



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 7 Issue: IX Month of publication: September 2019

DOI: <http://doi.org/10.22214/ijraset.2019.9130>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

A Survey of Web Mining and Various Web Mining Techniques

Ashish Panchal¹, Kaushlendra Verma², Ayush Maru³

^{1, 2, 3}Department of Computer Applications, National Institute of Technology, Haryana

Abstract: Mining is essence of valuable information from the huge set of raw information. Mining techniques in data mining is known as web mining. The rapidly increasing number of web contents including image, multimedia, and digital data. The knowledge gained from the web can be utilised for increase the performance for searching of data. In the internet there is so many duplicate data in present. Thus how we can utilise the useful data from the whole collection of data is a tough task. Web mining shows the past work on using different web algorithms. Today The huge amount of data is present on web. The web crawler plays the essential role in updating the current data. The search engine are depends on the ranking technology instead of the other vector based approaches. This paper will focuses on different web mining techniques and the various algorithm used in it.

Keywords: Web mining, crawlers, structured mining, HITS, page rank, usage mining.

I. INTRODUCTION

It is a process that searches and filter the data over web for the purpose of sorting and excluding that repeated information or data. Web mining targets the content of structure of world wide web document which has large dispense enterprising data. Web mining included the subtasks –

- 1) Select from where to extract the data.
- 2) Choose the specific portion of data which is useful for web pages.
 - a) *Generalisation:* Automatically find the same pattern on one or more web sites.
 - b) *Analysis:* Assure that the information is true or false.

Web pages increases rapidly day by day and it is up to 3 trillion[12] and development of web pages overlapping some information exist and misleading data in web. Ten year ago there was lack of data for home users and this is difficult to identify information to make collection or analysis of web content which help to solve the problem of uncontrolled data in lesser way. Like some system content information in the internet for specific group of users as well as some system could search illegal data or information in the internet to take some legal action. Now days data in web pages is presented in the non structured or semi-structured form. The formation of the web pages is in non-convenient for data analysis system.

The main barriers is to identify or understanding the contain which is on web pages. The important of web mining is presentation and definition of possible web mining categories not cleared, The service requires a different structured framework which has ability to provide a handing substitute for user[1].

Web mining play an important role in achieving the useful information role which we want. It refers to discover and analysis the useful content over the world wide web. It is basically obtaining knowledge from number of web page in websites[2]. The Area of research increasing day by day because of the interest various research communities. The splendid growth of knowledge resources accessible on internet, Now a days interested in E-commerce. The situation that is observed to exist partly create distraction. Although the constitutes web mining is the technique for fetching the information either online or offline from the text content which is present on the web like newsletter, newsgroup the text content html document achieve by deleting html tags and web resources a selected by manually[11]. The information selection is type of conversion process of initial data. We suggest decomposition web mining in to the sub task.

A. Resources Identify

Task of fetching intended web document.

B. Gaining Knowledge of Collection and Pre-Preparation

In this collection and pre-preparation or automatically selected individual data for retrievals resources.

C. Generalization

It is automatically find unexpectedly or during search general pattern at the single websites as well as several sites.

D. Survey

A validation or definition for the mine patterns.

Web mining have a too much attachment with software relationship or intelligent agent and many of this agent are performed data mining task to attain their goal.

Different classification of software agents—

- 1) UIA (User interface agents)
- 2) DA (Distributive agents)
- 3) MA (Mobile agents)

User interface agent attempt to increase the capacity of existing user interaction with the system by modify conduct.

User interface agent that can also be partitioned into the web mining agents categories filtering agents and individual assistant agents.

Distributive agents technology is study with problem solving by multiple of agents and applicable agents in this categories are distributed agents for knowledge detection.[4]

There are two routinely used methods for developing intelligent agent that would assist user detect as well as fetch useful information from web.

It presume that if some users rate any item or product then the other user with the same interest also give rating this item better also so this method mainly uses the data or user rating viewed in this light we could classified the content based methods web content mining and the collaborative approaches as web uses mining. However collaboration approaches might also combined with web contain.

II. WEB MINING

Raw information are mainly stored in a Data Warehouse, involving the focused on the establishment is to find certain data for analyse behaviour of seamanship. It is three types:

- 1) Information consisting in web page like text, video.
- 2) User ip address, login time, how much time he spend of web.
- 3) Set of details of user.

A. Five Classification Methods

Classified by decision tree, neural networks, svm, classification based on association, bayes.

B. Four Clustering Type

Partition method, density type method, grid-type method, model-type method.

C. Types of Union Rule

Multi type associated law, multi dimensional associated rules, quantitative associated rules.

Comparison in Data mining and Web mining—

Data Mining	Web Mining
Data mining involves using methods to search basic structure and relationship in huge amount of data	Web mining involves analysis of web server logs of a website
Most of the data mining apps search patterns in a structure data such as data based	Web mining used for search pattern in a semi-structured data example world wide web.
It can handle large amount of data	It can handle big data compare then traditional data mining.
When applied data mining of corporate knowledge, the information is secure and generally ask right to read.	Web mining, the information is public and really ask arrival right
Data mining fetch data by a database, which gives some levels of obvious framework.	web mining is processing of non framework or semi-framework Data from world wide web. Even the basic information of web pages arrive from a database, generally is hidden by html.

III. DIFFERENT WEB MINING TYPES

It is three types—

- 1) WUM (where WUM is Web usage mining)
- 2) WCM (where WCM is Web contents mining)
- 3) WSM (where WSM is Web Structured mining)

A. WUM (where WUM is web Usage is Mining)

In usage mining it follow the procedure of data mining methods for find knowledgeable data from user's activity. This files contains the user's Activities when the user search on the website.[9]

- 1) Determine the market importance for each customers.
- 2) Build effective market strategy.
- a) *Web Servers Data*: Customers data are automatically fetched in the web servers, which includes user IP address, page information, agent details and time of acquiring the website.
- b) *Apps Server Data*: E-commerce applications use online transaction application servers. It stores different kind of commercial functions and log files related to those users.

B. Web Content Mining (WCM)

It is slightly nonidentical as compare to the data mining technique we use in this type of mining. It is similar to text mining for the reason that web contents are also texts. It deals with data which is unstructured and semi structured data. Content mining applies methods to web documents. Web content mining used in multimedia data mining.

Tools for web content mining are—

- 1) Crawlera
- 2) Webql
- 3) Xml miner
- 4) Bixo

C. WSM (where WSM is Web Structure Mining)

Web structure mining known as linkage construction analysis, it is the area of web mining. Different ways on huge hyperspace linkage storeroom to prompt systematic way regarding websites as well as webpages by analysing the link association that can be used for increasing the page hits and to improve page ranking. The data come from web structure mining includes of textual webpages assembled by slowpoke from all above the connection of multiple server. It contains the four basics tide.

- 1) *Data Gathering*: Collecting the data for analysis from various interconnected links.
- 2) *Pre-Processing*: It contains the four procedure data cleans, data in corporation, data evaluation as well as variation, data depletion also includes the task link validation, link originality.[3]
- 3) *Knowledge Searching*: By trying many data mining methods used in process data which is statistical explication, clustering, format matching study.

Some of the web structure mining algo.:-

- a) Page Rank
- b) Hyper text induced topic search
- c) Weight page rank
- d) Time Rank.
- e) WPCR (where WPCR is Weighted page contents rank.)
- f) Hyperspace page ranking using Link attributes.
- g) Eigen Rumour.
- h) Distance rank.
- i) Query Dependent Ranking

Web structured mining mainly uses page rank and Hyperlink-Induced Topic Search for better and user relevant.

Tools used in web structured mining

- 4) *Majestic Tools*: It is a huge effective Market Analytic tool it provides service for Search websites to Optimized strategy, website developers and media analysts. from the help of this tool, you can get reliable data so that you can study the performance of your websites You can become completely clear about your site's ranking in terms of backlinks. And also compare your website with your competitors. HITS, page Rank algorithm.

IV. RELATED WORK

The Purpose Approached is to build an Information system that analysis the current website performance and avoid the data redundancy Data that can be achieved high performance of our websites. Many researchers find out the way of representing the web mining

V. ALGORITHMS USED IN WEB MINING.

A. Page Rank Algorithm

PageRank is an algorithm which is being used by Google for Searching for ranking web pages into search engine. It is the first algorithm that was used by the google. Its name based on Larry Page, [5] founders of Google. PageRank is a method which calculates significance for web pages. PageRank accepts that page has better rank if total of the rank of its out links is more.[10] It is known the base for all present time Search Engines. The key approval for crucial webpage for accept extra nodes from different web pages.

Google stated that PageRank works with the help of calculating the number and standard for nodes to web page to setup a dry idea for crucial web pages. Ranks pages depend on the number of outlinks indicating to them. The algorithm allocates pages a Total PageRank based on the PageRanks of the outlinks specify to the page. The links to a page can be arranged into the following types: Inbound links which is links in between given site from external source pages. Outbound links which are links from the given page to pages in the same site or other sites and The links which has no outgoing link is known as dangling link.

The Page rank of the web sites is evaluated as a total of the Page ranks for every pages incoming nodes and split up by the number for outgoing node for every web pages.

Page rank of $P = (1 - \text{damping factor}) + \text{damping factor} \sum_{i=1}^n (\text{page rank}(K_i) / O(K_i))$

Where Page rank (P) is the PageRank of page P

Page rank (K_i) is the PageRank of pages K_i which link to page P

$O(K_i)$ is the number of outbound links on page K_i

damping factor which can be set between 0 and 1. It depends on the number of clicks, usually set to 0.85

n is the number of inlinks of page P.

Mentioned is example for PR algorithm. Suppose Web Graph shown in Fig. 1

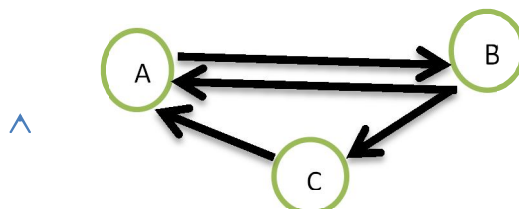


Figure 1: Example of web graph with links and outlinks

Page, A being referenced by pages B and C. C, B has 1,2 outlinks. Page rank value of the page A is given as:

$$PR(A) = 1 - d + d(PR(C)/1 + PR(B)/2)$$

The Algorithm will not rank the entire website, but it is negotiate for each page separately. and, $PR(A)$ is recursively identify from Page rank of which pages connect to page A.

- 1) *Iterative Approach of Page Rank*: Every pages is a allocated a initial page rank value for one. This rank values are repetition replaced page rank equations for found the final result. In more repetition would be follows to more desirable page rank .
- 2) *The Page Rank Algorithm Iteratively Stated Are*
 - a) Originally suppose one is the Page rank for all websites.
 - b) Evaluate page ranks for each pages by given method.
 - c) Do the evaluation until value of back to back Iterations Matches.

3) Advantages

- a) Since it pre evaluate the rank counts it appear in small time and so it is fast.
- b) It is extra achievable as calculate rank will be score at the time of indexing not at the time of querying.
- c) It gives crucial pages rank is evaluated on the basis of admired for page.

4) Disadvantages

- a) It give importance previous pages more, since a new page, even if it is better, Not having much links untill it not belong of an present websites.
- b) Contingency of consequence pages user query is not so much as this doesn't approve the data of web page.
- c) the problems exists in the type of Dangling links which arise pages which includes a link such that the hypertext objective to pages without outbound nodes.
- d) It leads to rank sinks trouble rise when n/w pages receive in endless link cycles.
- e) One more issue in Page rank is a set of pages in which if no links are there of with in the set of outside the another set.
- f) Whether there is circle denoting in your web pages, then it decrease home page's Page rank.
- g) page without outlinks is also there.

B. HITS Algorithm

Kleinberg developed algorithm based on WSM defined as HITS assume of particular query entered from user, It is collection of authority pages that are relevant and admitted focusing on the query and a releated pages consisting links to different authorities . a better hub page of subject points to different authorities pages at which content, and a better a page is pointed by various better hub pages on the simila subject. Authorities and hubs are shown in Fig. 2.

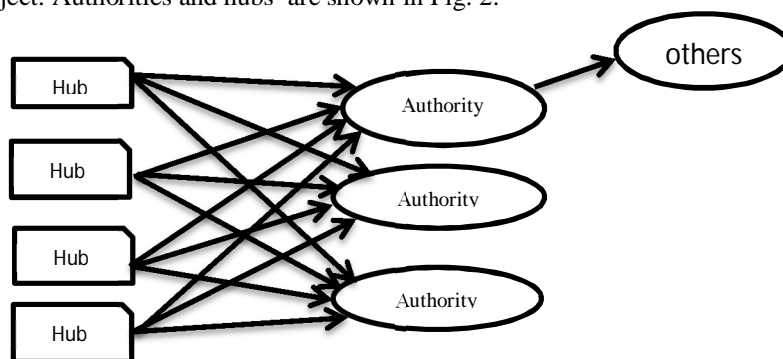


Fig. 2: Authority has backlink from many Hub

Kleinberg declare which page perhaps be a better hub and a better authority at the similar time. This spherical connection connects with the definition of an iterative algorithm known as HITS. The HITS algorithm behave towards world wide web as a directed graph $G(V,E)$,

where V = group for Vertices denoting pages

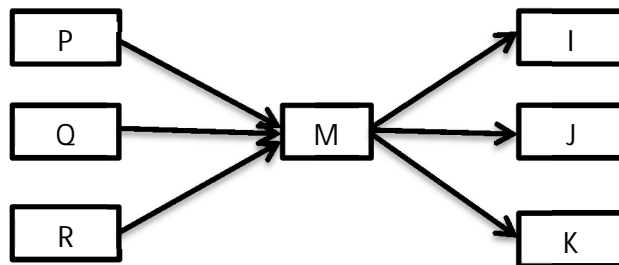
where E = group of edges which paired with links.

HITS algorithm steps contains one is sampling step and another is iterative step.. In the Sampling step, a group of releated pages for the available query are composed together i.e. a sub-graph S of G is fetch which is high in affecting pages. This algorithm starts with a root set R , a set of S is obtained, keeping in mind that S is comparatively small, rich in relevant pages regarding of the query and includes almost all of the better authorities. The second step, Iterative step, Identify hubs and authorities using the results of the sampling step using:[6]

$$H_y = \sum_{r \in J(q)} A_r$$

$$A_y = \sum_{r \in K(q)} H_r$$

Where hub weight is H_y , $A_y = A_w$, $J(q)$ and $K(q)$ notice as the collection of reference and referrer pages of page q . The page's A_w is proportional to the addition of the hub weights of pages that it connected to it; and also, a page's hub weight is proportional to the sum of the affected weights of pages that it links to. Fig. 3 shows an example of the evaluation of authority and hub counts.



$$X_M = Y_I + Y_J + Y_K \quad X_M = Y_P + Y_Q + Y_R$$

Fig. 3: hubs and Authorities counts.

1) Advantages of HITS

- The ranking may also be include with another data retrieval basis rankings.
- HITS is very quick to respond instructors queries (when compare with page rank)[7].
- Useful pages are found on when we observe authority and hubs data.
- HITS is a algorithm for evaluating authority and hubs in order to rank the retrieved data.
- Observations points to that HITS evaluate authority nodes and hub vertex in accurate way.

2) Disadvantages of HITS

- It is based on queries so evaluation time is too much.
- A condition can arise at the time a page that includes links to a huge quantity of other topics might be collect a better hub rank but it isn't belongs to the users query. So this page is not the better related source for any data, then also it has a very high hub rank if it points to better ranked authorities.
- HITS focus on mutual reinforcement in the middle of authority and hub web pages. A better authority is a page that points to lot of better hubs and abetter hub is a page that is pointed to by lot of better authorities.
- Topic drift can also arise when there are pages not related to the root set and pages are tightly connected with root. if the root set has non-identical pages, so it also affect on to the pages in the other set. and, the web graph created by the pages in the base set, will not have the most related nodes and it result the algorithm failed to find the top ranked hubs and authorities for a available query.

Year	Researchers	Algorithms/methods	Inputs	Results
2018	K. Sellamy, H. Jamil, Y. Lakhrissi	various web mining technique	-	proposed approached to renovate execution of employment for young graduate in morcco
2017	Kuber Mehan, Jitendra Kumar, Sanjay Kumar	investigate various web structure algorithms and a case study on page rank and weightage page rank	crawler4j	gives information about various web structures algorithms
2014	Rashmi Sharma, Kamayit Kaur	page rank weightage, page rank, hits	data on some url	weightage page rank is better then hits and page rank
2014	Anshul Bhargav, Munish Bhargav	web usage mining	web logs	increase the efficiency of web site.
2017	K Jayamalini, Dr. M Ponnaivoiko	different technology of web mining	-	gives knowledge about application of web mining
2011	Sanjay K Malik, Sam Rizvi	web usage mininag	server log	it analysis result of user interaction with the serves.
2012	Loraine Charlet Annie, Ashok Kumar	k-means	web page	its results that k apriori algorithms is better the apriori algo.
2013	A. Raiyani and prof. S.S Pandya	distict user identification	user logs	purpose approach fraud detection

VI. CONCLUSION AND FUTURE WORK

Till now we read out the research topic of web mining that focuses on content, usage and structure of web. Web structured mining deals with mainly hits and page rank. Web mining which is one of the Mining technique that extracts the information from web documents automatically. Page Rank algorithm is used in WSM to rank the related pages. In general web mining retrieve the data from websites for users in efficient manner but we catch some problems in hits and page rank algorithm that is our purposed work for future to solve the problem.

The future scope of the topic is to purpose a tool which check and act as the accuracy efficiency checker.

REFERENCES

- [1] K. Selamy, Y. Fakhri , S. Boualknadeeli , A. Momen, K.Haed(2016).Web mining methods and applications: Literature report and a proposed approach to improve performance of employment (2016)
- [2] Mrs. S. R. Kalailvi1, S. Maheshi2, V. Shobana. Web Mining: Data mining concepts applications and research directions.International Journal of Advanced Research in Computer and Communication Engineering (Vol. 4, Issue 11, November 2015)
- [3] Suvaraan Sharma, Department of Mathematics, Manit Bhopal. Data Pre processing Algorithm for Web Structure Mining. 5thInternational Conference on Eco-Friendl yCommunication Systems and Computing (ICECCS-2017)
- [4] Anshul Bhargav School of Computer Engineering,Pattern Discovery and Users Classification Through Web Usage Mining.(ICCICCT-2014)
- [5] S. Brin, and L. Page, The Anatomy of a Large Scale Hypertextual Web Search Engine, Computer Network and ISDN Systems (Vol. 30, Issue 1-7, pp. 107-117, 1998.)
- [6] K.R. Srinath Page Ranking Algorithms – A Comparison , Associate Professor, Department of Computer Science, Pragati Mahavidyalaya Degree and PG college, Telangana, e-ISSN: 2395-0056 Volume: 04 Issue: 12 |Dec-2017 www.irjet.net p-ISSN: 2395-0072
- [7] J. Kleinberg, “Authoritative Sources in a Hyper-Linked Environment”, Journal of the ACM 46(5), pp. 604-632,199
- [8] Nidhi grover mca scholer institute of IT & Managementcomperative analysis of page rank and HITS algorithm. ISSN 2278-0181 Volume:01 issue 8 oct-2002.
- [9] Ms. Gaikwad Surekha Naganath , Ms. Mali Supriya Pralhad, Web Mining-Types,Applications,Challenges and Tools, International Journal of Advanced Research in Computer Engineering & Technology (IJARCET) Volume 4 Issue 5, May 2015
- [10] Kuber Mohan M.Tech CS Department of Cs, A Survey on Web Structure Mining.IJARC Volume 8, No. 3, March – April 2017, ISSN No. 0976-5697
- [11] S. Jeyalatha CS dept. Design and Implementation of a Web Structure Mining Algorithm using bfs Strategy for Academic Search Application.ICITST 2002.
- [12] Web Mining Taxonomy Kiril Griahev,Department of IT, 978-1-5386-6737-8/2018 IEEE



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)