# Survey Paper on Random surfer model in Page Ranking Algorithm

Suthar Henilkumar D.
M-Tech.IT
U. V. Patel college of
Engineering,
Ganpat University, Mehsana

Rajendra J. Patel
Department of CE
U. V. Patel college of
Engineering,
Ganpat University, Mehsana

Nikhil Kumar Singh
Department of IT
U. V. Patel college of
Engineering,
Ganpat University, Mehsana

## ABSTRACT

The World Wide Web have a large amount of data, the user are find the some relevant information from WWW. If user want to search any information than numbers of URL or number of pages are opened on the web. The user want to be the show the important information are show on the top of the list than this page have a high ranking compare to other and this page is a most relevant page on the web site. The page rank are depends upon the user searching or clicking on the web pages, but some problem are created on given the page ranking about the page, in this paper we discussed the some problem about the page ranking on the world wide web.

## General Terms

Page Ranking Algorithm

## Keywords

Web log, Structure mining, Page rank

## 1. INTRODUCTION

Web mining is an application of data mining technique which uses data mining to extract some useful information from web document. The web is a lot of source of information and persists to increase in size and difficulty. Retrieving the necessary web page on the web, efficiently and effectively, is becoming a some challenge aspect nowadays [1].

World Wide Web is flurrying by numerous Web sites around the world, a global information system. Web servers can potentially host number of pages which make the number of web pages extremely difficult to track. Web page networks like the thousands of interconnected intertwined with the cells organized in a complex structure. Each Web site also contains a millions of Web pages the page contain hypertext mark-up and hyperlinks between Web pages. Web pages within a sitedo not exist in isolation, usually related to hyperlinks between documents, text which describes the hyperlink links between some Documents [1].

The World Wide Web is a collection of number of pages on the web and content and some resource.lot of source and content are available. In the World WideWeb the every Web pages are connected through the link.

The information can be discover by web mining through the web activity, from server log and web browser activity tracking, and web graph from link between pages people and other data. The web content is a data found on the web pagesand inside the document and inside the web pages someimportant information is stored. The Web mining is classifiedinto the three types: Web Structure Mining, Web Usage Mining and Web Content Mining.
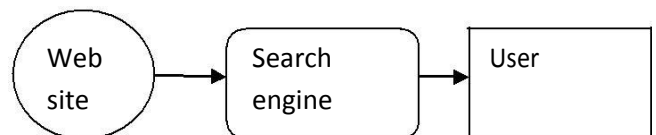


**Fig 1.1 Process of Web mining**

## 2. WEB MINING

Web mining is a one type of the data mining technique and the many type of the data are stored in the WWW. Now World Wide Web is a very popular on the world. The text, image and video, etc. Many types of the data are available on the web pages. The page rank technique is a one type of the web mining technique for given the Rank of the Web site pages or Web pages.

## 3. LITERATUTE REVIEW

3.1 The web data cleaning process is a most important as a research and 75% time are is spending in this process [2]. The many classification technique are used in to the data cleaning process. The major issue in web mining is a over lapping, scaling and high dimensional data are require in web mining.

The web mining is mostly depended upon the link structure of the web [2].

3.2 The page rank algorithm is used in Google search engine and page rank algorithm are generated only numeric value. In the page rank algorithm calculated based on the in link and out link of the pages [3]. Some important pages have a more of the back links are available, back link means in link of the pages. If any page have a many more back link are available than this page-page rank are high on to the search engine [3]. The high page rank are improved the efficiency of page on the search engine.

3.3 Increasing the efficiency of web pages are based on the page rank algorithm, now very most and successful search engine is a Google and we improve the efficiency of Google search engine pages than using the weighted page rank algorithm and increasing the efficiency of search engine [4]. The weighted page rank are calculated based on weight of the pages and this weight are depends upon some back link of the web pages. And this are calculating between the some in link and out link of the web pages on to the search engine [4].

3.4 The web structure mining is a generated the link structure of the web pages. The page ranking algorithm are depends upon the link structure of the web pages and the random surfer model are also fall in this class[5]. The structure mining have a one drawback of the web data, this are not generated structure for the vast amount of data [5]. The structure mining generated link structure by in link and out link on the web pages.

3.5 Now the WWW is a covered by wide area of network. The billions of pages are available on the web and user is search on these pages in any time so many in links and out link is available of this pages [6]. But some in relevant pages have a many more in link are available than this type of in relevant pages rank will be higher compare to other relevant pages so search engine are could not decided which pages rank are give higher [6]. So based on the page rank algorithm search engine are decided the higher rank of the web pages and some back link are connected between the relevant number of the pages user click on random number of pages than back link are connected between the this type of the pages and rank will be higher compare to other pages.

3.6 In the page rank algorithm the back link are very important by calculating the ranking of the web pages [7]. If any pages have a more number of in links are available than this page-page rank are will be higher compare to other pages and this type of page are first open in to the list. This back link and out link are calculated between the numerical calculating on the web [7]. The back link is improving the efficiency of the WebPages. So based on this link structure calculating the page rank of the web. Now widely link structure is

available on the web or on the search engine. And the random rank is one part of the page rank algorithm [7]. Calculating the random rank are based on to the user click randomly on the web pages and after that clicking event occur the some in link and out link and based on this in link and out link the random rank are generate on the search engine.

## 4. PAGE RANKING ALGORITHM

We assume page A has pages T1...TN which point to it (i.e., are citations). The parameter d is a damping factor, which can be set between 0 and 1. We usually set d to 0.85…….. C (A) is defined as the number of links going out of page A. The Page Rank of a page A is given as follows [6]

$$PR(A) = (1 - d) + d\left(PR(T1)/C(T1) + ... + PR(Tn)C(Tn)\right)$$

Where PR(A) = web page

d = damping factor = 0.85

PR(T1) = Number of in link

C(T1) = Number of out link

The Page rank algorithm is a analyze the some links and hyperlinks structure on the web. It assign the numerical weight of all the web pages and its denote the importance of the particular one pages to the some other relevance pages [3].

Example of Page ranking algorithm.



Two pages A and B are pointing to each other

Start with PR (A) = PR (B) = 1

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

        = (1-0.85) + 0.85 * (1/1)

        = 1

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

        = (1-0.85) + 0.85 (1/1))

        = 0.15 + 0.85 (1)

        = 1

Let's start with PR (A) = PR (B) = 10

After 1st iteration:

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

    = (1-0.85) + 0.85 * (10/1)

    = 0.15 + 0.85 + 10/1

    = 8.6

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

     = (1-0.85) + 0.85 * (8.65/1)

     = 0.15 + 0.85 * 8.65

     = 7.50

After 2nd iteration:

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

     = (1-0.85) + 0.85 (7.50/1)

     = 0.15 +0.85(7.50)

     = 6.52

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

     = (1-0.85) +0.85 (6.527/1)

     = 0.15 + 0.85(6.527)

     = 5.6

Till up to 1…

# 5. RANDOM SURFER MODEL WITH RANKSINK PROBLEM

The Random surfer model is a used into the Page ranking algorithm, in the random surfer model who is clicking on the page randomly. The randomly rank are generated in this model. Page Rank also has a second definition based upon the model of a random surfer navigating the Internet. In short, the model states that Page Rank models the behavior of someone who keeps clicking on successive links at random. However, occasionally the surfer gets bored and jumps to a random page chosen based on the distribution [8].The page rank and random surfer rank are generated based on the some in link and out link. Many forming loop are available in random surfer model algorithm and many effected link are generated in random surfer model algorithm[12].The rank sink problem should be more affected the page rank on search engine.

Equation for random surfer model [9] :

$$PR(J) = (1 - d) + d * \sum_{t \in B(J)} PR(t) / O(t)$$

Where PR (J) = which page counting the rank

     d   = damping factor = 0.85

     PR (t) = Number of in link
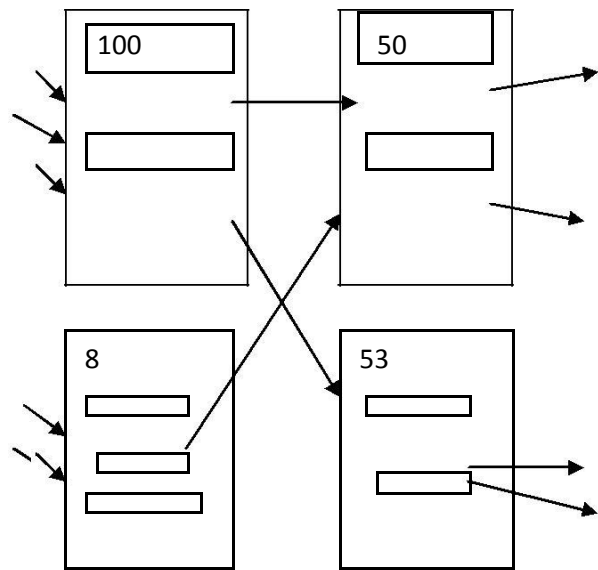
     O (t) = Number of out link



**Fig 2. Simplified Page Rank algorithm**

In the page rank algorithm every in link and out link are counting and based on this in link and out link the rank are generated about the web pages[13].
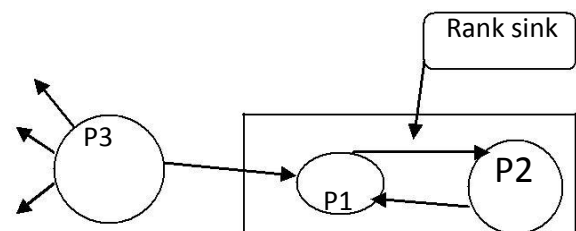
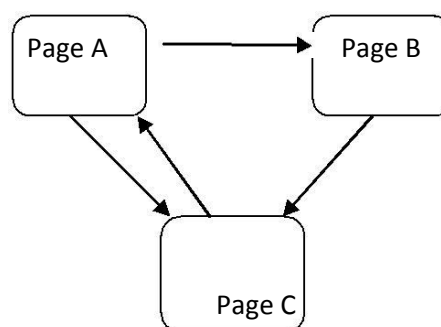Example of Rank sink problem



**Fig 3. Rank sink problem**



**Fig 4. Example of Rank sink problem**

Now we find the Page rank of A , B , C using the Random surfer model equation,

PR(A) = (1-d) + d $\sum$ PR(t)∕O(t)
t€B(J)

= (1-d) + d PR(C)∕O(C)

= (1-0.85) + 0.85 (2/1)

= 1.85

Same as Page B rank is

PR(B) = 1.075

Same as Page C is

PR(C)= 2.57
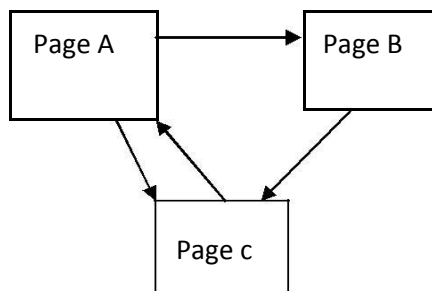
After three iteration

PR (A) = 3.77

PR (B) = 4.086

PR(C) = 6.72

## 6. RANDOM SURFER MODEL WITHOUT RANK SINK PROBLEM



Equation of Random surfer model without Rank sink problem
Equation of Random surfer model without Rank sink problem
PR(J) = ((1-d) + d * $\sum$ PR(t)∕O(t) + (1-d) + N)) / N

t€B(J)

Where PR (J) = Pointing of the pages d =
damping factor = 0.85

PR (t) = Number of the in link

O(t) = Number of Out link N = Total number of pages

B(J) = which page to the pointing After
three iteration

PR (A) = 3.67

PR (B) = 4.71

PR(C) = 7.5

## 7. COMPERISION BETWEEN RANDOM SURFER MODEL WITH AND WITHOUT RANK SINK PROBLEM

| Page | Random surfer model with Rank sink problem | Random surfer model without Rank sink problem |
|---|---|---|
| Page A | 3.7 | 3.67 |
| Page B | 4.086 | 4.71 |
| Page c | 6.72 | 7.5 |

## 8. CONCLUSION

This survey paper in removing the rank sink problem with the page ranking algorithm. Removing the rank sink problem because the rank sink problem are created the problem in given the page rank for some relevant pages on the searching url or affected the page rank for some relevant pages on the url. This model are work on the url clicking, if user clicking on some randomly pages than link are directly connected to the some clicking pages and after some relevant pages have not link connected than relevant pages are not considered between the any rank of the clicking event. So we needed the rank sink problem are remove in to the random surfer model in page rank algorithm. So we improve the result of the some relevance pages on the searching URL.

## 9. REFERENCES

[1] Brin, Sergey, and Lawrence Page. "The anatomy of a large-scale hyper textual Web search engine.[En línea]." Disponible en Web (1998).

[2] D. Jayalatchumy, Dr. P.Thambidurai, IOSR Journal of Engineering p- ISSN: 2278-8727Volume 14, Issue 3, Sep. - Oct. 2013, PP 20-27.

[3] Praveen Rani, International Journal of Advanced Researchin Engineering & Technology (IJARCET) Volume 2, Issue 3, and March 2013.

[4] A.M. Sote1, Dr. S. R. Pande2, IOSR Journal of Computer Science 2014 (ICAET-2014).

[5] Zakaria Suliman Zubi, ISBN: 978-1-61804-169-2, 2009.

[6] Ms.M.Sangeetha, Dr.K.Suresh Joseph, ISBN No.978-1-4799-3834-6/14/$31.00©2014 IEEE.

[7] Navadiya Hareshkumar, Dr. Deepak Garg, International Journal of Computer Applications (0975 – 8887) Volume 35– No.11, December 2011.

[8] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd, The Page Rank Citation Ranking: Bringing Order to the Web (1998).

[9] Richardson, M., Domingos P (2002). The Intelligent Surfer: Probabilistic Combin. Jaideep Srivastava, Prasanna Desikan, Vipin Kumar, "Web Mining - Concepts, Applications & ResearchDirections" Page Rank+ Citation Ranking: Bringing Order to theWeb",

[10] Ashish Jain, Rajeev Sharma, Gireesh Dixit and VarshaTomar "Page Ranking Algorithms in Web Mining,Limitations of Existing methods and a New Method for Indexing Web Pages" Conference - IEEE 2013.

[11] Page, Lawrence, et al. "The Page Rank citation ranking: bringing order to the Web."

[12] Wei Huang and Bin Li, "An Improved Method for the Computation of Page Rank", IEEE 2011.