

Chapter 1.2

Knowledge Structure and Data Mining Techniques

Rick L. Wilson

Oklahoma State University, USA

Peter A. Rosen

University of Evansville, USA

Mohammad Saad Al-Ahmadi

Oklahoma State University, USA

INTRODUCTION

Considerable research has been done in the recent past that compares the performance of different data mining techniques on various data sets (e.g., Lim, Low, & Shih, 2000). The goal of these studies is to try to determine which data mining technique performs best under what circumstances. Results are often conflicting—for instance, some articles find that neural networks (NN) outperform both traditional statistical techniques and inductive learning techniques, but then the opposite is found with other datasets (Sen & Gibbs, 1994; Sung, Chang, & Lee, 1999; Spangler, May, & Vargas, 1999). Most of these studies use publicly available datasets in their analysis, and because they are

not artificially created, it is difficult to control for possible data characteristics in the analysis. Another drawback of these datasets is that they are usually very small.

With conflicting empirical results in the knowledge discovery/data mining literature, there have been numerous calls for a more systematic study of different techniques using synthetic, well-understood data. The rationale for synthetic data is that various factors can be manipulated while others are controlled, which may lead to a better understanding of why technique X outperforms technique Y in some, but not all, circumstances (Scott & Wilkins, 1999).

This call for research dates back to Quinlan's seminal work in inductive learning algorithms.

In his 1994 study that analyzed the difference between neural networks and inductive decision trees, Quinlan conjectures the existence of what he called S-problems and P-problems. In his definition, S-problems are those that are unsuited for NN's, while P-problems are those unsuited for decision tree induction. More recently, the review work on neural networks by Tickle, Maine, Bologna, Andrews, and Diederich (2000) propose that determining whether a classification task belongs to the P-problem or S-problem set is a very important research question.

Recently, other researchers have proposed that the composition of the underlying knowledge in a dataset, or knowledge structure (KS), may be pertinent in understanding why knowledge discovery techniques perform well on one dataset and poorly on others. This term has been used by Hand, Mannila, and Smyth (2001), and Padmanabhan and Tuzhilin (2003) to refer to this phenomenon, while Scott and Wilkins (1999) used a similar term, structural regularities, to describe the same concept.

The goal of this article is to explore in more detail how the existence of a database's underlying *knowledge structure* might help explain past inconsistent results in the knowledge discovery literature. Management scholars will recognize the term knowledge structure, as Walsh (1995) refers to it as a "mental template...imposed on an information environment to give it form and meaning." Therefore, for the knowledge discovery context, we propose that knowledge structure is analogous to the form and meaning of the knowledge to be discovered in a database. Though we will not explore the concept too deeply, one also can define knowledge structure through the use of a parameter set P as proposed by Hand et al. (2001). The parameter set would be attribute-value pairs that detail the existence of a specific knowledge structure for a given knowledge concept/database pair.

This knowledge structure concept is an abstract concept, which may make it hard to visualize.

Typically, when a knowledge worker is using a technique to extract knowledge from a database, they will not have any idea about the underlying knowledge structure of the concept of interest. But, researchers have hypothesized that knowledge discovery in a database is optimized when the formalism of the tool matches this underlying structure of the knowledge (Hand et al., 2001). Based on this, we conjecture that if a knowledge worker did know the knowledge structure parameter values prior to exploring the data, he or she could find the optimal tool for the knowledge discovery process.

From a historical perspective, past knowledge discovery and data mining research results could be explained by whether a particular knowledge discovery tool was or was not a good "match" with the underlying knowledge structure. The idea of matching the tool to the structure is somewhat analogous to the concept of task-technology fit, studied in the MIS literature during the mid 1990s (Goodhue, 1995).

Recent research in other related areas has found that contradictory or difficult to explain results could be related to the concept of knowledge structure (Wilson & Rosen, 2003). In this study, the well-known IRIS and BUPA Liver datasets were used to examine the efficacy of knowledge discovery tools in protected (by data perturbation) confidential databases. The IRIS dataset is known to possess linearly separable classes, while the BUPA Liver dataset cases has been historically difficult to correctly classify for all knowledge discovery tools. An outcome of this research was the proposal that knowledge discovery tool effectiveness in a protected (perturbed) database could be impacted by both the database's underlying knowledge structure and the *noise* present in the database. The concept of *noise* is simply the degree to which the different classes can be separated or differentiated by the optimal tool, or, alternatively, a surrogate measure of how difficult cases are to classify (e.g., Li & Wang, 2004).

7 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/knowledge-structure-data-mining-techniques/7628

Related Content

Conceptual and Systematic Design Approach for XML Document Warehouses

Vicky Nassis, R. Rajagopalapillai, Tharam S. Dillon and Wenny Rahayu (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 485-508).

www.irma-international.org/chapter/conceptual-systematic-design-approach-xml/7661

Exploiting Captions for Web Data Mining

Neil C. Rowe (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 1461-1485).

www.irma-international.org/chapter/exploiting-captions-web-data-mining/7710

Bayesian Networks

Ahmad Bashir, Latifur Khan and Mamoun Awad (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 89-93).

www.irma-international.org/chapter/bayesian-networks/10572

Privacy in Data Mining Textbooks

James Lawler and John C. Molluzzo (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 2872-2879).

www.irma-international.org/chapter/privacy-data-mining-textbooks/7808

Discovering Frequent Embedded Subtree Patterns from Large Databases of Unordered Labeled Trees

Yongqiao Xiao, Jenq-Foung Yao and Guizhen Yang (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 3235-3251).

www.irma-international.org/chapter/discovering-frequent-embedded-subtree-patterns/7832