# Web Mining Techniques and its Applications

**P.S Thakur,**
**Associate Professor, Government Degree College,**
**Billawar, Kathua, Jammu and Kashmir, India**

**Abstract— Information on Cyberspace and especially on Web sites snowballing rapidly day by day, Web sites play an important role in this manner where a lot of Web users are always upload, download and brows a lot of stuffings based on their needs. The aim of Web Mining is to provide an algorithm or technique to make data accesses more efficient and convenient. In this paper Web Mining techniques and application are discussed briefly.**

**Keywords: web mining; Web Content Mining (WCM); Web Usage Mining (WUM); Web Structure Mining (WSM); data mining; log file.**

## I.    Introduction

Information mining is Information Revelation in Data set (KDD), it is for the most part characterized as finding cycle of fundamental examples from huge sum information that put away in enormous data sets. The huge development of information sum on Internet (WWW) where tremendous Website pages made regularly made the mining and breaking down of helpful data is a functional test. WWW comprise of billions of interconnected Website pages, which distributed by a huge number of creators on the world. Page is a report that is reasonable to see by WWW through utilizing Internet browsers. WCM is a course of extricating an important data or information from Page contents. Page might incorporate text, pictures, recordings, sounds, or other construction content like tables, and it intended to pass this substance on to the mentioned clients. WSM is a course of found and removing underlying data from Web records, where the mining system acted in two levels rely upon the sort of primary information that utilized in examining process, which are hyperlink level and site page level. WUM characterized as the most common way of applying Information Mining methods to figure out the examples of utilization information to grasp and better serve electronic application necessities. Web Mining are varying from one another in light of the sort of information to be mined and the sort of separated information. WCM zeroed in on Site items, for example, Website page text, pictures and other joined media, Information Mining methods applied on this class incorporates the pre-handling of Website page printed

contents as well as such other substance like pictures relies upon the application necessities. WCM strategies can be used in many Web applications that expect to find Web protests that having normal attributes or examples, for example, gathering of Site pages that discussing comparative points (subjects), or comparative Web pictures that incorporates specific articles like logos, watermarks, or appearances. Second Web Mining classification is WSM, which worry with the most common way of finding semantic highlights from Site pages, for example, the relationship among Pages that have a place either with comparative or different Sites. Thus, this class meant to find gathering of Pages or Sites that pertinent to one another in view of underlying connections geographies, examining of such an information.

## II.     Study Method

The primary objective behind WSM is to extricate already obscure connections among the Pages that have a place with single or many Sites. There are two distinct methodologies in WSM, which are connect geography mining and connection URL mining; the two methodologies utilized an alternate crude information and strategies. Geography mining percept the Internet as a diagram in such a manner the pages addressed as hubs and the edges are the hyperlinks among these pages. URL mining combination with connect geography for the source and targets pages to build more exact connection designs model. WSM can be utilized in later attempts to breaking down informal communities and recognizing clients' networks that having normal data sources, application upon WSM are represents as follow:

**A. Clustering Web**:

Incorporates a few items that exist for the most part with no coordinated design; the items in the WWW are the Pages, which are linkage to different pages through the connections. The principal objective of bunching strategy is to bunch comparative pages in light of the sort of design data utilized which incorporates: Hyperlinks, Archive construction and Connection examination. Accordingly, bunching empowers of associated Pages to lay out relationship of other related pages and permits clients to get to the ideal data through catchphrase affiliation and removed contents.

B. **Classification**:

Is a supervised Data Mining techniques that aim to assign class property from set of predefined set of classes, in Web data there two types of classification which are

**1) Link-Based Classification**: Connect (hyperlink) based grouping is the latest redesigning method in Web mining, the principal objective behind it is to foresee the classification of the Page in light of the connections ascribes (for example Joins among Pages, Anchor text and other HTML labels). Hyperlinks is an underlying piece that interface different area in Site pages to other area either inside the very page or different pages that has a place with single or many Sites. Connect (hyperlink) structure ordered into two sorts, which are between report hyperlinks and intra-record hyperlinks. Between archive hyperlink characterized as the hyperlink that interface different Site pages while intra-record hyperlink interfaces various parts inside a similar Page.

**2) Content-Based Classification:** Content based characterization means to grouping Page in view of the connection (hyperlink) contents (anchor text of the connection).

Subsequently, every Page alloted to class in light of the words that shows up in their connections, this approach is required iterative strategy for relegating the names, because of class of Pages may possibly changes through the connections that dispersed on the Website page itself.

## C. Retrieval Information

That containing in the hyperlinks assumes a significant part for recovering the outcomes in the web search tools, where the anchor text segments (text that show up in hyperlinks) of ancestor Site pages are as of now listed by Internet Worms. There are two sorts could be recovered in view of client question: specialists' pages and center points pages. Each Site page appoint two scores, one is called power score and the other its center points score by utilizing a calculation called HITS (hyperlink-actuated subject inquiry) which comes to take care of the issues of web indexes. Authority's pages are the pages that containing huge data about the question subject, while the pages that focuses to numerous power Website pages is called center pages and its helpful assets in the Internet. Consequently, the scores figured out which page is a decent center page in the event that it is containing pointers to numerous great specialists, and great power page assuming that numerous great center points pages point it and pages are positioned in light of score values in web crawlers.

## D. Web Usage Mining

WUM also called web log mining that endeavor to separate, find and examination fascinating clients' gets to, client exchange, clickstreams designs and other partner information produced from client cooperation with Web assets and put away it in standard text record design called log document that dwell Online server itself. Most web improvement applications utilized Web use information to examination and extricated data about client profile, access designs for the pages' items, activities, and others and it utilized by numerous World biggest organizations such Hurray, MSN, Amazon and others to gather information from Web log admittance to find their clients' entrance designs.

## III.    Data categorized into: Client level data, Proxy level data and Server level data.

**Client Level Log**: This sort of log document lives in client's program window which incorporate data connected with single client activities towards single or many Sites perusing conduct. Treats used to store status data of Site such passwords and some connected data like passwords, and this data put away in client machine and can perceive this client from others due it contains data about client's perusing working framework, this sort of log document in many Web applications isn't thought of as because of client treats might handicap by client for security and protection issues and gathering data from all clients is an essentially troublesome errand.

• **Proxy Level Log:** Proxy Log act as a transitional stage that take HTTP demands from entire clients and passing it to Web server, then, at that point, returns the outcomes by the Internet server and passing it again to the submitted clients. Web intermediary log document can record every one of the solicitations that made by clients' networks that admittance to the Web through Web access Suppliers (ISPs), it's feasible to distinguish have machine name making that solicitation alongside other subordinate data. The disadvantages of Intermediary level are the development Intermediary server is a troublesome errand and it required progressed organizing software engineers, and solicitation capture is restricted.

• **Server Level Log:** Server log file provide most exact and finish utilization information when clients associate with different Sites. Log information is essential utilized in WUM which contain access log information and application log information. Pre-handling of web information is a useful test in many Web applications because of the accompanying reasons: the size of Web information surpasses any ordinary data set, log information put away in Web server log record is in standard text document organization and comes in various configurations, surmising stored pages' references and clickstream designs and its relationship to other related information this large number of reasons make pre-handling is in many cases most tedious and computationally obtuse step.

## IV.    Future Direction

Web Mining become a fundamental field to many web improvement applications while a great deal of associations content dwell on the web are expanded immensely. Web mining is an original expansion of Information Mining methods that applied on the information types that live Online which exist in various structures and designs. Along these lines, there is relentless need to cover and gives thorough answers for some basic issues, for example, controlling, observing, discernment, and information portrayal of Web information created either by clients' networks or dispatching of data set framework to web administrations. Our future work zeroed in on blend more than Web mining strategy to give a superior discernment to such complex Web information and breaking down the strategies that follow to find the huge examples from the Internet, for example, examining Web utilization information and late methods utilized for finding use designs from it, and dissecting the critical parts in Pages' substance to be utilized for critical application, for example, suggestion framework or Site pages' grouping/bunching.

## V. Conclusions

In this paper, we survey various Web mining techniques that used by a lot of Web application recently. We had also reviewed a comparison among Web mining categories based on significant approaches used currently by most of research works. Since, Web environment is a huge area and there are a lot of work to do in future, we hope this paper could be providing a good starting point to knowing current Data Mining techniques that applied upon different Web data and also help to identifying opportunities for forthcoming research works by understanding the nature of the data that reside on different Web resources

## VI. References

[1] Joy Shalom Sona, Prof. Asha Ambhaikar" A Reconciling Website System to Enhance Efficiency withWeb Mining Techniques" International Journal of Scientific & Engineering Research Volume 3, Issue 2, February-2012 1 ISSN 2229-5518

[2] Aparna Ranade, Abhijit R. Joshi, Ph. D," Techniques for Understanding User Usage Behavior on theInternet" International Journal of Computer Applications (0975 – 8887) Volume 92 – No.7, April2014

[3] Karan Bhalla & Deepak Prasad," Data Preparation and Pattern Discovery for Web Usage Mining"

[4] Amit Pratap Singh1, Dr. R. C. Jain 2," A Survey on Different Phases of Web Usage Mining forAnomaly User Behavior Investigation" International Journal of Emerging Trends & Technology inComputer Science (IJETTCS)Volume 3, Issue 3, May – June 2014 ISSN 2278-6856

[5] R. Lokeshkumar1, R. Sindhuja2, Dr. P. Sengottuvelan, "A Survey on Pre-processing of Web Log Filein Web Usage Mining to Improve the Quality of Data" International Journal of Emerging Technologyand Advanced Engineering, ISSN 2250-2459, ISO 9001:2008 Certified Journal, Volume 4, Issue 8, August 2014

[6] Mitali Srivastava, Rakhi Garg, P. K. Mishra," Preprocessing Techniques in Web Usage Mining: ASurvey" International Journal of Computer Applications (0975 – 8887) Volume 97– No.18, July 2014

[7] http://www.slideshare.net/akhanna3/discovering-knowledge-using-web-structure-mining-27488978

[8] Ashish Kumar Garg, Mohammad Amir, Jarrar Ahmed, Man Singh, Sham Bansa," Implementation ofa Search Engine" International Journal of Science and Research (IJSR) ISSN (Online): 2319-7064

[9] C. Gomathi, M. Moorthi," Web Access Pattern Algorithms in Education Domain" Computer andinformation science journal vol. 1, No.4, November 2008

[10] Md. Zahid Hasan, Khawja Jakaria Ahmad Chisty and Nur-E-Zaman Ayshik, "Research Challengesin Web Data Mining", International Journal of Computer Science and Telecommunications Volume3, Issue 7, July 2012