# Association Rule Mining From Server Log File

[1]Kalpana Wani , [2]J.W.Bakal

PIIT New Panvel, Navi Mumbai, SSJCOE Dombivali, Mumbai
Email: [1]kjwani82@gmail.com, [2]bakaljw@gmail.com

**Abstract—** As we know World Wide Web plays vital role in serving the needs of the user's on web. The web log files are generated as a result of an interaction between the client and the service provider on web. Web log file contains the massive hidden valuable information pertaining to the visitors, if mined can be used for predicting the navigation behavior of the users. However the task of discovering frequent sequence patterns from the web log is challenging. This system focuses on adopting an intelligent technique that can provide personalized web service for accessing related web pages more efficiently and effectively, so that it can be determined which web pages are more likely to be accessed by the user in future. This system uses two intelligent algorithms for predicting the user behavior's namely FP Growth and Eclat.These algorithms save the time and space problem of existing system. Further from the frequent pages pattern Direct and Indirect Association Rules are generated and based on that Ranking is provided to pages which will help recommendation system to recommend similar search pages.

**Index Terms—** Association Rule, Indirect Association Rule, Parsing, Preprocessing.

## I. INTRODUCTION

The World Wide Web serves as a vast, widely distributed, global information service center for advertisement, consumer information, e-commerce, education, financial management, government, news and many other information services. So, it has become much more difficult to access relevant information from the web with the explosive growth of information available on the internet. Therefore, further research work needs to be carried out for extracting the appropriate content as per the user's needs. This can be performed using Web mining as it helps in extracting the common sequences of the user's accessed web pages. In the web environment, association rules are typically applied to HTTP server log data that contain historical user sessions. Web sessions are gathered without any user involvement and additionally, they reliably reflect user behavior while navigating throughout a web site. For that reason, web sessions can be regarded as an important source of information about users. Association rules that reveal similarities between web pages derived from user behavior can be simply utilized in recommender systems. The main goal of such a recommendation is to suggest to the current user some web pages that appear to be useful.

## II. LITERATURE SURVEY

In literature survey we discuss regarding various methods of Frequent Pattern generation and Mining direct and indirect association rule from weblog data.

### A. Frequent Pattern Generation:

Bina Kotiyalt, Ankit Kumar, Bhaskar Pant, R.H. Goudar,Shiv ali Chauhan and Sonam Junee(2013)[1] proposed a web page recommendation system by finding frequent pattern using Apriori and Eclat algorithm and compare both the algorithm for generating frequent pattern.

Dr. S. Vijayarani1, Ms. P. Sathya(2013)[2] Compare the Éclat and Rapid Association Rule mining algorithm for finding the frequent item sets in data streams. Éclat and Rapid Association Rule mining algorithm are used for finding the frequent item sets in data streams.

N. Venkatesan, Ramaraj. (2012) [4] proposed the novel bit search technique which is implemented in the existing association rule mining algorithms. Frequent item sets are generated with the help of apriori & Eclat based bit search technique.

Karthikeswaran D, Dinakar S. (2010)[8] Proposed recursive algorithm uses this doubly Linked tree to efficiently find all access patterns that satisfy user specified criteria.

### B. Association Rule Generation:

Ms Shweta1, Dr. Kanwal Garg (2013)[3] compares the Association rule algorithms to find out the best combination of different attributes in any data. In this paper author uses Apriori to find association rule. Here author consider three association rule algorithms: Apriori

Association Rule, Predictive Apriori Association Rule and Tertius Association Rule. Author compares the result of these three algorithms and presents the result. According to the result obtained using data mining tool author find that Apriori Association algorithm performs better than the Predictive Apriori Association Rule and Tertius Association Rule algorithms.

Rachna Somkunwar (2012) [6] proposed a new technique known as Hash-set technique that will help to increase the efficiency of algorithm by avoiding multiple scanning the database. Our Hash Techniques based algorithm is used to reduce the frequently scanning of item set also it helps to mine interesting association rules on the basis of lift and conviction technique.

PRZEMYSŁAW KAZIENKO (2009)[13] proposed algorithm to extracts complete indirect association rules with their important measure confidence, using pre-calculated direct rules. Both direct and indirect rules are joined into one set of complex association rules, which may be used for the recommendation of web pages. Performed experiments revealed the usefulness of indirect rules for the extension of a typical recommendation list. They also deliver new knowledge not available to direct ones. The relation between ranking lists created on the basis of direct association rules as well as hyperlinks existing on web pages is also examined.
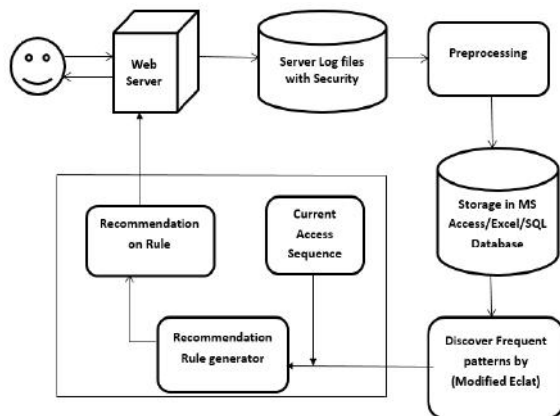
## III. PROPOSED SYSTEM



Fig.1. Architecture Diagram of Proposed System
Following steps are used for implementation:

Step 1: Data Cleaning & Preprocessing

Preprocessing is necessary, because Log file contain noisy &ambiguous data which may affect result of mining process. Some of web log file data are unnecessary for analysis process and could affect detection of web attack. Data preprocessing is an important steps to filter and organize only appropriate information before applying any web mining algorithm. Preprocessing reduce log file size also increase quality of available data. The purpose of data preprocessing is to improve data quality and increase mining accuracy. Preprocessing consists of field extraction, data cleansing, user identification, session identification.

Step 2: Finding Frequent Pattern

First count the occurrences of pages in database and then reorder the pages with higher count page first and so on. Draw the FP Growth tree and start mining FP tree by considering the suffix as page with lowest count first. Find the conditional pattern base for each page and if length of conditional pattern base is less than some threshold value then find the frequent pattern using Éclat algorithm and if length is greater than or equal to the threshold value then use FP Growth algorithm for finding the frequent pattern of visited pages.

Step 3: Direct Association Rule Generation between the Frequent Access Pages

After getting the frequently visited page sequence filter the pattern by applying some minimum confidence and support value and form the one page association rule between pages of frequent pattern and calculate its support and confidence values and filter it with respect to minimum confidence and support value.

Step 4: Indirect Association Rule Generation between the Frequent Access Pages

In direct association rule of pages there are many pages which are connected to each other indirectly via some transitive pages. Find such pages which are indirectly connected and filter it with respect to minimum confidence and support value.

Step 5: Complex Association Rule Generation

A complex association rule exists if at least one of two component rules exists, i.e., either direct or complete indirect or both of them. The main quality features of both direct and in indirect rules—confidences—are combined within complex association rules. A complex association rule is characterized by the complex confidence, $con*(d_i \rightarrow *d_j)$, as follows

$$Con*(d_i \rightarrow *d_j) = \alpha \cdot Con(d_i \rightarrow d_j) + (1 - \alpha) \cdot Con\#(d_i \rightarrow \#d_j)$$

Where $\alpha$ is the direct confidence reinforcing factor $\alpha \in [0, 1]$

Calculate the complex confidence value for different values of $\alpha$. Setting $\alpha$ we can emphasize or damp the direct confidence at the expense of the complete indirect one. The greater the value of $\alpha$, the closer the complex confidence to the direct one. Depending on highest value of $\alpha$ rank is

## IV. SYSTEM IMPLEMENTATION

In our system instead of using any readymade tool for data preprocessing we have developed our own processing method which process web log data present in W3c format i.e. in Extended Log Format. Home page of proposed system is shown below. It consist of the display of the Architecture Diagram of the system and a brief explanation of the working of the system.
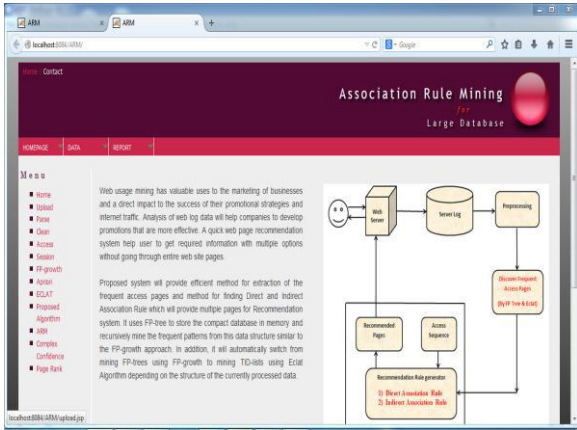


Fig.2 Home Page of Proposed System

On left hand side of the home page Menu options are given which are listed below

1) Home
2) Upload
3) Parse
4) Clean
5) Access
6) Session
7) FP-Growth
8) Apriori
9) Éclat
10) Proposed Algorithm(FEM)
11) ARM
12) Complex Confidence
13) Page rank

Upload

Second option is Upload, which facilitates to upload the weblog file in w3c format. We can browse the file from its location and click on upload button as shown below.



Fig.3 Web Log File Upload

After clicking on the upload button if file is uploaded successfully then message will be displayed about successful uploading of file and file contents are displayed as shown below. If file type is not matching then it display the error.



Fig.4 Web Log File ContentsParse

Once weblog file is uploaded successfully then we can start preprocessing of the file by selecting the Parse option from the Menu. We have to use Parse option for uploaded file only, if we try to use it without uploading a file it will give error. Following snapshot show that when we select Parse option for uploaded file it automatically display the file name.
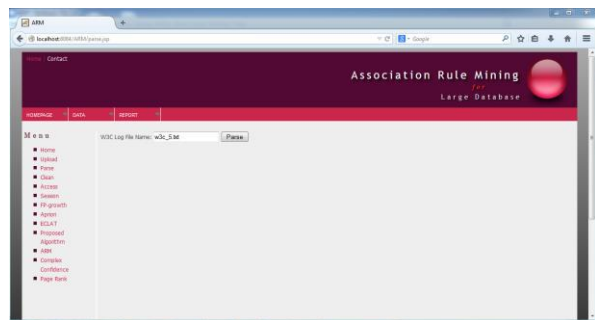


Fig.5 Web Log File Selection for Parsing of Data

After getting required file name click on parse button which will convert W3C format weblog file information in a tabular format as shown below.
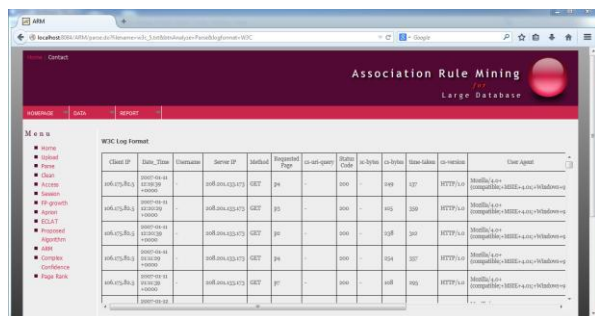


Fig.6 Tabular Representation of Web Log file data

Clean

After converting web log data into a tabular format we need to clean the data to know accurate information. In web log data many entries are not completed or

_____

sometimes we want information about some specific data format etc. Depending upon our requirement we can clean the data and retrieve the information. Different cleaning option provided in this system are given below.

a) Remove failed/invalid request.
b) Select only "GET" method
c) Remove multimedia object request.
d) Remove web 'Robots' request.

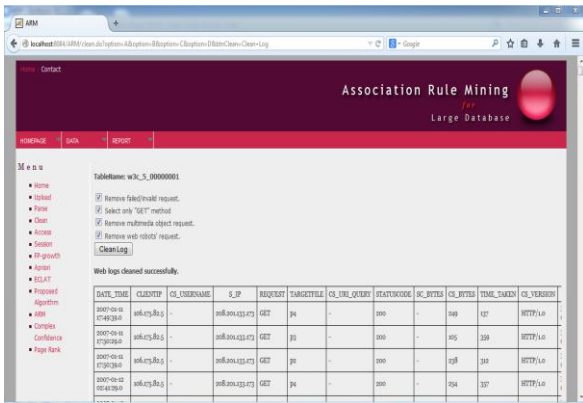We can select either one or multiple option depending upon our data requirements as shown below.


Fig.7 Web Log Data after cleaning Access

To know the user access pages list we can use this option. It will display list of users as well as their access pages as shown below.


Fig 8. User Access List Session

This option gives the details of visited users their sessions and the visited pages in respective sessions as shown in figure.


Fig 9 User Session

Next we can execute any of the given algorithms to find

frequent patterns from clean weblog data. Apriori, Éclat and FP-Growth algorithms implementation is not a part of proposed system. Only implementation of proposed algorithm is a part of this project and only for proposed FEM algorithm Direct and Indirect Association rules are generated. Apriori, Éclat are used in Existing system to generate frequent patterns. Apriori, Éclat and FP-Growth algorithms are implemented to compare performance of proposed system with respect to existing system.

Following snapshot shows the execution of all above algorithm for support value 20 except Apriori as it gives memory problem for low support value.
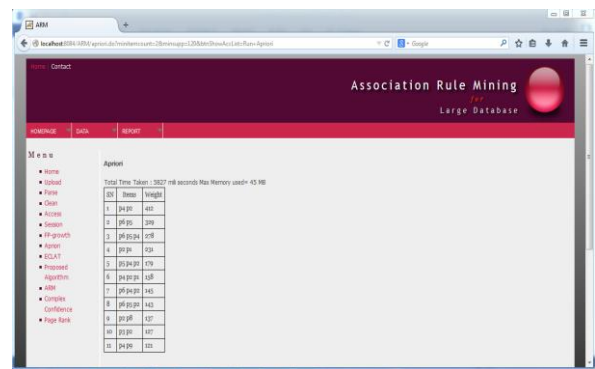
Apriori Algorithm
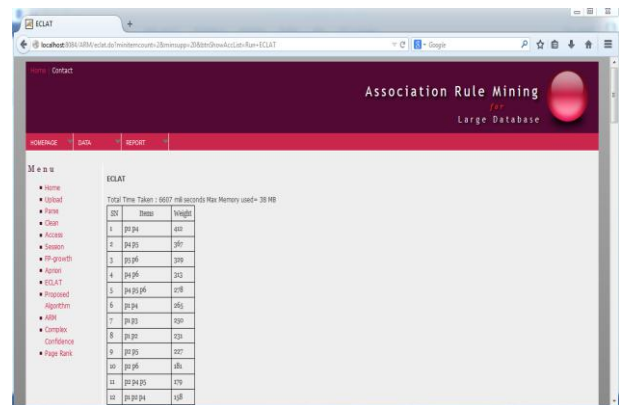

Fig .10 Apriori Algorithm

Éclat Algorithm


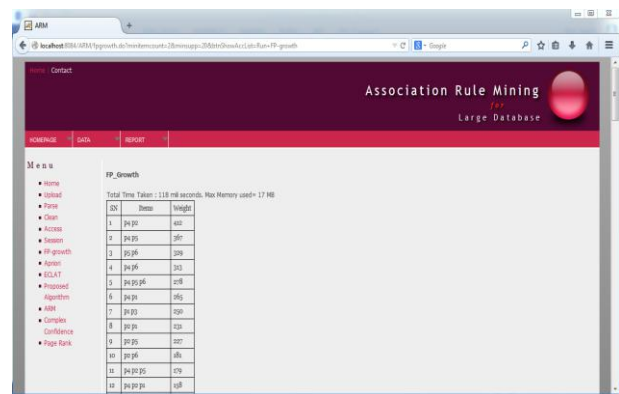Fig. 11 Éclat Algorithm

FP-Growth Algorithm


Fig .12 FP-Growth Algorithm
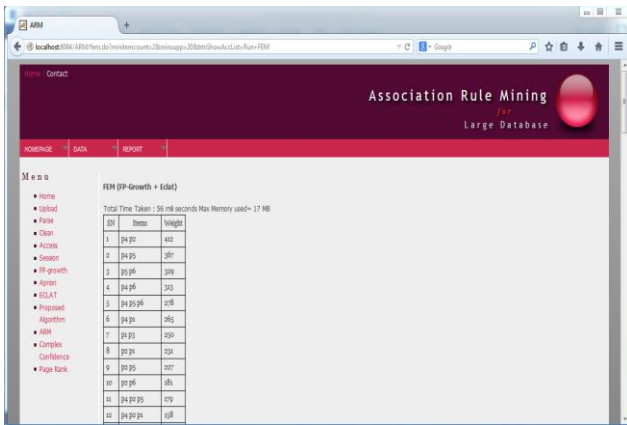
Proposed FEM Algorithm



Fig 13. Proposed FEM Algorithm

After executing FEM algorithm we can find direct and indirect Association rule snapshot for that given below. Depending on your selection it generate direct and indirect association rule as shown below.
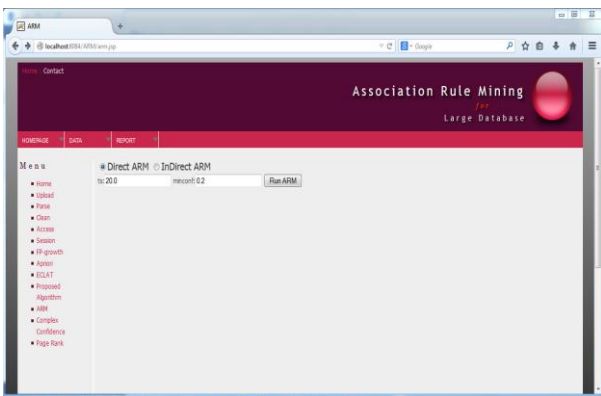


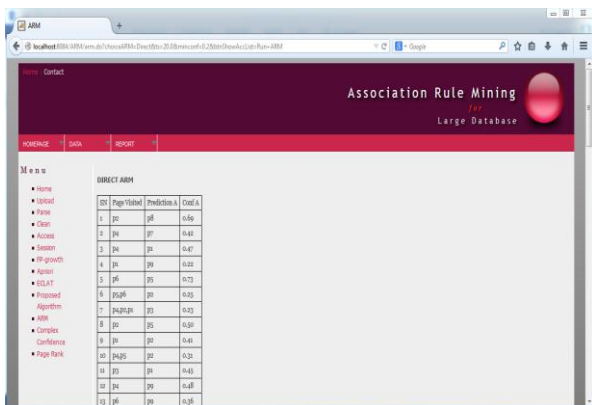Fig..14 Generation of Direct and Indirect Association Rule

Direct Association Rule



Fig..15 Generation of Direct Association Rule
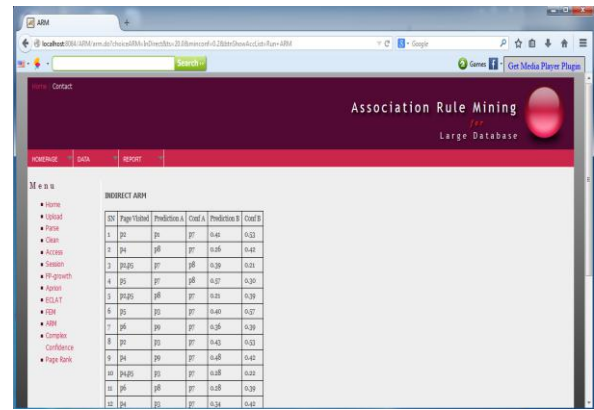
Indirect Association Rule



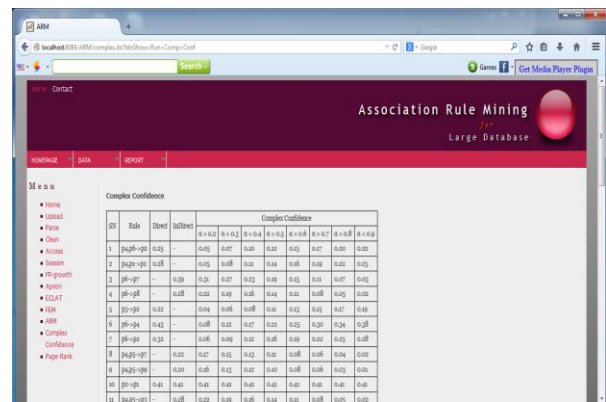Fig..16 Generation of Indirect Association Rule

Complex Confidence



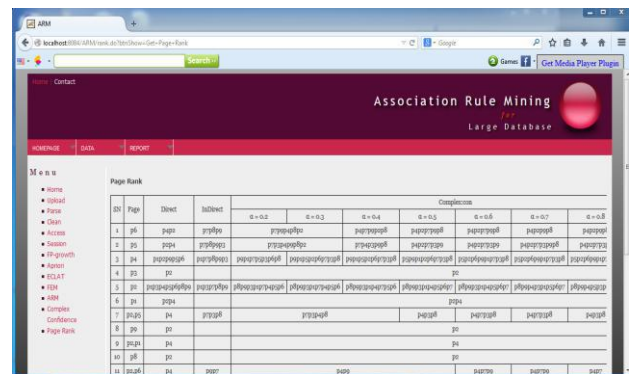Fig..17 Generation of Complex Confidence

Page Rank



Fig.18 Page ranking based on Complex Confidence

## IV. CONCLUSION

In today's world there is boom of online shopping and because of that if companies want to increase their business it's required to analyze the user behavior. This is possible only with help of web mining and its further processing to find frequent pattern and the association rule between them. This system will help to recommend appropriate page to user.

In future we are planning to complete remaining work by providing security to log file

## V. ACKNOWLEDGMENT

## REFERENCES

[1]   Bina Kotiyalt, Ankit Kumar, Bhaskar Pant, R.H. Goudar, Shiv ali Chauhan and Sonam Junee" User Behavior Analysis in Web Log through Comparative Study of Eclat and Apriori", Proceedings of7'h International Conference on Intelligent Systems and Control (ISCO 2013)

[2]   Dr. S. Vijayarani, Ms. P. Sathya   "An Efficient Algorithm for Mining Frequent Items in Data Streams",in International Journal of Innovative Research in Computer and Communication Engineering Vol. 1, Issue 3, May 2013

[3]   Ms Shweta, Dr. Kanwal Garg ,"Mining Efficient Association Rules Through Apriori Algorithm Using Attributes and Comparative Analysis of Various Association Rule Algorithms"in Volume 3, Issue 6, June 2013 ISSN: 2277 128X International Journal of Advanced Research in Computer Science and Software Engineering.

[4]   N. VENKATESAN, RAMARAJ," FREQUENT ITEMSET MINING WITH BIT SEARCH", Journal of Theoretical and Applied Information Technology, 15 July 2012. Vol. 41 No.1

[5]   Ravi Bhushan and Rajender Nath," Recommendation of Optimized Web Pages to Users Using Web Log Mining Techniques", 978-1-4673-4529-3/12/$31.00c 2012 IEEE

[6]   Rachna Somkunwar" Hash-Set Technique of Association Rules"in Volume 2, Issue 11, November 2012 ISSN: 2277 128X ,International Journal of Advanced Research in Computer Science and Software Engineering.

[7]   Chowdary Farha ahmed, Byeong-Soo Jeong, "Efficient mining of high utility patterns over data streams with a sliding window model".Springerlink.com, 2011

[8]   Yo unghee Kim, Won Young Kim and Ungmo Kim "Mining frequent item sets with normalized weight in continuous data streams". Journal of information processing systems. 2010.

[9]   L. Zhou, Z. Zhong, J. Chang, J. Li, J.Z. Huang, S. Feng,"Balanced parallel FP-Growth with MapReduce", in Conference on Information Computing and Telecommunications, IEEE, pp. 243 – 246, 2010.

[10]   Ding, Yau, S. S.: TCOM, "An innovative data structure for mining association rules among infrequent items", pp. 290-301. Comput. Math. Appl. 57, 2 (2009).

[11]   Ding, 1., Yau, S. S.: TCOM, an innovative data structure for mining association rules among infrequent items, pp. 290-301. Comput. Math. Appl. 57, 2 (2009).

[12]   Sang Lin ,Hu-yan Cui ,Ren Ying ,Zhou-lin Lin ," Algorithm Research for Mining Maximal Frequent Itemsets Based on Item Constraints", 2009 Second International Symposium on Information Science and Engineering, ISBN: 978-0-7695-3991-

[13]   PRZEMYSŁAW KAZIENKO "MINING INDIRECT ASSOCIATION RULES FOR WEB RECOMMENDATION", Int. J. Appl. Math. Comput. Sci., 2009, Vol. 19, No. 1, 165–186

❖ ❖ ❖