

# A Proficient Web Crawler Utilizing Page Rank Algorithm

Niraj Chandrabhan Prajapati

ASM Institute of Management & Computer Studies, Thane

**Abstract-** The Internet network of networks. The World Wide Web, via comparison, is a global database of records and other resources, related with URI and hyperlinks. The World Wide Web, generally consists of standard database locators, which are connected to hyperlinks and available via the Internet, describe the documents and other web services. However, the massive size of internet has unintentionally become an obstacle to the retrieval of information. For instance, the user has to turn through numerous pages in order to seek the required data. This obstacle has mainly led to exploration of Web Crawler. The core of search engines is Web crawlers. A Web crawler is a program or software that traverses across the web and systematically, aims to automate the retrieval of web documents. Web pages can also be checked with web crawlers to identify security bugs and evaluate if there is a leak. Within this article, I have also reviewed the functionalities of Web Crawler in the world of web, their classification and also its drawback experienced during its use in various search engines.

**Keywords-** web crawler, PageRank, HITS, and Link Based Search, solution for HITS algorithm.

## I. INTRODUCTION

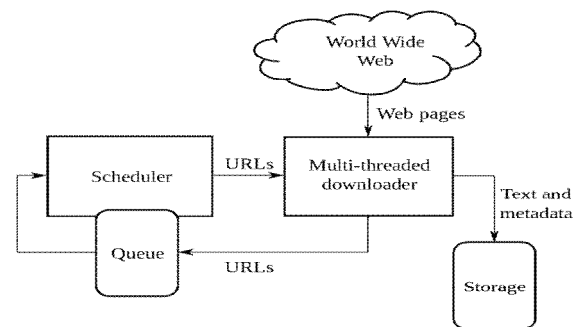
The Web contains an oversized volume of information on different topics. In contrast to traditional collections like libraries, the web world has no centrally organized content structure. This data is often downloaded using web crawler. So, Web crawler is software for downloading pages from the net automatically. It's also called web spider or web robot. Web crawling may be a vital method for collecting data on, and maintaining with, the rapidly expanding Internet.

**The general procedure that a crawler takes is as per the following: -**

- It checks for the subsequent page to download – the system keeps track of pages to be downloaded in a queue.
- Checks to see if the page is allowed to be downloaded - checking a robot's exclusion file and also reading the header of the page to work out if any exclusion

instructions were provided do that. Some people don't desire their pages archived by search engines.

- Download the entire page.
- Extract all links from the page and add those to the above-mentioned queue for later download.
- Extract all words & save them to a database related to this page, and save the order of the words in order that people can look for phrases, not just keywords.
- Alternatively filter for things like grown-up content, language type for the page, and so forth.
- Save the page summary and update the page's last processed date so the system will know when to recheck the page at a later stage.



## II. LITERATURE REVIEW

Title / Author	Descriptions	Drawbacks
Dr. Naresh Kumar, Shivank Awasthi and Devvrat Tyagi, "Web Crawler Challenges & Their Solutions", International Journal Scientific and Engineering Research, Volume 7, Issue 12, ISSN 2229-551	The HITS calculation in is out of date these days and subsequently effectiveness of results can be improved by utilizing a productive calculation.	Dr. Naresh has not determined any methodology for managing pages other than that of HTML type, which may prompt loss of significant data. Also, there is an extraordinary chance that the bookmarks of client may have a place with completely various areas prompting a condition of uncertainty for crawler.

**III. TRADITIONAL PAGERANK ALGORITHM**

PageRank works by checking the number and nature of connections to a page to decide good guess of how significant the website page is. At the point, when a web crawler experiences every site, it follows all the connections in the website, and checks for what number of connections are associated with each web page. And afterward it allots rate to every site page which speaks to the significance of the website page utilizing page rank calculation.

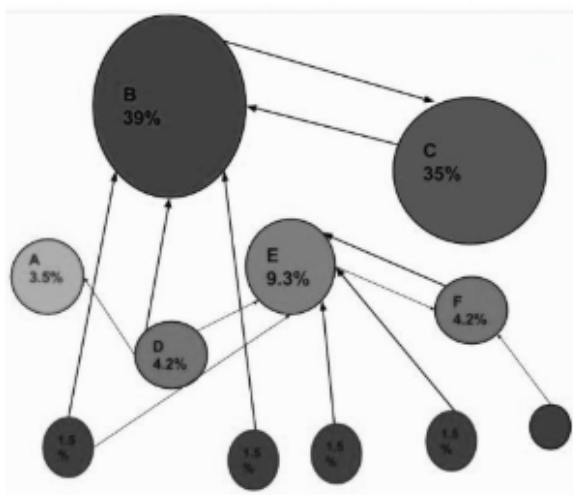
**Formula:**

$$PRGP(A) = (1 - d) + d ( PRGP ( T1 ) / C ( T1 ) + . . . . + PRGP( Tn ) / C ( Tn ) )$$

PRGP(A) : Page Rank of a given Page

d : Dump Factor

Ti : link's



**IV. DISADVANTAGES OF TRADITIONAL PAGERANK ALGORITHM**

- Rank Sinks: This issue happens when during a system pages get in limitless connection cycles.
- It could be a static calculation, in light of its total plan, well-known pages will normally remain mainstream for the foremost part. Fame of a site doesn't make sure the ideal data to the searcher.
- Dead Ends are basically pages with no friendly connections. PageRank doesn't cope with pages with no out edges well indeed, on the grounds that they refuse the PageRank generally speaking.
- Circular References: If you've got circle references in your website, then it will reduce your front page's PageRank.

**V. PROPOSED PAGE RANK ALGORITHM**

The proposed standardization strategy for Page Rank calculation relies upon mean estimation of page rank of all website pages with execution central focuses over the customary PageRank calculations.

No. Of iterations For Traditional PageRank algorithm	100
No. Of iteration for proposed PageRank algorithm	20

**Steps:**

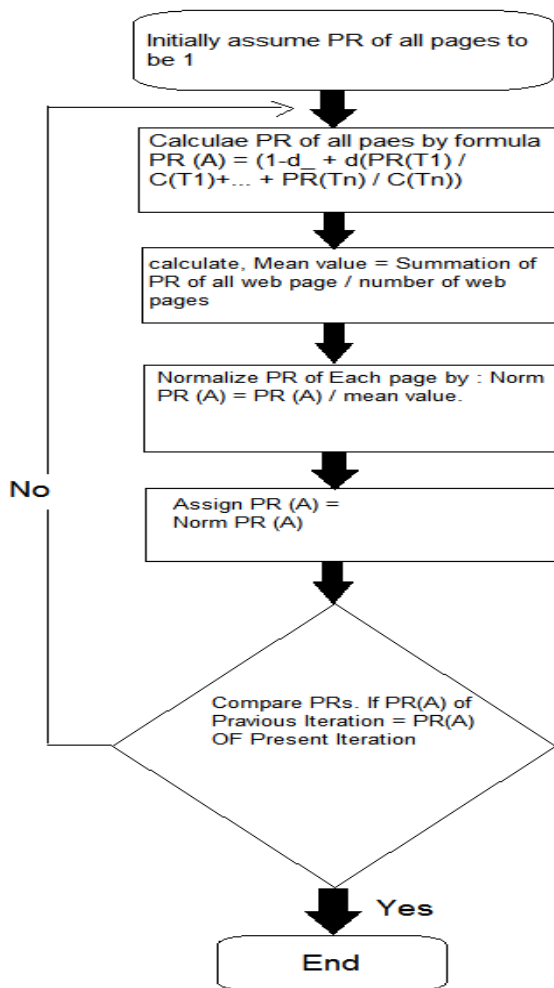
- 1) Initially accept PageRank of all sites pages to be any worth, let it be 1.
- 2) Calculate page rank of all pages by taking after equation  $PageRank(A) = .15 + .85 (PageRank(T1)/C(T1) + PageRank(T2)/C(T2) + \dots + PageRank(Tn)/C(Tn))$  Where (Default an incentive to damping factor which can be set somewhere in the range of 0 and 1.) T1 through Tn are pages give approaching connects to Page A PageRank(T1) is the PageRank of T1 PageRank (Tn) is the Page Rank of Tn C(Tn) is all out number of active connections on Tn
- 3) Calculate mean estimation of all pages positions by following equation: - Summation of PageRanks of all sites pages/number of website pages

4)Then standardize page rank of each page Norm PageRank (A) = PageRank (A)/mean worth Where standard PageRank (An) is Normalized PageRank of page An and PR (An) is page rank of page A

5)Assign PageRank(A)= Norm PageRank (A) 6) Repeat stage 2 to stage 4 until page rank estimations of two back to back cycles are same.

6) Now after applying this calculation the distinction will be between customary PageRank calculation and proposed page rank calculation.

**Flow chart for proposed efficient page rank algorithm:**



**VI. CONCLUSION**

Web crawlers are a crucial aspect of all the search engines. They're the essential component of all web services in order that they got to provide high performance. Whenever the user wants to retrieve any sort of data from entire world the Internet is the most effortless source accessible in present days. World Wide Web consists of a graphical visualization having web pages linked to each other through hyperlinks.

Web Crawlers generally use graphical visualization to navigate from one page to another. However, it fails somewhere do its task effectively. A number of crawling algorithms are employed by the search engines in order to present the required data to users. Therefore, a decent crawling algorithm should be implemented for better results and high performance.

**REFERENCES**

- [1] Christopher Olston & Marc Najork (2010), “Web Crawling”, now the essence of knowledge, Vol. 4, No. 3 (2010) 175–246.
- [2] DhirajKhurana, Satish Kumar (2012), “Web Crawler: A Review”, International Journal of Computer Science & Management Studies, Vol. 12, Issue 01, ISSN: 2231 – 5268.
- [3] Mohit Malhotra (2013), “Web Crawler And It’s Concepts”.
- [4] Subhendukumarpani, Deepak Mohapatra, BikramKeshariRatha (2010), “Integration of Web mining and web crawler: Relevance and State of Art”, International Journal on Computer Science and Engineering Vol. 02, No. 03, 772-776.
- [5] Nemeslaki, András; Pocsarovszky, Károly (2011), “Web crawler research methodology”, 22nd European Regional Conference of the International Telecommunications Society.
- [6] Salvatore A. Catanese, Pasquale De Meo, Emilio Ferrara, Giacomo Fiumara, Alessandro Provetti (2011), “Crawling Facebook for Social Network Analysis Purposes”, WIMS’11, May 25-27, 2011 Sogndal, Norway, ACM 978-1-4503-0148-0/11/05.
- [7] Ari Pirkola (2007), “Focused Crawling: A Means to Acquire Biological Data from the Web”, University of Tampere Finland, ACM 978-1-59593- 649-3/07/09.
- [8] Priyanka-Saxena (2012), “Mercator as a web crawler”, International Journal of Computer Science Issues, Vol. 9, Issue 1, No 1, ISSN: 1694- 0814.
- [9] Vladislav Shkapenyuk, Torsten Suel, “Design and Implementation of a High-Performance Distributed Web Crawler”, NSF CAREER Award, CCR-0093400.
- [10]Raja Iswary, Keshab Nath (2013), “Web Crawler”, International Journal of Advanced Research in Computer and Communication Engineering, Vol. 2, Issue 10, ISSN: 2278-1021.
- [11]Allan Heydon and Marc Najork, “Mercator: A Scalable, Extensible Web Crawler”, Compaq Systems Research Center
- [12]S. Chakrabarti. *Mining the Web*. Morgan Kaufmann, 2003.

- [13] P. Srinivasan, J. Mitchell, O. Bodenreider, G. Pant, and F. Menczer. Web crawling agents for retrieving biomedical information. In *ETTAB: Agents in Bioinformatics*, Bologna, Italy, 2002
- [14] S. Chakrabarti, M. van den Berg, and B. Dom. Focused Crawling: A new approach to topic-specific Web resource discovery. *Computer Networks*, 31(11–16):1623–1640, 1999.
- [15] K. M. and Michelsen, R. (2002). Search Engines and Web Dynamics. *Computer Networks*, vol. 39, pp. 289–302, June 2002.