# Java Data Mining: Strategy, Standard, and Practice

# The Morgan Kaufmann Series in Data Management Systems
Series Editor: Jim Gray, Microsoft Research

# Java Data Mining: Strategy, Standard, and Practice

## A Practical Guide for Architecture, Design, and Implementation

**Mark F. Hornick**
**Erik Marcadé**
**Sunil Venkayala**

This book is printed on acid-free paper.

# Contents

# Preface

The birth of a standard is an amazing event, highlighting the ability of individuals from vastly different and often competing companies to come together to design an interface for a domain such as data mining. For JSR-73, we drew on experts from data mining tool and application vendors, as well as users of data mining technology. Data mining, as a field, is remarkably diverse in scope, encompassing capabilities from a broad range of disciplines: artificial intelligence, machine learning, statistics, data analysis, and visualization. Producing a standard in such a space is a challenging and fascinating adventure.

Within a year or so of embarking on the JDM 1.0 standard, various expert group members suggested that we'd have to write a book about Java Data Mining (JDM) someday. And indeed, here we are. Our main motivation for writing this book is to introduce data mining to a much broader audience, one that may have never used or encountered data mining before. As such, we focus less on the technical and scholarly details of data mining than on its practical understanding and application. We have tried to include a reasonably broad set of references for individuals who want to dive down to the next level of detail. However, we have strived to make data mining concepts, process, and use through JDM more accessible to Java developers, who usually do not encounter data mining, and the colleagues they will work with to develop advanced analytic applications.

Advanced analytic applications—those augmented with advanced and predictive analytics such as data mining—provide greater business intelligence, yielding insight into business problems and guidance for improved decision making. Such applications are becoming most valuable to businesses, and hence can increase revenue and profits—both for the vendors who sell them and for the businesses that use them.

Readers of this book will find a somewhat unconventional approach to data mining. Other books on data mining provide much detail on algorithms and techniques. Although this information is important to those studying machine learning or wanting to become a data analyst, other potential users of data mining are left wondering how these algorithms or techniques will be applied to solve problems. As vendors of data mining technology strive to make data mining more accessible to a broader range of users, such as business analysts, information technology (IT) specialists, and database administrators (DBAs), it is no longer the details that users require, but the big picture. Users ask, "How can I use this powerful technology to provide value within my business?" In this book, we strive to approach this and other questions from several perspectives: the software developer, the software and systems architect, and the business and data analyst. We explore these perspectives in the following section, "Guide to Readers."

In this book, we provide insight into three key aspects of the Java Data Mining standard. The first aspect, covered in Part I, focuses on strategies for solving data mining–related business and scientific problems, and on the strategy the JDM Expert Group pursued in the design of the JDM standard. After an introduction to the data mining field, we discuss solving problems in various industries using data mining technology.

Although every industry has unique problems to solve, requiring custom and innovative solutions, each industry also shares many problems that can benefit from cross-industry solutions. For example, industries such as retail, financial services, and healthcare, as well as the public sector, all have customers. The cross-industry solution spaces include customer acquisition, customer retention, customer lifetime value, and targeted marketing.

Because data mining solutions typically do not take form or produce value in a vacuum, we then discuss the overall process, based on the industry standard data mining process CRISP-DM. Because users of data mining technology need to be minimally conversant in the terminology and concepts to problem solve with their colleagues,

we introduce the mining functions and algorithms defined in JDM. With this foundation, we explore the JDM strategy, answering questions such as: What drove the design of JDM? What is the role of standards? Lastly, before embarking on details of the Java Data Mining standard, we provide a "getting started" code example that follows the CRISP-DM process.

The second aspect, covered in Part II, focuses on the standard itself. This part introduces various concepts defined by or assimilated into the standard using examples based on business problems. After this, we explore the design of the JDM API and more detailed code examples to give readers a better understanding of how to use JDM to build applications and solve problems. Although JDM is foremost dedicated to being a standard Java language API, Java Data Mining also defines an XML schema representation for JDM objects, as well as a web services interface to enable the use of JDM functionality in a Services Oriented Architecture (SOA) environment. Part II also discusses these with specific examples of their use.

The third aspect, covered in Part III, focuses on using JDM in practice, building applications and tools that use the Java Data Mining API. We begin this part with several business scenarios (e.g., targeted marketing, key factor analysis, and customer segmentation). Because JDM is designed to be used by both application designers and data mining tool designers, we introduce code for building a simple tool graphical user interface (GUI), which manipulates JDM-persistent objects as well as enables the building and testing of a model. Having introduced web services in Part II, we give an example of a web services based application. Since data mining can impact the Information Technology (IT) infrastructure of most companies, we explore the impact of data mining along several dimensions, including hardware, software, data access, performance, and administration tools. Since the practice of using data mining often involves the use of commercial implementations, we introduce two such JDM implementations, from Oracle and KXEN. We also provide some guidelines or insights for implementers new to JDM.

Wrapping up in Part IV, we explore the evolution of data mining standards, which puts JDM in the broader context of other data mining standards. We also contrast the approaches taken by various data mining standards bodies. Since we note that no standard is ever complete, and JDM 1.1 itself covers only a subset of the possible data mining functions and algorithms, we highlight directions for JDM 2.0. We introduce features under consideration  such as transformations, time series, and apply for association models, among others.

## Acknowledgments

We first want to acknowledge the Java Data Mining expert group members who participated in the long process required to produce the JSR-73 standard. Their unwavering support through weekly conference calls and face-to-face meetings over the 4 years of the standards development is greatly appreciated. We also acknowledge the additional contributions of Hankil Yoon, Ka Kit Chan, Jim Dadashev, and Somesh Marepalli to the Technology Compatibility Kit (TCK) implementation, and Marwane Jai Lamimi to the Reference Implementation (RI).

We are very grateful for the general and specialist input provided by Frank Byrum, Jim Melton, and Osmar Zaiane on the developing manuscript. Over the past year, their detailed comments on both structure and content were a tremendous asset. We thank Jacek Myczkowski and Don Deutsch for their valuable comments on the final manuscript, as well as their support of the standards efforts for JSR-73 and JSR-247 at Oracle. We thank the JDM expert group members Michal Prussak, Alex Russakovsky, and Michael Smith who also provided valuable comments on the final manuscript, and David Urena and Samy Mechiri who contributed to the source code used in Part III of this book.

Of course, all remaining errors (which we expect exist despite careful review) are entirely our responsibility.

We offer our appreciation and gratitude to the wonderful people at Morgan Kaufmann Publishers as they guided us through the process of book writing and publishing. We thank Jim Melton, one of our reviewers, for putting us in contact with Diane Cerra, our talented and patient publisher, to begin this journey. We thank Diane, Asma Palmeiro, Misty Bergeron, Marilyn Rash, and Bruce Siebert who worked to make this book possible.

# Guide to Readers

Data mining is becoming a mainstream technology used in business intelligence applications supporting industries such as financial services, retail, healthcare, telecommunications, and higher education, and lines of business such as marketing, manufacturing, customer experiences, customer service, and sales. Many of the business problems that data mining can solve cut across industries such as customer retention and acquisition, cross-sell, and response modeling. Due to the cost, skillsets, and complexity required to bring data mining results into an established business process, early adopters were typically big companies and research labs with correspondingly large budgets and access to statisticians and machine learning experts. In recent years, however, data mining products have simplified data mining considerably by automating the process—making the fruits of the technology more widely accessible. New algorithms and heuristics have evolved to provide good results with little or no experimentation or data preparation. In addition, the availability of data mining has increased with in-database data mining capabilities.

Java Data Mining (JDM) furthers the adoption of data mining by providing a standard Java and web services Application Programming Interface (API) for data mining. This book introduces data mining to software developers and application architects who may have heard of the benefits of data mining but are unsure how to realize these benefits. This book is also targeted at business and data analysts

who want to learn how JDM helps in developing vendor-neutral data mining solutions. It does not require a reader to be familiar with data mining, statistics, or machine learning technologies.

We have organized this book into three main parts: *strategy, standard*, and *practice*. In Part I, JDM Strategy, we introduce data mining in general, uses of data mining in solving industry problems, data mining processes and techniques, the role of data mining standards, and a high-level introduction to the JDM Application Programming Interface (API). Most of this part doesn't require the reader to know the Java language.

In Part II, JDM Standard, we explain the concepts used in JDM by example, explore the JDM API design and its usage, and introduce the Java Data Mining XML schema and web services. This part requires readers to know the Java language, XML, and XML schema. It gives a brief introduction to web services in Chapter 11 before discussing the JDM web services.

In Part III, JDM Practice, we illustrate practical problem solving using the JDM API. We begin by developing a sample data mining tool using JDM and a sample data mining web service using JDM. We then introduce two JDM vendor implementations, exploring their functionality, architecture, and design tradeoffs before giving some guidance to others interested in implementing a JDM-compliant system.

In Part IV, Wrapping Up, we discuss the evolution of data mining standards, where they have been and where they might go. We give a preview of some of the features proposed for JDM 2.0.

## For the Software Developer

For software developers, and in particular Java and web services developers, this book introduces data mining and how to use JDM to develop data mining solutions. Part I introduces data mining and various types of business problems that can be solved using data mining, illustrates a standard process used to conduct a data mining project, describes data mining techniques used to solve business problems, explains the JDM standard strategy and why the JDM standard is necessary, and provides an overview of the JDM API. Even though software developers are not typically involved in the initial solving of a data mining problem, it is important to know

about concepts to understand the JDM API and how to develop data mining solutions.

Part II will familiarize developers with JDM concepts and the API. Readers of this part are required to know the Java language, Object Oriented Programming, the Unified Modeling Language and XML to understand the Java examples, API design concepts, JDM XML schema, and web services. This part introduces JDM concepts using examples, describes the design and usage of the Java API, and illustrates the Java Data Mining XML schema and web services interfaces.

Part III describes the use of the JDM API in practice with sample applications and detailed code examples both for the Java and web services API. It also provides JDM vendor implementation details and explains the process for other data mining vendors in adopting the JDM standard.

After reading this book, we expect the data mining knowledge gap between developers and data analysts will be greatly reduced to help them communicate more effectively when developing a data mining solution.

## For the Software Architect

Data mining is often integrated with existing software applications and business processes. Understanding of data mining processes provides greater insight for architects to enable this technology in existing or new applications. For example, an architect needs to understand how data mining works to add intelligent customer offers using data mining to an existing call center application.

For architects who want to be hands-on with the JDM API (e.g., to develop prototypes), all parts of this book are useful. Part I and Part III are particularly useful for architects. Part I introduces data mining in general and provides examples of how it is currently being applied to solve business problems. Most important, it introduces the data mining process and the role of the information technology department in implementing a data mining project.

Part II will be useful to understand the API-level concepts for the architects who want to be hands-on with the API, to develop prototypes, or to mentor the developers about the use of the API.

In Part III, we provide deeper insight into how JDM can be used in practice. Chapter 16, which discusses vendor implementations, is

particularly useful for data mining software architects who are interested in developing JDM compatible API's and extensions.

After reading this book, architects should be comfortable in integrating JDM-based data mining solutions with their applications and be able to develop a strategy to operationalize data mining results with their existing applications.

## For the Business/Data Analyst

For business and data analysts who want to extract actionable information from corporate data, this book provides an introduction to data mining and how it is used to solve various business problems across industries. In Part I, the data mining usage scenarios and process of implementing a data mining project will be particularly useful for the analyst unfamiliar with data mining. Chapter 5, JDM Strategy, enables analysts to understand why the JDM standard is important in implementing a data mining solution. Typically, analysts are not involved in the software implementation of the solution, yet Part II may be useful for understanding the data mining concepts used by JDM to facilitate communication with developers and data mining experts, and for using tools based on JDM.

For an analyst who is already familiar with data mining and who has expertise in data mining and statistics, this book gives details of Java Data Mining and its usage in developing data mining solutions. Data mining tools can often generate JDM-compatible code to easily deploy a solution to a JDM-compatible Data Mining Engine (DME).

After reading this book, an analyst previously unfamiliar with data mining should be able to better understand how data mining can help in solving business problems. A data mining expert analyst will be able to understand the supported data mining features in JDM and be able to communicate easily with the software architects/ developers to implement a data mining solution.