



# Web Page Ranking using Web Usage Mining

Asfiya Khatoon<sup>1</sup>, Kuldeep Jaiswal<sup>2</sup>

M.Tech Scholar, Department of Computer Science and Engineering, BBD University, Lucknow, India<sup>1</sup>

Assistant Professor, Department of Computer Science and Engineering, BBD University, Lucknow, India<sup>2</sup>

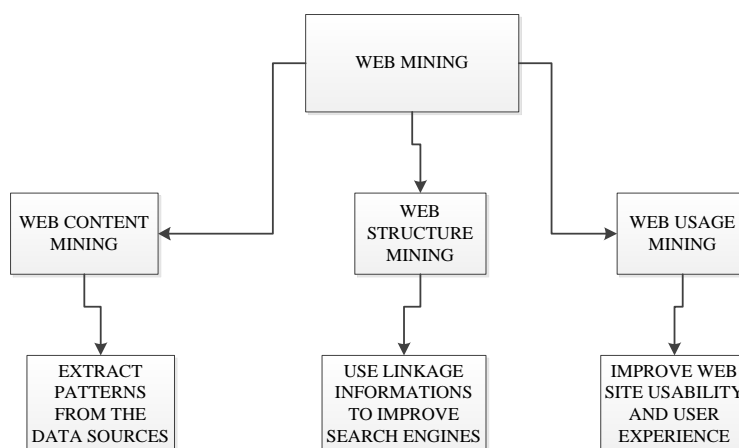
**Abstract:** Nowadays World Wide Web is a popular and interactive medium to distribute the information which develops increasingly. Web Mining is the extraction of interesting and potentially useful patterns and implicit information from artifacts or activity related to the World Wide Web. This extracted information can be further used to improve web usage mining, like next page prediction i.e. prediction of next page accessed by the user, improve web personalization, fraud detection and future prediction accessed by user, user profiling, and also to know about user browsing behavior. For predicting and identifying extracted information, it is needed to process Web Usage Mining. Web Usage Mining processing steps are explained in brief. Due to the difficulty in identifying unique sessions, additional prior knowledge is required. Different Web Usage Mining techniques are described for this knowledge and pattern discovery. Pattern analysis is needed to filter out uninterested rules or patterns from the set found in the pattern discovery phase.

**Keywords:** Web Mining, Web Usage mining, Web log files, fp growth, Association rules mining.

## I. INTRODUCTION

The World Wide Web is one of the main datasources for millions of people in the world in order to access the information from the huge amount of available data. While searching for the particular information on the web, it is important for the user to retrieve the information in less time. The web recommendation models provide access friendliness for users while browsing a website. The prediction models play a vital role in e-commerce, to give advertisement at specific pages of commercial website. Web access prediction is useful in personalization to send personalized web content to specific type of users. Web mining is a kind of data mining techniques to automatically discover and extract information through the analysis of Web contents, Web structure and Web usages. Predicting the normal users' browsing behavior is one of the web usage mining techniques. Web Usage Mining is the process of extracting useful patterns from the web log files. There are different techniques namely SVM, clustering, Page ranking, Markov model, Modified Markov model, Association rule mining, Markov model with clustering etc. has been used for web page prediction [1].

Web Mining is the utilization of information mining methods to consequently find and concentrate Information from web records and administrations. The World Wide Web, www or web is turning into an unpredictable universe. Actually, determining something significant out of it is focused on utilization of web mining [2]. Three sub classes:



- Web Content Mining
- Web Structure Mining
- Web Usage Mining



#### A. Web Content Mining

Web Content Mining alludes to the revelation of helpful data from the web content, here Content alludes to Text, Audio Video and so forth that various sites are holding. Web content mining gets to be entangled when it has to mine unstructured, organized, semi organized and mixed media information [2].

Case of Web Content Mining is:

Normal Google or Yahoo or Microsoft Bing seek that we do, and the resultant connections posting page we get is a case of substance mining. The way toward removing valuable data from the web content happens here.

#### B. Web Usage Mining

Web utilization mining process includes the log time of pages. The world's biggest gateway like hurray, msn and so forth., needs a considerable measure of bits of knowledge from the conduct of their users' web visits. Without this utilization reports, it will be hard to structure their adaptation endeavors. Use mining has direct effect on organizations [4]. The difficulties required in web utilization mining could be isolated in three stages:

##### A. Pre-handling.

The information accessible have a tendency to be loud, inadequate and conflicting. In this stage, the information accessible ought to be dealt with as indicated by the necessities of the following stage. It incorporates information cleaning, information joining, and information change and information diminishment.

##### B. Design revelation.

A few distinct strategies and calculations, for example, measurements, information mining, machine learning and example acknowledgment could be connected to recognize client designs.

##### C. Design Analysis.

This procedure focuses to comprehend, imagine and offer translation to these examples. Case of Web Usage Mining is: A specific component of site that is utilized by the guests oftentimes, that we need to upgrade and declare to build the utilization that can claim more to clients of the site Simply by comprehension the development of the visitors and the conduct of surfing the net, you can anticipate address the inclinations and the issues in a superior way and promote your site among the masses in the web field [3].

#### C. Web Structure Mining

Web structure mining is done at the hyper join level. This sort of mining tries to find the model basic the connection structure of the web. A pertinent illustration can be Google's Page rank. The objective of the Web Structure Mining is to produce the basic outline about the Web website and Web page. It tries to find the connection structure of the hyperlinks at the between report level. In view of the topology of the hyperlinks, Web Structure mining will arrange the Web pages and create the data like likeness and relationship between various Web locales. This sort of mining can be performed at the archive level (intra-page) or at the hyperlink level (between pages). It is essential to comprehend the Web information structure for Information Retrieval. The Web contains an assortment of articles with no bringing together structure, with contrasts in the composing style and substance much more noteworthy than in customary accumulations of content archives. The items in the WWW are site pages, and connections are in, out and co-reference i.e. two pages that are both connected to the same page. There are some conceivable undertakings of connection mining which are material in Web structure mining and are depicted as takes after: [5]

##### A. Join based Classification:

The latest redesign of an exemplary information mining undertaking to connect Domains. The undertaking is to concentrate on the expectation of the class of a website page, in light of words that happen on the page, joins between pages, stay content, html labels and other conceivable qualities found on the page.

##### B. Join based Cluster Analysis.

The objective in bunch investigation is to discover normally happening sub-classes. The information is sectioned into gatherings, where comparable items are assembled together, and unique articles are assembled into various gatherings. Not quite the same as the past errand, join based bunch investigation is unsupervised and can be utilized to find concealed examples from information.

##### C. Join Type.

There are an extensive variety of assignments concerning the forecast of the presence of connections, for example, foreseeing the kind of connection between two substances, or anticipating the motivation behind a connection.

**D. Join Strength**

Connections could be connected with weights.

**E. Join Cardinality**

The fundamental undertaking here is to foresee the quantity of connections between articles. There are a few employments of web structure mining as is it:

- Used to rank the user's question
- Deciding what page will be added to the accumulation
- Page arrangement
- Finding related pages
- Finding copied sites
- And likewise to discover similitude between them.

**II. WEB LOGES FILES**

Web Log Files are files that contain information about website visitor activity. Log files are created by web servers automatically. Each time a visitor requests any file (page, Image, etc.) from the site, information of his request is appended to a current log file. Most log files have text format and each log entry (hit) is saved as a line of text. Log file range 1KB to 100MB [7].

The information sources utilized as a part of Web Usage Mining may incorporate web information stores like:

**Web Server Logs:**

These are logs which keep up a background marked by page demands. The W3C keeps up a standard organization for web server log documents, however other exclusive configurations exist. Later sections are normally annexed to the end of the record. Data about the solicitation, including customer IP address, demand date/time, page asked for, HTTP code, bytes served, client specialist, and referrer are ordinarily included. This information can be consolidated into a solitary document, or isolated into particular logs, for example, an entrance log, blunder log, or referrer log. In any case, server logs regularly don't gather client particular data.

These documents are normally not open to general Internet clients, just to the website admin or other authoritative individual. A measurable examination of the server log might be utilized to look at activity designs by time of day, day of week, referrer, or client operator. Effective site organization, sufficient facilitating assets and the calibrating of offers endeavors can be supported by examination of the web server logs. Showcasing divisions of any association that possesses a site ought to be prepared to comprehend these effective apparatuses [6].

**Intermediary Server Logs:**

A Web intermediary is a storing component which lies between customer programs and Web servers. It diminishes the heap time of Web pages and also the system activity load at the server and customer side. Intermediary server logs contain the HTTP asks for from different customers to various Web servers. This may serve as an information source to find the utilization example of a gathering of mysterious clients, sharing a typical intermediary server.

**Program Logs:**

Various programs like Mozilla, Internet Explorer and so on can be adjusted or different JavaScript and Java applets can be utilized to gather customer side information. This usage of customer side information accumulation requires client participation, either in empowering the usefulness of the JavaScript and Java applets, or to deliberately utilize the altered program. Customer side gathering scores over server-side accumulation since it decreases both the boot and session recognizable proof issues [9].

**III. PHASES OF WEB USAGE MINING PROCESS**

The fundamental procedures in Web Usage Mining are:

**A. Preprocessing**

Data preprocessing depicts any sort of handling performed on crude information to set it up for another preparing system. Generally utilized as a preparatory information mining rehearse, information preprocessing changes the information into an organization that will be all the more effortlessly and viably handled with the end goal of the client. The distinctive sorts of preprocessing in Web Usage Mining are:

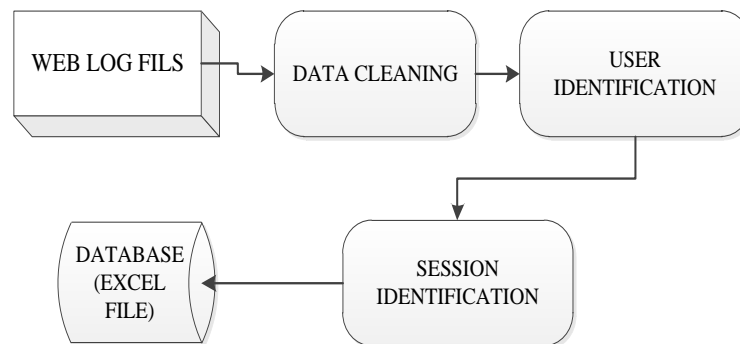


Fig. 1 Pre-Processing Phase of WUM

1. Utilization Pre-Processing: Pre-Processing identifying with Usage examples of clients.
2. Content Pre-Processing: Pre-Processing of substance got to.
3. Structure Pre-Processing: Pre-Processing identified with structure of the site.

Design Discovery: Web Usage mining can be utilized to reveal designs in server logs however is regularly completed just on tests of information. The mining procedure will be insufficient if the examples are not a decent representation of the bigger assortment of information. The accompanying are the example revelation strategies.

1. Factual Analysis
2. Affiliation Rules
3. Bunching
4. Grouping
5. Successive Patterns
6. Reliance Modeling

Design Analysis: This is the last stride in the Web Usage Mining process. After the preprocessing and example disclosure, the got use examples are examined to channel uninteresting data and concentrate the valuable data. The techniques like SQL (Structured Query Language) preparing and OLAP (Online Analytical Processing) can be utilized.

#### B. Pattern Discovery

In Pattern Discovery phase, data mining techniques like association rule mining and fp growth applied on web log files after preprocessing to discover the useful pattern.

#### Association Rule Mining:

It is usual relate pages that are most as often as possible referenced along amid a solitary server session. In connection of web Usage Mining, affiliation rules deliberate with set of pages that are gotten to close to a bolster esteem over some, for example, limit. Affiliation principle mining misuse it's usual relate pages that are most as often as possible referenced along amid a solitary server session. In setting of web Usage Mining, affiliation rules consult with set of pages that are gotten to next to a bolster cost over some, for example, limit. Affiliation principle mining abuse Apriori algorithmic system could see connection between's clients who went to a page containing electronic stock to those that get to a page with respect to don related stock. A common case of continuous itemset mining is business sector bushel examination. This technique investigates customer looking for propensities by discovering relationship between totally diverse things that clients place in their "shopping baskets (crate)". The innovation of such affiliations will encourage retailers create advancing courses by picking up understanding into that things are much of the time obtained along by clients.

#### FP Growth:

FP-development utilizes a blend of the vertical and even database design to store the database in primary memory. Rather than putting away the spread for each thing the database, it stores the real exchanges from the database in a tire structure and each thing has a connected rundown experiencing all exchanges that contain that thing. This new information structure is meant by FP-tree (Frequent Pattern tree).

#### C. Pattern Analysis

In Pattern Analysis phase, irrelevant pattern are remove from the pattern which identified during pattern discovery phase. The main purpose of pattern analysis is to analyze the pattern which is identified during pattern discovery phase.



#### IV. PATTERN DISCOVERY

##### A. Association Rules Method

Affiliation standard mining, a standout amongst the most critical and all around looked into strategies of information mining [9]. Are utilized to distinguish connections among an arrangement of things in a database. These connections are not in light of innate properties of the information themselves (as with useful conditions), but instead in view of co-event of the information items.[17] Association standards are generally utilized as a part of different territories, for example, telecom systems, market and hazard administration, and instruction [9]. Affiliation guideline mining is to discover affiliation decides that fulfill the predefined least backing and certainty from a given database. The issue is typically decayed into two sub issues. One is to discover those thing sets whose events surpass a predefined limit in the database; those thing sets are called successive or vast thing sets. The second issue is to create affiliation rules from those expansive thing sets with the imperatives of negligible certainty [9].

The primary sub-issue can be further partitioned into two sub-issues: hopeful substantial thing sets era prepare and visit thing sets era process. Those thing sets whose backing surpasses the bolster edge as vast or successive thing sets, those thing sets that are normal or have the would like to be huge or regular are called applicant thing sets. Much of the time, the calculations create a to a great degree huge number of affiliation standards, frequently in thousands or even millions. Further, the affiliation guidelines are now and then vast. It is about outlandish for the end clients to understand or approve such huge number of complex affiliation rules, along these lines restricting the convenience of the information mining results [12].

A few systems have been proposed to decrease the quantity of affiliation principles, for example, creating just —interestingl rules, producing just —non-redundantl controls, or creating just those tenets fulfilling certain other criteria, for example, scope, influence, lift or quality.

##### Association Rules Metrics:

Two essential measures for affiliation rules bolster (s) and certainty ( $\alpha$ ), can be characterized as takes after:

##### Support:

Definition 1: As the rate/division of records that contain to the aggregate number of records in the database. The mean everything is expanded by one each time the thing is experienced in various exchange T in database D amid the filtering procedure. It implies the bolster tally does not consider the amount of the thing [17].

Definition 2: Support of a tenet is a measure of how much of the time the things required in it happen together. Utilizing likelihood documentation: bolster (An infers B) =  $P(A, B)$  [12]. Support(s) is figured by the accompanying recipe [9]:

From the definition we can see, backing of a thing is a measurable importance of an affiliation guideline. Assume the backing of a thing is 0.1%, it implies just 0.1 percent of the exchange contains buying of this thing. Prior to the mining procedure, clients can indicate the base backing as an edge, which implies they are just keen on certain affiliation decides that are created from those thing sets whose backings surpass that limit [9].

##### Certainty:

Definition 1: Confidence of a principle is the contingent likelihood of B given A. Utilizing likelihood documentation: certainty (An infers B) =  $P(B \text{ given } A)$  [12].

Definition 2: the rate/division of the quantity of exchanges that contain to the aggregate number of records that contain X, where if the rate surpasses the edge of certainty a fascinating affiliation standard can be produced [9]. Support(s) is figured by the accompanying equation [9]:

##### B.FP-Growth

FP-development utilizes a blend of the vertical and even database design to store the database in primary memory. Rather than putting away the spread for each thing the database, it stores the real exchanges from the database in a tire structure and each thing has a connected rundown experiencing all exchanges that contain that thing. This new information structure is meant by FP-tree (Frequent Pattern tree).

FP-Growth regular example mining is utilized as a part of the advancement of affiliation principle mining. FP-Growth calculation conquers the issue found in Apriori calculation. The successive thing set era process requires just two ignores the database there is no requirement for competitor era process. By dodging the competitor era process and less ignores the database, FP-Growth observed to be speedier than the Apriori calculation [16].



A FP-Growth is a prefix tree for exchanges; each hub in the tree speaks to one thing and every way speaks to the arrangement of exchanges that include with the specific thing. All hubs alluding to the same thing are connected together in a rundown, so that every one of the exchanges that containing the same thing can be effectively found and tallied [16].

FP-Growth calculation includes the era of successive examples utilizing the continuous examples era process which incorporates two sub forms:

- Constructing the FP-Growth.
- Generation of regular examples from the FP-Growth.

The way toward developing the FP-Tree is as per the following [16]:

- (1) The database is checked interestingly, amid this examining the bolster tally of every things are gathered. Therefore the incessant 1 - thing sets are produced procedure is the same as in Apriori calculation. Those regular thing sets are sorted in a slipping request of their backings. Additionally the head table of requested incessant 1 - thing sets is made.
- (2) Create the root hub of the FP-Tree T with a name of Root. The database is filtered again to build the FP-Tree with the head table, for every exchange the request of frequent.

## V. PROPOSED APPROACH FOR PATTERN DISCOVER

The proposed approach for web usage mining is shown in figure below. This is a three step process.

1. The beginning is collecting of web log file and performing preprocessing operation on it. After preprocessing processed web log file is stored in database i.e. in excel file.
2. We can apply any data mining technique like association rule mining, fp growth for pattern discovery. Here combined approach of fp growth and association rule mining is used. Thus in step 2 partitioning based fp growth itemset is used to find user's having common behavior and access patterns
3. Finally in step 3 association rule mining technique will be used to find user's access patterns from this itemset group of data.

Due to this approach fp tree will be performed and every user will be assigned to specific itemset according to behavior & access patterns. Applying association rule mining technique on this itemset data will help to find results having less computing time and better accuracy.

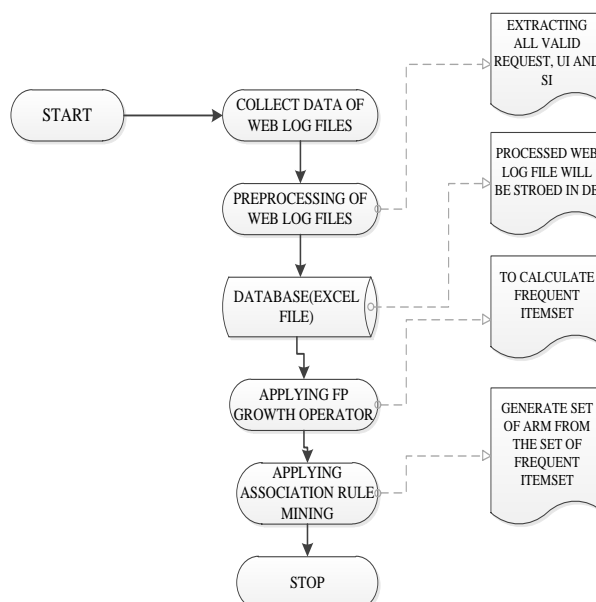


Fig 2: Proposed Approach for Pattern Discovery in WUM

## VI. IMPLEMENTATION

Web Usage Mining is implemented on web server log files as input. Then apply preprocessing on web log file and store into the database i.e. excel files. We can generate useful pattern from web log file by association rule mining and fp growth. The following figure shows step wise implementation:



Step 1: Raw web log files are choose from where it is stored.

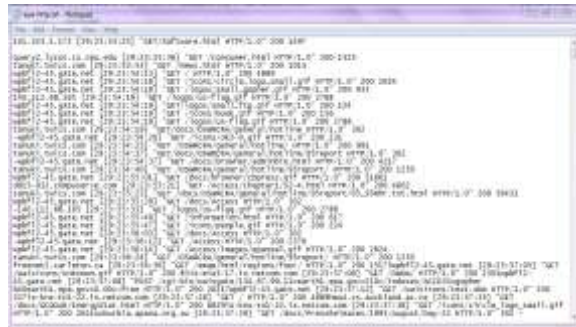


Fig 3: Choose raw web log files

Step 2: Apply the preprocessing on web log files and store them into the database.

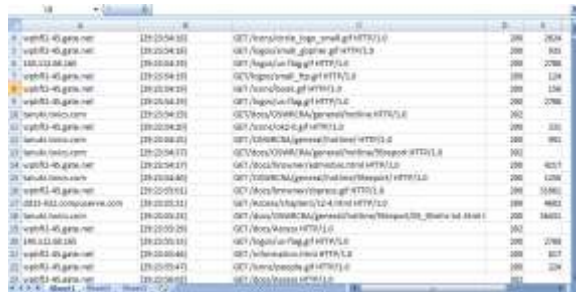


Fig 4: Raw web log files after preprocessing

Step 3: Unique users and webpages are identified from web log files.

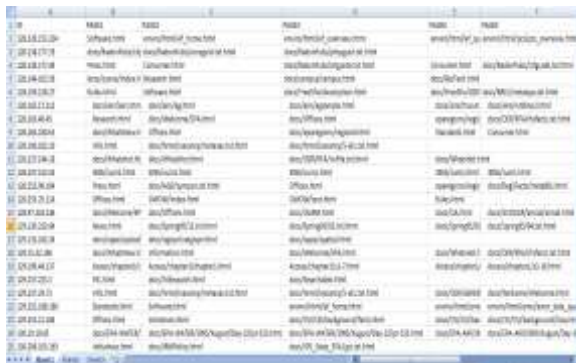


Fig 5: Unique users and webpages

Step 4: FP growth and association rule mining is applied on web log files and get frequently accessed webpages.

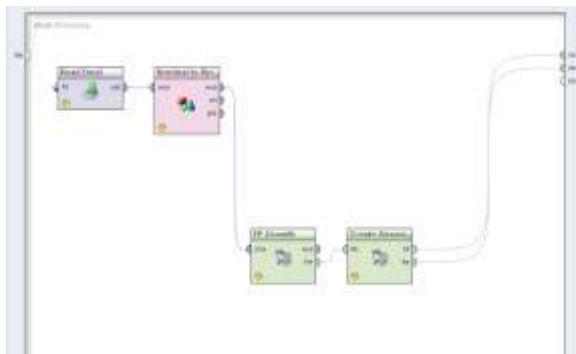


Fig 6: Process for carrying out Association Rule Mining



Item	Antecedent	Consequent	Confidence
1	PAGE1	PAGE2	0.75
2	PAGE2	PAGE3	0.75
3	PAGE1, PAGE2	PAGE4	0.75
4	PAGE1	PAGE5	0.75
5	PAGE2	PAGE6	0.75
6	PAGE1, PAGE2	PAGE7	0.75
7	PAGE1	PAGE8	0.75
8	PAGE2	PAGE9	0.75
9	PAGE1, PAGE2	PAGE10	0.75
10	PAGE1	PAGE11	0.75
11	PAGE2	PAGE12	0.75
12	PAGE1, PAGE2	PAGE13	0.75
13	PAGE1	PAGE14	0.75
14	PAGE2	PAGE15	0.75
15	PAGE1, PAGE2	PAGE16	0.75
16	PAGE1	PAGE17	0.75
17	PAGE2	PAGE18	0.75
18	PAGE1, PAGE2	PAGE19	0.75
19	PAGE1	PAGE20	0.75
20	PAGE2	PAGE21	0.75
21	PAGE1, PAGE2	PAGE22	0.75
22	PAGE1	PAGE23	0.75
23	PAGE2	PAGE24	0.75
24	PAGE1, PAGE2	PAGE25	0.75

Fig 7: showing the association rule of web pages confidence 0.75

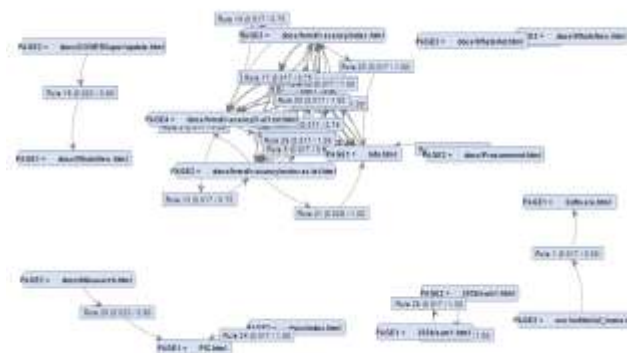


Fig 8: Graphical Plot of the Web Page Association using confidence 0.75

V. CONCLUSION

Web site is a group of web pages. Web pages may contain text, images, and videos. They are linked by hyperlinks through which navigation happens. Whenever user accesses any website, log files are created. Log file records entire information about each user’s website access. Due to increase in usage of web sites the size of log files is increasing day by day. Data stored in Web log files can exist in various formats such as the NCSAs Common log file format, the W3C Extended Log File format or the IIS log file format. There are different kinds of log files which includes Error logs, Referrer logs, and Access logs. Log files are created in various locations like web server, proxy server, and Client browser. For getting optimum results we need to extract data from all three log files. Analysis of the patterns of user’s habits and interests helps in increasing performance of web site, by improving web site design.

Web mining is the application of various data mining techniques to discover data from web documents and services. Web mining is divided into three types Web content mining, Web structure mining and Web usage mining. Web Content Mining deals with the discovery of information from the contents or data or documents or services of web. Web Structure Mining mines the structure of hyperlinks within the website. Web Usage Mining mines the usage data stored in the logs. Web usage mining analyzes information about web pages which were navigated by users. Analysis of such information helps us to discover the unknown and potentially interesting patterns.

REFERENCES

[1] R. Khosala, H. Blockeel, Web Mining Research: A Survey, ACM SIGKDD, 2000.  
 [2] A. Abraham, Business Intelligence from web usage mining, Journal of Information and Knowledge Management, 2003.  
 [3] Srivastava J., Cooley R., Deshpande M. and Tan. P.N., ‘Web Usage Mining: Discovery and applications of Usage Pattern from web data.’ SIGKDD Explorations 2000.  
 [4] J. Han and M. Kamber, Data Mining: Concepts and Techniques, 2<sup>nd</sup> Edition, Elsevier Morgan Kaufmann Publishers, San Francisco, USA, 2006.  
 [5] Jiang M.; Shyong S. and Liao T. ‘Data Types Generalization for Data Mining Algorithms’, National Science Council of the Republic of China, 1994.  
 [6] C. Ding, X. He, P. Husbands, H. Zha, and H. Simon, "Link Analysis: Hubs and Authorities on the World". Technical Report: 47847, 2001. [14] Wenpu Xing and Ali Ghorbani, “Weighted PageRank Algorithm”, Proceedings of the Second Annual Conference on Communication Networks and Services Research (CNSR '04), IEEE, 2004.  
 [7] Ali Mohammad ZarehBidoki and Nasser Yazdani, “DistanceRank: An Intelligent Ranking Algorithm for Web Pages”, Information Processing and Management, 2007.





- [8] Fabrizio Lamberti, Andrea Sanna and Claudio Demartini, "A Relation- Based Page Rank Algorithm for. Semantic Web Search Engines", In IEEE Transaction of KDE, Vol. 21, No. 1, Jan 2009.
- [9] Lian-Wang Lee, Jung-Yi Jiang, ChunDer Wu, Shie-Jue Lee, "A Query- Dependent Ranking Approach for Search Engines", Second International Workshop on Computer Science and Engineering, Vol. 1, PP. 259-263, 2009..
- [10] Su Cheng, PanYunTao, YuanJunPeng, GuoHong, YuZhengLu and Hu ZhiYu "PageRank, "HITS and Impact Factor for Journal Ranking", Inproceedings of the 2009 WRI World Congress on Computer Science and Information Engineering – Vol. 06, PP. 285-290, 2009 .
- [11] NeelamDuhan, A.K.Sharma and Komal Kumar Bhatia, Page Ranking Algorithms: In proceedings of the IEEE International Advanced Computing Conference (IACC), 2009.
- [12] Xiang Lian and Lei Chen, "Ranked Query Processing in Uncertain databases", In IEEE KDE, Vol. 22, No. 3, March 2010.
- [13] P Ravi Kumar, and Singh Ashutoshkumar, "Web Structure Mining Exploring Hyperlinks and Algorithms for Information Retrieval", American Journal of applied sciences, 7 (6) 840-845 2010.
- [14] Saeko Nomura, Tetsuo Hayamizu, "Analysis and Improvement of HITS Algorithm for Detecting Web Communities". Volume 11-No 08, 2011.
- [15] Rekha Jain, DrG.N.Purohit, "Page Ranking Algorithms for Web Mining", International Journal of Computer application, Vol 13, Jan 2011.
- [16] Tamanna Bhatia, "Link Analysis Algorithms For Web Mining", IJCST Vol. 2, Issue 2, June 2011. 43
- [17] W.Xing and Ali Ghorbani, "Weighted PageRank Algorithm", Proc. Of the Second Annual Conference on Communication Networks and Services Research, IEEE, 2013. 48
- [18] A.M. Sote, Dr. S. R. Pande "Application of Page Ranking Algorithm in Web Mining" International Conference on Advances in Engineering & Technology-2014.
- [19] Wenpu Xing and Ali Ghorbani, "Weighted PageRank Algorithm", Proceedings of the Second Annual Conference on Communication Networks and Services Research (CNSR '04), IEEE, 2004.
- [20] Ko Fujimura, Takafumi Inoue and Masayuki Sugisaki, "TheEigenRumor Algorithm for Ranking Blogs", In WWW 2005 2nd Annual Workshop on the Weblogging Ecosystem, 2005.