

## Exploratory Data Mining and Data Cleaning

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: *David J. Balding, Peter Bloomfield, Noel A. C. Cressie,  
Nicholas I. Fisher, Iain M. Johnstone, J. B. Kadane, Louise M. Ryan,  
David W. Scott, Adrian F. M. Smith, Jozef L. Teugels;*

Editors Emeriti: *Vic Barnett, J. Stuart Hunter, David G. Kendall*

A complete list of the titles in this series appears at the end of this volume.

# Exploratory Data Mining and Data Cleaning

TAMRAPARNI DASU

THEODORE JOHNSON

AT&T Labs, Research Division  
Florham Park, NJ



A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2003 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.  
Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400, fax 978-750-4470, or on the web at [www.copyright.com](http://www.copyright.com). Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, e-mail: [permreq@wiley.com](mailto:permreq@wiley.com).

**Limit of Liability/Disclaimer of Warranty:** While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993 or fax 317-572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print, however, may not be available in electronic format.

***Library of Congress Cataloging-in-Publication Data:***

Dasu, Tamraparni.

Exploratory data mining and data cleaning / Tamraparni Dasu, Theodor Johnson.

p. cm.

Includes bibliographical references and index.

ISBN 0-471-26851-8 (cloth)

1. Data mining. 2. Electronic data processing—Data preparation. 3. Electronic data processing—Quality control. I. Johnson, Theodore. II. Title.

QA76.9.D343 D34 2003

006.3—dc21

2002191085

Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

# Contents

<b>Preface</b>	<b>ix</b>
<b>1. Exploratory Data Mining and Data Cleaning: An Overview</b>	<b>1</b>
1.1 Introduction, 1	
1.2 Cautionary Tales, 2	
1.3 Taming the Data, 4	
1.4 Challenges, 4	
1.5 Methods, 6	
1.6 EDM, 7	
1.6.1 EDM Summaries—Parametric, 8	
1.6.2 EDM Summaries—Nonparametric, 9	
1.7 End-to-End Data Quality (DQ), 12	
1.7.1 DQ in Data Preparation, 13	
1.7.2 EDM and Data Glitches, 13	
1.7.3 Tools for DQ, 14	
1.7.4 End-to-End DQ: The Data Quality Continuum, 14	
1.7.5 Measuring Data Quality, 15	
1.8 Conclusion, 16	
<b>2. Exploratory Data Mining</b>	<b>17</b>
2.1 Introduction, 17	
2.2 Uncertainty, 19	
2.2.1 Annotated Bibliography, 23	
2.3 EDM: Exploratory Data Mining, 23	
2.4 EDM Summaries, 25	
2.4.1 Typical Values, 26	
2.4.2 Attribute Variation, 33	

- 2.4.3 Example, 41
- 2.4.4 Attribute Relationships, 42
- 2.4.5 Annotated Bibliography, 49
- 2.5 What Makes a Summary Useful?, 50
  - 2.5.1 Statistical Properties, 51
  - 2.5.2 Computational Criteria, 54
  - 2.5.3 Annotated Bibliography, 54
- 2.6 Data-Driven Approach—Nonparametric Analysis, 54
  - 2.6.1 The Joy of Counting, 55
  - 2.6.2 Empirical Cumulative Distribution Function (ECDF), 57
  - 2.6.3 Univariate Histograms, 59
  - 2.6.4 Annotated Bibliography, 61
- 2.7 EDM in Higher Dimensions, 62
- 2.8 Rectilinear Histograms, 62
- 2.9 Depth and Multivariate Binning, 64
  - 2.9.1 Data Depth, 65
  - 2.9.2 Aside: Depth-Related Topics, 66
  - 2.9.3 Annotated Bibliography, 68
- 2.10 Conclusion, 68

### 3. Partitions and Piecewise Models

69

- 3.1 Divide and Conquer, 69
  - 3.1.1 Why Do We Need Partitions?, 70
  - 3.1.2 Dividing Data, 71
  - 3.1.3 Applications of Partition-Based EDM Summaries, 73
- 3.2 Axis-Aligned Partitions and Data Cubes, 74
  - 3.2.1 Annotated Bibliography, 77
- 3.3 Nonlinear Partitions, 77
  - 3.3.1 Annotated Bibliography, 78
- 3.4 DataSpheres (DS), 78
  - 3.4.1 Layers, 79
  - 3.4.2 Data Pyramids, 81
  - 3.4.3 EDM Summaries, 82
  - 3.4.4 Annotated Bibliography, 82
- 3.5 Set Comparison Using EDM Summaries, 82
  - 3.5.1 Motivation, 83
  - 3.5.2 Comparison Strategy, 83
  - 3.5.3 Statistical Tests for Change, 84

3.5.4	Application—Two Case Studies, 85	
3.5.5	Annotated Bibliography, 88	
3.6	Discovering Complex Structure in Data with EDM Summaries, 89	
3.6.1	Exploratory Model Fitting in Interactive Response Time, 89	
3.6.2	Annotated Bibliography, 90	
3.7	Piecewise Linear Regression, 90	
3.7.1	An Application, 92	
3.7.2	Regression Coefficients, 92	
3.7.3	Improvement in Fit, 94	
3.7.4	Annotated Bibliography, 94	
3.8	One-Pass Classification, 95	
3.8.1	Quantile-Based Prediction with Piecewise Models, 95	
3.8.2	Simulation Study, 96	
3.8.3	Annotated Bibliography, 98	
3.9	Conclusion, 98	
<b>4.</b>	<b>Data Quality</b>	<b>99</b>
4.1	Introduction, 99	
4.2	The Meaning of Data Quality, 102	
4.2.1	An Example, 102	
4.2.2	Data Glitches, 103	
4.2.3	Conventional Definition of DQ, 105	
4.2.4	Times Have Changed, 106	
4.2.5	Annotated Bibliography, 108	
4.3	Updating DQ Metrics: Data Quality Continuum, 108	
4.3.1	Data Gathering, 109	
4.3.2	Data Delivery, 110	
4.3.3	Data Monitoring, 113	
4.3.4	Data Storage, 116	
4.3.5	Data Integration, 118	
4.3.6	Data Retrieval, 120	
4.3.7	Data Mining/Analysis, 121	
4.3.8	Annotated Bibliography, 123	
4.4	The Meaning of Data Quality Revisited, 123	
4.4.1	Data Interpretation, 124	
4.4.2	Data Suitability, 124	
4.4.3	Dataset Type, 124	

4.4.4	Attribute Type, 128	
4.4.5	Application Type, 129	
4.4.6	Data Quality—A Many Splendored Thing, 129	
4.4.7	Annotated Bibliography, 130	
4.5	Measuring Data Quality, 130	
4.5.1	DQ Components and Their Measurement, 131	
4.5.2	Combining DQ Metrics, 134	
4.6	The DQ Process, 134	
4.7	Conclusion, 136	
4.7.1	Four Complementary Approaches, 136	
4.7.2	Annotated Bibliography, 137	
<b>5.</b>	<b>Data Quality: Techniques and Algorithms</b>	<b>139</b>
5.1	Introduction, 139	
5.2	DQ Tools Based on Statistical Techniques, 140	
5.2.1	Missing Values, 141	
5.2.2	Incomplete Data, 144	
5.2.3	Outliers, 146	
5.2.4	Detecting Glitches Using Set Comparison, 151	
5.2.5	Time Series Outliers: A Case Study, 154	
5.2.6	Goodness-of-Fit, 160	
5.2.7	Annotated Bibliography, 161	
5.3	Database Techniques for DQ, 162	
5.3.1	What is a Relational Database?, 162	
5.3.2	Why Are Data Dirty?, 165	
5.3.3	Extraction, Transformation, and Loading (ETL), 166	
5.3.4	Approximate Matching, 168	
5.3.5	Database Profiling, 172	
5.3.6	Annotated Bibliography, 175	
5.4	Metadata and Domain Expertise, 176	
5.4.1	Lineage Tracing, 179	
5.4.2	Annotated Bibliography, 179	
5.5	Measuring Data Quality?, 180	
5.5.1	Inventory Building—A Case Study, 180	
5.5.2	Learning and Recommendations, 186	
5.6	Data Quality and Its Challenges, 188	
	<b>Bibliography</b>	<b>189</b>
	<b>Index</b>	<b>197</b>



# Preface

As data analysts at a large information-intensive business, we often have been asked to analyze new (to us) data sets. This experience was the original motivation for our interest in the topics of exploratory data mining and data quality. Most data mining and analysis techniques assume that the data have been joined into a single table and cleaned, and that the analyst already knows what she or he is looking for. Unfortunately, the data set is usually dirty, composed of many tables, and has unknown properties. Before any results can be produced, the data must be cleaned and explored—often a long and difficult task.

Current books on data mining and analysis usually focus on the last stage of the analysis process (getting the results) and spend little time on how data exploration and cleaning is done. Usually, their primary aim is to discuss the efficient implementation of the data mining algorithms and the interpretation of the results. However, the true challenges in the task of data mining are:

- Creating a data set that contains the relevant and accurate information, and
- Determining the appropriate analysis techniques.

In our experience, the tasks of exploratory data mining and data cleaning constitute 80% of the effort that determines 80% of the value of the ultimate data mining results. Data mining books (a good one is [56]) provide a great amount of detail about the analytical process and advanced data mining techniques. However they assume that the data has already been gathered, cleaned, explored, and understood.

As we gained experience with exploratory data mining and data quality issues, we became involved in projects in which data quality improvement was the goal of the project (i.e., for operational databases) rather than a prerequisite. Several books recently have been published on the topic of ensuring data quality (e.g., the books by Loshin [84], by Redman [107]), and by English [41]). However, these books are written for managers and take a

managerial viewpoint. While the problem of ensuring data quality requires a significant managerial support, there is also a need for technical and analytic tools. At the time of this writing, we have not seen any organized exposition of the technical aspects of data quality management. The most closely related book is Pyle [102], which discusses data preparation for data mining. However, this text has little discussion of data quality issues or of exploratory data mining—pre-requisites even to preparing data for data mining.

Our focus in this book is to develop a systematic process of data exploration and data quality management. We have found these seemingly unrelated topics to be inseparable. The exploratory phase of any data analysis project inevitably involves sorting out data quality problems, and any data quality improvement project inevitably involves data exploration. As a further benefit, data exploration sheds light on appropriate analytic strategies.

Data quality is a notoriously messy problem that refuses to be put into a neat container, and therefore is often viewed as technically intractable. We have found that data quality problems can be addressed, but doing so requires that we draw on methods from many disciplines: statistics, exploratory data mining (EDM), databases, management, and metadata. Our focus in this book is to present an integrated approach to EDM and data quality. Because of the very broad nature of the subject, the exposition tends to be a summarization of material discussed in great detail elsewhere (for which we provide references), with an emphasis on how the techniques relate to each other and to EDM and data quality. Some topics (such as data quality metrics and certain aspects of EDM) have no other good source, so we discuss them in greater detail.

## **EXPLORATORY DATA MINING (EDM)**

Data sets of the twenty-first century are different from the ones that motivated analytical techniques of statistics, machine learning and others. Earlier data sets were reasonably small and relatively homogeneous so that the structure in them could be captured with compact models that had large but a manageable number of parameters. Many researchers have focused on scaling the methods to run efficiently and quickly on the much larger data sets collected by automated devices. In addition, methods have been developed specifically for massive data (i.e., data mining techniques). However, there are two fundamental issues that need to be addressed before these methods can be applied.

- A “data set” is often a patchwork of data collected from many sources, which might not have been designed for integration. One example of this problem is when two corporate entities providing different services to a common customer base merge to become a single entity. Another is when different divisions of a “federation enterprise” need to merge their data

stores. In such situations, approximate matching heuristics are used to combine the data. The resulting patchwork data set will have many data quality issues that need to be addressed. The data are likely to contain many other data glitches, and these need to be treated as well.

- Data mining methods often do not focus on the “appropriateness of the model for the data,” namely, goodness-of-fit. While finding the best model in a given class of models is desirable, it is equally important to determine the class of models that best fits the data.

There is no simple or single method for analyzing a complex, unfamiliar data set. The task typically requires the sequential application of disparate techniques, leveraging the additional information acquired at each stage to converge to a powerful, accurate and fast method. The end-product is often a “piecewise technique” where at each stage we might have had to adapt or extend, to improvise on an existing method. The importance of such an approach has been emphasized by statisticians such as John Tukey [123] and more recently in the machine learning community, for instance, in the Auto-Class project [19].

## DATA QUALITY

A major confounding factor in EDM is the presence of data quality issues. These are often unearthed as “interesting patterns” but on closer examination prove to be artifacts. We emphasize this aspect in our case study, since typically data analysts spend a significant portion of their time weeding-out data quality problems. No matter how sophisticated the data mining techniques, bad data will lead to misleading findings.

While most practitioners of data analysis are aware of the pitfalls of data quality issues, it is only recently that there has been an emphasis on the systematic detection and removal of data problems. There have been efforts directed at managing processes that generate the data, at cleaning up databases (e.g. merging/purging of duplicates), and at finding tools and algorithms for the automatic detection of data glitches. Statistical methods for process control (predominantly univariate) that date back to quality control charts developed for detecting batches of poorly produced lots in industrial manufacturing are often adapted to monitor fluctuations in variables that populate databases.

For operations databases, data quality is an end in itself. Most business (and governmental, etc.) processes involve complex interactions between many databases. Data quality problems can have very expensive manifestations (e.g., “losing” a cross-country cable, forgetting to bill customers). In this electronic age, many businesses (and governmental organizations, etc.) would like to “e-enable” their customers—that is, let them examine the relevant parts of the

operational databases to manage their own accounts. Depending on the state of the underlying databases, this can be embarrassing or even impossible.

## SUMMARY

In this book, we intend to:

- Focus on developing a modeling strategy through an iterative data exploration loop and incorporation of domain knowledge;
- Address methods for dealing with data quality issues that can have a significant impact on findings and decisions, using commercially available tools as well as new algorithmic approaches;
- Emphasize application in real-life scenarios throughout the narrative with examples;
- Highlight new approaches and methodologies, such as the *DataSphere* space partitioning and summary-based analysis techniques, and approaches to developing data quality metrics.

The book is intended for serious data analysts everywhere that need to analyze large amounts of unfamiliar, potentially noisy data, and for managers of operations databases. It can also serve as a text on data quality to supplement an advanced undergraduate or graduate level course in large-scale data analysis and data mining. The book is especially appropriate for a cross-disciplinary course in statistics and computer science.

## ACKNOWLEDGMENTS

We wish to thank the following people who have contributed to the material in this book: Deepak Agarwal, Dave Belanger, Bob Bell, Simon Byers, Corinna Cortes, Ken Church, Christos Faloutsos, Mary Fernandez, Joel Gottlieb, Andrew Hume, Nick Koudas, Eleftheris Koutsofios, Bala Krishnamurthy, Ken Lyons, David Poole, Daryl Pregibon, Matthew Roughan, Gregg Vesonder, and Jon Wright.