively small data sets, data mining involves very large and sometimes enormous quantities of information. One talks about megabytes and terabytes of information. A megabyte represents a million bytes, with a byte being the number of bits needed to encode a single character of text. A typical English book in plain text format (500 pages with 2000

characters per page) amounts to about 1 MB. A terabyte is a million megabytes, and an exabyte is a million terabytes.

Data mining attempts to extract useful information from such large data sets. Data mining explores and analyzes large quantities of data in order to discover meaningful patterns. The *scale* of a typical data mining application, with its large number of cases and many variables, exceeds that of a standard statistical investigation. The analysis of millions of cases and thousands of variables also puts pressure on the *speed* that is needed to accomplish the search and modeling steps of the typical data mining application. This is why researchers refer to data mining as statistics at scale and speed. The large scale (lots of available data) and the requirements on speed (solutions are needed quickly) create a large demand for automation. Data mining uses a combination of pattern-recognition rules, statistical rules, as well as rules drawn from machine learning (an area of computer science).

Data mining has wide applicability, with applications in intelligence and security analysis, genetics, the social and natural sciences, and business. Studying which buyers are more likely to buy, respond to an advertisement, declare bankruptcy, commit fraud, or abandon subscription services are of vital importance to business.

Many data mining problems deal with categorical outcome data (e.g., no/yes outcomes), and this is what makes machine learning methods, which have their origins in the analysis of categorical data, so useful. Statistics, on the other hand, has its origins in the analysis of continuous data. This makes statistics especially useful for correlation-type analyses where one sifts through a large number of correlations to find the largest ones.

The analysis of large data sets requires an efficient way of storing the data so that it can be accessed easily for calculations. Issues of data warehousing and how to best organize the data are certainly very important, but they are not emphasized in this book. The book focuses on the analysis tools and targets their statistical foundation.

Because of the often enormous quantities of data (number of cases/replicates), the role of traditional statistical concepts such as confidence intervals and statistical significance tests is greatly reduced. With large data sets, almost any small difference becomes significant. It is the problem of overfitting models (i.e., using more explanatory variables than are actually needed to predict a certain phenomenon) that becomes of central importance. Parsimonious representations are important as simpler models tend to give more insight into a problem. Large models overfitted on training data sets usually turn out to be extremely poor predictors in new situations as unneeded predictor variables increase the prediction error variance. Furthermore, overparameterized models are of little use if it is difficult to collect data on predictor variables in the future. Methods that help avoid such overfitting are needed, and they are covered in this book. The partitioning of the data into training and evaluation (test) data sets is central to most data mining methods. One must always check whether the relationships found in the training data set will hold up in the future.

Many data mining tools deal with problems for which there is no designated response that one wants to predict. It is common to refer to such analysis as *unsupervised learning*. Cluster analysis is one example where one uses feature (variable) data on numerous objects to group the objects (i.e., the cases) into a

smaller number of groups (also called *clusters*). Dimension reduction applications are other examples for such type of problems; here one tries to reduce the many features on an object to a manageable few. Association rules also fall into this category of problems; here one studies whether the occurrence of one feature is related to the occurrence of others. Who would not want to know whether the sales of chips are being “lifted” to a higher level by the concurrent sales of beer?

Other data mining tools deal with problems for which there is a designated response, such as the volume of sales (a quantitative response) or whether someone buys a product (a categorical response). One refers to such analysis as *supervised learning*. The predictor variables that help explain (predict) the response can be quantitative (such as the income of the buyer or the price of a product) or categorical (such as the gender and profession of the buyer or the qualitative characteristics of the product such as new or old). Regression methods, regression trees, and nearest neighbor methods are well suited for problems that involve a continuous response. Logistic regression, classification trees, nearest neighbor methods, discriminant analysis (for continuous predictor variables) and naïve Bayes methods (mostly for categorical predictor variables) are well suited for problems that involve a categorical response.

Data mining should be viewed as a *process*. As with all good statistical analyses, one needs to be clear about the purpose of the analysis. Just to “mine data” without a clear purpose, without an appreciation of the subject area, and without a modeling strategy will usually not be successful. The data mining process involves several interrelated steps:

1. Efficient data storage and data preprocessing steps are very critical to the success of the analysis.
2. One needs to select appropriate response variables and decide on the number of variables that should be investigated.
3. The data needs to be screened for outliers, and missing values need to be addressed (with missing values either omitted or appropriately imputed through one of several available methods).
4. Data sets need to be partitioned into training and evaluation data sets. In very large data sets, which cannot be analyzed easily as a whole, data must be sampled for analysis.
5. Before applying sophisticated models and methods, the data need to be visualized and summarized. It is often said that a picture is worth a 1000 words. Basic graphs such as line graphs for time series, bar charts for categorical variables, scatter plots and matrix plots for continuous variables, box plots and histograms (often after stratification on useful covariates), maps for displaying correlation matrices, multidimensional graphs using color, trellis graphs, overlay plots, tree maps for visualizing network data, and geo maps for spatial data are just a few examples of the more useful graphical displays. In constructing good graphs, one needs to be careful about the right scaling, the correct labeling, and issues of stratification and aggregation.
6. Summary of the data involves the typical summary statistics such as mean, percentiles and median, standard deviation, and correlation, as well as more advanced summaries such as principal components.



7. Appropriate methods from the data mining tool bag need to be applied. Depending on the problem, this may involve regression, logistic regression, regression/classification trees, nearest neighbor methods,  $k$ -means clustering, and so on.
8. The findings from these models need to be confirmed, typically on an evaluation (test or holdout) data set.
9. Finally, the insights one gains from the analysis need to be implemented. One must act on the findings and spring to action. This is what W.E. Deming had in mind when he talked about process improvement and his Deming (Shewhart) wheel of “plan, do, check, and act” (Ledolter and Burrill, 1999).

Some data mining applications require an enormous amount of effort to just collect the relevant information. For example, an investigation of Pre-Civil War court cases of Missouri slaves seeking their freedom involves tedious study of handwritten court proceedings and Census records, electronic scanning of the records, and the use of character-recognition software to extract the relevant characteristics of the cases and the people involved. The process involves double and triple checking unclear information (such as different spellings, illegible entries, and missing information), selecting the appropriate number of variables, categorizing text information, and deciding on the most appropriate coding of the information. At the end, one will have created a fairly good master list of all available cases and their relevant characteristics. Despite all the diligent work, there will be plenty of missing information, information that is in error, and way too many variables and categories than are ultimately needed to tell the story behind the judicial process of gaining freedom.

Data preparation often takes a lot more time than the eventual modeling. The subsequent modeling is usually only a small component of the overall effort; quite often, relatively simple methods and a few well-constructed graphs can tell the whole story. It is the creation of the master list that is the most challenging task. The steps that are involved in the construction of the master list in such problems depend heavily on the subject area, and one can only give rough guidelines on how to proceed. It is also difficult to make this process automatic. Furthermore, even if some of the “data cleaning” steps can be made automatic, the investigator must constantly check and question any adjustments that are being made. Great care, lots of double and triple checking, and much common sense are needed to create a reliable master list. But without a reliable master list, the findings will be suspect, as we know that wrong data usually lead to wrong conclusions. The old saying “garbage in—garbage out” also applies to data mining.

Fortunately many large business data sets can be created almost automatically. Much of today’s business data is collected for transactional purposes, that is, for payment and for shipping. Examples of such data sets are transactions that originate from scanner sales in super markets, telephone records that are collected by mobile telephone providers, and sales and rental histories that are collected by companies such as Amazon and Netflix. In all these cases, the data collection effort is minimal,