# A Review of-Web Mining

K. Abirami and Dr.P. Mayilvahanan

**Abstract**--- This study presents the role of the alarming rate at which the World Wide Web (WWW) is growing in both the sheer volume of traffic and the complexity of websites, it has become very important to analyze this web traffic and the usage of the web sites by the users. This is the review paper which show deep and intense study of various technologies available for web mining and it is the application of data mining techniques to extract knowledge from web. present advances in each of the three different types of web mining are reviewed in the categories of web content mining, web usage mining, and web structure mining.

**Index Terms**--- Web Usage Mining, Web Content Mining, Web Mining, Web Structure Mining, Web Mining Software

## I. INTRODUCTION

Web may be the principle and prevalent sourball from claiming data available, reachable and receptive in low expense gives fast light of the clients besides diminishes load on the clients of physical developments. That information on the web is loud that complaint goes starting with two major wellsprings. in a significant Web page holds a large number bits of information, e.g. , the principle substance of the page, directing links, advertisements, copyright notices, protection policies, and so forth. Second, because of the way that those Web doesn't need incredibleness control of information, i.e. , you quit offering on that one might compose practically anything that person likes, an expansive amount about data on the Web may be about low quality, erroneous, or much misdirecting.

Retrieving of the needed web page on the web, effectively and effectively, will be getting to be a troublesome.

Web mining is a requisition of information mining which need get to be a noteworthy region for investigate because of inconceivable measure for planet totally Web administrations in late a long time. Those developing field for web mining plans in discovering and extracting important majority of the data that is stowed away over Web-related data, specifically for content documents distributed on the Web. Those review looking into information mining system is constructed with admiration to Clustering, Classification, arrangement design Mining, cooperation standard mining and Visualization [1]. Those examine fill in done by separate clients portraying those pros and cons are examined.

## II. WEB MINING

Web mining aims to discover helpful information or knowledge from the Web hyperlink structure, page content, and usage data. Although Web mining uses many data mining techniques, as mentioned above it is not purely an application of traditional data mining due to the heterogeneity and semi-structured or unstructured nature of the Web data. Many new mining tasks and algorithms were invented in the past decade. Web mining tasks can be categorized into three types: Web structure mining, Web content mining and Web usage mining.

*K. Abirami, Research Scholar, Vels University, Chennai–600117, India. E-mail:abiramidharmarajna@gmail.com*
*Dr.P. Mayilvahanan, Head, Department of MCA, Chennai, India. E-mail:hodmca@velsuniv.org*
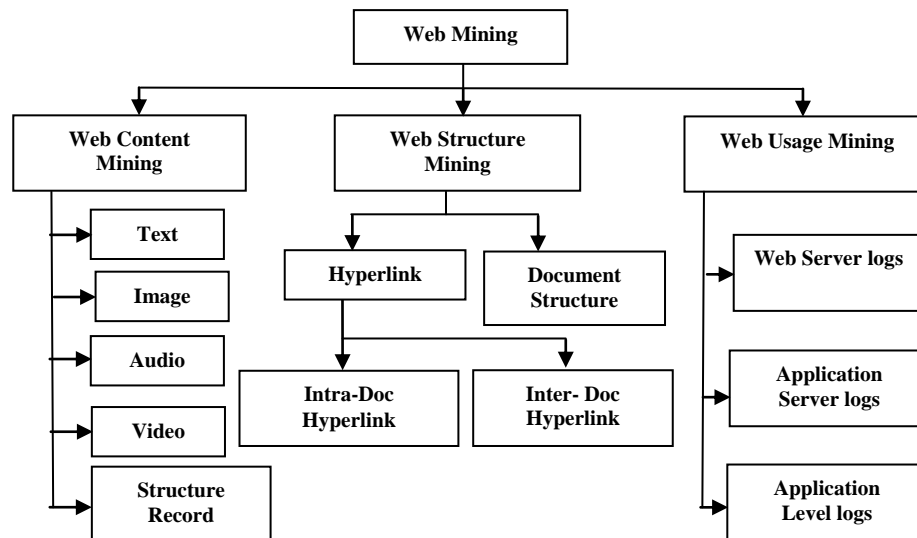
Fig. 1: Web Mining Categories

*Web structure mining:* Web structure mining uncovers helpful information starting with hyperlinks (or joins for short), which speak to the structure of the Web. For example, starting with those links, we came uncover significant Web pages, which, incidentally, will be a magic innovation organization utilized within quest engines. We came additionally uncover groups about clients who stake regular investment. Accepted information mining doesn't perform such errands on there is as a rule no join structure done a social table.

*Web content mining:* Web content mining extracts or mines suitable data alternately learning starting with Web page substance. To example, we might naturally arrange and group Web pages as stated by their topics. These assignments need aid comparative on the individuals clinched alongside accepted information mining. However, we came additionally uncover designs for Web pages with extricate functional information for example, portrayals of products, postings of forums, etc, to a large number purposes.Additionally, we might mine client reviews and also gathering postings to find shopper sentiments.

*Web usage mining:* Web use mining alludes all the of the disclosure for client right examples from web use logs, which record each click aggravated toward each client. Web utilization mining applies large portions information mining

calculations. A standout amongst the key issues over web use mining will be those pre-processing from claiming click stream information on utilization logs in place to prepare those correct information to mining.

## III. SURVEY ON WEB CONTENT MINING

Web content mining is the process of extracting useful information from the contents of web documents. Content data is the collection of facts a web page is designed to contain [6]. It may consist of text, images, audio, video, or structured records such as lists and tables [1].

Table 1: Web Content Mining Using Different Algorithms

| WEB CONTENT MINING | | |
|---|---|---|
| *Author* | *Representation* | *Method Used* |
| (Ahonen, 1998) | Bag of words and word positions | Episode rules |
| (Billsus & Pazzani, 1999) | Bag of words | TFIDF Naïve Bayes |
| (Cohen, 1995) | Relational | Propositional rule based system Inductive Logic Programming |
| (Dumais, 1998) | Bag of words - Phrases | - TFIDF - Decision trees - Naïve Bayes -Bayes nets |

| WEB CONTENT MINING | | |
|---|---|---|
| *Author* | *Representation* | *Method Used* |
| | | - Support Vector Machines |
| (Feldman & Dagan, 1995) | Concept categories | Relative entropy |
| (Feldman, 1998) | Terms | Association rules |
| (Frank, 1998) | Phrases and their positions | Naïve Bayes |
| (Freitag & McCallum, 1999) | Bag of words | Hidden Markov Models |
| (Hoffmann, 1999) | Bag of words | Unsupervised statistical Method |
| (Junker, 1999) | Relational | Inductive Logic Programming |
| (Kargupta, 1999) | Bag of words with n grams | - Unsupervised hierarchical clustering - Decision trees - Statistical analysis |
| (Nahm & Mooney, 2000) | Bag of words | Decision trees |
| (Nigam, 1999) | Bag of words | Maximum entropy |
| (Scott & Matwin, 1999) | - Bag of words - Phrases - Hyponyms and synonyms | Rule based system |
| (Witten, 1999) | Named entity | Text compression |
| (Yang, 1999) | Bag if words and phrases | -Clustering algorithms - K-Nearest Neighbor - Decision tree |
| (Genersereth and Nilsson, 1987) | set of objects | ontology |

Those web substance information comprise about unstructured information for example, allowed texts, semi-structured information for example, html documents, and a a greater amount organized information for example, information in the tables alternately database produced html pages. So, two fundamental methodologies for web substance mining arise, (1) unstructured content mining methodology Also (2) Semi-Structured Furthermore organized mining methodology. In this area we start Toward reviewing exactly of the critical issues that Web substance mining plans on take care of. We At that point rundown a portion of the separate methodologies in this field arranged rely on upon those diverse sorts of Web substance information. Previously, every methodology we rundown a few of the The greater part utilized systems.

Those different grouping procedure are takes after: content built grouping : the content built bunch methodologies describe each report as stated by its, i.e. those expressions confined to it. Those fundamental thought will be that whether two documents hold huge numbers basic expressions then it will be really could be allowed that the two archive need aid fundamentally the same. The methodologies in this class came make further sorted accounting of the grouping strategy utilized under those emulating categories: partition, Hierarchical, Chart Based, Probabilistic calculations.

## IV. SURVEY ON STRUCTURE MINING

The challenge for Web structure mining is to deal with the structure of the hyperlinks within the Web itself. Link analysis and Stochastic Approach for Link-Structure Analysis (SALSA) InDegree are an old area of research [4]. The Web contains a variety of objects with almost no unify structure, with differences in the authoring style and content much greater than in traditional collections of text documents. The link analysis algorithm contains page rank, weighted page rank and HITS [3].

Table 2: Web Structure Mining using Different Algorithms

| WEB STRUCTURE MINING | | |
|---|---|---|
| *Algorithms Used* | *Author* | *Year* |
| In Degree | Marchiori | 1997 |
| Page Rank | Brin and Page | 1998 |
| Link Analysis | Kleinberg | 1998 |
| HITS | Klienberg | 1999 |

| WEB STRUCTURE MINING | | |
|---|---|---|
| *Algorithms Used* | *Author* | *Year* |
| PHITS | Cohn and Chang | 2000 |
| SALSA | Lempel and Moran | 2000 |
| Weighted Page Rank | Wenpu Xing and Ali Ghorbani | 2004 |
| Page Rank based on visits of links | Gyanendra Kumar, Neelam Duhan, A. K. Sharma | 2011 |
| Weighted Page Rank based on visits of links(VOL) | Neelam Tyagi, Simple Sharma | 2012 |

### A.  HITS (Hyper-link Induced Topic Search)

A HIT is a purely link-based algorithm. It is used to rank pages that are retrieved from the Web, based on their textual contents to a given query. Once these pages have been assemble, the HITS algorithm ignores textual content and focuses itself on the structure of the Web only.

### B.  Weighted Page Rank (WPR):

The only tip of the iceberg mainstream webpages would the that's only the tip of the iceberg linkages that different webpages have a tendency should must them alternately need aid joined with Toward them. Those suggested broadened PageRank algorithm–a Weighted PageRank calculation assigns bigger rank values will that's only the tip of the iceberg imperative (popular) pages As opposed to separating those rank worth of a page uniformly "around its outlink pages. Every outlink page gets An esteem proportional will its Ubiquity (its amount of inlinks Furthermore outlinks).

### C.  Page Rank Algorithm

Pageranking algorithms are the heart of search engine and give result that suites best in user expectation. Need of best quality results are the main reason in innovation of different page ranking algorithms, HITS, PageRank, Weighted PageRank, DistanceRank, DirichletRank Algorithm , Page content ranking are different examples of page ranking used in different scenario. Since GOOGLE search engine has great importance now days and this affect many web users now days, so page rank algorithm used by GOOGLE become very significant to researches [2].

### D.  Page Rank Based on VOL

We have seen that original Page Rank algorithm, the rank score of a page p, is equally divided among its outgoing links or we can say for a page, an inbound links brings rank value from base page, p( rank value of page p divided by number of links on that page)[3]. Which more rank value is assigned to the outgoing links which is most visited by users. In this manner a page rank value is calculate based on visits of inbound links.

### E.  Result Analysis

This section compares the page rank of web pages using standard Weighted PageRank (WPR), Weighted PageRank using VOL ($WPR_{VOL}$) and the proposed algorithm. We have calculated rank value of each page based on WPR, $WPR_{VOL}$ and proposed algorithm i.e. $EWPR_{VOL}$ for a web graph shown in Table2.

Table 3: Comparison of Different Algorithms

| *Algorithm* | *PageRank* | *Weighted PageRank* | *PageRank with VOL* | *Weighted PageRank with VOL* |
|---|---|---|---|---|
| Web mining technique used | Web Structure mining | Web Structure mining | Web structure mining, web usage mining | Web structure mining, web usage mining |
| Input Parameters | Backlinks | Backlinks, Forward links | Backlinks and VOL | Backlinks and VOL |
| Importance | More | More | More | More |
| Relevancy | Less | Less | More | More |

The values of page rank using WPR, $WPR_{VOL}$ and $EWPR_{VOL}$ have been compared [5]. The values retrieved by EWPRVOL are better than original WPR and $WPR_{VOL}$. The WPR uses only web structure mining to calculate the value of page rank, $WPR_{VOL}$ uses both web structure mining and

web usage mining to calculate value of page rank but it uses popularity only from the number of inlinks not from the number of outlinks. The proposed algorithm EWPR$_{VOL}$ method uses number of visits of inlinks and outlinks to calculate values of page rank and gives more rank to important pages.

## V.    SURVEY ON WEB USAGE MINING

Web Usage Mining is the application of Data Mining to discover and analyze patterns from clickstreams, user transactions and other logged user interactions with a website. The goal is to capture, model and analyze the behavior of users, and define patterns and profiles from the captured behaviors. There are three phases: data collection and pre-processing, pattern discovery, and pattern psychiatry.

Data collection and pre-processing: this concerns the generating and cleaning of web data and transforming it to a set of user transactions representing activities of each user during his/her website visit. This step will influence the quality and result of the pattern discovery and analysis following therefore it needs to be done very carefully [8].

Pattern discovery: during pattern discovery, information is analyzed using methods and algorithms to identify patterns. Patterns can be found using various techniques such as statistics, data mining, machine learning and pattern recognition

Pattern analysis: it describes the filtering of uninteresting and misleading patterns. The content and structure information of a website can be used to filter patterns.
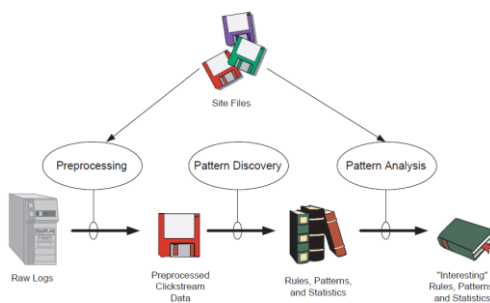


Fig. 2: Web Usage Mining Process

They are web server data, application server data and application level data. Web server data correspond to the user logs that are collected at Web server. Some of the typical data collected at a Web server include IP addresses, page references, and access time of the users and is the main input to the present Research. This Research work concentrates on web usage mining and in particular focuses on discovering the web usage patterns of websites from the server log files.

The result of Web Usage Mining process is usually an aggregated user model, which describes the behavior of user groups or pinpoints a trend in user behavior. The discovery of user access patterns from the user access logs, referrer logs, user registration logs etc is the main purpose of the Web Usage Mining activity. We then list some of the different approaches in this field classified depend on the different types of Web usage data. In each approach we list some of the most used techniques.

Table 4: Web Usage Mining using Different Algorithms

| WEB USAGE MINING | | |
|---|---|---|
| *Algorithms Used* | *Author* | *Year* |
| fuzzy clustering | Bezdek | 1981 |
| Self-Organizing Map | Kohonen | 1982 |
| Association Rules | Agrawal | 1993 |
| Ontologies | Gruber | 1993 |
| Apriori or FP Growth Module | Agrawal and R. Srikant | 1994 |
| Direct Hashing and Pruning | J. S. Park, M. Chen, P.S. Yu | 1995 |
| Sequential Patterns | R. Agrawal and R. Srikant | 1995 |
| Generalized Sequential Pattern | R. Srikant and R. Agrawal | 1996 |
| Parameter Space Partition | Shiffrin & Nobel | 1998 |
| FP-GROWTH | Jiawei Han, Jian Pei, Yiwen Yin | 2000 |
| Vertical data format | Zaki | 2000 |
| TREE-PROJECTION | Ramesh C. Agarwal, Charu C. Aggarwal, V.V.V. Prasad | 2000 |
| Baraglia and Palmerini | SUGGEST | 2002 |

**WEB USAGE MINING**

| Algorithms Used | Author | Year |
|---|---|---|
| An average linear time algorithm | José Borges , Mark Levene | 2004 |
| Harmony | Wang et al | 2005 |
| semantic web mining | Berendt | 2005 |
| Frequent pattern-based classification | Cheng et al | 2007 |
| Lee and Fu | pattern-growth principl | 2008 |
| Tree-based frequent patterns | Fan et al | 2008 |
| Zhihua Zhang | intelligent algorithm | 2009 |
| Sequential pattern mining with $K^{th}$ order  Markov model clustering | A. Anitha | 2010 |
| Mehrdad, Norwati Ali, Md Nasir | LCS Algorithm, clustering | 2010 |
| Bing Liu's | tools & technology | 2011 |
| Nicolas Poggi, Vinod Muthusamy, David Carrera, and Rania Khalaf | process mining techniques | 2013 |

## VI.    CONCLUSION

The most important task of the Web Usage Mining process is data preparation. Now a day Web mining become very popular, interactive and innovative technique and it is application of the Data Mining technique that mechanically discovers or extracts the information from web documents. In this paper have provided a more current evaluation study research papers the various algorithms methods, techniques, phases that are used for web mining and its three categories. This paper has provided the well-organized algorithms of web mining to have an idea about in their application and effectiveness. Weighted Page Content Rank user can get relevant and important pages easily as it employs web structure mining. The new approach uses different technique using Genetic Algorithm (GA) for web content mining. Web usage can combine FP-Tree with Apriori candidate generation method to solve the disadvantages of both apriori and FP-growth. Web usage mining is used by e-commerce sites to organize their sites and to increase profits.  Since this

is a broad area, and there a lot of work to do, we wish this paper could be a useful for identifying opportunities for further research.

## REFERENCE

[1] Dushyant Rathod, "A Review on Web Mining," IJERT, vol. 1, Issue 2, April – 2012.
[2] KaushalKumar, Abhaya and Fungayi Donewell Mukoko, "PageRank algorithm and its variations: A Survey report", IOSR-JCE., Vol 14, Issue 1, Sep. - Oct. 2013, PP 38-45.
[3] Sonal Tuteja," Enhancement in Weighted PageRank Algorithm Using VOL," in IOSR-JCE, Volume 14, Issue 5 Sep. - Oct. 2013), PP 135-141.
[4] Tamanna Bhatia, "Link Analysis Algorithms For Web Mining," in IJCST Vol. 2, Issue 2, June 2011.
[5] ALLAN BORODIN, GARETH O. ROBERTS , JEFFREY S. ROSENTHAL , and PANAYIOTIS TSAPARAS "Link Analysis Ranking: Algorithms, Theory,and Experiments",  ACM Transactions on Internet Technology, Vol. 5, No. 1, February 2005, Pages 231–297.
[6] D. Jayalatchumy, and  P.Thambidurai, "Web Mining Research Issues and Future Directions – A Survey," IOSR-JCE, Vol 14, Issue 3 ,Sep. - Oct. 2013, PP 20-27.
[7] Michael Azmy, "Web Content Mining Research: A Survey", DRAFT Version 1, - Nov. 2005.
[8] J Vellingiri, and  S.Chenthur Pandian, "A Survey on Web Usage Mining," in Global Journals Inc. (USA), Vol 1, Issue 4 Version 1.0 March 2011.