

# Template Guided Association Rule Mining from XML Documents

Rahman AliMohammadzadeh

Database Research Group  
Faculty of ECE,  
School of Engineering  
University of Tehran, Iran,  
+98-912-5068928

r.mohammadzadeh@ece.ut.ac.ir

Sadegh Soltan

Database Research Group  
Faculty of ECE,  
School of Engineering  
University of Tehran, Iran,  
+98-912-3092377

s.soltan@ece.ut.ac.ir

Masoud Rahgozar

Control and Intelligent Processing  
Center of Excellence,  
Faculty of ECE, School of Engineering,  
University of Tehran, Iran,  
+98-21-82084304

rahgozar@ut.ac.ir

## ABSTRACT

Compared with traditional association rule mining in the structured world (e.g. Relational Databases), mining from XML data is confronted with more challenges due to the inherent flexibilities of XML in both structure and semantics. The major challenges include 1) a more complicated hierarchical data structure; 2) an ordered data context; and 3) a much bigger size for each data element. In order to make XML-enabled association rule mining truly practical and computationally tractable, we propose a practical model for mining association rules from XML documents and demonstrate the usability and effectiveness of model through a set of experiments on real-life data.

## Categories and Subject Descriptors

H.2.8 [Database Applications]: Data Mining

## General Terms

Design, Experimentation, Performance

## Keywords

XML, Data Mining, Association Rule Mining.

## 1. INTRODUCTION

Data mining is usually used to extract interesting knowledge from large amounts of data stored in databases or data warehouses. This knowledge can be represented in many different ways such as clusters, decision trees, decision rules, etc. Among them, association rules have been proved effective in discovering interesting relations in massive amounts of data.

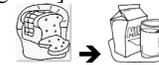
Currently, XML is penetrating virtually all areas of Internet application programming and is bringing about huge amount of data encoded in XML. With the continuous growth in XML data sources, the ability to extract knowledge from them for decision support becomes increasingly important and desirable [3]. Due to the inherent flexibilities of XML, in both structure and semantics, mining knowledge in the XML Era is faced with more challenges than in the traditional structured world.

In this paper, we propose a practical model for mining association rules from XML documents. Our model is based on XML-enabled association rule framework that was introduced by Feng [2]. XML-AR framework extends the notion of associated items to XML fragments to present associations among trees rather than simple-structured items of atomic values.

Although this framework is flexible and powerful enough to represent simple and complex structured association rules in XML documents [4] but to our best knowledge no implementation model has been proposed yet.

## 2. XML Association Rules

Association rules were first introduced by Agrawal et al. to analyze customer habits in retail databases. Association rule is an implication of the form  $X \Rightarrow Y$ , where the rule *body*  $X$  and *head*  $Y$  are subsets of the set  $I$  of items ( $I = \{I_1, I_2, \dots, I_n\}$ ) within a set of transactions  $D$  and  $X \cap Y = \emptyset$ . A rule  $X \Rightarrow Y$  states that the transactions  $T$  that contain the items in  $X$  are *likely* to contain also the items in  $Y$ . Association rules are characterized by two measures: the *support*, which measures the percentage of transactions in  $D$  that contain both items  $X$  and  $Y$ ; the *confidence*, which measures the percentage of transactions in  $D$  containing the items  $X$  that also contain the items  $Y$  [Figure 1]. In XML context, both  $D$  and  $I$  are collections of trees [1], in the same way  $X$  and  $Y$  are XML fragments [Figure 2].



[Support = 2%, Confidence = 95%]

Figure 1. Association rule between bread and milk

<Author> Rakesh Agrawal </author>

<Keyword> Data Mining </keyword>

Figure 2. XML Association rule

## 3. XML Association Rule Mining (Practical Model)

We consider the problem of mining XML association rules from content [3] of XML documents based on user provided rule template. We suggest an implementation model for the XML-AR framework that was introduced by Feng [2]. Our practical model consists of 5 steps (see Figure 6): Filtering, Generating Virtual Transactions, Finding Association Rules, Converting extracted rules to XML AR rules and Visualizing.

Filtering and Generating virtual transactions are most important steps in this model so we describe these two steps in more details. Filtering step uses the XML-AR template and extracts only those parts of XML that are interesting for the user. In the next step, we define a transaction context, based on tag nesting in XML document and use it to generate virtual transactions that can be used as input format by association rule mining algorithms (e.g. Apriori). As an example, consider the problem of mining frequent associations among people who appear as coauthors, with our

