

CHAPTER 1

INTRODUCTION

1.1. AN OVERVIEW

The World Wide Web (also referred to as the WWW) is the repository of information that has acquired explosive growth. WWW has a unique combination of flexibility, portability, and user-friendly features that distinguish it from other services provided by the Internet [21]. Since its evolution in the late 1980's, WWW has already grown up to a great extent and is still growing rapidly [11]. The billions of pages are added every week. The information present on web is very large, unorganized and unstructured. So it is very difficult to search on web without the help of some external agents. There exist some "search engines" as such, which do this required task for us. So a search engine is designed to search information on the World Wide Web and the search results are generally presented in a list format. Information that we are searching on the web may consist of web pages, photos, graphs or any other type of files or data. A search engine make a local store for the web by downloading the web pages with the help of an agent called crawler. An indexer module builds an index by indexing the information brought by the crawler in the local store. Later Searcher gets the query from the user and searches the keywords of query in the index. A standard search engine consists of the following components:

Crawler

Crawling is the real undertaking performed by a web search engine in the field of web technology. A search engine should have a highly scalable crawler because the crowd of web pages is increasing in an exponential rate. A web crawler also known as spider is a program that browses the WWW to retrieve web pages. It has to constantly monitor whether the web page has changed and refresh the downloaded pages. A page is fresh when it is similar to the page present at the web server. It interacts with several web servers and visits the web pages' URLs provided by the URL servers and goes through the hyperlinks given on that page. It adds these URLs in its list and makes its own DNS cache [1]. It also gives unique id to each web page called doc id. Figure 1.1 shows how a crawler works.

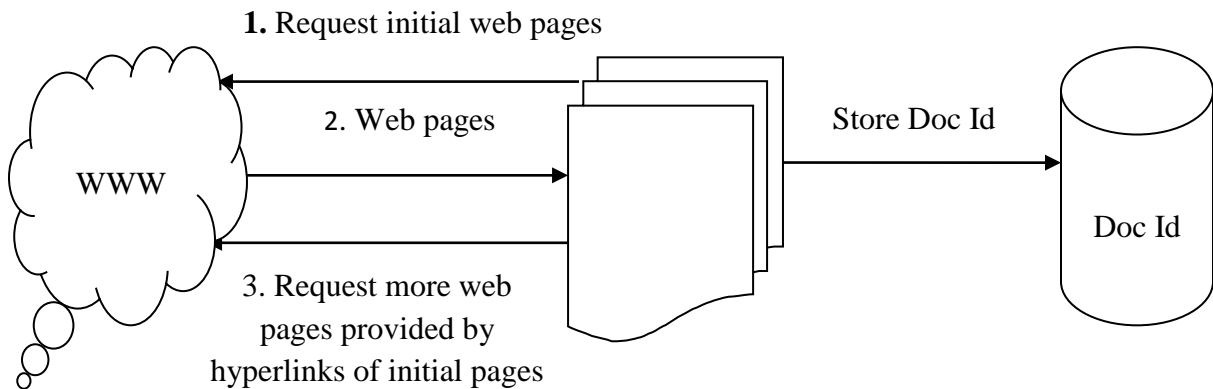


Figure 1.1: Working of a Crawler

Indexer

Indexing is the process of collecting, parsing and storing of the data for the use of searcher module. Indexing associates the keywords to the documents in which keywords are present. Without indexing, search engine would take large amount of time to reply for a query as search engine would have to search all the web pages in the repository at the time of searching. It firstly parses the web pages to extract keywords and generates inverted index. Inverted index consists of keywords and the doc id in which keyword is present. The working of an indexer has been shown in Figure 1.2.

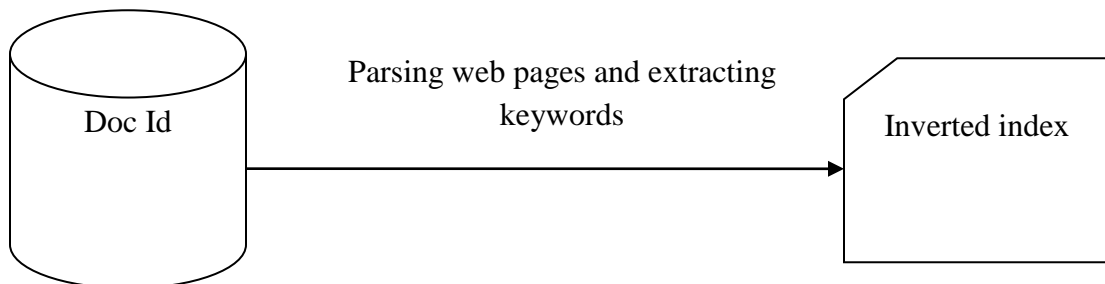


Figure 1.2: Working of an Indexer

Searcher

Searcher gets the query from user and searches the keywords of the query in the inverted index. Then it returns the best matching web pages associated with the query. As the size of web is very large, number of documents retrieved for a particular query is also large. So a search engine has to use some ranking algorithm to prioritize the retrieved web pages. Relevant documents are appeared at the top of the results. Boolean logic operators are used in query to narrow or expand the searches. Logic operator 'OR' is used to expand and 'AND' is used to narrow the results. Figure 1.3 shows the working of searching component.

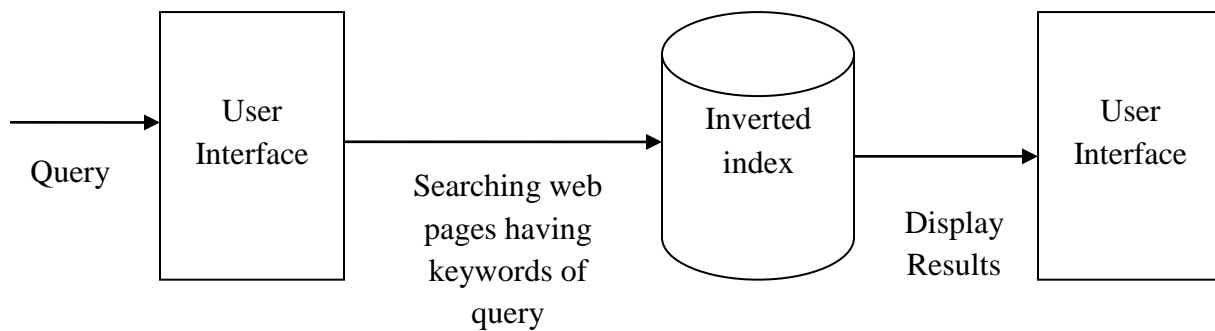


Figure 1.3: Working of a Searcher

Due to enormous amount of data present on the web, the number of web pages returned by a searcher is very large out of which only few are relevant. In a standard search engine, firstly the crawler downloads pages from the WWW. Then the web pages are sent to indexer module that builds the index on the basis of keywords present in the web pages. The query processor module accepts the query from the user and returns the list of web pages which matches the keywords present in the query. But before presenting the resultant web pages to the user, query processor module has to sort the results so that relevant pages are displayed at the top. For sorting the web pages, we need some kind of page ranking algorithms. Web mining plays a vital role in the development of page ranking algorithms. Figure 1.4 shows the working of a standard search engine [12].

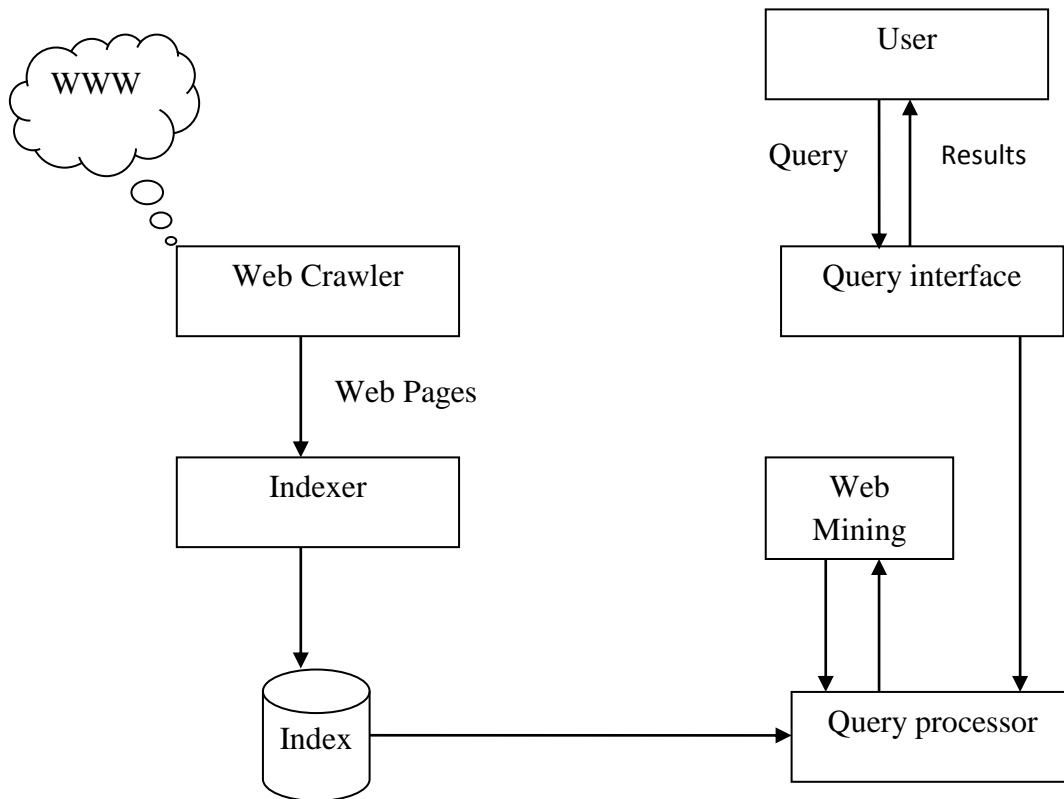


Figure 1.4: Working of a Search Engine [12]

The need is to sort the search results leads to the development of first page ranking algorithm called Page Rank by Surgey Brin and Larry Page. It uses web graph to calculate the value of page rank of each page. Later on, other variation of page ranking algorithms based on web structure mining and web content mining were also developed like Weighted PageRank, Page Content Ranking, SALSA etc. These ranking algorithms sort the search results returned by a search engine so that most relevant web pages are displayed at the top. These algorithms are based on web structure mining or web content mining or combination of both. But web structure mining only considers link structure of the web and web content mining is not able to cope up with multimedia such as images, mp3 and videos [4]. These algorithms do not take web usage behavior into account. The relevancy of a web page for a user can be determined by how many users click on the link (VOL), or recent uses of link or time spent on the link. So our aim is to design an algorithm which takes user relevance ranking into account.

1.2. MOTIVATION OF THE WORK

Due to enormous amount of data present on the web, information retrieval has been a great issue. It is becoming harder day by day to get the relevant documents in response of a query. This is the reason why we need ranking algorithms to prioritize the search results. Ranking algorithms are driving force behind search engine's working. The first ranking algorithm was introduced by Sergey Brin and Lawrence Page which was based on web structure mining. Later on, other algorithms were developed to rank the web pages such as Weighted PageRank, Page Content Ranking, and HITS etc. Most of the page ranking algorithms are based on web structure mining (WSM) or web content mining (WCM) or combination of both. WSM is extracting information from the structure of the web. WCM is extracting information from the contents of the web. A web graph contains nodes representing web pages and links representing hyperlinks between the connected pages.

The aim of the research work is to first survey the page ranking algorithms and then to propose an algorithm by combining the web structure mining and web usage mining(WUM) to calculate the value of page rank. WUM intends to extract what users are looking on the web. The relevancy of a web page can be estimated by usage trends of web pages. There are different criteria to check the relevancy of web page like frequency of visit of a web page, recent visit of a web page, time spent on a web page etc. The thesis uses number of visit of links (VOL) and combines with WSM to calculate the value of page rank.

1.3. GOALS OF THE THESIS

The overall goal of the thesis is to survey the existing page ranking algorithms and to propose a new ranking algorithm based on web structure mining and web usage mining. The goals of the thesis are to:

- To survey the existing page ranking algorithms and discussing their advantages and disadvantages.

- To propose a new ranking algorithm based on web structure mining and web usage mining as web usage mining can help in determining the relevance of web pages from users' point of view.
- Apply the proposed technique on a web graph to validate the proposed algorithm.

1.4. ORGANIZATION OF THESIS

The remainder of the thesis is structured as follows:

Chapter 2 discusses the previous work done in the field of page ranking algorithms. This includes the extensive study of various page ranking algorithms that have been proposed in the literature so far. It also highlights some of the most relevant works in the direction of field of work presented in the thesis.

Chapter 3, Classification of web mining has been introduced and how it has been used in ranking of web pages. We also survey existing page ranking algorithms with their advantages and disadvantages.

Chapter 4 introduces the proposed algorithm which uses combination of WSM and WUM to calculate the value of page rank. The formulae have been derived to calculate the value of page rank. It also discusses the benefits of proposed algorithm.

Chapter 5 describes the results and compares the results with an existing algorithm.

The final section concludes the thesis with a summary of the results, and a discussion on possible future directions along with the references used.

CHAPTER 2

LITERATURE REVIEW

WWW has ample number of hyperlinked documents and these documents contain heterogeneous information including text, image, audio, video, and metadata. So there are lots of search results corresponding to a user's query out of which only some are relevant. It is becoming harder day by day to get the relevant documents in response of a query. Ranking algorithms are driving force behind search engines' working. The first ranking algorithm was introduced by Sergey Brin and Lawrence Page which was based on web structure mining. Later on, other algorithms were developed to rank the web pages such as Weighted PageRank, Page Content Ranking, and HITS etc. Most of the page ranking algorithms are based on web structure mining (WSM) or web content mining (WCM) or combination of both. WSM is extracting information from the structure of the web. WCM is extracting information from the contents of the web. The relevancy of a web page is calculated by search engines using page ranking algorithms. Ranking algorithms are required to sort the results so that more relevant documents are displayed at the top.

Brin and Page [1] came up with an idea at Stanford University to use link structure of the web to calculate page rank of web pages. The algorithm was named PageRank after Larry Page (Cofounder of Google Search Engine). PageRank was the first ranking algorithm which was used to prioritize the results produced by keyword based search.

Wenpu Xing and Ali Ghorbani [18] proposed an algorithm called Weighted PageRank algorithm by extending standard PageRank. The working principle behind the algorithm was that an important page has more linkages from other web pages have to it or is linked to by it. Unlike standard PageRank, it did not evenly distribute the page rank of a page among its outgoing linked pages but the page rank of a web page is divided among its outgoing linked pages in proportional to the importance or popularity (its number of inlinks and outlinks). The popularity from the number of inlinks and the popularity from the number of outlinks were used to calculate the values of page rank.

Hyperlink-Induced topic search (HITS) algorithm based on WSM was developed by Kleinberg [3]. It works on the principle that for a given query, there is a set of authority pages that are relevant for a given query and set of hub pages that contain links to relevant pages which includes links to many authority pages also. The algorithm finds the set of authority pages relevant for a query using sampling.

Gyanendra Kumar et. al. [6] came up with a new idea to incorporate user's browsing behavior in calculating page rank. Previous algorithms were either based on web structure mining or web content mining but none of them took web usage mining into consideration. A new page ranking algorithm called Page Ranking based on Visits of Links (VOL) was proposed for search engines. It modifies the basic page ranking algorithm by taking into consideration the number of visits of inbound links of web pages. It helps to prioritize the web pages on the basis of user's browsing behavior.

Neelam Tyagi and Simple Sharma [13] incorporated user browsing behavior in Weighted PageRank algorithm to develop a new algorithm called Weighted PageRank based on number of visits of links (VOL). The algorithm assigns more rank to the outgoing links having high VOL. It only considers the popularity from the number of inlinks and ignores the popularity from the number of outlinks which was incorporated in Weighted PageRank algorithm.

CHAPTER 3

PAGE RANKING ALGORITHMS: A SURVEY

3.1. WEB MINING

Data mining can be defined as the process of extracting useful information from large amount of data. The application of data mining techniques to extract relevant information from the web is called as web mining [14] [2]. The process of web mining [2] can be shown in Figure 3.1.



Figure 3.1: Process of Web Mining [2]

The data is retrieved from the web, some data mining and machine learning techniques are applied on the data and useful pattern can be discovered. Web mining is categorized into three techniques [14] as given in Figure 3.2.

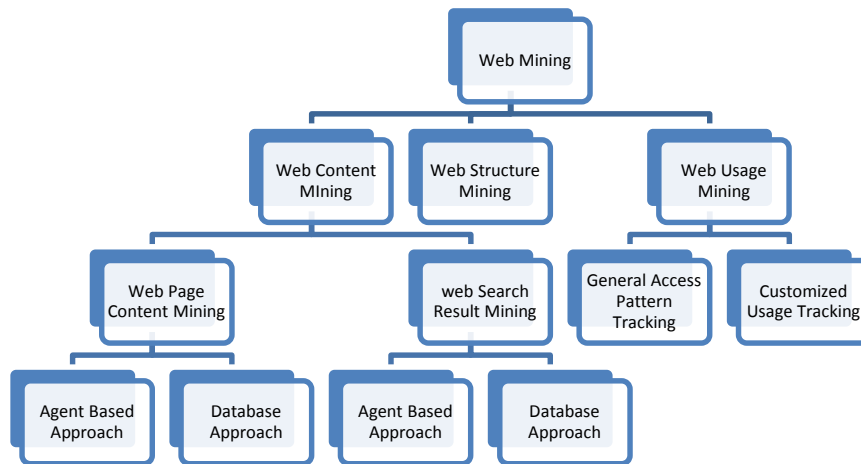


Figure 3.2: Categories of Web Mining [12]

Web Content Mining

Web Content Mining is used to mine the content of web pages. This technique can be applied either on the web pages or on the result pages obtained by the query processor of a search engine. WCM can be differentiated from two different views: Information Retrieval (IR) View and Database (DB) View. IR view works for unstructured and semi-structured data [9]. To represent unstructured data, bag of words is used and HTML structure inside the documents is used to represent semi-structured data. In DB view, a web site can be transformed to represent a multi-level database and web mining tries to infer the structure of the web site from this database.

Web Structure Mining

WSM can be defined as extracting information from the structure of the web [15]. A web graph contains nodes representing web pages and links representing hyperlinks between the connected pages as shown in Figure 3.3. A hyperlink is a structural unit that is used to move from a location in a web page to a different location, either within the same web page or on a different web page. A hyperlink can be differentiated from two different views:

Inlinks: The links which point to a web page are called as inlinks of that page. e.g.; the links from web pages A, B and C are inlinks of web page P.

Outlinks: The links from a web page to others web page are called outlink of that page. e.g.; the links from web page P to Y and Z are outlinks of page P.

The number of outlinks from a page and the number of inlinks to a page are very important parameter in the area of web mining [15]. The importance of the web page is measured by the number of inlinks in PageRank algorithm and by the number of inlinks and outlinks in Weighted Page Rank algorithm. So WSM is a very important area in the field of page ranking algorithm. Structure of a web graph is shown in Figure 3.3.

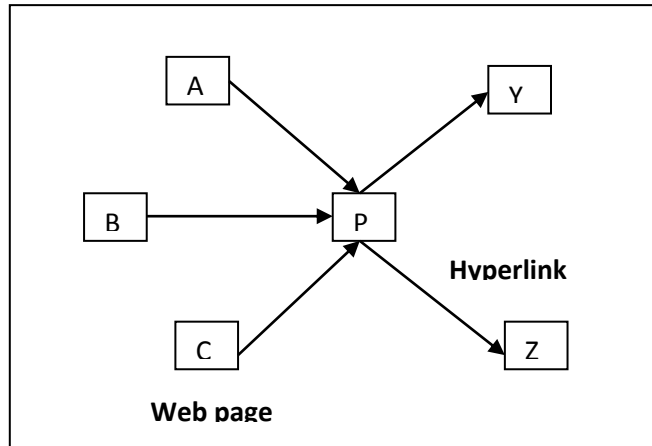


Figure 3.3: Structure of a Web Graph

Web Usage Mining

Web Usage Mining (WUM) is used to extract information from the server logs which are maintained during interaction with the web. WUM analyze what are the people looking for on the web. Server logs provide information about identity of web users, access time and their browsing behavior [22] which may help in understanding of web based application. The server logs can provide information like

- the frequency of visits per document
- most recent visit per document
- who is visiting which documents
- frequency of use of each hyperlink
- most recent use of each hyperlink

which can be used in calculating the value of page rank [8]. It can be further categorized in finding the general access patterns or finding the patterns matching the specified parameters. These web mining techniques are used in page ranking algorithm. In the next section, the page ranking techniques have been explained.

3.2. VARIOUS PAGE RANKING ALGORITHMS

WWW has ample number of hyperlinked documents and these documents contain heterogeneous information including text, image, audio, video, and metadata. So there are lots of search results corresponding to a user's query out of which only some are relevant. The relevancy of a web page is calculated by search engines using page ranking algorithms. Ranking algorithms are required to sort the results so that more relevant documents are displayed at the top. Various ranking algorithms have been developed such as PageRank, Weighted PageRank, WPR using link visits etc.

3.2.1. PAGERANK ALGORITHM

Brin and Page [1] came up with an idea at Stanford University to use link structure of the web to calculate page rank of web pages. The algorithm was named PageRank after Larry Page (Cofounder of Google Search Engine). The algorithm is used to prioritize the results produced by keyword based search. It works on the principle that if a web page has important

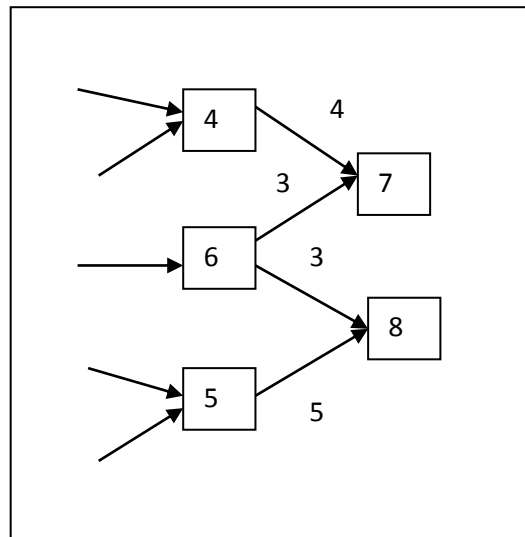


Figure 3.4: Working of PageRank Algorithm

links towards it then the links of this page to other pages are also considered important. PageRank algorithm considers only backlinks into account and propagates the ranking through links. i.e., A web page divides equally its rank to its outgoing links which is illustrated in Figure 3.4. In Figure 3.4, the web page having rank six divides its rank equally to its two outgoing links

and the rank of a web page is calculated by summing up the ranks from its incoming links. So the overall rank of a page using PageRank algorithm is calculated by the formula given in equation 3.1:

$$PR(u) = c \sum_{v \in B(u)} \frac{PR(v)}{N_v} \quad \dots (3.1)$$

Notations are:

- u represents a web page,
- $PR(u)$ and $PR(v)$ represents the page rank of web pages u and v respectively,
- $B(u)$ is the set of web pages pointing to u ,
- N_v represents the total numbers of outlinks of web page v ,
- c is a factor used for normalization.

Original PageRank algorithm was modified by taking into consideration that not all users follow direct links on WWW. The modified formula for calculating page rank is given in equation 3.2.

$$PR(u) = (1 - d) + d \sum_{v \in B(u)} \frac{PR(v)}{N_v} \quad \dots (3.2)$$

Where d is a dampening factor which represents the probability of user using direct links and it can be set between 0 and 1.

Algorithm to calculate Standard Page Rank

To calculate the value of page rank using SPR, following steps are required.

1. *Finding a web graph:* A web graph with rich hyperlinks is required because the SPR algorithm makes use of the web structure to calculate values of page rank.
2. *Building a link matrix:* The link matrix consisting of binary values is created in which 1 means hyperlink present between web pages and 0 means no hyperlink present between web pages.
3. *Calculate intermediate values:* For each web page v , calculate the total number of outlinks of web page v .

4. *Calculate page rank of each web page:* The intermediate values are substituted in equation 3.2 to calculate values of page rank.
5. *Repetition of step 5:* The step 4 is used recursively until a stable value of page rank is obtained.

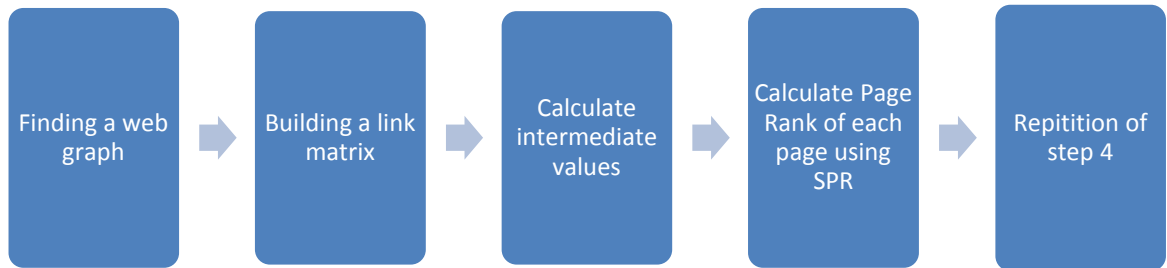


Figure 3.5: Algorithm of Standard PageRank Algorithm

Illustration of PageRank algorithm

The working of Pagerank algorithm can be illustrated by taking an example shown in Figure 3.5.

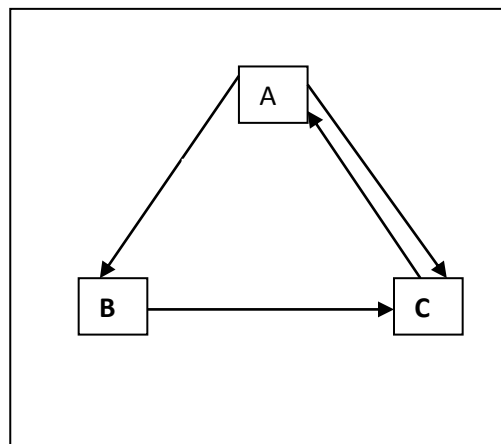


Figure 3.6: A web graph [13]

It has three web pages A, B and C and hyperlinks between them. The page rank of A, B and C can be calculated using equation (3.2) as:

$$PR(A) = (1 - d) + d \left(\frac{PR(C)}{1} \right) \quad \dots (3.2.1)$$

$$PR(B) = (1 - d) + d \left(\frac{PR(A)}{2} \right) \quad \dots (3.2.2)$$

$$PR(C) = (1 - d) + d \left(\frac{PR(A)}{2} + PR(B) \right) \quad \dots (3.2.3)$$

Initially the rank of web page C is considered 1 and value of d is taken as 0.5. Then page ranks are calculated using equation 3.2.1, 3.2.2 and 3.2.3 and the values are A=1, B=0.75000 and C=1.125. These values are used iteratively until the values get stabilized. The intermediate values of A, B and C have been shown in Table 3.1.

A	B	C
1	0.75000	1.12500
1.06250	0.76563	1.14844
1.07422	0.76855	1.15283
1.07642	0.76910	1.15366

Table 3.1: Values of Page Ranks using PageRank Algorithm

The PageRank algorithm was the first ranking algorithm to use link structure of web to sort the web pages and it helped the users of search engines to find the web pages of their interest. But the method only considered the links of web pages and relevancy of a web page was totally ignored. The presence of query terms in web pages did not affect the rank of web page.

3.2.2. WEIGHTED PAGERANK ALGORITHM

Wenpu Xing and Ali Ghorbani [18] proposed an algorithm called Weighted PageRank algorithm by extending standard PageRank. It works on the principle that if a page is important, more linkages from other web pages have to it or are linked to by it. Unlike standard PageRank, it does not evenly distribute the page rank of a page among its outgoing linked pages. The page rank of a web page is divided among its outgoing linked pages in proportional to the importance or popularity (its number of inlinks and outlinks). The popularity from the number of inlinks and the popularity from the number of outlinks are used to calculate the values of page rank.

$W^{\text{in}}(v, u)$, the popularity from the number of inlinks, is calculated based on the number of inlinks of page u and the number of inlinks of all reference pages of page v as given in equation 3.3.

$$W^{\text{in}}(v, u) = \frac{I_u}{\sum_{p \in R(v)} I_p} \quad \dots (3.3)$$

Notations are:

- I_u and I_p are the number of inlinks of page u and p respectively,
- $R(v)$ represents the set of web pages pointed by v .

$W^{\text{out}}(v, u)$, the popularity from the number of outlinks, is calculated based on the number of outlinks of page u and the number of outlinks of all reference pages of page v as given in equation 3.4.

$$W^{\text{out}}(v, u) = \frac{O_u}{\sum_{p \in R(v)} O_p} \quad \dots (3.4)$$

Notations are:

- O_u and O_p are the number of outlinks of page u and p respectively,
- $R(v)$ represents the set of webpages pointed by v .

The page rank using Weighted PageRank algorithm is calculated by the formula given in equation 3.5.

$$WPR(u) = (1 - d) + d \sum_{v \in B(u)} WPR(v)W^{in}(v, u)W^{out}(v, u) \quad \dots (3.5)$$

Notations are:

- $WPR(u)$ and $WPR(v)$ are page rank of web page u and v respectively,
- $W^{in}(v, u)$ represents the popularity from the number of inlinks of web page u ,
- $W^{out}(v, u)$ represents the popularity from the number of outlinks of web page u ,
- $B(u)$ is the set of web pages that point to u ,
- d is the dampening factor.

Algorithm to calculate WPR

To calculate the value of page rank using WPR, following steps are required [18].

1. *Finding a web graph:* A web graph with rich hyperlinks is required because the WPR algorithm makes use of the web structure to calculate values of page rank.
2. *Building a link matrix:* The link matrix consisting of binary values is created in which 1 means hyperlink present between web pages and 0 means no hyperlink present between web pages.
3. *Calculate intermediate values:* For each web page u , the value of popularity from the number of inlinks ($W^{in}(v, u)$) and the value of popularity from number of outlinks ($W^{out}(v, u)$) is calculated by the formula given in equation 3.3 and 3.4 respectively.
4. *Calculate page rank of each web page:* The intermediate values are substituted in equation 3.5 to calculate values of page rank.
5. *Repetition of step 4:* The step 4 is used recursively until a stable value of page rank is obtained.

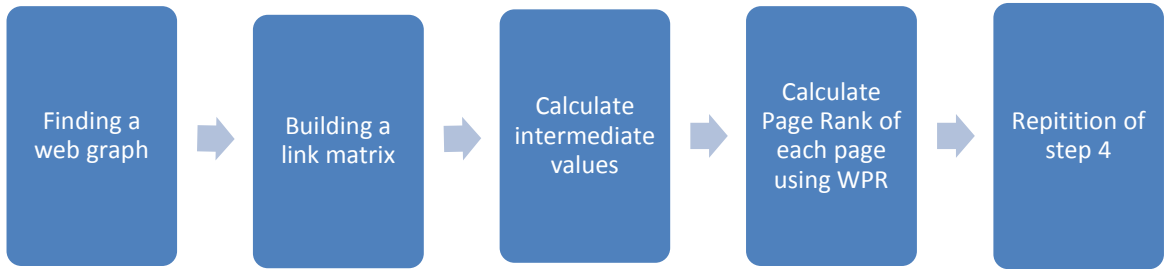


Figure 3.7: Algorithm of Weighted PageRank Algorithm

Illustration of Weighted Page Rank algorithm

The working of Weighted Pagerank algorithm can be illustrated by taking the example shown in Figure 3.5. It has three web pages A, B and C and hyperlinks between them. The page rank of A, B and C can be calculated using equation 3.5 as:

$$WPR(A) = (1 - d) + d \left(WPR(C) * W^{in}(C, A) * W^{out}(C, A) \right) \quad \dots (3.5.1)$$

$$WPR(B) = (1 - d) + d(WPR(A) * W^{in}(A, B) * W^{out}(A, B)) \quad \dots (3.5.2)$$

$$WPR(C) = (1 - d) + d(WPR(A) * W^{in}(A, C) * W^{out}(A, C) + WPR(B) * W^{in}(B, C) * W^{out}(B, C)) \quad \dots (3.5.3)$$

The values of $W^{in}(v, u)$ and $W^{out}(v, u)$ can be calculated using equation 3.3 and 3.4.

$$W^{in}(C, A) = \frac{I_A}{I_A} = \frac{1}{1} = 1$$

$$W^{out}(C, A) = \frac{O_A}{O_A} = \frac{2}{2} = 1$$

$$W^{in}(A, B) = \frac{I_B}{I_B + I_C} = \frac{1}{1 + 2} = \frac{1}{3}$$

$$W^{out}(A, B) = \frac{O_B}{O_B + O_C} = \frac{1}{1 + 1} = \frac{1}{2}$$

$$W^{in}(A, C) = \frac{I_C}{I_C + I_B} = \frac{2}{2 + 1} = \frac{2}{3}$$

$$W^{out}(A, C) = \frac{O_C}{O_C + O_B} = \frac{1}{1 + 1} = \frac{1}{2}$$

$$W^{in}(B, C) = \frac{I_C}{I_C} = \frac{2}{2} = 1$$

$$W^{out}(B, C) = \frac{O_C}{O_C} = \frac{1}{1} = 1$$

These values are put in equation 3.5.1, 3.5.2 and 3.5.3 and the resultant equations are:

$$WPR(A) = (1 - d) + d(WPR(C) * 1 * 1) \quad \dots (3.5.1)$$

$$WPR(B) = (1 - d) + d\left(WPR(A) * \frac{1}{3} * \frac{1}{2}\right) \quad \dots (3.5.2)$$

$$WPR(C) = (1 - d) + d\left(WPR(A) * \frac{1}{3} * \frac{1}{2} + WPR(B) * 1 * 1\right) \quad \dots (3.5.3)$$

Initially the rank of web page C is considered 1 and value of d is taken as 0.5. Then page ranks are calculated using equation 3.5.1, 3.5.2 and 3.5.3 and the values are A=1, B=0.58333 and C=0.95833. These values are used iteratively until the values get stabilized. The intermediate values of A, B and C have been shown in Table 3.2.

A	B	C
1	0.58333	0.95833
0.97917	0.58160	0.95399
0.97701	0.58142	0.95354
0.97677	0.58142	0.95351

Table 3.2: Values of Page Ranks using Weighted PageRank Algorithm

Unlike PageRank algorithm, Weighted PageRank algorithm used both inlinks as well as outlinks to calculate the values of page rank. The page rank of a web page was not divided equally among

its outlinked pages but in proportion to the popularity of outlinked pages. But this method also considered only the links of web pages and relevancy of a web page was totally ignored. The presence of query terms in web pages did not affect the rank of web page.

3.2.5. PAGERANK WITH NUMBER OF VISITS OF LINKS

Gyanendra Kumar et. al. [6] came up with a new idea to incorporate user's browsing behavior in calculating page rank. Previous algorithms were either based on web structure mining or web content mining but none of them took web usage mining into consideration. A new page ranking algorithm called Page Ranking based on Visits of Links (VOL) was proposed for search engines. It modifies the basic page ranking algorithm by taking into consideration the number of visits of inbound links of web pages. It helps to prioritize the web pages on the basis of user's browsing behavior.

In the original PageRank algorithm, the rank of a page p is evenly distributed among its outgoing links but in this algorithm, rank values are assigned in proportional to the number of visits of links. The more rank value is assigned to the link which is most visited by user. The Page Ranking based on Visits of Links (VOL) can be calculated by the formula given in equation 3.6.

$$PR(u) = (1 - d) + d \sum_{v \in B(u)} L_u \frac{PR(v)}{TL(v)} \quad \dots (3.6)$$

Notations are:

- $PR(u)$ and $PR(v)$ represent page rank of web pages u and v respectively,
- d is dampening factor,
- $B(u)$ is the set of web pages pointing to u ,
- L_u is number of visits of links pointing from v to u ,
- $TL(v)$ is the total number of visits of all links from v .

Algorithm for PageRank using Visits of Links

The algorithm to calculate the value of page rank is:

1. *Finding a web graph:* A web graph with rich hyperlinks is to be selected because the algorithm depends on the hyper structure of website.
2. *Building a link matrix:* The link matrix consisting of binary values is created in which 1 means hyperlink present between web pages and 0 means no hyperlink present between web pages.
3. *Calculate intermediate values:* For each web page u , total number of visits of all links from u ($TL(u)$) and the number of visits of links pointing from v to u (L_u), is calculated where $v \in B(u)$ which is set of web pages pointing to u .
4. *Calculate page rank of each web page:* The intermediate values are substituted in equation 6 to calculate values of page rank.
5. *Repetition of step 4:* The step 4 is used recursively until a stable value of page rank is obtained.

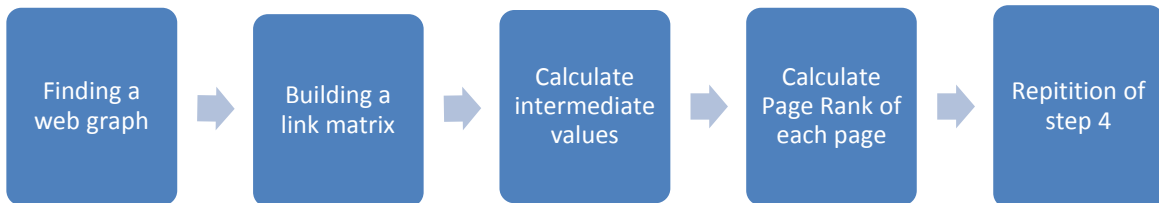


Figure 3.8: Algorithm of PageRank using VOL

Figure 3.12 shown above explains the steps required to calculate page rank using visit of links algorithm.

Illustration of PageRank using VOL

The working of PAGERANK algorithm using visits of links can be illustrated by taking the example shown in Figure 3.13. It has three web pages A, B and C and hyperlinks between them shows

number of visits of links. i.e., how many users have accessed that link. The page rank of A, B and C can be calculated using equation (3.6.1) as:

$$PR(A) = (1 - d) + d \left(\frac{2}{2} * PR(C) \right) \quad \dots (3.6.1)$$

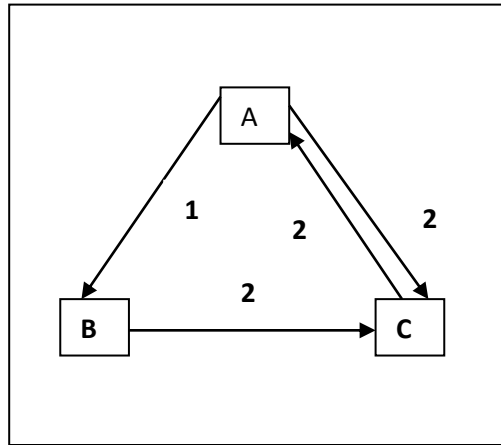


Figure 3.9: A Web Graph with VOL [13]

$$PR(B) = (1 - d) + d \left(\frac{1}{3} * PR(A) \right) \quad \dots (3.6.2)$$

$$PR(C) = (1 - d) + d \left(\frac{2}{3} * PR(A) + \frac{2}{2} * PR(B) \right) \quad \dots (3.6.3)$$

Then page ranks are calculated using equation 3.6.1, 3.6.2 and 3.6.3 and the values are A=1, B=0.66667 and C=1.6667. These values are used iteratively until the values get stabilized. The intermediate values of A, B and C have been shown in Table 3.3.

A	B	C
1	0.66667	1.16667
1.08334	0.68056	1.20139
1.10071	0.68345	1.20863
1.10432	0.68405	1.21013

Table 3.3: Values of Page Ranks using PageRank using VOL Algorithm

Unlike PageRank and Weighted PageRank algorithm, PageRank using visits of links makes use of web structure mining and web usage mining to calculate the value of page rank. So the web pages obtained by this algorithm are more relevant to the users as compared to the web pages obtained from PageRank and Weighted PageRank algorithm. In this algorithm too, the presence of query terms in web pages did not affect the rank of web page.

3.2.6. WEIGHTED PAGERANK ALGORITHM USING VISITS OF LINKS

Neelam Tyagi and Simple Sharma [13] incorporated user browsing behavior in Weighted PageRank algorithm to develop a new algorithm called Weighted PageRank based on number of visits of links (VOL). The algorithm assigns more rank to the outgoing links having high VOL .It only considers the popularity from the number of inlinks and ignores the popularity from the number of outlinks which was incorporated in Weighted PageRank algorithm.

In the original Weighted PageRank algorithm, the page rank of a web page is divided among its outgoing linked pages in proportional to the importance or popularity (its number of inlinks and outlinks) but in this algorithm, number of visits of inbound links of web pages are also taken into consideration.

$W^{in}(v, u)$, the popularity from the number of inlinks, is calculated based on the number of inlinks of page u and the number of inlinks of all reference pages of page v as given in equation 3.7.

$$W^{in}(v, u) = \frac{I_u}{\sum_{p \in R(v)} I_p} \quad \dots (3.7)$$

Notations are:

- I_u and I_p are the number of inlinks of page u and p respectively,
- $R(v)$ represents the set of webpages pointed by v .

The rank of web page using this algorithm can be calculated as given in equation 3.8.

$$WPR_{VOL}(u) = (1 - d) + d \sum_{v \in B(u)} \frac{L_u WPR_{VOL}(v) W^{in}(v, u)}{TL(v)} \quad \dots (3.8)$$

Notations are:

- $WPR_{VOL}(u)$ and $WPR_{VOL}(v)$ represent page rank of web page u and v respectively,
- d is the dampening factor,
- $B(u)$ is the set of web pages pointing to u ,
- L_u is number of visits of links pointing from v to u ,
- $TL(v)$ is the total number of visits of all links from v ,
- $W^{in}(v, u)$ represents the popularity from the number of inlinks of u .

Algorithm to calculate Weighted PageRank using VOL

The algorithm to calculate the value of page rank is:

1. *Finding a web graph:* The web graph with rich hyperlinks is to be selected because the algorithm depends on the hyper structure of web graph.
2. *Building a link matrix:* The link matrix consisting of binary values is created in which 1 means hyperlink present between web pages and 0 means no hyperlink present between web pages.
3. *Calculate intermediate values:* For each web page u , total number of visits of all links from u ($TL(u)$) and the number of visits of links pointing from v to u (L_u) and the value of popularity from the number of inlinks ($W^{in}(v, u)$) is calculated.
4. *Calculate page rank of each web page:* The intermediate values are substituted in equation 3.22 to calculate values of page rank.
5. *Repetition of step 4:* The step 4 is used recursively until a stable value of page rank is obtained.

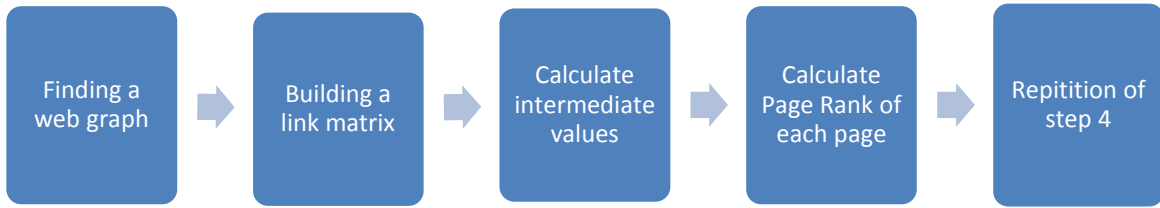


Figure 3.10: Algorithm of WPR_{VOL}

Figure 3.14 shown below explains the steps required to calculate page rank using $WPR_{VOL}(u)$.

Illustration of Weighed PageRank using VOL

The working of Pagerank algorithm using visits of links can be illustrated by taking the example shown in Figure 3.14.

$$WPR_{VOL}(A) = (1 - d) + d \left(\frac{2 * WPR_{VOL}(C) * W^{in}(C, A)}{2} \right) \quad \dots (3.8.1)$$

$$WPR_{VOL}(B) = (1 - d) + d \left(\frac{1 * WPR_{VOL}(A) * W^{in}(A, B)}{3} \right) \quad \dots (3.8.2)$$

$$WPR_{VOL}(C) = (1 - d) + d \left(\frac{2 * WPR_{VOL}(A) * W^{in}(A, C)}{3} \right) + d \left(\frac{2 * WPR_{VOL}(B) * W^{in}(B, C)}{2} \right) \quad \dots (3.8.3)$$

The values of $W^{in}(v, u)$ can be calculated using equation 3.3.

$$W^{in}(C, A) = \frac{I_A}{I_A} = \frac{1}{1} = 1$$

$$W^{in}(A, B) = \frac{I_B}{I_B + I_C} = \frac{1}{1 + 2} = \frac{1}{3}$$

$$W^{in}(A, C) = \frac{I_C}{I_C + I_B} = \frac{2}{2 + 1} = \frac{2}{3}$$

$$W^{in}(B, C) = \frac{I_C}{I_C} = \frac{2}{2} = 1$$

These values are put in equation 3.8.1, 3.8.2 and 3.8.3 and the resultant equations are:

$$WPR_{VOL}(A) = (1 - d) + d \left(\frac{2 * WPR_{VOL}(C)}{2} * \frac{1}{1} \right) \quad \dots (3.8.1)$$

$$WPR_{VOL}(B) = (1 - d) + d \left(\frac{1 * WPR_{VOL}(A)}{3} * \frac{1}{3} \right) \quad \dots (3.8.2)$$

$$WPR_{VOL}(C) = (1 - d) + d \left(\frac{2 * WPR_{VOL}(A)}{3} * \frac{2}{3} + \frac{2 * WPR_{VOL}(B)}{2} * \frac{2}{2} \right) \quad \dots (3.8.3)$$

Initially the rank of web page C is considered 1 and value of d is taken as 0.5. Then page ranks

A	B	C
1	0.55556	1
1	0.55556	1

Table 3.4: Value of Page Ranks using Weighted PageRank using VOL Algorithm

are calculated using equation 3.22.1, 3.22.2 and 3.22.3 and the values are A=1, B=0.55556 and C=1. These values are used iteratively until the values get stabilized. The intermediate values of A, B and C have been shown in Table 3.4. Weighted PageRank using visits of links makes use of web structure mining and web usage mining to calculate the value of page rank. So the web pages obtained by this algorithm are more relevant to the users as compared to the web pages obtained from PageRank and Weighted PageRank algorithm. In this algorithm too, the presence of query terms in web pages did not affect the rank of web page and it ignores the popularity from the number of outlinks, $W^{out}(v, u)$ which was used in Weighted PageRank algorithm. Table 3.5 gives a brief description of above algorithm using some parameters from [12].

PageRank algorithm and Weighted PageRank algorithms are based on web structure mining only. PageRank algorithm using visits of links and Weighted PageRank algorithm using visits of links are based on combination of web structure mining and web usage mining. PageRank algorithm relies only on the backlinks to calculate the value of page rank.

ALGORITHM	WEB MINING TECHNIQUE USED	INPUT PARAMETERS	IMPORTANCE	RELEVANCE
PageRank	Web structure mining	Backlinks	More	Less
Weighted PageRank	Web structure mining	Backlinks, Forward links	More	Less
PageRank with VOL	Web structure mining, Web usage mining	Content	More	More
Weighted PageRank with VOL	Web structure mining, Web usage mining	Backlinks and VOL	More	More

Table 3.5: Comparison of Page Ranking Algorithms

Weighted PageRank algorithm relies on the backlinks and forward links to calculate the value of page rank. PageRank using VOL and Weighted PageRank using VOL relies on the backlinks and visits of links to calculate the value of page rank. PageRank and Weighted PageRank algorithms do not consider relevancy of web pages into account. PageRank using VOL and Weighted PageRank using VOL consider relevancy from users' point of view as number of visits of links give the relevancy of a web page.

CHAPTER 4

PROPOSED WORK

The original Weighted PageRank algorithm distributes the rank of a web page among its outgoing linked pages in proportional to their importance or popularity. The algorithm is purely based on web structure mining and uses web graph to calculate the rank of web pages. $W^{in}(v, u)$, the popularity from the number of inlinks and $W^{out}(v, u)$, the popularity from the number of outlinks do not include usage trends. The rank of web page remains constant whether it has been visited by users or not. It does not give more popularity to the links most visited by the users. i.e.; the relevancy of a web page from user point of view is ignored and the weighted PageRank using VOL makes use of web structure mining and web usage mining to calculate the value of page rank.

In proposed algorithm, $W_{VOL}^{in}(v, u)$, the popularity from the number of visits of inlinks is used to calculate the values of page rank.

$W_{VOL}^{in}(v, u)$ is the weight of link(v, u) which is calculated based on the number of visits of inlinks of page u and the number of visits of inlinks of all reference pages of page v as shown in equation 4.1.

$$W_{VOL}^{in}(v, u) = \frac{I_{u(VOL)}}{\sum_{p \in R(v)} I_{p(VOL)}} \quad \dots (4.1)$$

Notations are:

- $I_{u(VOL)}$ and $I_{p(VOL)}$ represent the incoming visits of links of page u and p respectively.
- $R(v)$ represents the set of reference pages of page v.

Then $W_{VOL}^{in}(v, u)$, the popularity from the number of visits of inlinks are used to calculate page rank using equation 4.2.

$$EWPR_{VOL}(u) = (1 - d) + d \sum_{v \in B(u)} \frac{L_u EWPR_{VOL}(v) W_{VOL}^{in}(v, u)}{TL(v)} \quad \dots (4.2)$$

Notations are:

- d is a dampening factor.
- $B(u)$ is the set of pages that point to u .
- $EWPR_{VOL}(u)$ and $EWPR_{VOL}(v)$ are the rank scores of page u and v respectively.
- $W_{VOL}^{in}(v, u)$ represents the popularity from the number of visits of inlinks.

4.1. ALGORITHM TO CALCULATE $EWPR_{VOL}$

The algorithm depicts the steps required to calculate page rank of web pages using proposed algorithm.

1. *Finding a web graph*: This step requires finding a web graph which has rich hyperlinks because the algorithm depends on the hyper structure of web graph. The web graph having rich hyperlinks will help in better distribution of page rank.
2. *Building a link matrix*: The link matrix consisting of binary values is created in which 1 means hyperlink present between web pages and 0 means no hyperlink present between web pages .
3. *Calculate intermediate values*: For each web page u , total number of visits of all links from u ($TL(u)$) and the number of visits of links pointing from v to u (L_u) and the value of popularity from the number of visit of inlinks ($W_{VOL}^{in}(v, u)$) is calculated.
4. *Calculate page rank of each web page*: The values of $W_{VOL}^{in}(v, u)$, the popularity from the number of visits of inlinks, total number of visits of all links from u ($TL(u)$) and the number of visits of links pointing from v to u (L_u) calculated as intermediate values are substituted in equation 9 to calculate values of page rank.
5. *Repetition of step 4*: The step 4 is used recursively until a stable value of page rank is obtained.

Figure 4.1 shown below explains the steps required to calculate page rank using $EWPR_{VOL}$.

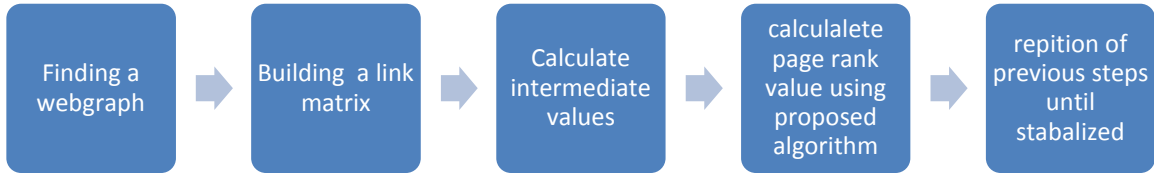


Figure 4.1: Algorithm to calculate $EWPR_{VOL}$

4.2. BENEFITS OF PROPOSED ALGORITHM

The original Weighted PageRank algorithm distributes the rank of a web page among its outgoing linked pages in proportional to their importance or popularity. $W^{in}(v, u)$, the popularity from the number of inlinks and $W^{out}(v, u)$, the popularity from the number of outlinks does not include usage trends. It does not give more popularity to the links most visited by the users. The weighted PageRank using VOL makes use of web structure mining and web usage mining. In proposed algorithm, $W_{VOL}^{in}(v, u)$, the popularity from the number of visits of inlinks is used to calculate the value of page rank. In this way the algorithm helps in sorting the resultant web pages in accordance of users need and it has following advantages.

1. The page rank using original WPR remains unaffected whether the page has been accessed by the users or not. i.e.; the relevancy of a web page is ignored. But the page rank using proposed method $EWPR_{VOL}$ assigns high rank to web pages having more visits of links.
2. The page rank using original WPR depends only on the link structure of the web and remains same whether the web page has been accessed by the user or not. Although the algorithm WPR_{VOL} makes use of web structure mining and web usage mining to calculate the value of page rank.

On the other side, our proposed method $EWPR_{VOL}$ makes use of $W_{VOL}^{in}(v, u)$, the popularity from the number of visits of inlinks is used to calculate page rank.

3. The proposed method uses number of visits of links to calculate the rank of web pages.
So the resultant pages are popular and more relevant to the users need.

CHAPTER 5

RESULTS AND IMPLEMENTATION

In this section, we compare the page rank of web pages using original WPR algorithm, WPR_{VOL} and the proposed algorithm. The original WPR makes use of web structure mining only to calculate the value of page rank. $W^{in}(v,u)$, the popularity from number of inlinks and $W^{out}(v,u)$, the popularity from the number of outlinks use the hyperlinks of graph and it do not change if the numbers of users accessing a link change. Although the WPR_{VOL} algorithm makes use of both web structure mining and web usage mining to calculate the value of page rank but it does not incorporate the popularity from the number of outlinks. But the proposed algorithm calculates $W_{VOL}^{in}(v,u)$, the popularity from number of visits of inlinks by analyzing the user behavior. When the numbers of visits of links change, the rank of web pages also changes. So this technique considers the relevancy from the users' point of view and gives high rank to those web pages which are frequently accessed by users. In this way, the proposed method gives improved results than standard WPR and WPR_{VOL} .

The values of page rank using WPR , WPR_{VOL} and $EWPR_{VOL}$ have been compared using a bar chart. The values retrieved by $EWPR_{VOL}$ are better than original WPR and WPR_{VOL} . The WPR uses only web structure mining to calculate the value of page rank, WPR_{VOL} uses both web structure mining and web usage mining to calculate value of page rank but it uses popularity only from the number of inlinks not from the number of outlinks. The proposed algorithm $EWPR_{VOL}$ method uses number of visits of inlinks to calculate values of page rank and gives more rank to important pages. Figure 5.1,5.2,5.3 shows the page ranks of A, B and C using WPR , WPR_{VOL} and $EWPR_{VOL}$ for $d=0.85$ and 5.4 compares the page rank values as shown in fig given below.

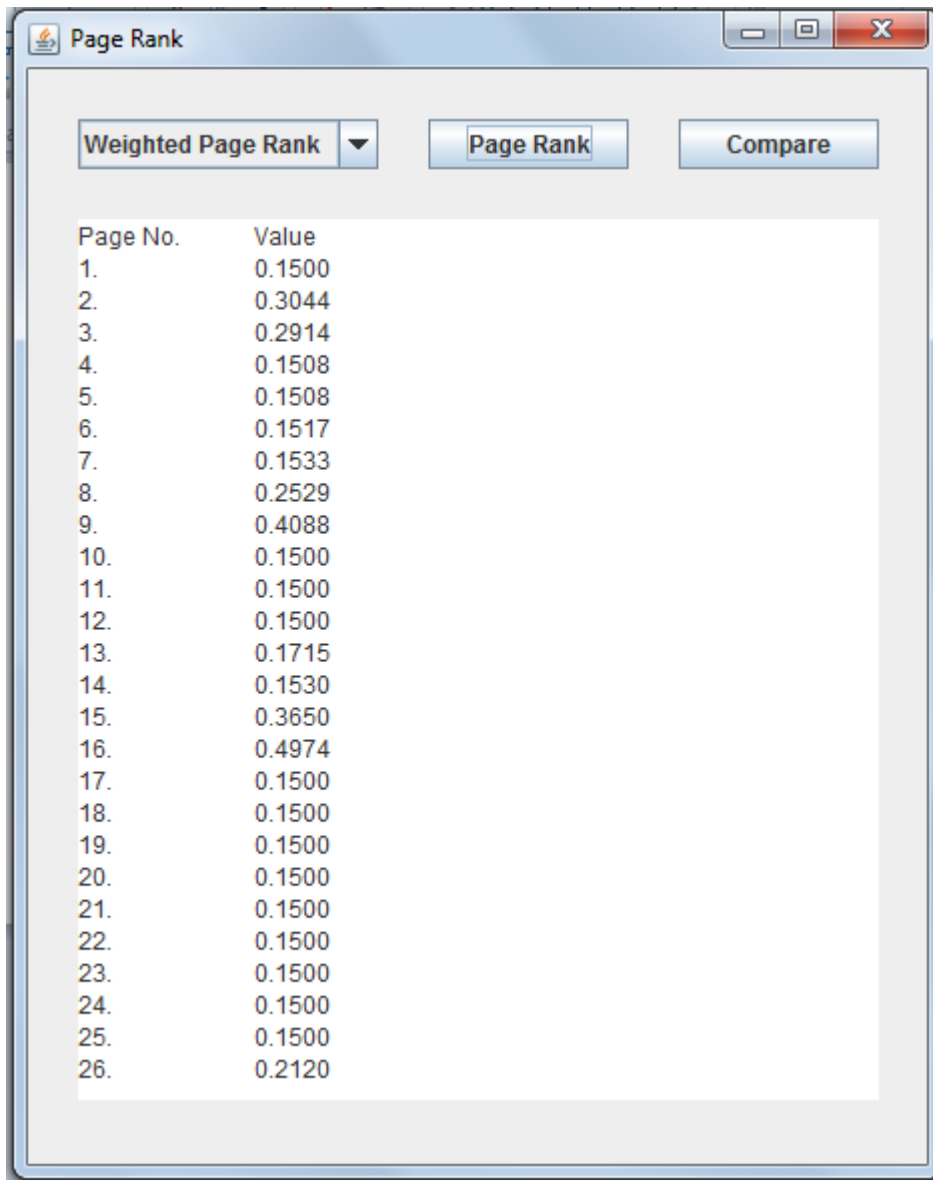


Figure 5.1 Weighted Page Rank Values

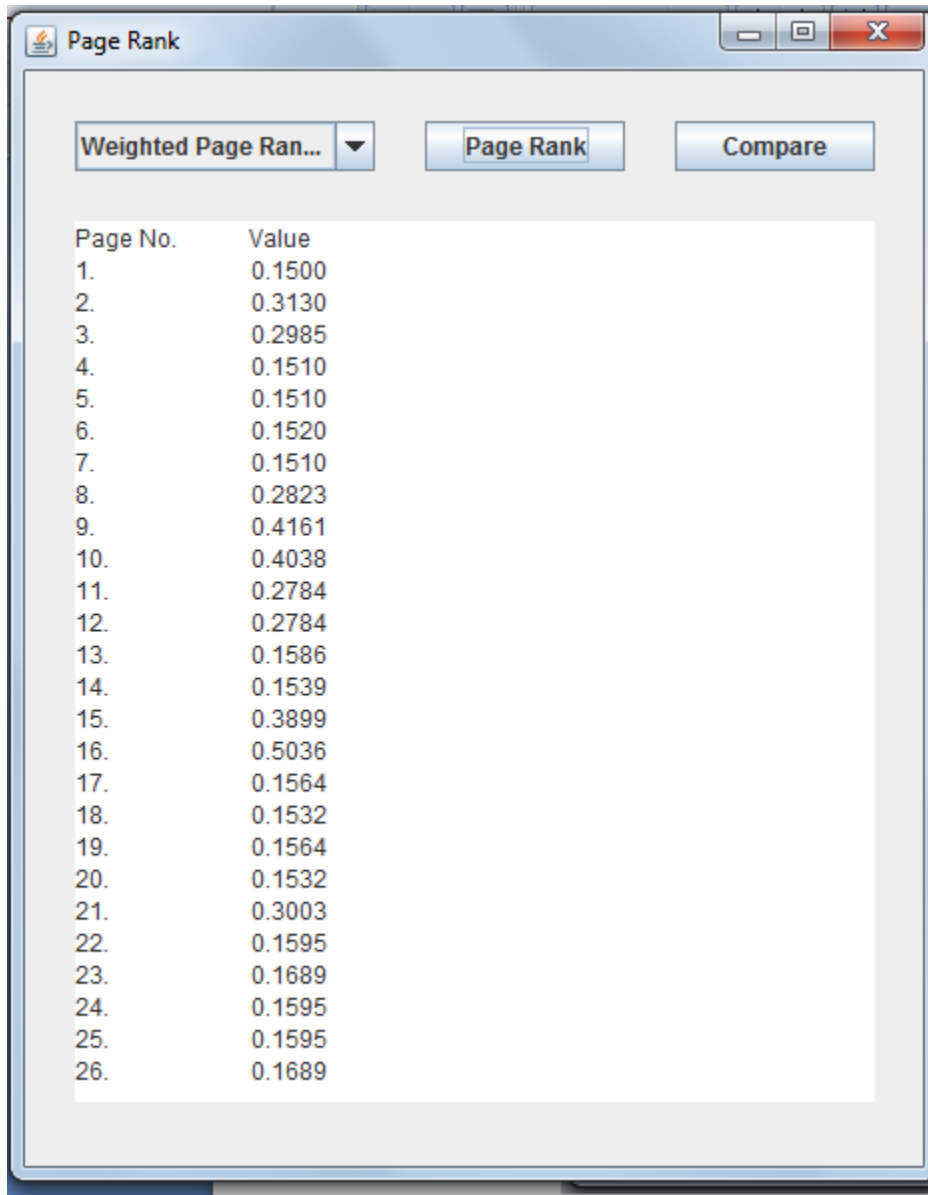


Figure 5.2 Weighted Page Rank using link visit Values

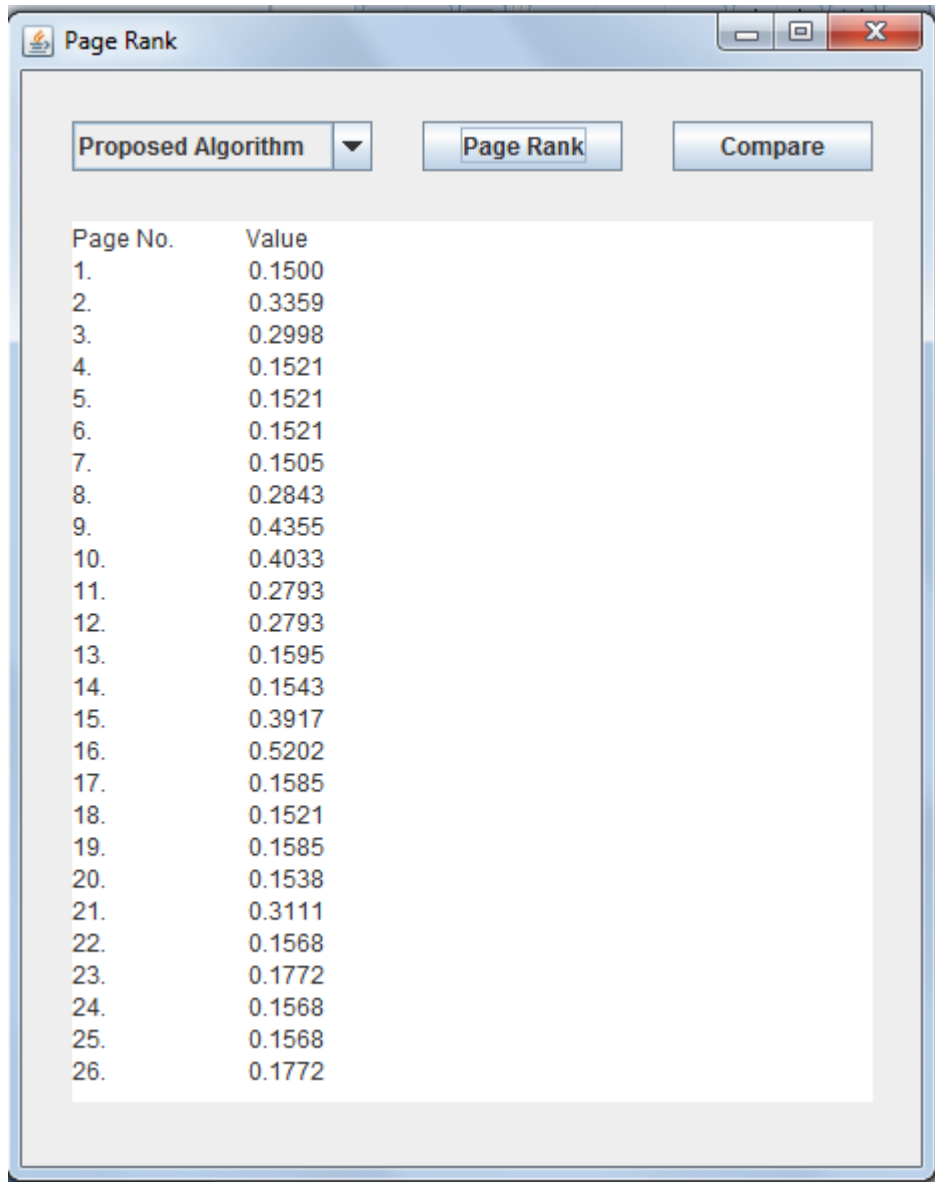


Figure 5.3 Proposed algorithm Values

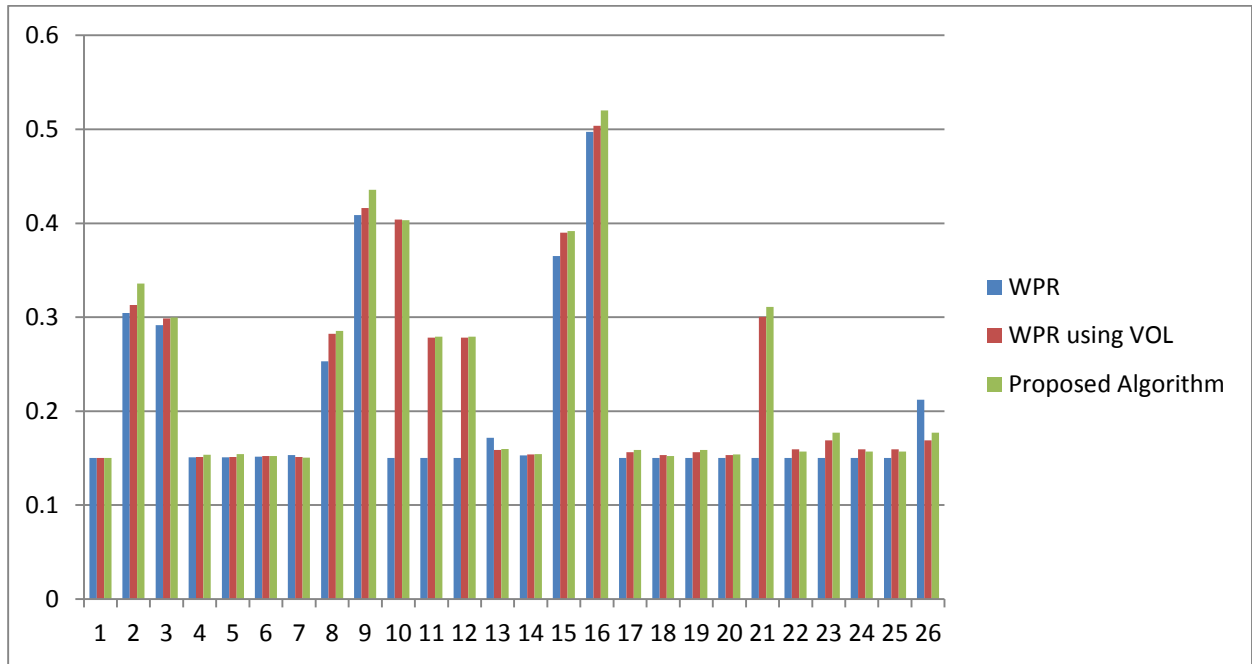


Figure 5.4 Comparison of WPR, WPR using VOL, Proposed Algorithm

CONCLUSION AND FUTURE WORK

Due to enormous amount of information present on the web, the users have to spend lot of time to get pages relevant to them. So it is necessary for search engine to sort the resultant web pages before presenting to the user. The original *WPR* algorithm calculates the page rank by using only the web graph and it ignores the relevancy of web pages from user's point of view. $W^{\text{in}}(v, u)$, the popularity from the number of inlinks and $W^{\text{out}}(v, u)$, the popularity from the number of outlinks does not include usage trends. It does not give more popularity to the links most visited by the users. The weighted PageRank using VOL makes use of web structure mining and web usage mining. The proposed algorithm $EWPR_{VOL}$ makes use of number of visits of links (VOL) to calculate the values of page rank so that more relevant results are retrieved first. In this way, it may help users to get the relevant information quickly. Some of the future works for the proposed algorithm are:

1. The values of page rank have been calculated on a small web graph only. A web graph with large number of websites and hyperlinks should be used to check the accuracy and importance of method.
2. The graph and number of visits of links (VOL) used in the validation of the algorithm are done on small graph. The algorithm needs to be validated with big set of data so that the actual relevancy of the method could be evaluated.
3. The number of visits of links (VOL) only is not strong enough to determine the value of page rank. Some other measures like most recent use of link, information about the user and time spent on web page corresponding to a link can also be used to calculate the value of page rank. So the future work includes deriving a formula for page rank using these parameters also.

REFERENCES

- [1] Brin, Sergey and Page, Lawrence, "The Anatomy of a Large-Scale Hypertextual Web Search Engine. In Seventh International World-Wide Web Conference (WWW 1998), 14-18 April, 1998, Brisbane, Australia.
- [2] Companion slides for the text by Dr. M. H. Dunham, "Data Mining: Introductory and Advanced Topics", Prentice Hall, 2002.
- [3] C. Ding, X. He, P. Husbands, H. Zha, and H. Simon, "Link Analysis: Hubs and Authorities on the World". Technical report:47847, 2001.
- [4] Dell Zhang, Yisheng Dong, "A novel Web usage mining approach for search engines", *Computer Networks* 39 (2002) 303–310
- [5] D. Cohn and H. Chang, "Learning to Probabilistically identify Authoritative Documents". In *Proceedings of 17th International Conf. on Machine Learning*, pages 167-174. Morgan Kaufmann, San Francisco, CA, 2000.
- [6] Gyanendra Kumar, Neelam Duhan, A. K. Sharma, "Page Ranking Based on Number of Visits of Links of Web Page", Department of Computer Engineering, YMCA University of Science & Technology, Faridabad, India.
- [7] Jaroslav Pokorny, Jozef Smizansky, "Page Content Rank: An Approach to the Web Content Mining".
- [8] Jinguang Liu & Roopa Datla, "Web Usage Mining – Pattern Discovery and its application"
- [9] Kosala, Raymond; Hendrik Blockeel, "Web Mining Research: A Survey". *SIGKDD Explorations* 2 (1) July 2000).
- [10] Longzhuang Li, Yi Shang, and Wei Zhang, "Improvement of HITS- based Algorithms on Web Documents", WWW2002, May 7-11, 2002, Honolulu, Hawaii, USA. ACM 1-58113-449-5/02/0005.

- [11] Naresh Barsagade, "Web Usage Mining And Pattern Discovery: A Survey Paper", CSE 8331, Dec.8, 2003.
- [12] Neelam Duhan, A. K. Sharma, Komal Kumar Bhatia, "Page Ranking Algorithms: A Survey" Advance Computing Conference, 2009. IACC 2009 IEEE International.
- [13] Neelam Tyagi, Simple Sharma, "Weighted Page Rank Algorithm Based on Number of Visits of Links of Web Page", International Journal of Soft Computing and Engineering (IJSCE) ISSN: 2231-2307, Volume-2, Issue-3, July 2012.
- [14] R.Cooley, B.Mobasher and J.Srivastava,"Web Mining: Information and Pattern Discovery on the World Wide Web". In Proceedings of the 9th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'97), 1997.
- [15] Raymond Kosala, Hendrik Blockee, "Web Mining Research: A Survey", ACM Sigkdd Explorations Newsletter, June 2000, Volume 2.
- [16] Saeko Nomura, Satoshi Oyama, Tetsuo Hayamizu, "Analysis and Improvement of HITS Algorithm for Detecting Web Communities".
- [17] Salton G. and Buckley, C., "Weighting Approaches in Automatic Text Retrieval". In Information Processing and Management, 1998, Vol. 24, pp. 513-523.
- [18] Wenpu Xing and Ali Ghorbani, "Weighted PageRank Algorithm", Faculty of Computer Science, University of New Brunswick, Fredericton, NB, E3B 5A3, Canada.
- [19] Zdravko Markov and Daniel T. Larose, "Mining the Web: Uncovering Patterns in Web Content, Structure, and Usage Data". Copyright 2007 John Wiley & Sons, Inc.
- [20] http://en.wikipedia.org/wiki/Web_search_engine
- [21] http://en.wikipedia.org/wiki/World_Wide_Web
- [22] http://en.wikipedia.org/wiki/Web_mining