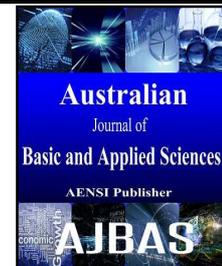




ISSN:1991-8178

## Australian Journal of Basic and Applied Sciences

Journal home page: www.ajbasweb.com



### A Novel Method to Managing Semi Structured Data in Distributed Environment using Modified Tree based Association Rules(TAR)

<sup>1</sup>E. Seshatheri and <sup>2</sup>Dr.T. Bhuvaneshwari

<sup>1</sup>Research Scholar Manonmaniam Sundaranar University Tirunelveli.

<sup>2</sup>Assistant Professor Department of Computer Science and Application L.N. Government Arts College Ponneri

#### ARTICLE INFO

##### Article history:

Received 10 October 2015

Accepted 30 November 2015

Available online 24 December 2015

##### Keywords:

XML document, Query-answering, Tree-based Association Rule, Data Mining

#### ABSTRACT

The large amount of internet causes enormous amount of data to be stored and processed which contains structured and semi structured data formats. Due to extreme size of semi structured document or data, results of a specific query may be massive which leads to retrieving interpretable knowledge a deadly process. The query-answering system makes it attainable to make queries and recuperate results for XML documents and reduced processing of big sized files of data. In this research, it is proposed a method for retrieving more efficient more accurate results for the queries made by the users on the XML document. The original XML document is interpreted to Modified Tree based Association Rules(TAR) files which were shaped by frequent patterns on the original document. The Tree based Association Rules files provide considered information on the structure and content of XML document. Also here it's using CMTree Miner algorithm to mine most frequent set of XML sub-trees, and also comparing Bit cube, Range Cube and Proposed(i.e., MTAR) method with several existing approaches.

© 2015 AENSI Publisher All rights reserved.

**ToCite This Article:** E. Seshatheri and Dr.T. Bhuvaneshwari., A Novel Method to Managing Semi Structured Data in Distributed Environment using Modified Tree based Association Rules(TAR). *Aust. J. Basic & Appl. Sci.*, 9(35): 277-286, 2015

#### INTRODUCTION

The goal of data mining is to extract or mine knowledge from large amounts of data. Consequently, data mining consists of more than gathering and managing data, it also includes analysis and prediction. The eXtensible Markup Language (XML) (Agrawal, R. and R. Srikant, 1994) has become a standard language for data depiction and exchange XML is a Standard, flexible syntax for data exchanging Regular, Structured and Semi Structured data. Mining of XML documents significantly differs from structured data mining and text mining. The XML allows the representation of semi-structured and hierarchical data containing not only the values of individual items, but also the relationships between data items. Due to the characteristic flexibility of XML, in both structure and semantics, discovering knowledge from XML data is faced with new challenges as well as assistances. Mining of structure along with content provides new insights and means into the process of knowledge discovery.

As for query-answering, since query languages for semi structured data depend on the document structure in a file to carry its semantics, in order for query formulation to be actual users need to know this structure in advance, which is often not the case.

This limitation is a crucial problem which did not arise in the context of relational database management systems (RDBMS). As a consequence, when accessing for the first time a large dataset, gaining some general information about its main structural and semantic characteristics helps study on more specific details. This research work addresses the need of getting the essence of the document earlier querying it, both in terms of content and structure. Determining regular patterns inside XML documents provides high-quality knowledge about the document content: frequent patterns are in fact intentional information about the data contained in the document itself, that is, they specify the document in terms of a set of properties rather than by means of data. As different to detailed and precise information conveyed by the data information, this information is partial and often approximate, but synthetic, and concerns both the document structure and its content.

#### 1.1 Tree-Based Association Rules from XML Document:

Association rules describe the co-occurrence of data items in a large amount of collected data and are usually represented as implications in the form  $X \Rightarrow Y$ , where X and Y are two arbitrary sets of data items,

**Corresponding Author:** E. Seshatheri, Research Scholar Manonmaniam Sundaranar University Tirunelveli.

E-mail: seshathriphd@gmail.com

such that  $X \cup Y = \Phi$ ; The quality of an association rule is usually measured by means of support and confidence. Support corresponds to the frequency of the set  $X \cup Y$  in the dataset, while self-assurance corresponds to the conditional probability of finding  $Y$ , having found  $X$  and is given by  $\text{sup}(X \cup Y) = \text{sup}(X)$ . In this work here it is extend the notion of association rule originally introduced in the context of relational databases, in order to adapt it to the hierarchical nature of XML documents. In particular, it is consider the element-only Info set content model, which allows an XML nonterminal tag to include only other elements and/or attributes, while the text is confined to terminal elements. Furthermore, without loss of generalization, here it is do not consider some features of the Info set that are not relevant to the present work, such as names IDREF attributes, URIs, and Links.

Following the Info set conventions, here it is represent an XML document by a labeled tree  $(N, E, \gamma)$  where  $N$  is the set of nodes,  $\gamma \in N$  is the root of the tree (i.e. the root of the XML document),  $E$  is the set of edges. Moreover, the following properties on nodes and edges hold:

1) Each node  $n_i$  has a tuple of labels  $NL_i = \{Ntag_i; Ntype_i; Ncontent_i\}$ ; the type label  $Ntype_i$  indicates whether the node is the root, an element, text, or attribute, whereas the label  $Ncontent_i$  can

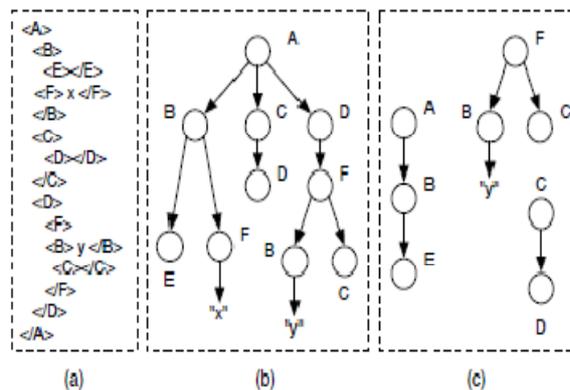
assume as value a PCDATA or  $\perp$  (undefined, for no terminals).

2) Each edge  $e_j = h(n_h, n_k, EL_j)$ , with  $n_h$  and  $n_k$  in  $N$ , has a label  $EL_j = (Etype_j)$ ,  $Etype_j \in \{\text{attribute of, sub-element of}\}$ . Note that edges represent the "containment" relationship between different items of an XML document, thus edges do not have names. Here in this work interested in finding relationships among subtrees of XML documents. Thus, here do not distinguish between textual content of leaf elements and value of attributes. As a consequence, in order to draw graphical concepts in a more readable way, here do not report the edge label and the node type label. Attributes and elements are characterized by empty circles, whereas the textual content of elements, or the value of attributes, is reported under the outgoing edge of the element or attribute it refers to.

**1.2 Fundamental concepts:**

Given two labeled trees  $T = \langle NT; E_T; r_T \rangle$  and  $S = \langle N_S; E_S; r_S \rangle$ ,  $S$  is said to be an induced subtree of  $T$  if and only if there exists a mapping  $\theta : N_S \rightarrow N_T$  such that  $\forall n_i \in N_S; NL_i = NL_j$ , where  $\theta(n_i) = n_j$  and for each edge  $e_j = h(n_1; n_2); EL_{ji} \in E_S; h(\theta(n_1); \theta(n_2)); EL_{ji} \in E_T$ .

Figure 1 shows an example of an XML document (Figure 1(a)), its tree-based representation (Figure 1(b)) and three induced subtrees of the document (Figure 1(c)).



**Fig.1:** a) an example of XML document, b) its tree-based representation, and c) three induced sub-trees

A Tree-based Association Rule (TAR) is a tuple of the form  $T_r = \langle S_B; S_H; s_{Tr}; c_{Tr} \rangle$ , where  $S_B = \langle N_B; E_B; r_B \rangle$  and  $S_H = \langle N_H; E_H; r_H \rangle$  are trees and  $s_{Tr}$  and  $c_{Tr}$  are real numbers representing the support and confidence of the rule respectively. Furthermore,  $S_B$  is an induced subtree of  $S_H$  with an additional property

on the node labels. Indeed, the following properties must hold:

- $N_B \subseteq N_H$
- $E_B \subseteq E_H$  and  $\forall n; m \in N_B; \langle (n; m); EL \rangle \in E_B \text{ iff } \langle (n; m); EL \rangle \in E_H$
- the set of tags of  $S_B$  is equal to the set of tags of  $S_H$  with the addition of the empty label "C".

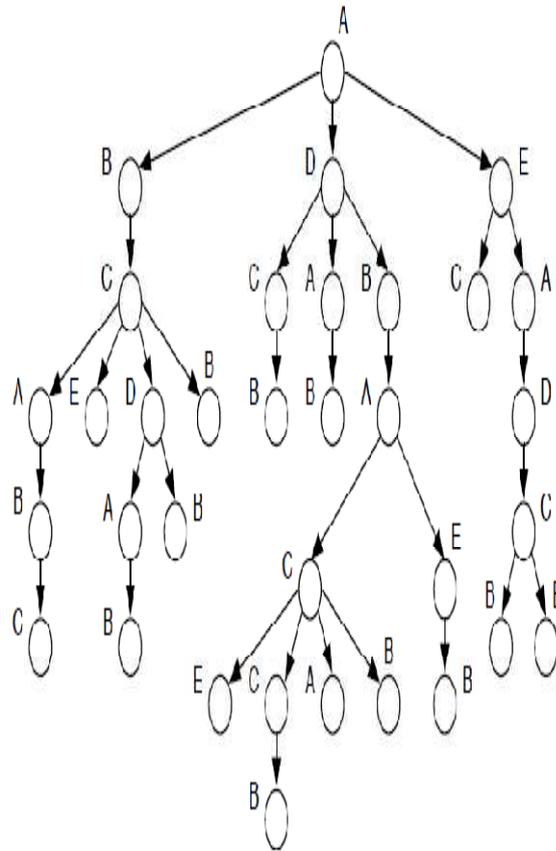


Fig.2: Sample dataset

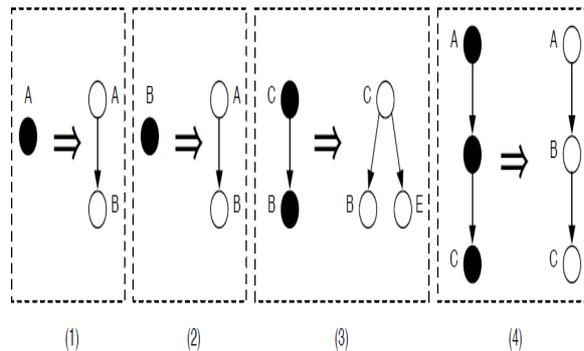


Fig.3: Sample sTARs (structure Tree-based Association Rules)

The empty label is introduced because the body of a rule may contain nodes with unspecified tags (a sort of placeholder nodes) and these tags will be declared in the head part of the rule (see the rule (2) of Figure 3). For the sake of clarity the label  $\_$  is omitted in the figures and all nodes with empty labels do not present any label at all. It is worth noticing that Tree-based association rules are different from XML association rules, because the first requires that  $(X \not\subseteq Y) \wedge (Y \not\subseteq X)$ , i.e. the two trees  $X$  and  $Y$  have to be disjoint, while Tree-based association rules require that  $X$  is an induced subtree of  $Y$ .

A TAR represents intensional knowledge in the form  $S_B \Rightarrow S_H$ , where  $S_B$  is the body tree and  $S_H$  the head tree of the rule. Indeed, the rule  $S_B \Rightarrow S_H$  states that if the tree  $S_B$  appears in an XML document  $D$ , it

is likely that the bigger", or more specific", tree  $S_H$  also appears. In our graphical representation, we will render the nodes of the body of a rule by black circles, and the nodes of the head by empty circles.

Thus, every tree-based association rule is characterized by two measures:  $s_{Tr}$  support, measures the frequency of the tree  $S_H$  in the XML document  $c_{Tr}$  confidence, measures the reliability of a rule, that is the frequency of the tree  $S_H$ , once  $S_B$  has already been found. Given function  $count(S;D)$  denoting the number of occurrences of a subtree  $S$  in the tree  $D$  and function  $cardinality(D)$  denoting the number of nodes of  $D$ , it is possible to define formally the two measures as:

2. Literature Survey:

The problem of association rule mining was initially proposed and many implementations of the algorithms were developed in the database literature. Some methods use XQuery to extract association rules from simple XML documents. They propose a set of functions, written in XQuery, which implement the Apriori algorithm.

The problem related to XML context was produced in the year 2003 by J.W. Won & G. Dobbie which uses XQuery to extract association rules. This approach performs well on simple XML documents but not on complex XML document with irregular structure also the tool called XQuery is a language which is used for finding & extracting element, attributes from XML documents.

An algorithm PATHJOIN is Proposed by Y. Xiao et.al to discover all maximal frequent sub trees given some minimum support threshold. Mohammed J. Zaki, has introduced the mining of embedded sub trees in a (forest) database of trees and introduced a novel algorithm, TREEMINER, for tree mining. TREEMINER uses depth-first search; it also uses the scope-list vertical representation of trees. The limitation of this algorithm is that it cannot find structure of XML document.

Another algorithm is presented by Yun Chi, Yirong Yang, Yi Xia, and Richard R. Muntz, called CMTreeMiner, a computationally efficient algorithm that discovers all closed and maximal frequent sub trees in a database of rooted unordered trees.

Termier *et al.* show that DRYADEPARENT is currently the fastest tree mining algorithm. However, Dryadeparent extracts embedded sub trees which are trees that maintain the ancestor relationship between nodes but do not distinguish, among the ancestor-descendant pairs and the parent-child ones.

Wan and Dobbie show that their approach performs well on simple XML documents but it is very difficult to apply to complex XML documents with an irregular structure.

This limitation is overcome by where Braga *et al.*, introduced proposal to enrich XQuery with data mining and knowledge discovery capabilities, by introducing XMINERULE, an operator for mining association rules for XML documents. They formalize the syntax and semantics for the operator and propose some examples of complex association rules. However, XMINE is based on the MINERULE operator, which works on relational data only. This resource that, after a step of pruning of unnecessary information, the XML document is translated into the relational format. Moreover, in many techniques, the designer is forced to specify the structure of the rule to be extracted and then to mine it, if possible. This means that the designer has to specify what should be contained in the body and head of the rule, i.e., the designer has to know the structure of the XML document in advance, and this is an irrational requirement when the document does not have a Document Type Definition. A document type

definition (DTD) is a set of markup declarations that define a document type for an SGML-family markup language (SGML, XML and HTML). A DTD uses a terse formal syntax that declares precisely which elements and references may appear wherein the document of the particular type, and what the elements' contents and attributes are.

#### Disadvantages of Existing System

- Search intention for a keyword query is not easy to determine.
- It returns low result quality in term of query relevance.
- Rank the individual matches of all these queries are challenging.

In our research, the above all disadvantages is also incorporated and discover the new technique.

### 3. Problem Statement:

- The goal of data mining is to extract or mine knowledge from huge amounts of data. Thus, data mining consists of more than collecting and managing data. Sometimes, it also includes analysis and prediction of the data. The XML allows the representation of hierarchical and semi-structured data containing not only the values of individual items but also the relationships between data items from the different type of database. Due to the intrinsic flexibility of XML, discovering knowledge from XML data is faced with new challenges as well as benefits.

Usually people have attempted to handle the XML processing problems using one of four methods as follows:

(1) focus on techniques for "shredding" XML into tables, followed by combining the tables, and later re-joining the results to produce required XML output;

(2) make a few modifications to object oriented or semi-structured databases, which also inherently support hierarchy, hence they support the XML;

(3) use a top-down tree-traversal strategy for executing queries;

(4) use a custom wrapper at the source end for index like retrieval of only the necessary content. Before here it is describe the Tukwila architecture, it is useful to briefly examine these previous approaches, including their relative strengths and weaknesses.

### 4. Objectives of the Proposed System:

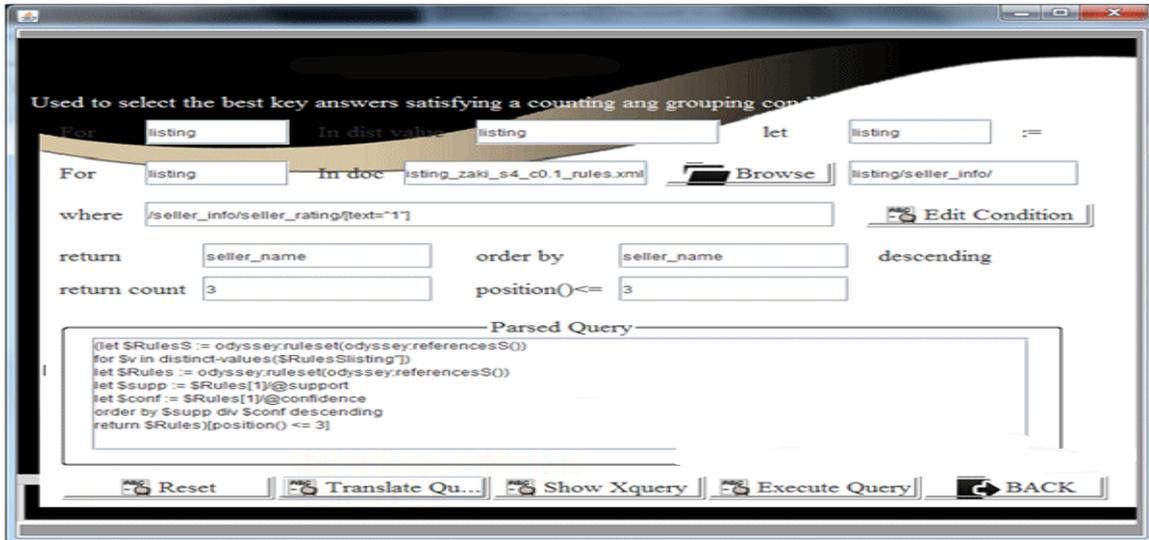
1. Query answering system using Tree Based Association Rules on XML document to extract the most relevant feeds from the large file directly.

2. To provide an approach which actually find frequent pattern and Tree Based Association Rules (TAR) from the XML file from various sources/documents which will help to improve performance and retrieve relevant information from the XML document.

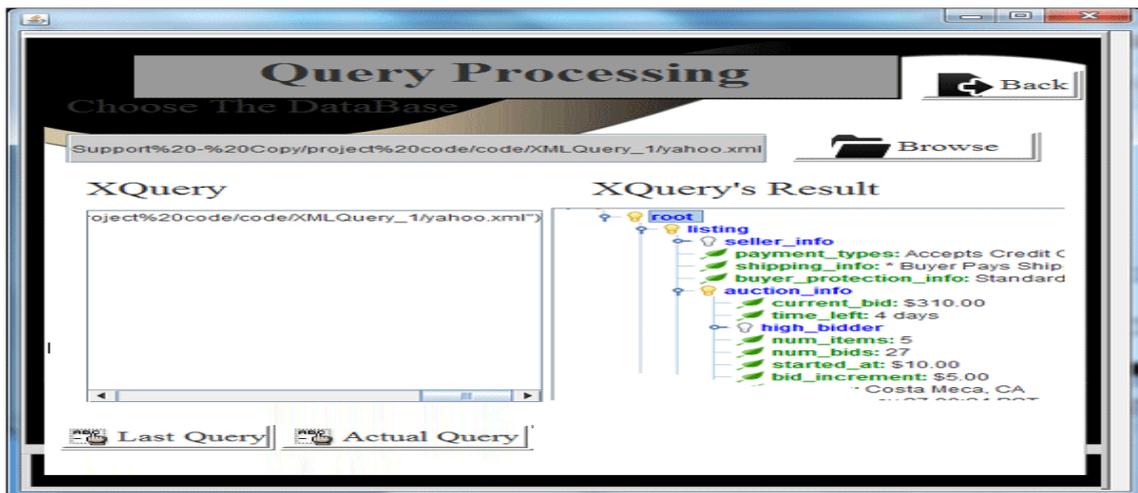
**5. Proposed System – Experimental Design:**

The proposed method for query-answering system is based on the association mining. The XML document (rss, XML or doc) is given as input. The XML tree is constructed using XML parser

which is extracted for the user defined values of support and confidence. This system provides the intentional knowledge in the form of XML document for the entire XML documents.



**Fig.4:** XML Query Processing using WHERE Clause



**Fig.5:** XQuery processing results

Then the query will be applied for XML and appropriate clause, document which was set for by the users. This gives the result related to the query. Then the XML document is modified and the process is repeated until goal reached. The modifications are

stored in the same XML document which was already formed by the values of support and confidence combine all this document parsed is evaluate as again when required.

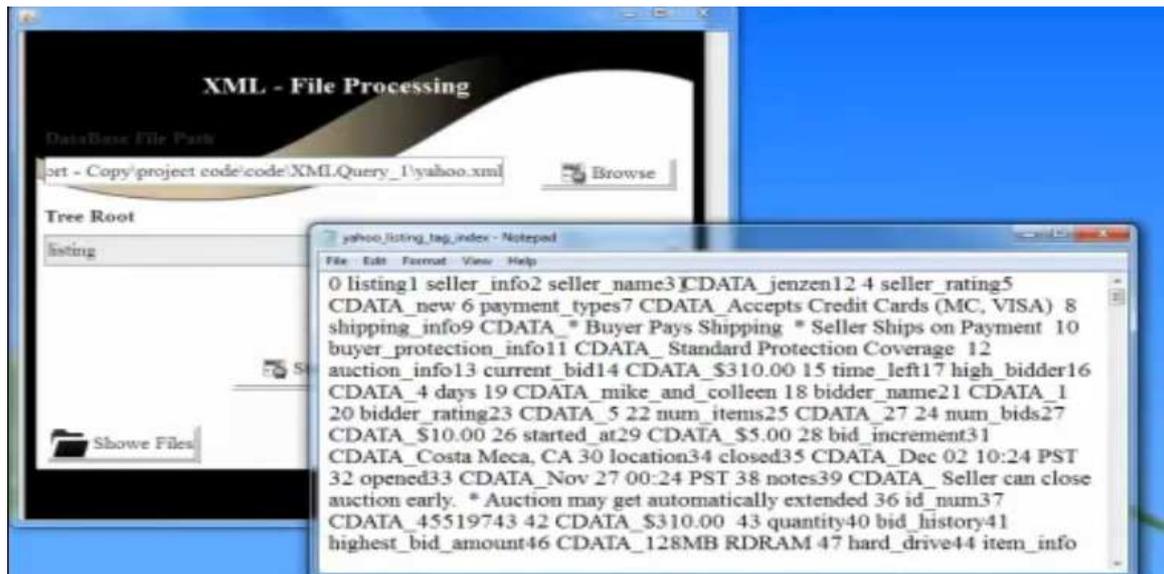


Fig.6:XML File Processing using Data File

### 5.1 Construction of Trees:

The XML documents have a flexible architecture. XML documents can be pre-processed. XML pre-processing is done by XML parser. The DOM (Document Object Model) parser is used here which is issued to construct the tree from the XML document. Before a XML document can be accessed, it must be loaded into an XML DOM object. The DOM creates a tree structure in the internal memory from the given XML document. It stores the all the document into memory before processing those documents. It allows the users to traverse the document using top-down approach XML trees, access, insert, update the content, style and structure of the document and also to delete the nodes from the tree. Therefore XML document forms a tree structure. Also the XML document should be validated (i.e) the tags should be started and ended correctly without leaving any tag without its pair.

### 5.2 Extraction of Frequent Pattern:

Association rules describe the frequent occurrence of data items in a large amount of data collected. X and Y are the two considered data items. Now they are represented in the form of X $\cup$ Y. Here Association rule is measured by means of Support and Confidence. In this the Support represents the frequency of the set (X and Y) found in the data set and Confidence represents the conditional probability of finding Y, having got X. The interesting patterns among the subtrees of the given XML document can be identified. A Tree-based Association Rule (TAR) (Suganya, I., et al., 2013), the form of XML document defines pattern of subtrees in the XML document. The entire XML document can be accessed by providing the values of

support and confidence. Mining TAR is a two step process (Suganya, I., et al., 2013).

This includes:

1. Mining frequent sub trees
2. Computing interesting rules

After obtaining the set of files the Proposed model will merge them to obtain one XML document which contains data from all the included files. The step next to this is to obtain the TAR of all the files. Once it is done, the new novel Proposed model will give the most frequent feeds of all the files. After that the feed search will be performed by the Proposed model which will then give us the filtered result. The obtained results are better than the other methods. The Comparative results are analyzed in the last section.

### 6. Experimental Results:

Here in this research used Oracle 9i (Enterprise Edition Release 9.2.0.8.0) and Stylus Studio 2008 XML Enterprise Suite Release 2, to evaluate the different query results in the case of a centralized database system. Here in this research used the Visual C++ language using Borland C compiler (32 bit that supports up to 4 GB RAM) to implement our Proposed technique. It is used an Intel Processor with 2.13 GHz, 4 GB of RAM under the Windows 7 Ultimate Operating system.

To support the Oracle 9i database is used the Windows 7 Ultimate operating system. Here in this work used the (well structured) XML datasets in (Mazuran, M., et al., 2012) to run comparisons of the XQuery language and our Proposed technique using file sizes of 2.5MB, 5.6MB, 10.8MB, 23.48MB and 110.35MB.

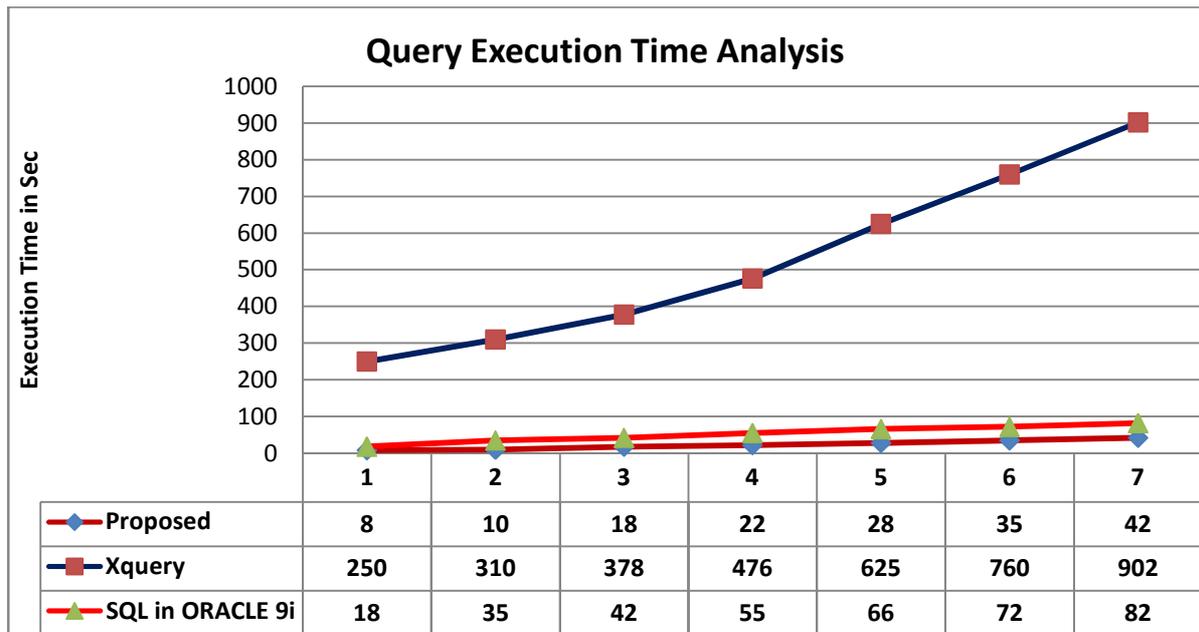


Fig.6: Query execution time analysis.( Number of predicates using AND Condition WHERE Clause)

The experimental results for XQuery execution time, execution time in Oracle and proposed are presented in Figure 4, which for the purpose of clarity only depicts the results for 110.35MB files, but which show comparative performance that is typical of all the tested file sizes. It can be seen that proposed method's execution times were substantially lower than those of XQuery, and also that proposed method's outperforms the highly regarded Oracle execution times across the range of predicates tested. The results show that the improvement achieved by Proposed

method increases with the number of predicates. A second set of experiments was performed to compare the query execution times of Proposed Method, Range Cube and Bitcube (Paik, J., et al., 2005). These experiments used 5, 10, 15, 20 and 25 element paths (ePaths) per document, and a variety of words per element path. For all numbers of ePaths per document, and all numbers of documents, proposed method outperformed Bitcube as shown in Figure 5, which for the purpose of clarity only depicts the results for 25 ePaths per document, and 20 words per ePath.

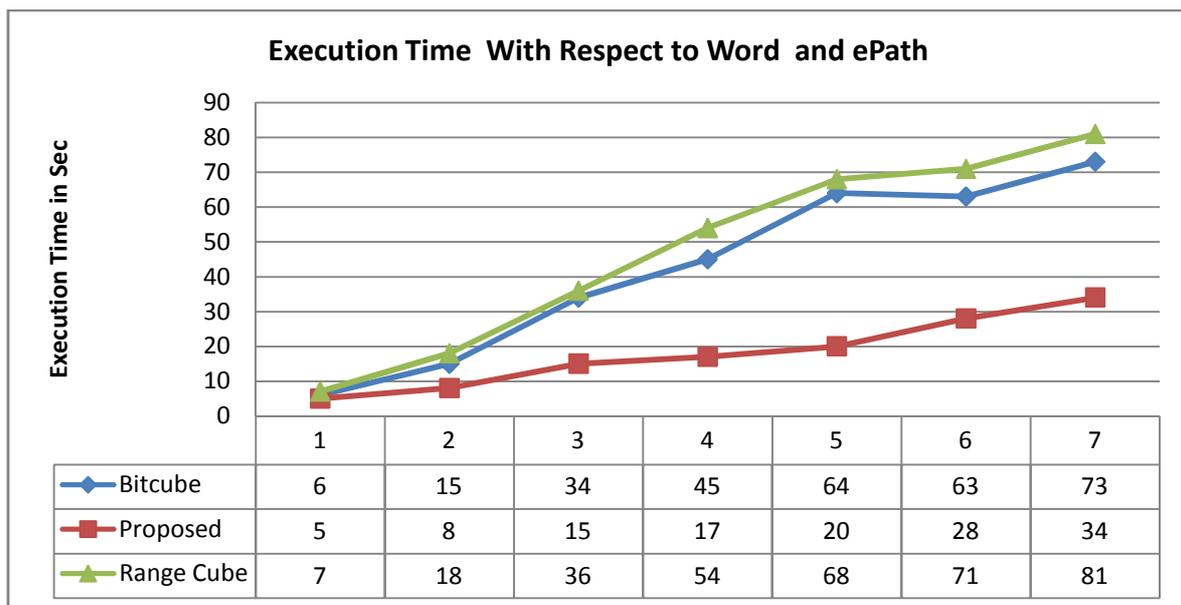


Fig.5: Execution time wrt word/ePath.

### Implementation Details:

In particular, the idea of mining association rules to provide summarized representations of XML documents has been investigated in many proposals either by using languages (e.g., XQuery) and techniques developed in the XML context, or by implementing graph-or-tree-based algorithms. This work proposes to develop an application to mine frequently occurring subtrees from the generated XML data set by RSS Feed links of different newspaper links. Then storing them in another XML document. These rules are extracted only if support and confidence of nodes is greater than the provided threshold values. Then apply our Frequent Pattern Miner algorithm for XML. The idea of using association rules as summarized representations of XML documents was also introduced in where the XML summary is based on the extraction of rules both on the structure (schema patterns) and on content (instance patterns) of XML data sets. The limitations of this approach are: 1) the root of the rule is established a-priori and 2) the patterns, used to describe general properties of the schema applying to all instances, are not mined, but derived as an abstraction of similar instance patterns and are less precise and reliable.

### A. Fundamental Concepts:

The idea of the project basically revolves around the concept of Tree base Association rule. Mining

tree-based association rules is a process composed of two steps:

1. Mining frequent subtrees from the XML document;
2. Computing interesting rules from the previously mined frequent subtrees.

Association rule (Agrawal, R. and R. Srikant, 1994) is an hint of the form  $A \cup B$ , where the rule  $A$  and  $B$  are subset of the set  $C$  of element in a set of transactions  $D$  and  $A \cap B = \emptyset$ . A rule  $X \rightarrow Y$  states that the transaction  $T$  that has the elements in  $A$  are likely to contain also the elements in  $B$ . Association rules are distinguished by two survey: the support, which gives the percentage of transactions present in  $D$  that contain both items  $A$  and  $B$  ( $A \cup B$ ); the confidence, which calculates the percentage of transactions in  $D$  containing the element  $A$  that also contain the elements  $B$  ( $\text{support}(A \cup B) / \text{support}(A)$ ). In XML context, both  $D$  and  $C$  are group of trees. In this work TAR is generated using Natural Language Processing. The project works for RSS news links of multiple newspapers. RSS stands for "Really Simple Syndication". It is a way to easily distribute a list of headlines, update notices, and sometimes content to a wide number of people. It is used by computer programs that organize those headlines and notices for easy reading. RSS works by having the website author maintain a list of notifications on their website in a standard way. This list of

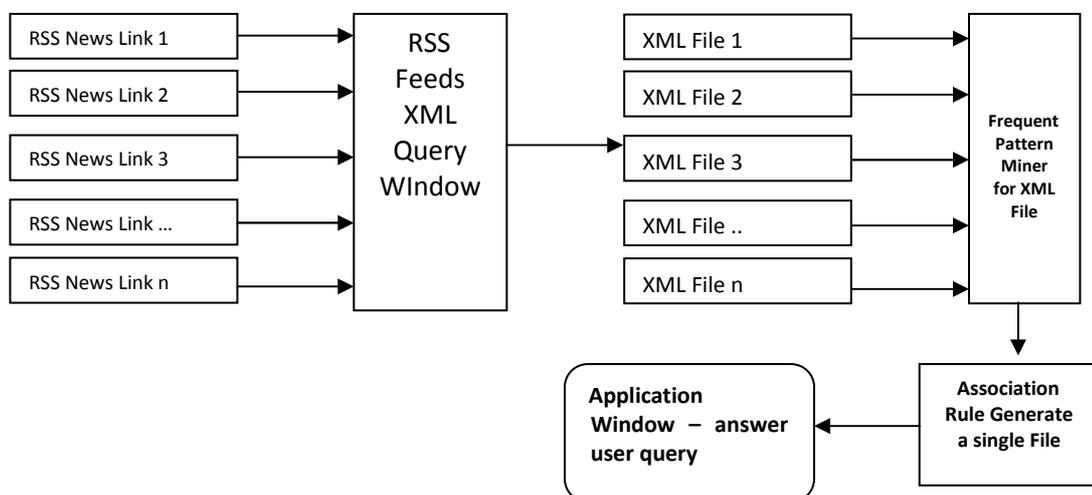


Fig. 6: Operation flow for frequent pattern Miner for XML.

notifications is called an "RSS Feed". People who are interested in finding out the latest headlines or changes can check this list. Special computer programs called "RSS aggregators" have been developed that automatically access the RSS feeds of websites you care about on your behalf and organize the results for you.

Producing an RSS feed is very simple and hundreds of thousands of websites now provide this feature, including major news organizations like the New York Times, the BBC, and Reuters, as well as many weblogs. RSS provides very basic information to do its notification. It is made up of a list of items presented in order from newest to oldest. Each item usually consists of a simple title describing the

item along with a more complete description and a link to a webpage with the actual information being described. Sometimes this description is the full information you want to read (such as the content of a weblog post) and sometimes it is just a summary.

**B. Architecture Details:**

This research work introduces a proposal for mining and storing Modified Tree-Based Association Rules (MTARs) as a means to represent intensional knowledge as an alternative, synthetic data set would be queried for providing quick and summarized answers. For calculating the MTAR we make use of Natural Language Processing. The algorithm for finding MTAR is given below:

Algorithm: Extracting Tree Based Association Rules from XML Files

Input: RSS Feed Files  $R_n$ , Similarity Threshold  $T_h$ , Support  $S$ , Confidence  $C$ .

Output: XML Tree Based Association Rules

1. Read all  $R_n$  by using XML Reader.
2. Extract Inner Element <Title> Tags within all  $R_n$  and insert them into a new XML File  $x_f$ .
3. Initialize Node List  $L_{st}$  to contain Visited  $T_{sub}$  in  $x_f$ .
4. Initialize ItemSet, a Key Value Pair structure to contain Element  $e$  and its support  $f_s$ .

5. For Each Subtree  $T_{sub}$  in  $X_f$
6. Extract <Description> Tag  $I_{desc}$  from  $T_{sub}$
7. If (Count( $L_{st}$ ) == 0)
8. Then Add  $I_{desc}$  into  $L_{st}$
9. Else For Each element  $e_i$  in  $L_{st}$
10. SimilarityScore = Compare( $I_{desc}$ ,  $e_i$ )
11. If (SimilarityScore >=  $T_h$ )
12. Then If ItemSet contains  $e_i$
13. Then Increment support of  $e_i$  by 1 from ItemSet
14. Else Add  $e_i$  to ItemSet and increment support of  $e_i$  by 2
15. Else Add  $e_i$  to ItemSet and increment support of  $e_i$  by 1
16. Add  $I_{desc}$  into  $L_{st}$
17. Using ItemSet as reference data structure containing Item  $e_i$  and its support  $f_s$ , desired support  $C$  and confidence  $C$ , generate frequent ItemSet  $F$ .
18. For Each ItemSet  $I_{fr}$  in  $F$
19. Insert  $I_{fr}$  in XML Document  $X_{tar}$
20. Return  $X_{tar}$

**Conclusion:**

The Comparative Analysis of Implemented Approach with existing Document summarization and Information Retrieval Methods as mentioned below in the Table 1;

**Table 1:** Comparative Analysis of Implemented Approach with existing Document summarization and Information Retrieval Methods

Methods	IR Ranking	Type of Input information	Generate Frequent Patterns?	Use of Dictionary	Query Independent	Applicable for Smart Phone Devices?	Applicable on Web Pages?
Automatic Document Summarization By Sentence Extraction	Statistical	Structured Only	No	No	No	No	No
Query – Specific Document Summarization	Statistical	Structured Only	No	No	No	No	No
A Language Independent Algorithm for Single and Multiple Document summarization	NLP	Structured Only	No	Yes	Yes	No	No
OLEDOT	NLP	Structured and Non structured data	Yes	Yes	Yes	No	Yes
Implemented Approach (Modified TAR)	NLP	Structured and Non structured data	Yes	Yes	Yes	Yes	Yes

The main goals of the work are to achieve:

1. Mine all frequent association rules restriction on the structure and the content of the rules;
2. Store mined information in XML format;
3. Use the extracted knowledge to gain information about the original data sets;
4. Modified mined Tree based Association Rules when the original XML datasets change.

**REFERENCES**

Agrawal, R. and R. Srikant, 1994. "Fast Algorithms for Mining Association Rules in Large Databases," Proc. 20th Int'l Conf. Very Large Data Bases, pp: 478-499.

Combi, C., B. Oliboni and R. Rossato, 2005. "Querying XML Documents by Using Association Rules," Proc. 16th Int'l Conf. Database and Expert Systems Applications, pp: 1020-1024.

Paik, J., H.Y. Youn and U.M. Kim, 2005. "A New Method for Mining Association Rules from a

Collection of XML Documents," Proc. Int'l Conf. Computational Science and Its Applications, pp: 936-945.

Yogesh R. Rochlani, Prof. A.R. Itkikar, 2012. "Integrating Heterogeneous Data Sources Using XML Mediator", *ijcsn*, 1: 3.

Arundhati Birari, Prof. Ranjit Gawande, 2013. "Mining Tree-Based Association Rules for XML Query Answering", *ijettcs* vol., 2: 3.

Wan, J.W.W. and G. Dobbie, 2003. "Extracting Association Rules from XML Documents Using XQuery," Proc. Fifth ACM Int'l Workshop Web Information and Data Management, pp: 94-97.

Barbosa, D., L. Mignet and P. Veltri, 2005. "Studying the XML Web: Gathering Statistics from an XML Sample," *World Wide Web*, 8(4): 413-438.

Suganya, I., N. Velmurugan, Dr. P. Ganeshkumar, 2013. "XML Query-Answering Support System using Association Mining Technique" IEEE conference on ICT.

Marouane Hachicha and Jérôme Darmont, Member, 2013. IEEE Computer Society, "A Survey of XML Tree Patterns", IEEE 25.

Yun Chi, Yirong Yang, Yi Xia and Richard R. Muntz, 2010. "CMTreeMiner: Mining Both Closed and Maximal Frequent Subtrees", International Conference on Knowledge Discovery and Data Mining.

Wang Lian, Nikos Mamoulis, David Wai-lok Cheung, 2005. "Indexing Useful Structural Patterns for XML Query Processing", *IEEE Transactions On Knowledge And Data Engineering*, 17: 7.

Neoklis Polyzotis, Minos Garofalakis, Yannis Ioannidis, 2003. "Approximate XML Query Answers", Intl. Conf. on Very Large Data Bases.

Devi Mahalakshmi, S., Dr. K. Vijayalakshmi, Dr. K. Muneeswaran, G. Priyanka, 2010. "Mining Intensional Information for answering XML-Queries using Tree-based Association Rules Approach", IEEE.

Chandra Sekhar, K. Dhanasree, 2012. "Extracting TARs from XML for Efficient Query Answering", *International Journal of Computer Science and Network (IJCSN)* 1: 6.

Agarwal, R. and R. Srikant, 1994. "Fast Algorithms for Mining Association Rules in Large

Databases," Proc. 20th Int'l Conf. Very Large Data Bases, pp: 478-499.

Borgwardt, K.M., H.P. Kriegel, P. Wackersreuther, 2006. "Pattern Mining in Frequent Dynamic Subgraphs," Proc. ICDM.

Braga, D., A. Campi, S. Ceri, M. Klemettinen and P. Lanzi, 2003. "Discovering Interesting Information in XML Data with Association Rules," Proc. ACM Symp. Applied Computing, pp: 450-454.

Gasparini, S. and E. Quintarelli, 2005. "Intensional Query Answering to XQuery Expressions," Proc. 16<sup>th</sup> Int'l Conf. Database and Expert Systems Applications, pp: 544-553.

Goldman, R. and J. Widom, 1997. "Data Guides: Enabling Query Formulation and Optimization in Semistructured Databases," Proc. 23rd Int'l Conf. Very Large Data Bases, pp: 436-445.

Mazuran, M., E. Quintarelli and L. Tanca, 2012. "Data Mining for XML Query-Answering Support," *IEEE Trans. Knowledge and Data Eng.*, 17(8): 1021-1035.

Paik, J., H.Y. Youn and U.M. Kim, 2005. "A New Method for Mining Association Rules from a Collection of XML Documents," Proc. Int'l Conf. Computational Science and Its Applications, pp: 936-945.

Asai, T., K. Abe, S. Kawasoe, H. Arimura, H. Sakamoto and S. Arikawa, 2002. Efficient Substructure Discovery from Large Semi-Structured Data, Proc. SIAM Int Conf. Data Mining.

Xiao, Y., J.F. Yao, Z. Li and M.H. Dunham, 2003. Efficient Data Mining for Maximal Frequent Subtrees, Proc. IEEE Third Int. Conf. Data Mining, pp: 379-386.

Zaki, M.J., 2005. Efficiently Mining Frequent Trees in a Forest: Algorithms and Applications, *IEEE Trans. Knowledge and Data Eng.*, 17(8): 1021-1035.

Chi, Y., Y. Yang, Y. Xia and R.R. Muntz, 2004. CMTreeMiner: Mining both Closed and Maximal Frequent Subtrees, Proc. Eighth Pacific-Asia Conf. Knowledge Discovery and Data Mining, pp: 63-73.

Termier, A., M. Rousset, M. Sebag, K. Ohara, T. Washio and H. Motoda, 2008. Dryad Parent, an Efficient and Robust Closed Attribute Tree Mining Algorithm, *IEEE Trans. Knowledge and Data Eng.*, 20(3): 300-320, Mar.