

Chapter I

Introduction to Data Mining Techniques via Multiple Criteria Optimization Approaches and Applications

Yong Shi

*University of the Chinese Academy of Sciences, China
and University of Nebraska at Omaha, USA*

Yi Peng

University of Nebraska at Omaha, USA

Gang Kou

University of Nebraska at Omaha, USA

Zhengxin Chen

University of Nebraska at Omaha, USA

ABSTRACT

This chapter provides an overview of a series of multiple criteria optimization-based data mining methods, which utilize multiple criteria programming (MCP) to solve data mining problems, and outlines some research challenges and opportunities for the data mining community. To achieve these goals, this chapter first introduces the basic notions and mathematical formulations for multiple criteria optimization-based classification models, including the multiple criteria linear programming model, multiple criteria quadratic programming model, and multiple criteria fuzzy linear programming model. Then it presents the real-life applications of these models in credit card scoring management, HIV-1 associated dementia (HAD) neuronal damage and dropout, and network intrusion detection. Finally, the chapter discusses research challenges and opportunities.

INTRODUCTION

Data mining has become a powerful information technology tool in today's competitive business world. As the sizes and varieties of electronic data-sets grow, the interest in data mining is increasing rapidly. Data mining is established on the basis of many disciplines, such as machine learning, databases, statistics, computer science, and operations research. Each field comprehends data mining from its own perspective and makes its distinct contributions. It is this multidisciplinary nature that brings vitality to data mining. One of the application roots of data mining can be regarded as statistical data analysis in the pharmaceutical industry. Nowadays the financial industry, including commercial banks, has benefited from the use of data mining. In addition to statistics, decision trees, neural networks, rough sets, fuzzy sets, and vector support machines have gradually become popular data mining methods over the last 10 years. Due to the difficulty of accessing the accuracy of hidden data and increasing the predicting rate in a complex large-scale database, researchers and practitioners have always desired to seek new or alternative data mining techniques. This is a key motivation for the proposed multiple criteria optimization-based data mining methods.

The objective of this chapter is to provide an overview of a series of multiple criteria optimization-based methods, which utilize the multiple criteria programming (MCP) to solve classification problems. In addition to giving an overview, this chapter lists some data mining research challenges and opportunities for the data mining community. To achieve these goals, the next section introduces the basic notions and mathematical formulations for three multiple criteria optimization-based classification models: the multiple criteria linear programming model, multiple criteria quadratic programming model, and multiple criteria fuzzy linear programming model. The third section presents some real-life applications of these models, including credit card

scoring management, classifications on HIV-1 associated dementia (HAD) neuronal damage and dropout, and network intrusion detection. The chapter then outlines research challenges and opportunities, and the conclusion is presented.

MULTIPLE CRITERIA OPTIMIZATION-BASED CLASSIFICATION MODELS

This section explores solving classification problems, one of the major areas of data mining, through the use of multiple criteria mathematical programming-based methods (Shi, Wise, Luo, & Lin, 2001; Shi, Peng, Kou, & Chen, 2005). Such methods have shown its strong applicability in solving a variety of classification problems (e.g., Kou et al., 2005; Zheng et al., 2004).

Classification

Although the definition of classification in data mining varies, the basic idea of classification can be generally described as to "predicate the most likely state of a categorical variable (the class) given the values of other variables" (Bradley, Fayyad, & Mangasarian, 1999, p. 6). Classification is a two-step process. The first step constructs a predictive model based on training dataset. The second step applies the predictive model constructed from the first step to testing dataset. If the classification accuracy of testing dataset is acceptable, the model can be used to predicate unknown data (Han & Kamber, 2000; Olson & Shi, 2005).

Using the multiple criteria programming, the classification task can be defined as follows: *for a given set of variables in the database, the boundaries between the classes are represented by scalars in the constraint availabilities*. Then, the standards of classification are measured by minimizing the total overlapping of data and maximizing the distances of every data to its class boundary

23 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/introduction-data-mining-techniques-via/7543

Related Content

Evaluating NoSQL Databases for Big Data Processing within the Brazilian Ministry of Planning, Budget, and Management

Ruben C. Huacarpuma, Daniel da C. Rodrigues, Antonio M. Rubio Serrano, João Paulo C. Lustosa da Costa, Rafael T. de Sousa Júnior, Lizane Leite, Edward Ribeiro, Maristela Holandaand Aleteia P. F. Araujo (2016). *Big Data: Concepts, Methodologies, Tools, and Applications* (pp. 1110-1128). www.irma-international.org/chapter/evaluating-nosql-databases-for-big-data-processing-within-the-brazilian-ministry-of-planning-budget-and-management/150208

Combining BPSO and ELM Models for Inferring Novel lncRNA-Disease Associations

Wenqing Yang, Xianghan Zheng, QiongXia Huang, Yu Liu, Yimi Chenand ZhiGang Song (2023). *International Journal of Data Warehousing and Mining* (pp. 1-18). www.irma-international.org/article/combining-bpso-and-elm-models-for-inferring-novel-lncrna-disease-associations/317092

A Survey of Multidimensional Modeling Methodologies

Oscar Romeroand Alberto Abelló (2009). *International Journal of Data Warehousing and Mining* (pp. 1-23). www.irma-international.org/article/survey-multidimensional-modeling-methodologies/1824

Research on Data Mining and Investment Recommendation of Individual Users Based on Financial Time Series Analysis

Shiya Wang (2020). *International Journal of Data Warehousing and Mining* (pp. 64-80). www.irma-international.org/article/research-on-data-mining-and-investment-recommendation-of-individual-users-based-on-financial-time-series-analysis/247921

A BPMN-Based Design and Maintenance Framework for ETL Processes

Zineb El Akkaoui, Esteban Zimányi, Jose-Norberto Mazónand Juan Trujillo (2013). *International Journal of Data Warehousing and Mining* (pp. 46-72). www.irma-international.org/article/bpmn-based-design-maintenance-framework/78375