



ISSN: 2319-5967

ISO 9001:2008 Certified

International Journal of Engineering Science and Innovative Technology (IJESIT)

Volume 3, Issue 2, March 2014

Application of Association Mining Technique in XML QA Support System (AMTXSS)

Ujwal A.Bodke, Santosh Kumar, Sandip Darade

Abstract— Withdrawal information from semi structured documents is a very difficult task, and is going to become unfavorable as the amount of digital information of Internet is growing. Actually, documents are generally so large that the data set returned as answer to a query may be very large to fetch interpretable knowledge. In this paper, we explain an application based on Tree-Based Association Rules (TARs): mined rules, which provide almost accurate, purposive information on both the structure and the contents of Extensible Markup Language (XML) documents, and can be stored in XML format. The aim of our proposal is to provide a way to use intentional knowledge as a substitute of the original document during querying and not to enhance the execution time of the queries over the original XML dataset. A prototype system and experimental results determine the effectiveness of the approach.

Index Terms— XML dataset, Query Answering, TARs, Data Mining.

I. INTRODUCTION

Extract or mine knowledge from large amounts of data is the goal of data mining. Data mining not only collecting and managing data, it also includes analysis and prediction. The XML is an Extensible Markup Language. It has become a standard language for data representation and exchange XML is a flexible syntax for data exchanging. Mining of XML documents different from structured data mining and text mining. XML allows the relationships between data items with the representation of semi-structured and hierarchal data containing not only the values of individual items. Because of inherent flexibility of XML, in both structure and semantics, discovering knowledge from XML data is faced with modern challenges and advantages. XML provides new insights and means into the process of knowledge discovery by mining of structure at the same time as content. In Query answering system, languages for semi structured data depend on one document structure to transport its semantics, for query effective formulation but for this users need to know that structure in advance, which is often not the case but this limitation, is a crucial problem which did not emerge in the surroundings of relational database management systems. When retrieving for the first time a large dataset, gaining some common information about its main structural and semantic characteristics helps analysis on more individual component. The need of getting the general meaning of the document before querying it, in terms of content and structure. Xml finds frequent patterns inside XML documents which provide high-grade knowledge about the document content. Recurrent patterns are in fact intended information about the data contained in the document itself, it means, they specify the document in terms of a set of properties instead of by means of data. As opposed to detailed and precise information fetched by the data, this information is limited and often relative, but not genuine, and concerns both the document structure and its content.

II. LITERATURE SERVEY

In the recent years the database research held has concentrated on XML as an expressive and flexible hierarchical model suitable to represent huge amounts of data with no absolute and fixed schema, and with a possibly irregular and incomplete structure. More recently the problem has been investigated also in the XML context “Discovering interesting information in xml data with association rules”, “Extracting association rules from xml documents using XQuery” and “A new method for mining association rules from a collection of xml documents”. In “Discovering interesting information in xml data with association rules” the authors use XQuery to extract association rules from simple XML documents. They propose a set of functions written only in XQuery which implement together the A-priori algorithm. It is show that their approach performs well on simple XML documents; however it is very difficult to apply this proposal to complex XML documents with an irregular structure. This limitation has been overcome in “Extracting association rules from xml documents using XQuery”, where the authors introduce a proposal to enrich XQuery with data mining and knowledge discovery capabilities, by introducing XMINE RULE,



ISSN: 2319-5967

ISO 9001:2008 Certified

International Journal of Engineering Science and Innovative Technology (IJESIT)

Volume 3, Issue 2, March 2014

a specific operator for mining association rules for native XML documents. They formalize the syntax and an intuitive semantics for the operator and propose some examples of complex association rules. Our idea is to take a more general approach to the problem of extracting association rules from XML documents, i.e. to mine all frequent rules, without having any a-priori knowledge of the XML dataset. A similar idea was presented in "A new method for mining association rules from a collection of xml documents" where the authors introduced HoPS, an algorithm for extracting association rules in a set of XML documents. Such rules are called XML association rules and are implications of the form $X \rightarrow Y$, where X and Y are fragments of an XML document. In particular the two trees X and Y have to be disjoint. In our work, XML association rules are mined starting from frequent sub trees of the tree-based representation of a document. In the database literature it is possible to and many proposals of algorithms to extract frequent structures from tree/graph-based data structures. Just to cite some of them, Tree Miner, Path Join, Close Graph propose algorithms to directly mine frequent item sets-not association rules-from XML documents. Tree Miner and Close Graph do not preserve the exact structure of the item sets extracted -only the "descendant-of" (and not the "child-of") relationship between nodes is preserved -whereas Path Join does. In this work we propose an algorithm that extends Path Join to mine generic tree-based association rules directly from XML documents.

III. PROPOSED SYSTEM

- 1] Extract all recurrent association rules without forcing any algorithm rules on the structure and the content of the rules.
- 2] Store extracted knowledge in XML form.
- 3] Use extracted data to get information about the original XML datasets.

In existing system textual content of the XML document is the main data retrieval and because of that there is no any benefit is derived from the semantics fetched by the document structure. As for QA, After all query languages for semi structured data depend on the a document which is most important to transform its semantics for query computation to be successful users want to realize this structure in advance, but this is not happened.

IV. METHODOLOGY

We increasing productivity and potency for data mining in user friendly manner from an xml dataset, but the recommended system will be useful for mined data by user query. Here different types of mining techniques such as:

- A) TARS
- B) TAR mining
- c) Intended answer
- D) Tree ruler model

(a) TAR (TREE BASED ASSOCIATION RULES)

Association rule [1] 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 (support ($A \cup B$)/support (B)). In XML context, both D and C are group of trees. In this work we extend the concept of association rule found in the surroundings of relational databases to modify it to the hierarchical essence of XML documents. Succeeding the Info set traditions, we represent an XML document as a tree $[n, e, r, l, c]$ where n - the set of nodes, r -is the root of the tree, e - is the set of edges, l - label function which returns the tag of nodes. In following figure no.1 shows XML dataset/document, in 2nd section we are getting tree structured document. In 3rd part we are getting combined pictorial representation of two subtrees and last fourth parts showing final rooted subtree.

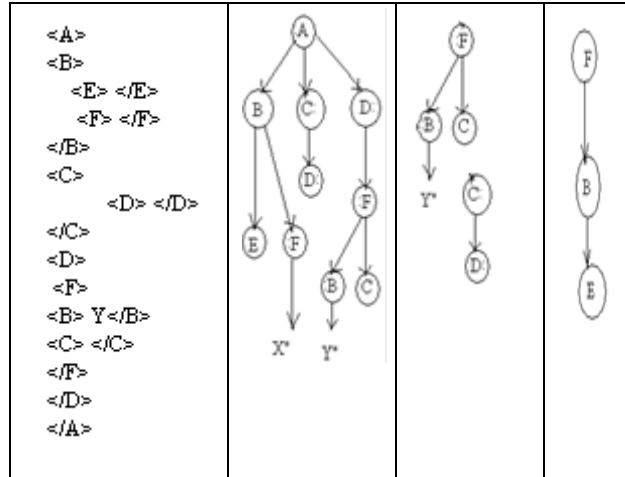


ISSN: 2319-5967

ISO 9001:2008 Certified

International Journal of Engineering Science and Innovative Technology (IJESIT)

Volume 3, Issue 2, March 2014



(b) TAR Mining

Extracting TARs through data mining is a two steps process. In the first step recurrent sub trees that fulfill given support are mined. In the second step fascinating rules that have confidence above given threshold are calculated from the recurrent sub trees. Finding recurrent sub trees is described in. Algorithm 1 which finds recurrent sub trees and computes interesting rules.

Algorithm 1. Get –Interesting-Rules (D, minsupp minconf)

- 1: FS1=FindFrequentSubtrees (D, minsupp)
- 2: ruleSet = \emptyset
- 3: for all st \in FS1 do
- 4: tempSet = Compute-Rules (st, minconf)
- 5: ruleSet = ruleSet \cup tempSet
- 6: end for
- 7: return ruleSet.

Function 1 Compute-Rules(st, minconf)

```

ruleSet , blacklist =  $\emptyset$ ;
for all CS, subtrees of st do
if CS is not a sub tree of any item in blacklist then conf = sup (st) / sup(CS)

if conf  $\cdot$  minconf then
newRule = (CS, st , conf, supp(st))
ruleSet = ruleSet  $\cup$  (newRule)
else
blackList = blackList  $\cup$  CS
end if end if end for
return ruleSet

```

Algorithm 1 finds recurrent sub trees and then hands each of them over to a function which calculates all the possible rules. Depending on the number of recurrent sub trees and their cardinality, the amount of rules created by a simple Compute Rules function may be very high. Given a sub tree with n nodes, we could generate rules, making the algorithm exponential. This surge occurs also in the relational context thus, based on similar inspection, it is possible to state the following property that allows us to suggest the enhanced version of compute rules shown in function 2.



ISSN: 2319-5967

ISO 9001:2008 Certified

International Journal of Engineering Science and Innovative Technology (IJESIT)

Volume 3, Issue 2, March 2014

Algorithm 2 Create-Index (D)

```

1: for all Di D do
2: for all dj
3: referral (root (d2)) = refferal (root (d2)) U references
  (root (dj))
4: SumChildren (d2, dj)
5: end for
6: end for
7: return D.

```

Function 2 SumChildern(T1,T2)

```

for all A children (root(T2)) do
if ch children (root(T2)) | ch= A then
referral (root(ch)) = referral (root(ch)) U referral
(root(A))
ch = sumchildern (ch, A) else add child (root(T2),A)
end if
end for return

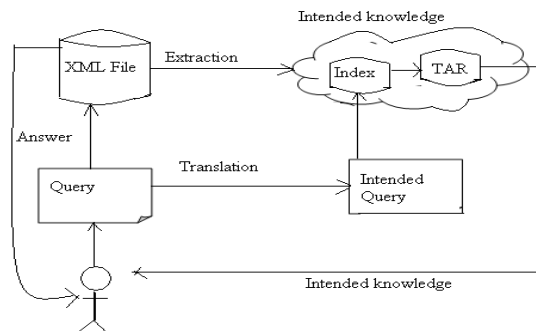
```

(c) INTENDED ANSWERS

Raw intended view of the XML document is express by TARs, which is not specific because, it explains the data in terms of its properties. A user query over the original dataset can be automatically converted into a query over the extracted TARs. The answer will be intended, because, instead of providing the set of data satisfying the query, the system will answer with a set of properties. Following are the two major benefits: i) querying TARs requires small time than querying the original XML document; ii) Intended answers are in some cases more useful than the extensional one.

d) THE TREERULER PROTOTYPE

Tree Ruler is a prototype tool that integrates all the functionalities proposed in our approach. Given an XML document, the tool is able to extract intentional knowledge and allows the user to compose traditional queries as well as queries over the intentional knowledge. Figure 1 shows the architecture of the tool. In particular, given an XML document, it is possible to extract Tree-based rules and the corresponding indexed. The user formulates XQuery expressions on the data and these queries are automatically translated in order to be executed on the intentional knowledge. The answer is given in terms of the set of Tree-based rules which re etc. the search criteria. It is com posed by several tabs for performing different tasks. In particular, there are three tabs:





ISSN: 2319-5967

ISO 9001:2008 Certified

International Journal of Engineering Science and Innovative Technology (IJESIT)
Volume 3, Issue 2, March 2014

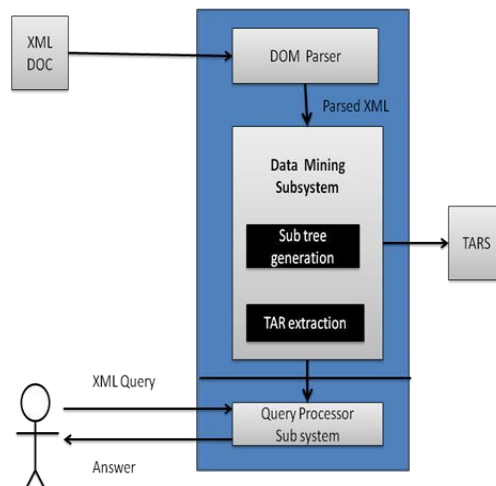
Proposed Scheme

I] mine all frequent association rules without forcing any algorithm's restriction on the structure and the content.

II] Store extracted data in XML format.

III] Use extracted knowledge to gain information about the original datasets.

Textual content of the document; this means that no any benefit is acquire from the semantics transformed by the document structure. As for query-answering, after all query languages for semi structured data depend on the one document structure to pass its semantics, for better query computation to be fruitfull,so users require to know this structure in advance but always situation is different. User puts up XML query to the query processor. Before query processor subsystem comes up with an answer, it goes through the below process. The XML DOM contains methods to convert XML trees, and thus the XML document becomes flexible on which we can perform no. of operations like access, insert, and delete nodes. An XML parser reads XML, and converts it into an XML DOM object that can be accessed with JavaScript. This parsed XML send to the Data mining subsystem. Data mining subsystem incorporates two parts: sub tree generation & TAR extraction. Data mining subsystem generates TARS. TAR extraction section provides the intended data to the query processor subsystem. And finally query processor subsystem provides the answer to the user.



Benefits of proposed scheme

- Resolve keyword inexactness Problems.
- To effectively recognize the type of target node.
- To effectively search via node.
- Rank the individual matches of all possible search purpose.

Class 1: Node or child node queries: This query is used to bring down simple and complex to count the number of elements having operators with restrictions on them.

Class 2: count-queries: These queries are used to retrieve number of elements or records from the dataset.

Class 3: Top-k queries: These queries are used to select the top k records from dataset which satisfy a grouping condition. These queries will retrieve the highest degree elements from the XML file.



ISSN: 2319-5967

ISO 9001:2008 Certified

International Journal of Engineering Science and Innovative Technology (IJESIT)

Volume 3, Issue 2, March 2014

Class 4: Aggregate Queries: Aggregate query provides aggregate functions like sum, minimum, maximum, average. For example if we want to retrieve minimum salary from employees then we can use this function min (salary).

Class 5: DML Queries: These queries used to manipulate the XML document by user, like Delete, Insert, Update without having any knowledge of XML. Dynamically updating the TAR files when the dataset is changes.

This project work has been divided in 4 modules:

1. XML Dataset Selection and Validation:- In this module, the system will facilitates the service to select xml dataset and it validates the file is correct or not.
2. XML Parser creates XML tree:- In this module, the system will transformed given xml dataset into tree structure for using xml parser.
3. Implement available 3 classes (Query Type):- In this module, the system will create a template or framework for all available classes to provide answer for user query.
4. Create XML QA System:- In this module, the system will supply the answer in tree form for all available classes , depends on user query.

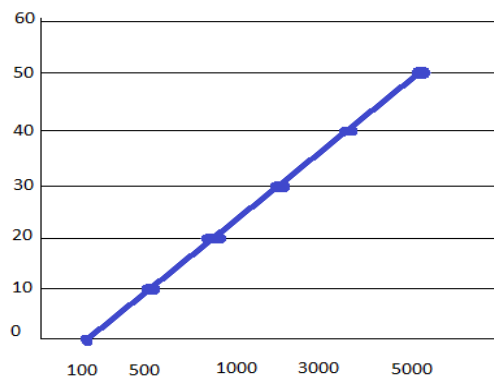
V. RESULTS

Tree Ruler is a tool used in our approach. When the XML document is given, it makes users to retrieve the intentional information for the queries. We use the sample dataset yahoo which is in XML format. The TAR files and its index files created with the references were stored in the XML format.

Types of Class Queries

1. Retrieve all listing where current bid is greater than 10
2. Retrieve the number of payment type
3. Retrieve the best k seller rating

These class queries were in the form of XQuery. These TARs provides the quick and approximate intentional answers more benefit than the extensional answers of the XML document. This querying over TAR files needs less time when compares to querying the original XML document. The extraction time was calculated for processing the intentional knowledge over the various numbers of nodes in the XML document. Precision and recall values were used to evaluate the accuracy of the results retrieved. Fig.4 shows the extraction time calculated for the nodes when queries are made.





ISSN: 2319-5967

ISO 9001:2008 Certified

International Journal of Engineering Science and Innovative Technology (IJESIT)
Volume 3, Issue 2, March 2014

Future work: In future we can have most of the provisions of SQL with XML.

VI. CONCLUSION

Extraction of frequent association rules from the structure and the content is the main objectives we are going to achieve in this work. Store extracted information in XML form. Use extracted data to get information about the original XML datasets. A framework or template will be developed to test the effectiveness of our proposal.

ACKNOWLEDGMENT

The writer would like to thank Prof. Santosh Kumar for his help. This work was supported in part by the Sandip Foundation.

REFERENCES

- [1] G. Seshadri Sekhar¹, Dr.S. Murali Krishna, Efficient Data Mining for XML QASS/IOSR Journal of Computer Engineering (IOSRJCE) ISSN: 2278-0661 Volume 4, Issue 6 (Sep.-Oct. 2012), PP 13-22.
- [2] KC. Ravi Kumar¹, E. Krishnaveni Reddy², Ramadevi.G³,\Data Mining for XML QASSt ,IOSR Journal of Computer Engineering (IOSR-JCE) ISSN:2278-0661, ISBN: 2278-8727 Volume 5, Issue 6 (Sep- Oct. 2012), PP25-29.
- [3] Chandra Sekhar.K, 2Dhanasree, Extracting TARs from XML for Efficient QA International Journal of Computer Science and Network (IJCSN) Volume 1, Issue 6, December 2012.
- [4] Mining Association Rules from XML Document using Modified Index Table, 2013 International Conference on Computer Communication and Informatics (ICCCI -2013), Jan. 04 06, 2013.
- [5] Anam V Bhaskara Reddy, Answering Xml Query Using Tree Based Association Rules, International Journal of Latest Trends in Engineering and Technology (IJLTET).
- [6] D.Karthiga¹, S.Gunasekaran, Optimization of Query Processing in XML Document Using Association and Path Based Indexing, International Journal of Innovative Research in Computer and Communication Engineering Vol. 1, Issue 2, April 2013.
- [7] Saranya T.J.,\MINING TREE-BASED ASSOCIATION RULES FROM XML DOCUMENT, International Journal of Advanced Technology and Engineering Research (IJATER).
- [8] Mrs.Mopuri Sujatha,\XML Query Answering using Tree based Association Rules\,International Journal of Advanced Research in Computer Engineering and Technology (IJARCET)Volume 1, Issue 8, October 2012.XML Query Answer System REFERANCES 21 [9] P. B. Vikhe¹, B. L. Gunjal, Extracting Tree Based Association Rules from XML Document ,International Journal, Volume 3, Issue 6, June 2013.
- [9] Arundhati Birari¹, Prof. Ranjit Gawande, Mining Tree-Based Association Rules for XML Query Answering, International Journal, Volume 2, Issue 3, May June 2013.

AUTHOR BIOGRAPHY

Miss.Ujwal A. Bodke She is post graduate student of computer engineering at SITRC Nashik under University of Pune. Her area of interest includes Web mining.

Mr. Santosh Kumar is working as Assistant Professor in Computer department at SITRC, Nashik, and Maharashtra, India.

Mr. Sandip Darade is working as IT Analyst in TCS.