



International Journal of Advanced Research in Computer Science and Software Engineering

Research Paper

Available online at: www.ijarcsse.com

Mining Web Log Files for Web Analytics and Usage Patterns to Improve Web Organization

Sana Siddiqui
CSE , RGPV
Bhopal, India

Imran Qadri
CSE, RGPV
Bhopal, India

Abstract: *The web access log is the best repositories for the information source; it keeps the whole record of even a tiny low event. One can simply determine internet usage patterns for numerous internet users. The web usage pattern analysis is a method of distinguishing browsing patterns by analyzing the user's navigation and behaviour. The internet server log files that store the knowledge concerning the guests of internet sites is employed as input for the web usage pattern analysis method. Most of existing websites have a stratified content organisation. This manner of organizing could also be slightly different from the visitor's expectation of organizing the web site. Particularly, it's typically unclear to the visitant at which location a specific document is found. First, these log files are pre-processed and regenerate into specified formats therefore web usage mining techniques will apply on these web logs. This paper reviews the method of discovering helpful patterns from the online server log file of an educational institute. The obtained results is employed in totally different applications like net traffic analysis, economical web site administration, website modifications, system improvement and personalization and business intelligence etc.*

Keywords: NCSA, W3C, WHOWEDA, HTTP, DOS

I. INTRODUCTION

With the technological advancements, businesses have gone online. The World Wide Web has, since then, been the ultimate and vast source of information. For example, people today can buy desired things by just clicking on a button in the computer. Because of the growing popularity of the World Wide Web, many websites typically experience thousands of visitor's every day. Analysis of who browsed what can give important insight into, for example, what are the buying patterns of existing customers. Interesting information extracted from the

visitors browsing data help analysts to predict, for example, what will be the buying trends of potential customers. Correct and timely decisions made based on this knowledge have helped organizations in reaching new heights in the market. The massive data growth provides several challenges and opportunities to the user and web miners. A data stream is an ordered sequence of items that arrives in timely order. Mining steam data is a significant challenge in web data mining. Web data mining is the process to extract the interesting (nontrivial, implicit, previously unknown and potentially useful) knowledge from huge amount of data. Stream data grows rapidly, so there is an augmented need to perform pre-processing on stream data.

II. WEB LOG MINING

Web Usage Mining addresses the problem of extracting behavioural patterns from one or more web access logs [1]. the entire process can be divided into three major steps. The first step, pre-processing, is the task of accurately identifying pages accessed by web visitors. This is a very difficult task because of page caching and accesses by web crawlers. The second step, pattern discovery, involves applications of data mining algorithms to the pre-processed data to discover patterns. The last step, pattern analysis, involves analysis of patterns discovered to judge their interestingness.

Web server records all users' activities of the web site as web servers Logs. Most log files have text format and each log entry is saved as a line of text. There are many types of web logs such as NCSA format, W3C format and IIS format, but they share the same basic information. Log data represented in W3C extended format is shown in figure1.

These log data can be used in web site designing, modifying and also to improve the overall performance of web site. After identifying the different web server log data files there is a need to merge the log files.

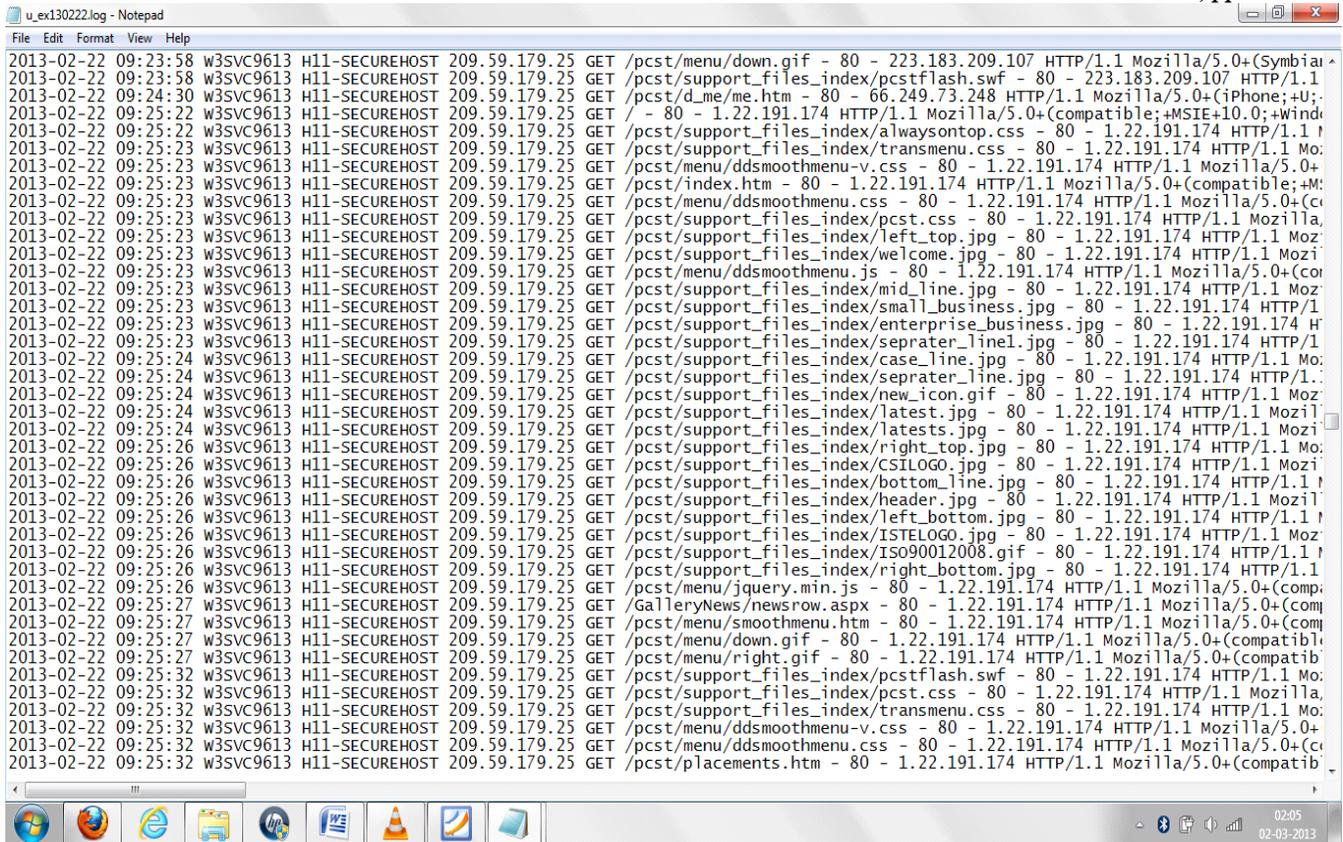


Figure 1: Web Log File

III. RELATED WORK

There has been considerable work on mining web logs; however, none of them include the idea of using backtracks to find expected locations of web pages.

Web usage mining [3] is referred to the discovery of user access patterns from web usage logs, which records every click made by the users. This information is frequently gathered and automatically stored into access logs through Web server. Web usage mining process is similar to data mining process. The difference is in data collection phase. The data are collected from databases for data mining whereas it is collected from web log files in web usage mining. In conventional data mining techniques information pre-process includes data cleaning, integration, transformation and reduction. But web mining pre-processing categorize into Content pre-processing, Structure pre-processing, Usage pre-processing. Once the data is collected from log files, a three-step process is performed in web usage mining namely data preparation, pattern discovery and pattern analysis.

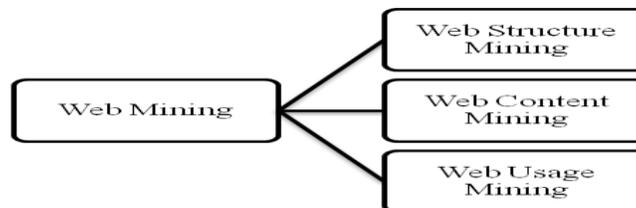


Figure 2: Taxonomy of Web Mining

Maheswara Rao [2] introduced a new framework to separate human user and search engine access intelligently with less time span. And also Data Cleaning, User Identification, Sessionization and Path Completion are designed correctly. The framework reduces the error rate and improves significant learning performance of the algorithm.

Perkowitz et al. [4] [5] investigate the problem of index page synthesis, which is the automatic creation of pages that facilitate a visitor's navigation of a website. By analyzing the web log, their cluster mining algorithm finds collections of pages that tend to co-occur in visits and puts them under one topic. They then generate index pages consisting of links to pages pertaining to a particular topic.

Spiliopoulou et al. [6] [7] propose a “web utilization miner” (WUM) to find interesting navigation patterns. The interestingness criteria for navigation patterns are dynamically specified by the human expert using WUM’s mining language which supports the specification of statistical, structural and textual criteria.

A model called WHOWEDA (Warehouse of Web Data) has been proposed by Sanjay Madria, Sourav S Bhowmick [8] in which a discussion has been performed on various issues in web mining area. Various experiments have been performed for implementing web data as a web personalization tool [9] in which they have categorized the process of web mining in five phases i.e. i) data gathering, ii) data preparation, iii) navigation pattern discovery, iv) pattern analysis and visualization, and v) pattern applications. A model has been proposed to get the benefit of combining the Semantic Web and Web Mining [10]. Accurate Web usage information could help to attract new customers, retain current customers, improve cross marketing/sales, effectiveness of promotional campaigns, track leaving customers and find the most effective logical structure for their Web space [11]. A very good model has been proposed using decision trees which analyses the hyper links of the pages and their hierarchies of arrangements to analyse the page and their structure [08]. Some of the researchers have analysed the pattern using different algorithms like Apriori, Hash tree and Fuzzy and then we used enhanced Apriori algorithm to give the solution for Crisp Boundary problem with higher optimized efficiency while comparing to other algorithm [13]. Few have given the detailed review of web mining as another form of data mining [14]. Another aspect of web mining has been also given using two different views i.e. process-centric view which defined web mining as a sequence of tasks, and data-centric- view which defined web mining in terms of the types of web data that was being used in the mining process [15].

IV. PRE-PROCESSING OF LOG

A Web log [16] is a listing of page reference data sometimes it is referred to as click stream data. Raw web log data is not a suitable format usable by mining applications. Therefore, it is necessary to apply pre-processing techniques that may reformat and cleansed the data for mining application. The process of Web Usage Mining [17] includes three phases namely pre-processing, pattern discovery and pattern analysis. Preprocessing is a primary work in web data mining. Pre-processing [18] consists of data integration, data cleaning, user identification and session identification. It eliminates unnecessary records and validates the important records that are saved into the database, which facilitates effective data mining.

V. PROPOSED APPROACH

The proposed method consists of several phases such as file integration or merging, pre-processing, pattern discovery and pattern analysis. This paper focuses only on pre-processing phase that deals with three major issues such as data cleaning, user identification and session identification. The figure 2 shows the complete web data mining work of this paper.

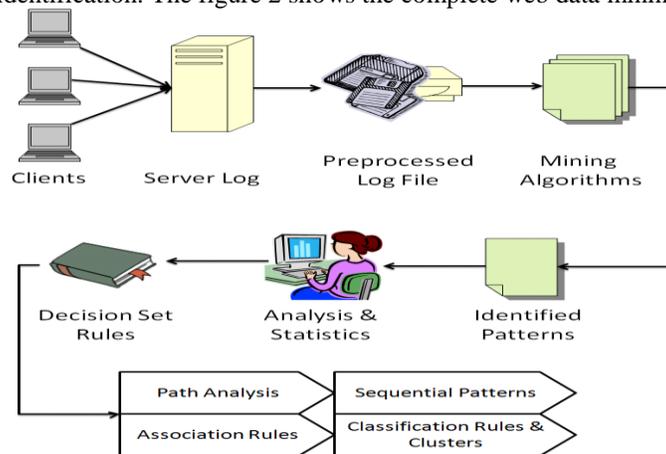


Figure 3: Proposed Methodology

Log File Integration: The integration of content, structure, and user data in other phases of the Web usage mining may also be essential in providing the ability to further analyze and reason about the discovered patterns. The integration of data mining approaches can contribute to create better and more effective intrusion detection system.

Web Usage Mining: Web usage mining attempts to discover knowledge for the data generated by the Web surfer’s sessions or behaviors. Web site servers generate a large volume of data from user accesses. This data help us to determine life time value of users, to improve Web site structure design, to evaluate the efficiency of Web services. Clustering and classification on Web server log file is a process that group the users, Web pages, or user requests on the basis of the access request similarities. Association rule mining task is to discover the correlation among a variety of issues like accesses to files, time of accesses, and identities who requested the accesses.

The proposed work in this paper will be carried out on different log records, which is taken from the multiple web servers. These web log records are then combined together and undergo for the pre-processing phase. Where the unwanted data and non-relevant entries of log will be removed and then user identification will be done.

Algorithm: Data Cleaning

```

begin
    while(!EOF)
    begin
        readLine();
        Check for keywords (bot, slurp, spider)
        if the line contains keyword
            begin
                botflag=true;
                botcounter++;
            end
        else
            botflag=false;
    end
end
    
```

After the pre-processing phase, the cleaned log file will be send to the analysis phase, where data mining algorithm will be applied to this log data to identify the log patterns using association rule mining. These association rules and extracted knowledge are then used by the website administrator to organize the web data and pages according to the popularity and usefulness.

VI. EXPERIMENTAL RESULTS

This section demonstrates a simple walk-through of proposed work. In user session, the browsed pages will be recorded in the log file according to transactional sequences. Web usage mining intelligent system retrieves the useful information from web access log which stored at backend, apart from home page there is others links are used by the visitors. The projected system enforced in the Microsoft .Net environment, that aims to visualizing the x-y plots, bar charts, line graphs, etc., with that a variety of subtle visualization techniques going to be used in information visualization. Information typically consist of variety of dimension and variables, relying on data different visualization are doable. The proposed scheme demonstrated by developing a web application, using Microsoft .Net Framework 4.0, Visual Studio 2008. That is tested on windows environment with IIS 6.0. That is one of the best combinations for making data-driven sites. Since each effort is cooperative in nature, there is continually lots of support from documentation and mailing lists. Bugs are rectified quickly, and requests for features are continually detected, evaluated, and if possible, implemented.

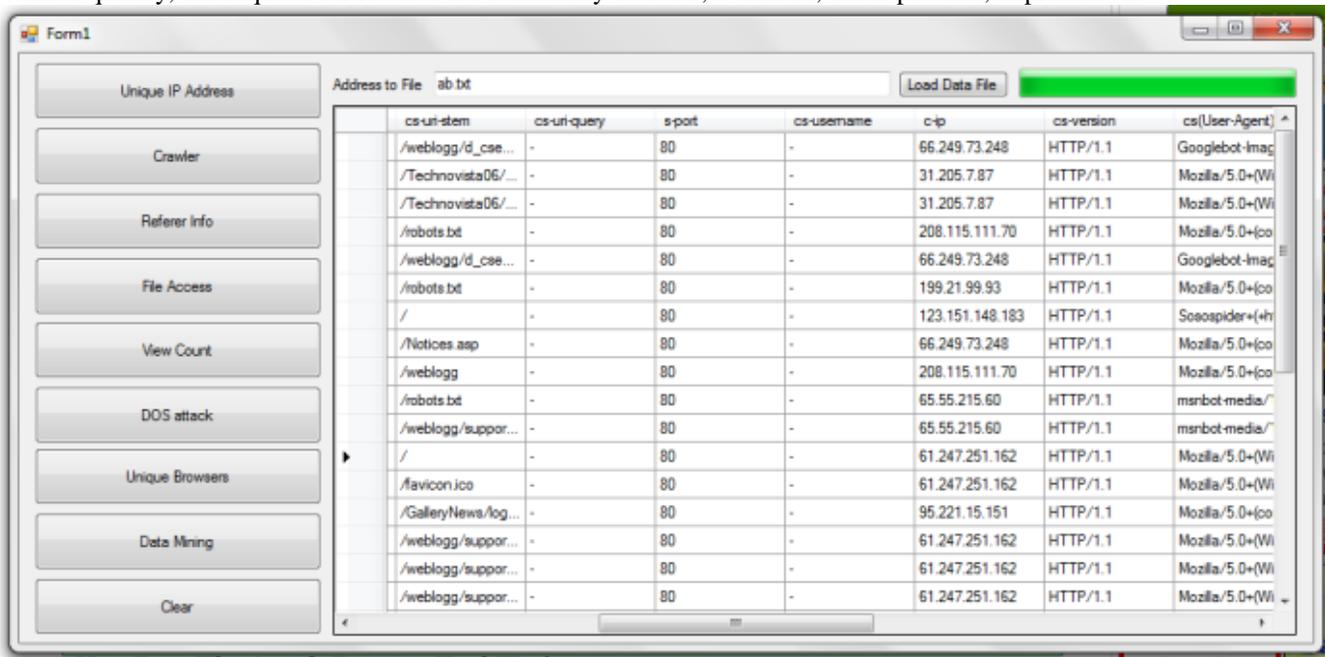


Figure 1: Load Dataset

In fig 5. HTTP codes are shown, that provides an information of DOS attacks, an HTTP-based threat, aims to cause an unseen Denial of Service (DoS) by exploiting a transmission control protocol (TCP) persist timer vulnerability. The attack sends a legitimate protocol request to a server, on the other hand reads the response terribly slowly, forcing connections to remain open. information fills the server's send buffer and holds it there because the server continues to poll the client for write-space accessibility till a DoS happens. Slow Read's response-based behavior is probably going to go unseen because it should apparently mimic a legitimate slow client or because the server isn't sufficiently tuned to shield against slow protocol attacks.

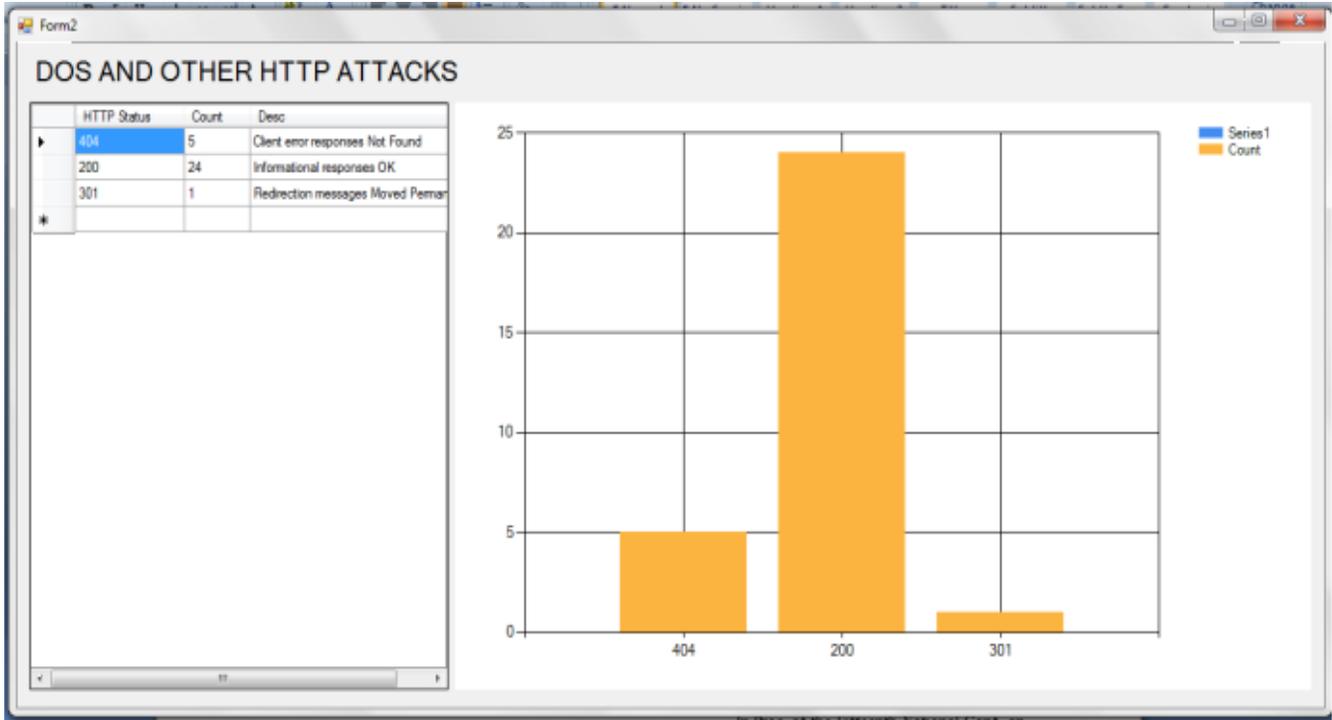


Figure 2: Visualization of DOS attacks

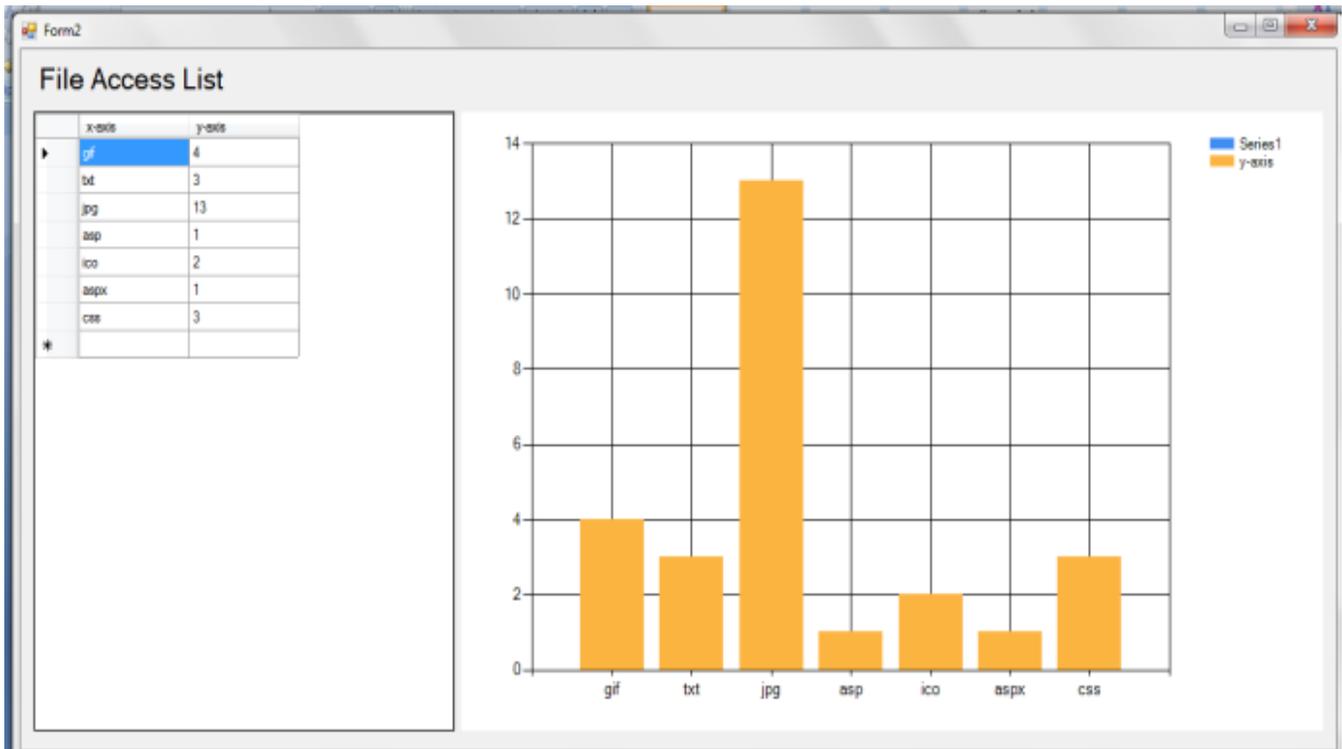


Figure 3: Visualization of file access records

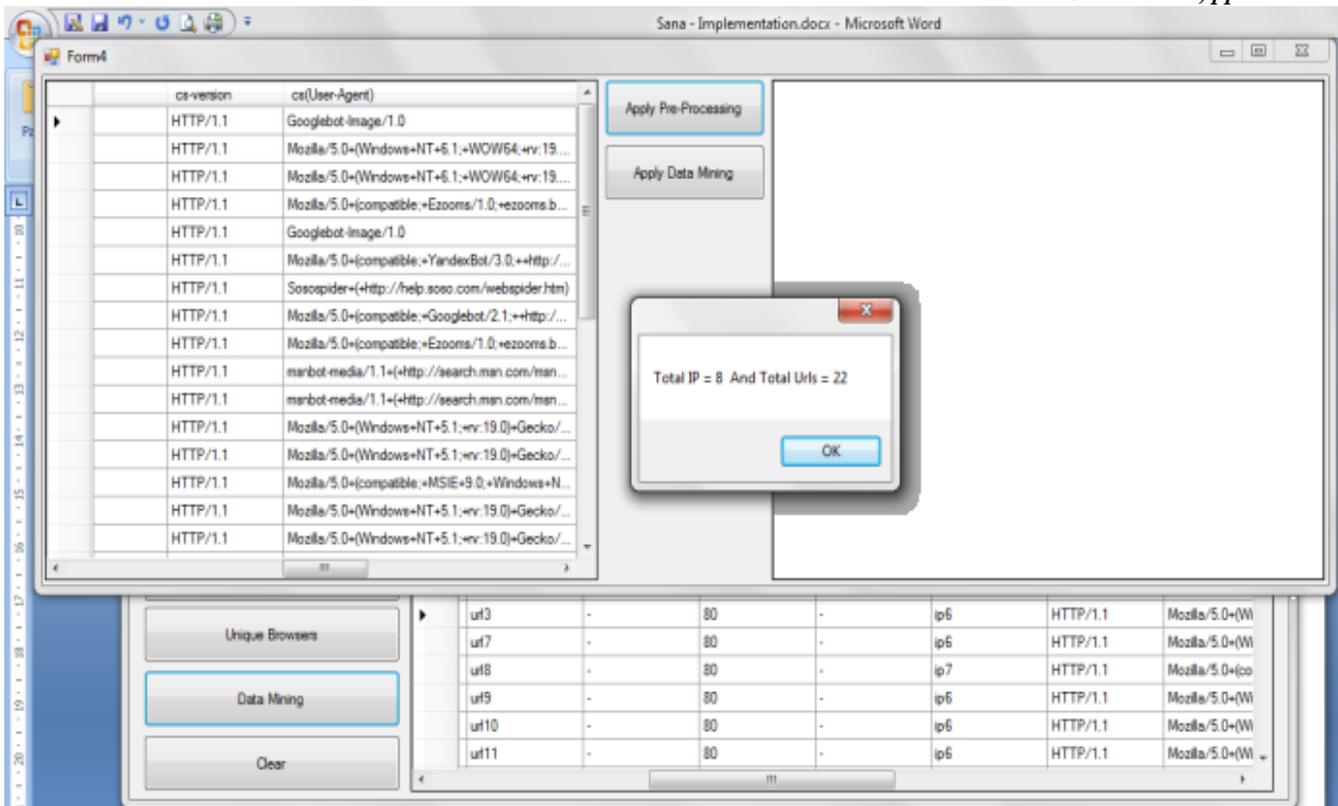


Figure 4: Pre-processing step for web usage mining

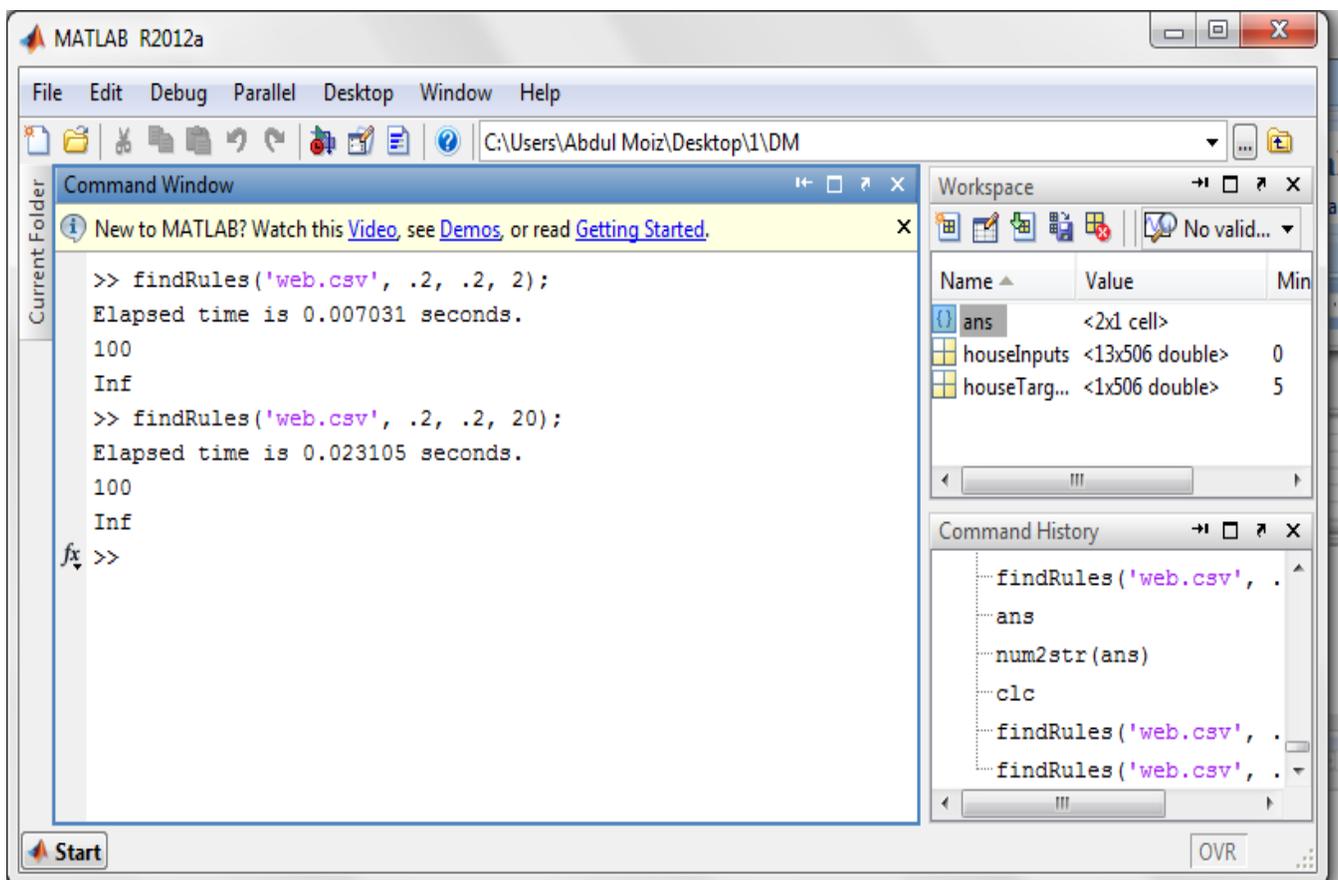


Figure 5: applying efficient pattern mining algorithm based on frequent item sets

```

Rule (Support, Confidence)
1 -> 2 (100%, Inf%)
2 -> 1 (100%, Inf%)
1 -> 3 (100%, Inf%)
3 -> 1 (100%, Inf%)
1 -> 4 (100%, Inf%)
4 -> 1 (100%, Inf%)
1 -> 5 (100%, Inf%)
5 -> 1 (100%, Inf%)
1 -> 6 (100%, Inf%)
6 -> 1 (100%, Inf%)
1 -> 7 (100%, Inf%)
7 -> 1 (100%, Inf%)
2 -> 3 (100%, Inf%)
3 -> 2 (100%, Inf%)

```

Figure 6: Results obtained from rules

In the above fig.9, the output of mining patterns are shown, in which 1,2,3,4,5,6,..... etc specifying the url's. The implemented modern system for creating visualisations have evolved to the extent that non-experts can create meaningful representations of their data. However, the process is still not easy enough, mainly because the visual effects of processing, realising and rendering data are well-understood by the user.

As we used an institute log file for web analytics and mining. It was thought that analysis of long transactions may provide attention-grabbing insights into the patterns of access. Hence, the long transactions from the transaction set were manually analyzed. A remarkable pattern discovered was that some guests tend to seem for various obtainable programs and towards the end of their visit they try to find pages associated with the information E-booklet for the courses they require to pursue. Different discovered pattern unconcealed that some guests attempt to notice contact information of employees members.

In Analysis of internet visitors usage habits may provide vital clues concerning current market trends and facilitate organizations to predict the longer term trends of potential customers. Analysis of long visit-paths of users could indicate the necessity of restructuring of the web site to assist visitors reach desired information quickly. Also, the mined data are often accustomed to most well-liked web content to visitors.

Comparative study of existing systems with the proposed system for web usage mining and visualization of events.

Table 1: Comparison of the proposed system with the existing systems

| | | Proposed Tool | Web-Miner [19] | Web Data Mining [20] | Web Personalization [9] |
|----|------------------------------------|--|----------------|----------------------|---------------------------|
| 1. | Pre-Processing | Yes | Yes | Yes | Yes |
| 2. | Log Parser | Yes | Yes | Yes | No |
| 3. | Visualization of low level events. | Yes | No | No | No |
| 4. | Visualization of High Level Events | Yes | No | No | No |
| 5. | Attacks Detection | Yes | No | No | No |
| 6. | Apply Mining | Yes | Partially | Yes | Yes |
| 7. | Useful Patterns | Yes, Output are shown in results section | No | No | Yes, Output are not shown |

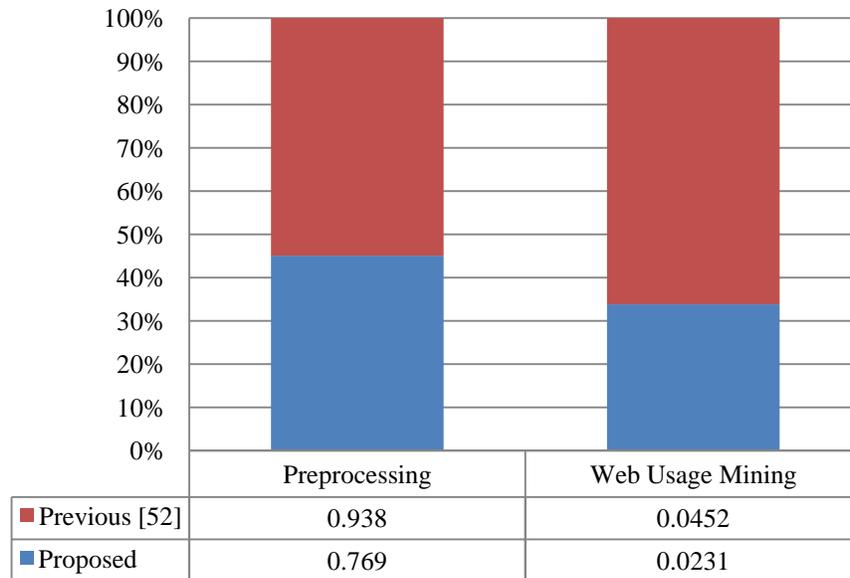


Figure 7: Result comparison for a raw log file of size 610KB

VII. CONCLUSION

With these main concerns, we decide to work for organisation website management and prepare a new reactive approach, which uses the web usage data information, site topology, academic calendar of university, in order to produce more specific process and results for organisation environment. In general this concept may of use to any website organization, A part of this various things may be of help to the system administrator like, analysis of errors helps to know the problems while accessing the website, analysis of references to website during special event will help administrator to know and balance load, analysis of navigational patterns and duration will help the administrator with the knowledge about how to decrease the duration of the user by providing layout change, decrease in duration helps to usage of less bandwidth. While applying pre-processing technique in the sample data, it is found that only 20% of the data contains useful information and the remaining 80% data is not useful for mining. From this 20% data, several frequent patterns are generated using pattern discovery techniques like association rule mining. The result is mainly useful for web page prediction and web site modifications. This work can be extended for massive web log data that is useful for strategic plan.

REFERENCES

- [1] F. Masseglia, P. Poncelet, and M. Teisseire. Using data mining techniques on web access logs to dynamically improve hypertext structure. In *ACM SigWeb Letters*,8(3): 13-19, 1999.
- [2] Maheswara Rao.V.V.R and Valli Kumari.V, 2011. "An Enhanced Pre-Processing Research Framework for Web Log Data Using a Learning Algorithm", *Computer Science and Information Technology*, DOI: 10.5121/csit.2011.1101, pp.01-15.
- [3] Anitha.A, 2010. "A New Web Usage Mining Approach for Next Page Access", *International Journal of Computer Applications (0975-8887)*, Vol. 8, No.11, pp.7-10.
- [4] M. Perkowitz and O. Etzioni. Adaptive web sites: Automatically synthesizing web pages. In *Proc. of the Fifteenth National Conf. on Artificial Intelligence (AAAI)*, pages 727-732, 1998.
- [5] M. Perkowitz and O. Etzioni. Towards adaptive sites:Conceptual framework and case study. In *Proc. of the Eighth Int'l World Wide Web Conf*, Toronto, Canada, May 1999.
- [6] M. Spiliopoulou and L. C. Faulstich. Wum: A web utilization miner. In *Proc. of EDBT WorkshopWebDB98*, Valencia, Spain, March 1998.
- [7] M. Spiliopoulou, L. C. Faulstich, and K. Wilkler. A data miner analyzing the navigational behaviour of web users. In *Proc. of the Workshop on Machine Learning in User Modelling of the ACAI99*, Greece, July 1999.
- [8] Sanjay Madria, Sourav s Bhowmick, w. -k ng, e. P. Lim, "Research Issues in Web Data Mining"
- [9] A.Jebaraj Ratnakumar,"An Implementation of Web Personalization Using Web Mining Techniques", *Journal Of Theoretical And Applied Information Technology*", 2005 - 2010 *JATIT*
- [10] *Data Engineering and Automated Learning (IDEAL'07)*, 16-19 December, 2007, Birmingham, UK.
- [11] J. I. Hong, , J. Heer, S. Waterson, and J. A. Landay,WebQuilt: A proxy-based approach to remote web usability testing, *ACM Transactions*
- [12] Naresh Barsagade, "Web Usage Mining and Pattern Discovery: A Survey Paper ", December 8, 2003

- [13] Thales Sehn Korting, “C4.5 algorithm and Multivariate Decision Trees”, Image Processing Division, National Institute for Space Research – INPES” ao Jos’e dos Campos – SP, Brazil S.Veeramalai , N.Jaisankar and A.Kannan, “Efficient Web Log Mining Using Enhanced Apriori Algorithm with Hash Tree and Fuzzy”, International journal of computer science & information Technology (IJCSIT) Vol.2, No.4, August 2010
- [14] Mr. Dushyant Rathod, “A Review On Web Mining “,International Journal of Engineering Research and Technology (IJERT) Vol. 1 Issue 2, April – 2012 , SSN: 2278-0181
- [15] Jaideep Srivastava, Prasanna Desikan, Vipin Kumar, “Web Mining— Concepts, Applications, and Research Directions”, Page 400-417
- [16] Castellano.G, A. M. Fanelli and M. A. Torsello, 2007. “Log Data Preparation For Mining Web Usage Patterns”, International Conference Applied Computing, pp.371-378.
- [17] Doru Tanasa and Trousse B, 2004. “Advanced Data Preprocessing for Intersites Web Usage Mining”, IEEE Intell Syst, Vol.19, No.2, pp.59-65, DOI: 10.1109 /MIS.2004.1274912
- [18] Dipa Dixit et. al., 2010. “Preprocessing Of Web Logs”, International Journal on Computer Science and Engineering, Vol. 02, No. 07, pp.2447-2452.
- [19] Roop Ranjan, Sameena Naaz, Neeraj Kaushik, “Web Miner: A Tool for Discovery of Usage Patterns From Web Data”, International Journal on Computer Science and Engineering (IJCSSE), Vol. 5 No. 05 May 2013, pp. 286-293, ISSN : 0975-3397.
- [20] M. Malarvizhi, S. A. Sahaaya Arul Mary, “Preprocessing of Educational Institution Web Log Data for Finding Frequent Patterns using Weighted Association Rule Mining Technique”, European Journal of Scientific Research ISSN 1450-216X Vol.74 No.4 (2012), pp. 617-633.