

Chapter 5.4

Interesting Knowledge Patterns in Databases

Rajesh Natarajan

Indian Institute of Management Lucknow (IIML), India

B. Shekar

Indian Institute of Management Bangalore (IIMB), India

INTRODUCTION

Knowledge management (KM) transforms a firm's knowledge-based resources into a source of competitive advantage. Knowledge creation, a KM process, deals with the conversion of tacit knowledge to explicit knowledge and moving knowledge from the individual level to the group, organizational, and interorganizational levels (Alavi & Leidner, 2001). Four modes³—namely, socialization, externalization, combination, and internalization⁴—create knowledge through the interaction and interplay between tacit and explicit knowledge. The “combination” mode consists of combining or reconfiguring disparate bodies of existing explicit knowledge (like documents) that lead to the production of new explicit knowledge (Choo, 1998). Transactional databases are a source of rich information about a firm's processes and its business environment. Knowl-

edge Discovery in Databases (KDD), or data mining, aims at uncovering trends and patterns that would otherwise remain buried in a firm's operational databases. KDD is “the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data.” (Fayyad, Piatetsky-Shapiro, & Smyth, 1996). KDD is a typical example of IT-enabled combination mode of knowledge creation (Alavi & Leidner, 2001).

An important issue in KDD concerns the glut of patterns generated by any knowledge discovery system. The sheer number of these patterns makes manual inspection infeasible. In addition, one cannot obtain a good overview of the domain. Most of the discovered patterns are uninteresting since they represent well-known domain facts. The two problems—namely, rule quality and rule quantity—are interdependent. Knowledge of a rule's quality can help in reduc-

ing the number of rules. End-users of data mining outputs are typically managers, hard pressed for time. Hence, the need for automated methods to identify interesting, relevant, and significant patterns. This article discusses the interestingness of KDD patterns. We use the association rule (AR) (Agrawal, Imielinski, & Swami, 1993) in a market-basket context as an example of a typical KDD pattern. However, the discussions are also applicable to patterns like classification rules.

BACKGROUND

The Rule Quantity Problem: Solution Perspectives

The rule quantity problem may be a result of the automated nature of many KDD methods, such as AR mining methods. In one study, Brin, Motwani, Ullman, and Tsur (1997) discovered 23,712 rules on mining a census database. Approaches to alleviate this problem aim at reducing the number of rules required for examination while preserving relevant information present in the original set. Redundancy reduction, rule templates, incorporation of additional constraints, ranking, grouping, and visualization are some of the techniques that address the rule quantity problem.

In AR mining, additional constraints in conjunction with support and confidence thresholds can reveal specific relationships between items. These constraints reduce the search space and bring out fewer, relevant, and focused rules. Rule templates (Klemettinen, Mannila, Ronkainen, Toivonen, & Verkamo, 1994) help in selecting interesting rules by allowing a user to pre-specify the structure of interesting and uninteresting class of rules in inclusive and restrictive templates, respectively. Rules matching an inclusive template are interesting. Such templates are typical post-processing filters. Constraint-based mining (Bayardo, Agrawal, & Gunopulos, 2000) embeds user-specified rule constraints in the mining

process. These constraints eliminate any rule that can be simplified to yield a rule of equal or higher predictive ability. Association patterns like negative ARs (Savasere, Omiecinski, & Navathe, 1998; Subramanian, Ananthanarayana, & Narasimha Murty, 2003), cyclic ARs (Ozden, Sridhar, & Silberschatz, 1998), inter-transactional ARs (Lu, Feng, & Han, 2000), ratio rules (Korn, Labrinidis, Kotidis, & Faloutsos, 1998), and substitution rules (Teng, Hsieh, & Chen, 2002) bring out particular relationships between items. In the market-basket context, negative ARs reveal the set of items a customer is unlikely to purchase with another set. Cyclic association rules reveal purchases that display periodicity over time. Thus, imposition of additional constraints offers insight into the domain by discovering focused and tighter relationships. However, each method discovers a specific kind of behaviour. A large number of mined patterns might necessitate the use of other pruning methods. Except for rule templates, methods that enforce constraints are characterized by low user-involvement.

Redundancy reduction methods remove rules that do not convey new information. If many rules refer to the same feature of the data, then the most general rule may be retained. "Rule covers" (Toivonen, Klemettinen, Ronkainen, Hatonen, & Mannila, 1995) is a method that retains a subset of the original set of rules. This subset refers to all rows (in a relational database) that the original ruleset covered. Another strategy in AR mining (Zaki, 2000) is to determine a subset of frequently occurring closed item sets from their supersets. The magnitude of cardinality of the subset is several orders less than that of the superset. This implies fewer rules. This is done without any loss of information. Sometimes, one rule can be generated from another using a certain inference system. Retaining the basic rules may reduce the cardinality of the original rule set (Cristofor & Simovici, 2002). This process being reversible can generate the original ruleset if required. Care is taken to retain the information content

9 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/interesting-knowledge-patterns-databases/7997

Related Content

Biological Data Mining

George Tzani, Christos Berberidis and Ioannis Vlahavas (2005). *Encyclopedia of Database Technologies and Applications* (pp. 35-41).

www.irma-international.org/chapter/biological-data-mining/11119

The Quality of Online Privacy Policies: A Resource-Dependency Perspective

Veda C. Storey, Gerald C. Kane and Kathy Stewart Schwaig (2009). *Journal of Database Management* (pp. 19-37).

www.irma-international.org/article/quality-online-privacy-policies/3402

Methodology Evaluation Framework for Component-Based System Development

Ajantha Dahanayake, Henk Soland Zoran Stojanovic (2003). *Journal of Database Management* (pp. 1-26).

www.irma-international.org/article/methodology-evaluation-framework-component-based/3288

An Extended Relational Model & SQL for Fuzzy Multidatabases

Awadhesh Kumar Sharma, A. Goswami and D.K. Gupta (2011). *Advanced Database Query Systems: Techniques, Applications and Technologies* (pp. 185-219).

www.irma-international.org/chapter/extended-relational-model-sql-fuzzy/52302

A Measurement Ontology Generalizable for Emerging Domain Applications on the Semantic Web

Henry M. Kim, Arijit Sengupta, Mark S. Fox and Mehmet Dalkilic (2007). *Journal of Database Management* (pp. 20-42).

www.irma-international.org/article/measurement-ontology-generalizable-emerging-domain/3365