



SURVEY ON WEB STRUCTURE MINING

B. L. Shivakumar¹ and T. Mylsami²

¹Department of Computer Applications, Sri Ramakrishna Engineering College, Coimbatore, Tamil Nadu, India

²Department of Computer Science and Information Technology, Dr. G.R. Damodaran College of Science, Coimbatore, Tamil Nadu, India

E-Mail: blshiva@yahoo.com

ABSTRACT

In recent days the data generation is enormous in all the fields. Same as in Internet the data generation is high and there is no control over the data generation. To retrieve the exact data required by the online consumer is a tedious task. To achieve the same is done by data mining methods and its techniques. The data mining concept consist of web mining methods. The term web mining extracts the required information to user and to reach the necessary goal in the website. To attain the goal, use the concept of web mining. Web mining divides into web content, web structure and usage mining. Web structure mining plays very significant role in web mining process. The future algorithms for web structure mining such as Pagerank Algorithm, HITS, Weighted Pagerank Algorithm, Weighted page content rank Algorithm (WPCR) and soon. In this paper, identify their strengths and limitations of different algorithms used in web mining.

Keywords: web mining, web structure, pagerank, hyper link inducted topic search, weighted pagerank.

INTRODUCTION

Data mining is defined as computer based method of finding and analyzing large amount of data and the process of discovering interesting useful format and its relationships. The key role is to extract the unfamiliar, useful and logical patterns from large database. Data mining concentrate on automatic discovery of patterns, creation of action statements and expected output. The development of internet data in website, make the user difficult to browse webpage effectively. To increase the performance of web sites design and the web server activities are changed as per users view and requirement. Various applications like commerce, personalization, web site designing, recommender systems are built efficiently by knowing users routing through website. Web mining is the application of data mining techniques to automatically retrieve, extract and evaluate information for knowledge discovery from web documents. The structure of Web consists of graph structured by documents and hyperlinks, the mining results may be on Web contents or Web structures. Web mining focused on three areas like Web content mining, Web structure mining and Web usage mining.

Web mining

The application of data mining techniques used on a website to discover interesting patterns. The application of data mining techniques to extract knowledge from Web content, structure, and usage. Web mining have an interface with data mining, is the process discovers knowledge from vast amount of information without users intervention. The technologies are that fulfill the possible of extracting valuable knowledge from the World Wide Web and its usage guide. Web mining

techniques specifically used to extract knowledge from Web data, and the data in the forms like web documents, hyperlinks between documents, and usage logs of web sites. The extracted data from internet can be filtered by content based filtering and Collaborative filtering. This work provides an overview of web structure and its uses. [1, 2, 3, 4].

Data mining aims at discovering valuable information that is hidden in conventional databases, the emerging field of Web mining aims at finding and extracting relevant information that is hidden in Web-related data; particular in text documents that are published on the web. Figure-1 represents the complete process of Web mining process and extracting knowledge from web data.

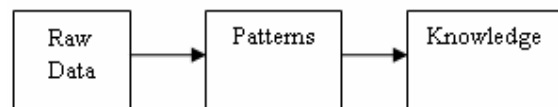


Figure-1. Web mining process and extracting knowledge.

The above concept works in the order of raw data converts into patterns. The patterns have been generated from the smaller unit data using processing steps. Then followed by using methods, the semi processed data called patterns have transferred in to knowledge. Finally the knowledge is utilized for extracting data effective from the website for users' requirements.

Figure-2 illustrates the complete layout of web mining.



www.arpnjournals.com

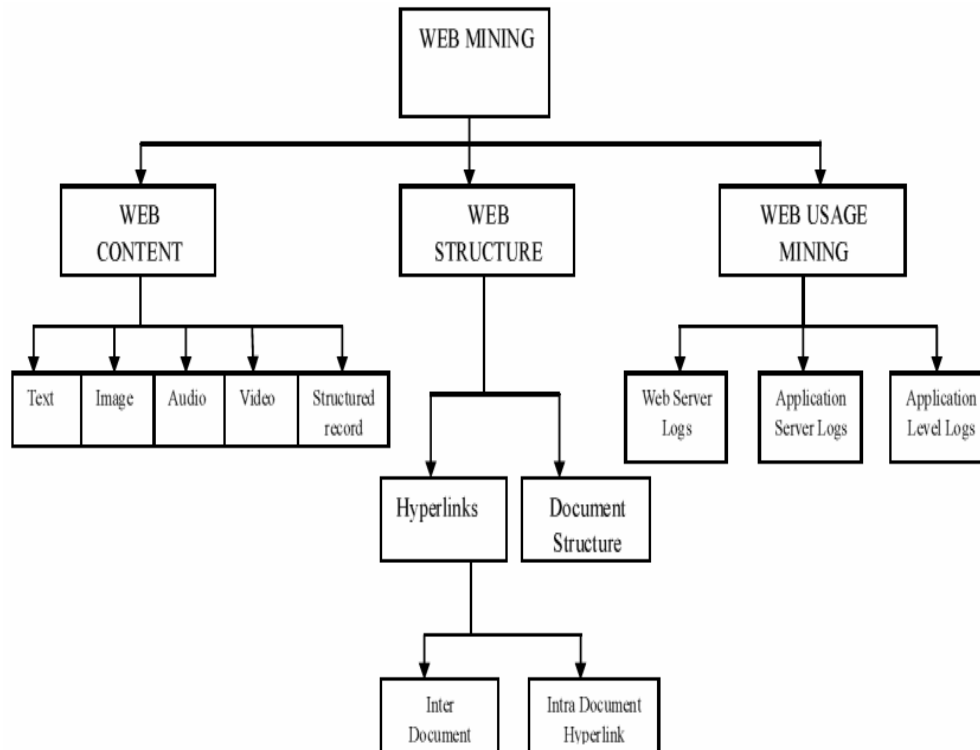


Figure-2. Layout of Web Mining.

From nature of the data, user can differentiate three main areas of research within the Web mining community [5, 6, 7].

a. Web content mining: Application of data mining techniques to unstructured or semi-structured data, usually HTML-documents. It is the automatic search of information resources available online. It aims to extract knowledge from web page contents. The web content mining is different from text mining and data mining. Web content mining completely deals with semi structured data, but data mining deals with structured data. Web content mining requires creative applications of data mining or text mining techniques. The technologies that are normally used in web content mining are natural language processing (NLP) and Information retrieval (IR). [1, 8, 9, 10].

b. Web structure mining: It uses the hyperlink structure of the Web as an (additional) information source. It generates the structural summary for the web sites and web pages. The aim of the web structure mining is to generate the structural abstract about the websites and webpage. It establish the link construction of the

hyperlinks at the inter text level. The topology used in web structure mining is that will categorize the web pages and spawn the information like similarity and relationship between the different websites. [1, 8, 9, 11].

c. Web usage mining: The process of extracting useful information from server logs. The key role is that finding user's requirement on Internet. Among the users level they differ in searching the data, one may be text data and other category of user search data related to multimedia and soon. It is the discovery of user access patterns from Web servers. Web usage mining is the process to identify the browsing patterns by analyzing the procedure followed by the user. It involves the automatic discovery of used access pattern from one or more web servers. Through the available methodologies in data mining user can discover the data need on Internet. It consists of three stages namely preprocessing, pattern discovery and pattern analysis. The above three steps play an eminent role in web mining tasks. Each component has distinct functions in the whole system. In among the three levels the web structure mining concentrate on hyperlink between the webpage and its functions. [1, 8, 9].

**Table-1.** Web mining categories.

	Web mining			
	Web content mining		Web structure mining	Web usage mining
View of data	Unstructured Structured	Semi Structured Web site as DB	Link structure	Interactivity
Main data	Text documents Hypertext documents	Hypertext documents	Link structure	Server logs Browser logs
Representation	Bag of words, n- gram terms, phrases, relational	Edge labeled graphs. Relational	Graph	Relational table Graph
Method	Machine learning Statistical (Include NLP)	Proprietary algorithms Association rules	Proprietary algorithms	Machine learning Statistical Association rule
Application categories	Categorization Clustering Finding extract rules Findings patterns in text.	Finding frequent sub structures. Web site schema discovery.	Categorization Clustering	Site construction Adaption and management. Marketing. User modeling.

Web structure mining

The WWW is one of the most important resources for information generation and the retrievals of data also an eminent step in web. The knowledge is discovered with the help of stable increasing of the amount of data generated in online. Considering the web aspect, the online users get easily lost in the web's loaded hyper structure. Through the available application of data mining methods leads to the perfect solution for knowledge discovery on the Web. [12, 13, 14].

The knowledge extracted from the Web can be used to raise the performances for Web information retrievals, question answering and Web based data warehousing. Web structure mining, one of three categories of web mining for data, is a tool used to identify the relationship between Web pages linked by information or direct link connection. It offers information about how different pages are linked together to form this huge web. Web Structure Mining finds hidden basic structures and uses hyperlinks for more web applications such as web search.

The continuous growth and spread of the internet using Web Mining to improve the quality of different services has become a necessity. Web Mining is nothing else than applying data mining techniques and algorithms on web data. The work concentrates on functions of various algorithms used in Mining namely Page Rank and HITS. The algorithms draw their origin from social networks analysis and they are modeled based on the Theory of Markov Chains. Page Rank algorithms have been used in many search engine application and among the popular search engine like GOOGLE and as the same HITS concept has been utilized by the search engine like

CLEVER. We present their strengths, weakness and other areas of applicability.

Web structure mining is the process of analyzing the hyperlink and mine important information from it and steps to achieve the information is tedious one. Likewise the remaining also used to mine the structure of document, analyze the structure of page and to describe the HTML format or XML usage. The primary objective of the Web Structure Mining is to generate the structural synopsis about the Web site and Web page. Web Structure mining will sort out the Web pages in different category and from the category to generate the information like the similarity and relationship between different Web sites. The type of mining can be either performed at document level called as Intra-page and at the same time another level is performed at hyperlink level called inter-page mining.

The challenge for web structure mining is to deal with the structure of the hyperlinks within the web itself. Link analysis is an old and traditional method in research areas. In growing interest in web mining areas, the research of structure analysis has given a new path for research area called Link mining. Link mining had produced some demonstration on some of the traditional data mining tasks. Some of the possible tasks of link mining which are applicable in web structure mining.

- a) Link based Classification - Link based classification is the most recent upgrade of a classic data mining task to linked domains. The task is to focus on the prediction of the category of a web page, based on words that occur on the page, links between web pages.



- b) Link based Cluster analysis - The primary objective of cluster analysis is to find logically taking place sub classes. Then the data is categories into groups with similar objects are group and contradictory objects as another groups.
- c) Link Type - Extensive range of tasks relating to guess of the existence of the links, such guess the type of link between two entities.
- d) Link Strength - Links could be related with weights.
- e) Link Cardinality - The key task is to expect the number of associations between objects.

Functions of web structure mining

Figure-3 represents the hierarchical structure for Web Structure Mining [15].

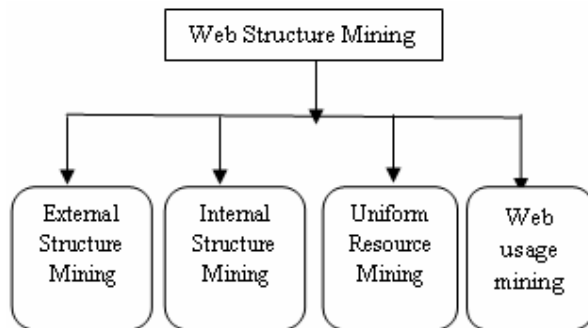


Figure-3. Hierarchical structure for Web Structure Mining.

From the above Figure the Hierarchical structure for Web Structure Mining is also known as “Link Analysis” process. The research of structure analysis had increased in value and focuses on future research concept with scope and we named called as Link mining. The Web contains a variety of stuff with almost no identical or standard structure, and it differs in the design / style and as well the content to be better than in usual collections of text documents.

The uses of web structure mining in among the online users.

- (i) It used to rank the online users queries.
- (ii) Used to finding the related web pages from the website.
- (iii) Finding the similarity between the websites and its category.
- (iv) Improving navigation of web pages on business websites.
- (v) It used to mine the previously unidentified link between web pages.
- (vi) Discovering the structure of web document.
- (vii) Structure mining can be used to reveal the structure (Schema) of web pages.

Some of the popular web mining algorithms

The frequently used algorithms in web structure mining, to access the webpage or website effectively for user in online.

- a) Page rank algorithm.
- b) HITS algorithms(Hyperlink-Induced Topic Search)
- c) Weighted page rank algorithm.
- d) Distance rank algorithm.
- e) Weighted page content rank algorithm.
- f) Webpage ranking using Link attributes.
- g) Eigen Rumor Algorithm
- h) Time Rank Algorithm
- i) Tag Rank Algorithm
- j) Query Dependent Ranking Algorithm

a. Page rank algorithms

The famous author Brin and Page (1998) has developed Page Rank algorithm during their Ph.D. research work at Stanford University based on the extract analysis. They suggested that the world fame search engine called Google has been created with help of Page Rank algorithm. They applied the extract analysis in Web search by treating the incoming links as credentials to the Web pages. Page rank algorithm is the most frequently used algorithm for ranking the various pages. Functioning of the page rank [2, 6, 12].

The Page Rank leads a better approach that can calculate the importance of web page by simply including the number of pages that are linking to it. The above calculated links are called as backlinks. In case a backlink generates from a key page and then this link is given higher weightage than those which are coming from non-important pages. The link from one page to another is measured as a vote. Not only the number of votes that a page receives is important but the importance of pages that casts the vote is also important. Figure-4 shows the backlinks tree structure.

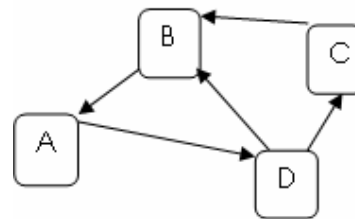


Figure-4. Sample tree structure for back link.

Page and Brin proposed a formula to calculate the Pagerank of a page A as stated below-

$$PR(A) = (1-d) + d(PR(T_1)/C(T_1) + \dots + PR(T_n)/C(T_n))$$

here $PR(T_i)$ is the Pagerank of the Pages T_i which links to page A, $C(T_i)$ is number of outlinks on page T_i and d is



damping factor. It is used to prevent other pages having too much authority.

The Pagerank forms a likelihood distribution over the web pages and the sum of Pagerank of all web pages will be one. The Pagerank of a page can be calculated without knowing the final value of Pagerank of other pages. It is an iterative algorithm which follows the principle of normalized link matrix of web. Pagerank of a page depends on the number of pages pointing to a page. Page rank algorithms require a minimum of little hours to calculate the rank millions of pages. The likelihood distribution algorithm used to represent the person randomly clicking on links will appear at any particular page. A probability is expressed as a numeric value between 0 and 1.

b. Hyper link induced topic search (HITS) algorithm

Jon Kleinberg suggested highlights that the Hyperlink - Induced Topic Search is defined with two attributes like hubs and authorities. It uses the link analysis algorithm rates the web pages and invented developed by He developed an algorithm that made use of the link structure of the web in order to discover and rank pages relevant for a particular topic. HITS (*hyperlink-induced topic search*) are now part of the Ask search engine. The HITS algorithm uses the Sampling and Iterative steps. In the Sampling step, a set of related pages for the given query are collected i.e. a sub-graph S of G is retrieved which is high in influence pages. This algorithm starts with a root set B , a set of S is obtained, keeping in mind that S is comparatively small, rich in relevant pages about the query and contains most of the good authorities. The second step, Iterative step, finds hubs and authorities using the output of the sampling. [2, 6, 15, 16].

The key objective of the algorithm is that by viewing the one mode web graph actually comprising two modes called hubs and authorities. A hub is a node primarily with edges to authorities and so a good hub has links to many authorities. A good authority is a page that is linked to by many hubs. Starting with a specific search objective, HITS algorithm performs a text based search to seed an initial set of results. An iterative relaxation algorithm will assign hubs and authority weights using matrix power iteration. The CLEVER search engine is built primarily using the basics of HITS algorithm [17].

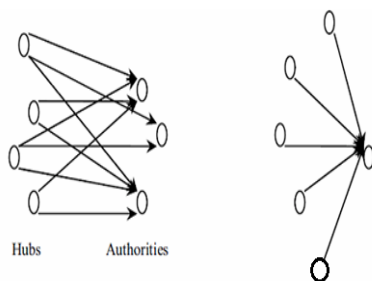


Figure-5. Sample tree structure for Hub and Authority.

Hyperlink-induced topic search (HITS) is an iterative algorithm for mining the Web graph to categorize the topic hubs and it symbolize the pages with good sources of content/data and authorities symbolize pages with good sources of links / hyperlink. Authorities are highly ranked pages for a given topic; hubs are pages with links to authorities. This algorithm takes as input search results returned by usual text indexing techniques and it reduces the results to identify hubs and authorities. The number value and weight of hubs pointing to a web page decide the page's authority. And also the algorithm assigns weight to a hub based on the validity of the pages it points to.

Figure-6 shows the webpage layout structure which exists in HITS algorithm.

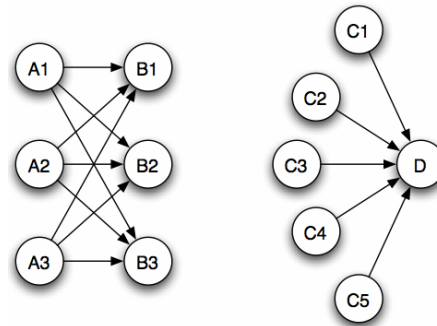


Figure-6. Layout structure for webpage using hub and authority.

The HITS algorithm performs a chain of iterations, each consisting of two fundamental steps.

Authority update - This is to update each nodes authority score to be equal to the sum of the hub scores of each node that points to it. A node with high authority score is being linked to by pages that are recognized as Hubs for information.

Hub Update - This is to update each nodes Hub score to be equal to the sum of the authority scores of each node that it points to it. A node with high hub score is linking to nodes that are considered to be authorities on the subject.

The constraints of HITS algorithm

Hubs and authorities: This is not simple to differentiate between hubs and authorities and since many sites are hubs as well as authorities.

Topic drift: In fewer cases the HITS algorithm may not bring into being the most appropriate documents to the user queries because of equivalent weights.

Automatically generated links: HITS gives equivalent importance for repeatedly generated links that might not produce relevant topics for the user requirements. HITS algorithm is not efficient in real time and as well in many web related applications.

HITS algorithm provides better results for online results and did not work well in all suitcases due to few reasons.



In quite common difficult situation like common relationship between hosts - The exact meaning is that a set of documents on one host point to a single document to the second node. As the same way in another possible way the single document on one host point to a set of document to the second node. The above reason will leads to wrong information about the hub and authority in the website. i.e. referring the same link.

Automatically generated links - In a web document the data generated through tools and have links were between the webpage is inserted by the tool.

Non relevant nodes - In website, few cases one page point to another with no significance to the uncertainty topic.

c. Weighted page rank algorithm

This algorithm was proposed by Wenpu Xing and Ali Ghorbani which is an extension of Page Rank algorithm. This Algorithm assigns rank values to pages according to their importance rather than dividing it evenly. Wenpu Xing and Ali Ghorbani expected a Weighted Pagerank (WPR) algorithm is similar and updated over the Page Rank algorithm. The strategy of this algorithm assigns a larger rank values to the more significant pages, relatively than separating the rank value of a page evenly among its outgoing linked pages. [2, 6].

This is an extension of the page rank algorithm, is assign to both back link and forward link. The incoming link is treated as the number of link points to that particular page and outgoing link is treated as the number links going out from the source. The algorithm is more efficient than page rank algorithm due to using of two parameters called back link and forward link. The status from the number of in links and out links is recorded and easy assign the label as W_{in} and W_{out} . The meaning is assigned in terms of weight values to incoming and out coming links. This is denoted as $W_{in}(m, n)$ and $W_{out}(m, n)$ respectively. $W_{in}(m, n)$ is the weight of links (m, n) as given in the equation. Finally the calculation is based on the number of incoming links to page n and the number of incoming links to all references pages of page m [12].

$$W_{in}(m, n) = \frac{I_n}{\sum P E R(m)}$$

I_n - is denoted as the number of incoming links of page n, I_p - is denoted as the number of incoming links of page p, $R(m)$ is the reference page list page m. $W_{out}(m, n)$ is the weight of links (m, n) as given the equation. The final value is calculated on the basis of outgoing links of page n and the number of outgoing links of all the reference pages of page m.

$$W_{out}(m, n) = \frac{O_n}{\sum P E R(m)}$$

O_n - is number of outgoing links of page n, O_p is number of outgoing links of page p. Then the weighted Page rank is calculated as follow as

$$WPR(n) = (1-d) + d \sum WPR(m) W_{in}(m, n) W_{out}(m, n)$$

Weighted page rank algorithm consider the rank score based on the recognition of the pages by taking into consideration and the importance of both the inlinks and outlinks of the pages. The algorithm provides high value of rank to the more accepted pages and it does not use the method of equally divide the rank of a page among it's outlink pages. But every out-link page is given a rank value based on its recognition. The recognition of a page is decided by observing its number of in links and out links of the particular website. As suggested, the presentation of WPR is to be experienced by using different websites and future work of the algorithm include to calculate the rank score by using more than one level of reference page list and increasing the number of user to categorize the web pages.

The key difference between the WPR from the Page Rank, is that grouping the resultant pages of a query into four categories based on their relevancy to the given condition.

- **Very Relevant Pages (VR):** The very relevant pages contain very important message associated to a given query.
- **Relevant Pages (R):** The Relevant pages do not have important message about given query.
- **Weak Relevant Pages (WR):** The Weak Relevant Pages do not have the relevant information and might have the query with keywords.
- **Irrelevant Pages (IR):** The Irrelevant Pages do not have both appropriate information and query based keywords.

The page rank and weighted page rank algorithm provide ranked pages in the sorting order to users' based on the given query. In the resultant list, the number of relevant pages and their order are very important for all user level. Relevance rule is used to calculate the relevancy value of each page in the list of pages. The relevancy of a page to a given query depends on its category and its position in the page list. The larger relevancy value the better is the result.

d. Distance rank algorithm

The distance between the two web page treated as penalty, called "Distance Rank" to compute the ranks of web pages. It is denoted by the number of average clicks between the two web pages. The idea is to minimize the penalty or distance, and finally a web page with less distance value to be considered as higher rank value to be used.

Like Page Rank algorithm, the rank of each page is denoted as the weighted sum of ranks and all the pages having links to the page. A page has a high rank and it has more back links or pages having links to one page have higher ranks. Specifically a web page is having as many as input links should have low distance and if pages pointing to this page have low distance then this page should have a low distance.



Two definitions has been applied to reach the better outcome of the webpage.

Definition-1. If page x points to page y then the weight of link between x and y is equal to $\log_{10} O(x)$ where $O(x)$ shows x 's out degree (number of forward links).

Definition-2. The distance between two pages x and y is the weight of the shortest path i.e. the path is having the minimum value from x to y . We call this logarithmic distance and denote it with d_{xy} .

e. Weighted page content rank algorithm

The Weighted Page Content Rank WPCR is the combination of structure mining and content mining. Maximum determination of the relevancy of the pages to the given query WPCR algorithms provide important information and relevancy about a given Query [2, 6].

Weighted Page Content Rank Algorithm (WPCR) is a projected page ranking algorithm which is used to give a sorted order to the web pages returned by a search engine in response to a user query and requirements. The WPCR uses the arithmetical value based on which the web pages are given an order. The algorithm design applied for both the web structure mining and web content mining techniques.

The role of Web structure mining logic has been used to calculate the weight of the page and same as the web content mining used to calculate the reputation of the page. i.e. how many pages are pointing to or are referred by this particular page. It can be calculated based on the number of inlinks and outlinks of the page. The key role of WPCR is that weight of web page calculated on the basis of input and outgoing links and also the weight of page is decided on importance. All the pages are sorted on basis of importance of it. The relevancy factor is not considered completely.

f. Webpage ranking using link attributes

In page algorithm concentrate on three factors like qualified position in the web page, tag where the link is contained and the length if the linked text. The same condition to be applied to the page rank concept with link attributes will leads to better result in performance. The page rank with link algorithm uses the page rank formula as base, it follows as $PR(p) = Q/T + (1-q) \sum PR(ri) / L(ri)$, is used to calculate the probability of the page.

But the page rank algorithm with link attributes uses the above modified expression as follows $R(i) = q/T + (1-q) \sum W(j,i)R(j) / \sum W(j,k)$. From the expression $R(i)$ corresponds to the probability to reach the page while searching the website.

g. Eigen rumor algorithm

The Eigen rumor algorithm gives a rank score to every blog in the webpage and it weights the scores of the hub and authority of the blogger based on the calculation of Eigen value. The Eigen algorithm concentrates on hub and authority value in each web page. This algorithm technique is applicable in web content mining and it uses

the methodology like adjacency matrix and it creates an agent to object for linkage, not like web page to web page. The algorithm applied in the search engine in blog community model.

The key objective of the algorithm is applicable in blog ranking concept. The performance in web site search is high and the limitation of the algorithm is not so better like other ranking algorithms.

h. Time rank algorithm

The time rank algorithm basically works on time based visiting model. This algorithm technique is applied in the field of web usage mining. The key function of the algorithm is that the visiting time is additional to the calculated score of the actual page rank value of that web page.

It uses the actual value of the page as given as input for processing and leads the average result as output. Relevancy of using this algorithm is high value when it compared to other ranking algorithms. It is more dynamic in nature with the attributes of duration of used by user and the basic structure of the links used in the webpage. The drawback of the algorithm is that omits the most relevant pages at the time search in web site. i.e. Avoid the recent accessed webpage by the user and need more steps to reach out the same webpage.

i. Tag rank algorithm

The rag rank algorithm consist of two factors are initial probabilistic tag relevance estimation and random walk refinement. Collaborative tagging systems allow users to assign keywords so called tags.

Tags are used for routing, finding resources and unexpected browsing and thus provide an immediate benefit for users. These systems usually include tag recommendation mechanisms reduction the process of finding good tags for a resource, but also consolidating the tag vocabulary across users.

This algorithm consists of an adaptation of user-based collaborative filtering and a graph-based recommender built on top of Folk Rank.

j. Query dependent ranking algorithm

The query-dependent ranking analyzes the relationship between the query results and the tuples in the database. The role of query-dependent ranking by analyzing the user's browsing choices and comparing different queries in terms of their similarity with each other without requiring knowledge of the Web database. The query dependent uses the method of learning to rank based on a distributional similarity measure for gauging the similar data between the queries. Each webpage is accessed through query statement and its nature.

In general the queries describe the user's information need and play an eminent role in the context of ranking for information retrieval and webpage search. The users search intention is based on navigation, informational and transactional queries. In such required, the query categorization has high correlation with users



different expectation on the result achieved through query dependent ranking.

Comparison of various algorithms

Table-2. Comparison state of page rank algorithms.

Algorithms	Page rank	HITS	Weighted page rank
Main Technique	Web Structure mining	Web Structure mining Web content mining	Web Structure mining
Methodology	It score for pages at the time of Indexing	It uses hubs and authority of the relevant pages.	It performs on the basis of Input and output links.
Input parameter	Back links	Content, back and forward links	Back link and forward link
Quality	Medium	Low	Higher than PR
Mining technique used	Web Structure Mining	Web Structure Mining and web content mining.	Web Structure Mining
Search engine	Google	IBM search engine Clever.	Research model
Limitations	Query independent	Efficiency problem	Query independent.

CONCLUSIONS

The term data mining concept focus on retrieving the data from any sources effectively with better outcome. In web mining is nothing but deals with retrieving the data from internet with best output. The web structure mining also deals with many algorithms that lead to fetch the data from any website. In general web structure mining is that retrieves the data from website for online user in effective manner. The unfilled requirements of user in online will be compensated through web structure mining concept and its techniques. The key role of web structure mining plays an eminent role in web access for all level of users in possible way.

Web mining is powerful technique used to extract the information from past behavior of users. The algorithms used in web structure mining to rank the relevant pages. Page rank and weighted page rank algorithms deals with Web structure mining. HITS deal with structure mining and web content mining. The main spotlight of Web structure mining is on link information; Page rank is a better approach for calculating the page value which is a numeric value that represents the importance of a page on the web. The purpose and the important of page ranking based algorithm used for information retrieval and compare those algorithms. An efficient page rank algorithm will meet the challenges efficiently with global technologies in web. Thus the above different algorithm in Page ranking gives a better result in web access for all levels of users in Internet. The scope of page ranking algorithm gives better performance in all possible way.

REFERENCES

[1] Neelam Tyagi and Simple Sharama. 2012. Comparative study of various page ranking

Algorithms in Web Structure Mining. International Journal of Innovative Technology and Exploring Engineering (IJTEE). 1(1).

- [2] Preeti Chopra and Md. Ataulah. 2013. A Survey on improving the Efficiency of different Web Structure Mining Algorithms. International Journal of Engineering and Advanced Technology. 2(3): 1-3.
- [3] Gurpreet Kaur. A Survey - Link Algorithms for Web Mining. International Journal of Computer Science and Communication Networks. 3(2): 105-110.
- [4] Ashish Gupta and Anil Khandekar. 2013. The Study of Web Mining. International Journal of Science, Engineering and technology research. 2(12): 2157-2161.
- [5] Kaur and Rinkle Rani Aggarwal. 2012. Web Mining Tasks and Types: A Survey. IJRIM. 2(2).
- [6] Arun Kumar Singh, Avinav Pathak and Dheeraj Sharma. 2013. A Survey on Enhancing the Efficiency of various web structure mining algorithms. International Journal of Computer Applications Technology and Research. 2(6): 771 -774.
- [7] J. Just. 2013. A Short Survey of Web Mining. WDS13 Proceedings of Contributed Papers Part I. pp. 59-62.
- [8] Rekha Jain and Dr. G.N. Purohit. 2011. Page Ranking Algorithms for Web mining. International Journal of Computer Applications (0975 - 8887). 13(5).
- [9] T. Nithya. 2013. Link Analysis Alogirthm for Web Structure mining. International Journal of Advanced



www.arnjournals.com

- Research in computer and Communication Engineering. 2(8).
- [10] Dashna Navadiya. 2012. Web Content Mining techniques - A Comprehensive Survey. International Journal of Engineering research and technology. 1(10): 1-6.
- [11] Claudia Elena Dinuca and Dumtru Ciobanu. 2011. On two Algorithms used in Web Structure Mining. Annals of University of Craiova - Economic Sciences Series. 3(39): 186-193.
- [12] S. Sathya Bama, M.S. Irfan Ahmed and A. Saravanan. 2013. Improved Page Rank Algorithm for Web structure mining. International Journal of Computer and Technology. 10(9): 1969-1976.
- [13] N.V. Pardakhe and R.R. Keole. 2013. Analysis of various Web Page Ranking Algorithms in Web Structure Mining. International Journal of Advanced Research in Computer and Communication Engineering. 2(12), December.
- [14] Raymond Kosala and Hendrik Blockeel. 2000. Web Mining Research: A Survey. ACM SIGKDD. 2(1): 1-15.
- [15] Miguel Gomes da Costa Junior and Zhiguo Gong. 2005. Web Structure Mining: An Introduction. International Conference on Information Acquisition, June 27- July 3 2005, China.
- [16] Ramesh Prajapati. 2012. A Survey paper on Hyperlink - Inducted Topic Search (HITS) Algorithms for Web Mining. International Journal of Engineering Research and technology. 1(2): 1-8.
- [17] T. Munibalaji and C. Balamurugan. 2012. Analysis of Link Algorithms for Web Mining. International Journal of Engineering and Innovative Technology (IJEIT). 1(1): 81-86.
- [18] Saeko Nomura, Satoshi Oyama, Tetsuo Hayamizu and Toru Ishida. 2004. Analysis and Improvement of HITS for detecting Web communities. Systems and computers in Japan. 35(13): 32-42.
- [19] Lamberti, Sanna. A and Demartini. C. 2008. A relation Based Page rank Algorithm for Semantic Web Search Engine. Knowledge and Data Engineering, IEEE Transaction. 21(1): 123-136.
- [20] T.A. Runkler and J.C. Bezdek. 2003. Web Mining with relational clustering. International Journal of Approximate reasoning. 32(2-3): 217-236.
- [21] Bamshad Mobasher and Robert Cooley, Jaideep Srivastava. 2000. Automatic personalization based on web usage mining. Communications of ACM. 43(8): 142-151.
- [22] Lihui Chen and Wai Lian Chue. 2005. Using Web structure and summarisation techniques for Web content mining. Information Processing and Management. 41: 1225-1242.
- [23] J. Hou and Y. Zhang. 2003. Effectively Finding Relevant Web Pages from Linkage Information. IEEE Transactions on Knowledge and Data Engineering. 15(4).
- [24] N. Duhan, A. K. Sharma and K. K. Bhatia. 2009. Pageranking Algorithms: A Survey. Proceedings. The IEEE International Conference on Advance Computing.
- [25] J. Kleinberg. 1999. Authoritative Sources in a Hyper-Linked Environment. Journal of the ACM. 46(5): 604-632.
- [26] P. Ravi, Singh and Ashutosh Kumar. 2010. Web Structure Mining: Exploring hyperlinks and algorithms for Information retrieval. 7(6): 840-845.
- [27] M. G. da Gomes Jr. and Z. Gong. 2005. Web Structure Mining: An Introduction. Proceedings of the IEEE International Conference on Information Acquisition.
- [28] Haveliwala. T.H. 2003. Topic sensitive Page Rank: A Context sensitive ranking algorithm for Web Search. Knowledge and data engineering, IEEE Transactions. 15(4): 784-796.
- [29] S. Chakrabarti, B. Dom, D. Gibson, J. Kleinberg, R. Kumar, P. Raghavan, S. Rajagopalan and A. Tomkins. 1999. Mining the Link Structure of the World Wide Web. IEEE Computer. 32: 60-67.
- [30] Kun Yu, Xiaobing Chen and Jianhong Chen. 2012. A Multidimensional Pagerank algorithm of literatures. Journal of theoretical and applied information technology. 44(2): 308-315.
- [31] Fabrizio Lamberti, Andrea Sanna and Claudio Demartini. 2009. A Relatin based page rank Algorithm for semantic web search engines. IEEE Transactions on Knowledge and data engineering. 21(1): 123-136.
- [32] Gaurav Agarwal. 2012. A Novel Ranking Algorithm for Ordering Web Search Results. The Second International Conference on Computer Applications.
- [33] Maurice D Mulvenna, Sarabjot S. Anand and Alex G. Buchner. 2000. Personalization on the Net using Web



- Mining: Introduction. Communication of the ACM. 43(8): 122-125.
- [34] Federico Michele Facca and Pier Luca Lanzi. 2005. Mining interesting knowledge from weblogs: a Survey. Data and Knowledge Engineering. 53(3): 225-241.
- [35] Yuefeng Li and Ning Zhong. 2004. Web mining model and its applications for information gathering. Knowledge Based Systems. 17(5-6): 207-217.
- [36] Chen Lihui and Chue Wai Lian. 2005. Using Web structure and summarization techniques for Web content mining. Information Processing and Management. 41(5): 1225-1242.
- [37] Magdalini Eirnaki and Michalis Vazirgiannis. 2003. Web mining for Web personalization. ACM Transactions on Internet Technology. 3(1).
- [38] Qinbao Song and Martin Shepperd. 2006. Mining web browsing patterns for E-commerce. Computer in Industry. 57(7): 622-630.
- [39] Kyung - Joong Kim and Sung-Bae Cho. 2007. Personalized mining of web documents using link structures and fuzzy concept networks. Applied Soft Computing. 7(1): 398-410.
- [40] Ji-Hyun Lee and Wei-Kum Shiu. 2004. An adaptive website system to improve efficiency with web mining techniques. Advanced Engineering Informatics. 18(3): 129-142.