# Evolution Of Web Log Mining Projected On Improved
# Fuzzy C-Means Clustering Algorithm

P.Pranitha*, M.A.H Farquad# and Dr.G.Narsimha$
*Reseach Scholar, CSE , JNTUH Hyderabad, Telangana State
#Associate Professor, CSE, SRIIT Hyderabad, Telangana State
$Associate Professor, CSE, NTUH Hyderabad, Telangana State

## Abstract

Web usage mining is the method of mining valuable usage patterns as of the web data. Web personalization uses web usage mining method for the progression of knowledge attainment done by scrutinizing the user directional patterns awareness. Nowadays, the Web is an imperative source of information retrieval, and the users accessing the Web are from different families. The usage information about users is verified in web logs. Studying web log files to extract useful patterns is named Web Usage Mining. Web usage mining approaches consist of clustering, association rule mining, sequential pattern mining etc. The web usage mining approaches can be pragmatic to forecast next page access. As the size of cluster increases due to the rise in web users, it will become predictable need to augment the clusters. This paper proposes a cluster optimization methodology based on fuzzy logic and is used to reduce the redundancy for clustering Fuzzy C-Means (FCM) algorithm is used. Fuzzy cluster hunting algorithm for cluster optimization is used to personalize web page clusters of end users. Clustering is a data mining technique of grouping set of data objects into multiple groups or clusters so that objects inside the cluster have high similarity, but are very disparate to the objects inthe other clusters. Fuzzy C-Means is the most commonly used method where an element may have partial membership ratings in more than one fuzzy cluster. This investigation work makes use of MATLAB language to yield a fuzzy clustering algorithm for a URL database into several numbers of clusters. The clusters as well as the membership function has been implemented using MATLAB. The results obtained from the database detect n-clusters to handle the inaccurate and abstruse result. Future research work deliver a relative analysis of K-Means, Fuzzy C-Means and Improved Fuzzy C-Means clustering techniques that provide appropriate and accurate data analysis in the field of web log mining.

Index terms : WebLog Mining, K-Means algorithm, Fuzzy C-Meansalgorithm,ImprovedFuzzyC-Meansalgorithm

## 1. Introduction

Data mining is the process of analysing data from different angles and summarizing it in useful information. Technically, data mining is the process of creating relationships or patterns between dozens of fields in large relational databases. Mining on the internet is the application of mining data technologies to discover the patterns or trends that the user follows from the web. It is required because only a small fraction of the information on the web is convenient and gives the user what he wants from the web mining is required because the information stored on the web around the world grows fast and gives the user what he wants is very important. There are three main thrust areas of web mining. The styles that users follow through these three technologies are evaluated in the Web Ming, and these styles are analysed to get the desired output from the user. The desired output is then fed into the understandable graphical user interface.

The World Wide Web is a repository of information. It is used by the user to obtain the required information required through inquiries. Sometimes the user may not be satisfied with the response provided. This may be as user-requested pages that are not indexed because they are not indexed, and are not returned in response to the user's search request. In order to increase users ' satisfaction with requests made on the Internet, we need a new technology that enables the user to obtain the required information easily, efficiently and correctly, which removes the required information easily within a few seconds. This is called extraction of information on the Internet or the World Wide Web mining. It is a data mining technology on the World Wide Web. Mining on the web has three main areas of payment: Web usage Mining, Web content Mining, web structure mining.

Web-based mining is the extraction of web logs to detect access patterns for the pages that the user enters. Analysing the regularity of web logs can help us identify potential e-commerce customers, help personalize Web pages, and improve server performance. The Web server saves all the page entries that are accessed in the web logs. It includes the required URL, IP address, and timestamp. These log files can also be created on the client and proxy. Web history databases provide rich information about Web dynamics, which is why it is important to develop a technology that helps us remove Web log databases.

Web Log Mining is part of the Web Mining, which in turn is part of Data Mining. Since data extraction involves the concept of extracting valuable and valuable information from a large volume of data, the mining of Web Log involves extracting the usage characteristics of Web application users. This extracted information can then be used in a variety of ways, such as improved application, verification of fraudulent elements, etc. The use of Web Log Mining is often considered part of business intelligence in an organization rather than in the technical side. It is used to report business strategies through the effective use of Web applications. It is also critical for customer relationship management (CRM) as it ensures customer satisfaction as far as client-organization interaction is concerned. The main problem with Web Mining in general and Web Log Mining in particular is the nature of the data they are dealing with. With the mounting Internet in this millennium, Web data has become huge in nature and many transactions and uses occur in seconds. Regardless of the size of the data, the data is not fully organized. It's in a semi-structured format so that it needs a lot.

In recent days, the World Wide Web has become the most popular communication medium for Internet users. In this, a huge amount of web pages, links and data size are added daily [1, 7]. Such things make it difficult for a Web user to know the target information on the web search engines. The reason is that the web search engine gives thousands of results to a particular search query [2]. The traditional search engine uses complex algorithms to index documents for the Web user to retrieve relevant information, but has something difficult [9]. Personalized search is a type of search technology that delivers good results and retrieves information for each user according to his/her interests [10]. It differs from general web search and provides identical results for all users of identical search requests. Custom search on the Web can be categorized into two types, such as record-based clicks and personal file-based methods. In the clicks Log-based methods, which enforce bias on the pages clicked in the user's search request log. In profile-based methods improves the search experience with the Tex models.

When you give the search engine return results for the user query and the most appropriate Web page priorities in the screen that is performed by the rating function [12]. In the general search engine involves three processing steps such as crawling, index, and ranking. A crawler visits a webpage that creates graphs on the web. Crawler is used to retrieve Web pages and Web content. After the process is assembled, the content of each page is analysed to determine how it is indexed. The indexer is used to store and index information on retrieved pages. In the ranking stage, many Web page matching cases are arranged together according to user requests or preferences. Measures the importance of webpages returned [4, 8]. In 1990, text Search is used only from a user query to display the result by the search engine. However, because of the growing growth of the Web, this method is not satisfied with the user's needs. So the Web page order has been presented. In 1998, the founder of Google Search, Larry Page and Sergey Brin, invented the page-ordering algorithm for millions of web pages. Evaluates the quality and quantity of links to a Web page.

Web page rankings are associated with many challenges, such as some Web pages set up for navigation and other Web pages only, with a self-tagging quality. Some common pagination algorithms are

page rank algorithm [2, 3], weighted page rank algorithm and Hyperlink induced object Search algorithm [5]. For the results of related Web pages, it is necessary to provide effective algorithms for page rank [13].

## 2.  Web Log Mining

In terms of data mining, a Web history miner is the task of applying data extraction techniques to discover the usage patterns of web data in order to better identify and deliver user requirements on the web. As a data mining task, the Web-use mining process also consists of three main steps: (1) pre-processing (2) pattern detection and (3) pattern analysis. In discovering this business style, it means applying the declared recurring style detection methods to the record data. For this reason, data must be converted in pre-processing stages so that the conversion result can be used as an input for algorithms. Analysing the pattern means understanding the results obtained from the algorithms and imaging conclusions.

As can be noted, the input of the procedure is the registry data. The data must be processed in advance in order to obtain appropriate inputs for the mining algorithms. Different methods require different input formats, so the pre-processing phase can provide three types of output data. The discovery phase of recurring discoveries only needs Web pages visited by a particular user. In this case, the page sequences are irrelevant. Duplicates are also deleted from the same pages, and the pages are arranged in a predefined order. However, in the case of sequential mining, the unique order of the pages is also important, and if a page is visited more than once by a particular user in a time period specified by the user, it will also be relevant. This is why the full system pre-processing module provides Web page ratings by users or user sessions. For sub-tree mining, it is not only classifications, but also web page structure that a particular user visits. In this case, the back-direction discovery is deleted. Only related front-end moves, which make up a tree for each user.  After the discovery is achieved, the analysis of the patterns follows. Depending on the results of the analysis, you can either adjust the pre-processing step parameters (that is, by choosing another time period to determine user sessions) or only the parameters of the mining algorithms. The results obtained by the application can be used to form an acceptable portal structure for advertising reasons and to provide a more personalized web portal.

## 2.1 Clustering Analysis

Grouping analysis is the identification of the clusters that are confined to the data, where the group is a collection of data Objects  "Similar" to each other. Similarity can be expressed by distance functions, specified by users or experts. A good aggregation method produces high intelligence sets to ensure that the similarity between the groups is low and that the similarity within the group is high.

For example, one can group groups in a region according to their community and geographic locations. The closest nearby prediction algorithm that is simply mentioned is as follows:  "Nearby Objects " will also include similar forecast values. So, if you know the prediction value of an object, you can predict it from the nearest neighbour.
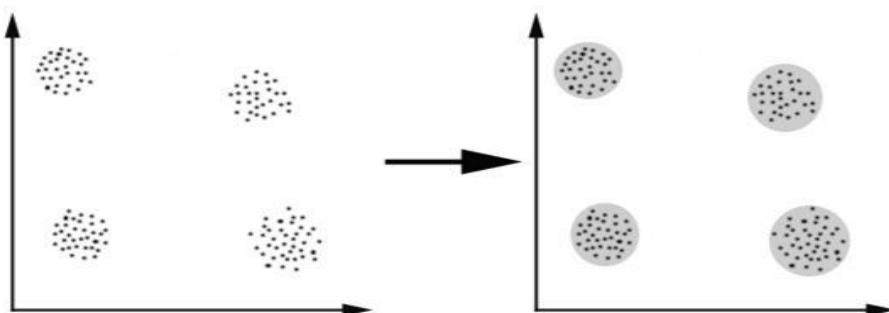


**Fig 2.1: showingfour clusters formed from the set of unlabelled data**

## 2.2 K-MeansAlgorithm

Given the database of the N and K elements, the number of clusters to form the partitioning algorithm organizes objects in the K partitions where each section refers to a group. These clusters are designed to improve the target partitioning standard, often called the similarity function, such as distance, so that objects within the group are  "symmetric ", while non-related group objects are  "differentiated " In terms of the attributes of the databases.

**Algorithm:**

The K-means algorithm is a split based on the average value of objects in the cluster.

Input: Number of clusters K and database include n objects.

Output: A set of k sets that reduce the standard error box.

**Method:**

1. The dataset is split into the K group and the dataset is set arbitrarily and manually or based on some indicative methods.

2. Each data point calculates the distance from the data point to each group. If Datapoint is adjacent to its own collection, leave it. If the data point is not close to its collection, move it to the nearby group.

3. Repeat the above step so that wide scrolling through all data points does not move any data point from one group to another. At this point, the cluster is stable and the assembly process ends.

$$J(V)=\sum_{j=1}^{C} \sum_{j=1}^{C_i}(\|m_i - v_j\|)^2$$

where,

$'\|m_i - v_j\|'$ is the Euclidean distance between $x_i$ and $v_j$.

$'C_i'$ is the number of data points in $i^{th}$ cluster.

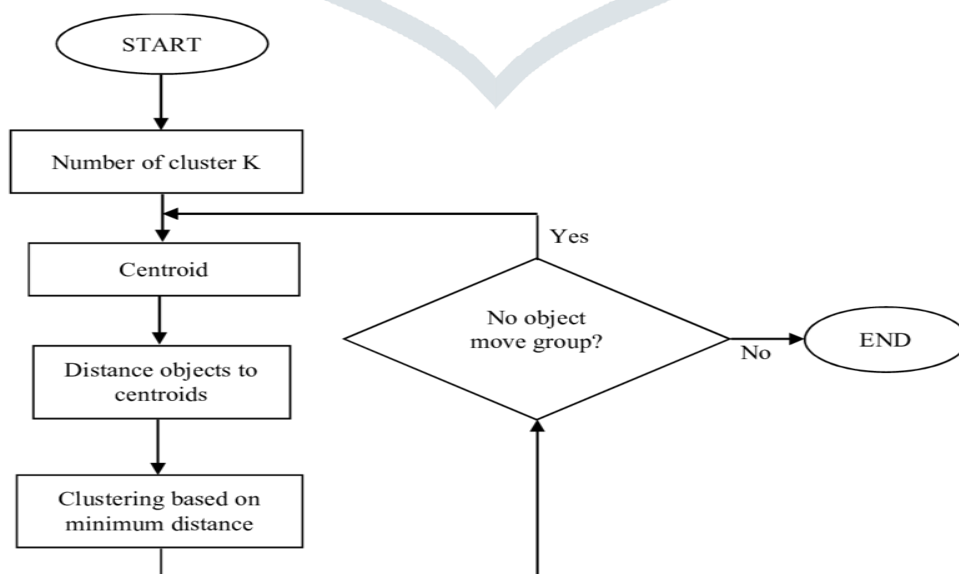$'C'$ is the number of cluster centers.



**Fig 2.2 Flow Diagram of K-Means algorithm**

The K-means algorithm takes the input parameter, K, and partitions a set of N objects into k groups so that they result in high intra-cluster similarity and low inter-cluster entanglement. The similarity of the group is measured in relation to the average value of the objects in a group, which can see the elimination of the gravitational centre of the mass.

## 2.3 Clustering features using Weighted Fuzzy C-Means

The W-FCM algorithm will display an iterative improvement similar to FCM, and as such it is amplified by some of its strengths, such as its convergence in a limited number of iterations, and its vulnerabilities, such as the algorithm, however, setting the c center randomly, does not guarantee the optimal solution. Also, the same applies to weights, which creates the possibility to be far from representing the link structures.

By seeing feature importance WFCM algorithm is specified as follows,

$$J(W,U,V) = \sum_{k=1}^{c} \sum_{i=1}^{n} u_{k,i}^{2} (d_{k,i}^{(W)})^{2} \qquad (1)$$

Where $d_{k,i}^{(W)}$ is calculated by:

$$d_{k,i}^{(W)} = \sqrt{\sum_{j=1}^{d} W^{\beta} (x_{i,j} - v_{k,j})^{2}} \qquad (2)$$

Minimizing eqn (3) $v_k$ and $u_{k,i}$ is specified as follows,

$$v_k = \sum_{i=1}^{n} (u_{k,i})^{m} x_i / \sum_{i=1}^{n} (u_{k,i})^{m}, \ \forall k = 1,\ldots c \qquad (3)$$

$$u_{k,i} = 1 / \sum_{k=1}^{c} (d_{k,i}^{(w)} / d_{k,\bar{k}}^{(w)})^{2/m-1} \qquad (4)$$

The main steps of weighted fuzzy c means clustering algorithm (W-FCM) is shown as follows :

- Fix the number of clusters $c$, where $2 \le c \le n$ and initialize the fuzzy partition matrix $U$ with a random value such that it satisfies situations.

$$\sum_{k=1}^{c} u_{k,j} = 1 \forall j \ and \ 0 < \sum_{i=1}^{n} u_{k,i} < n \forall k \qquad (5)$$

- Initialize the weighting vector $W$ with a random value such that it ones (2) and (3)
- Calculate the fuzzy centers $v_k$ using (6)
- Update the fuzzy partition matrix $U$ with (7)
- Update the weighting vector $W$

$$W_j = \{0 \ if \ D_j = 0 \ and \ 1 / \sum_{t=1}^{h} [D_j / D_t]^{1/B-1} \ if \ D_j \neq 0 \qquad (6)$$

$$D_j = \sum_{k=1}^{k} \sum_{i=1}^{n} u_{k,i} (x_{i,j} - v_{k,j})^{2} \qquad (7)$$

- Repeat the steps until one of the termination criterions is satisfied
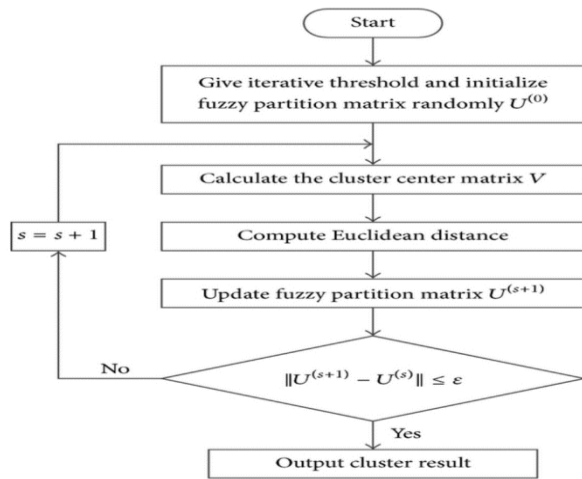
**Fig 2.3:  Flow Diagram of Fuzzy C-Means algorithm**

In order to find the data quality in the feature value for each group that is applied to the assembly algorithm to create a vector for the user query and find a space between the vector users, and then to specify the minimum raster space. Finally, WOA is included to optimize the best number of pages.

# 3. Proposed Solution

## 3.1 Proposed Cluster Analysis Using Improved FCM

A improved FCM algorithm is set up, in which the Data Set and group are set up, and the variable will be selected in general. Enhanced FCM is used to form the cluster group by capturing some dataset. The tolerance of mass is calculated by finding the distance measurement. The enhanced FCM method area is applied in a large volume of data and gives real prediction. The target function is minimized in an enhanced fuzzy c-mean algorithm. The covered function is less than the other by using the distance formula. The results are optimized by a Hermitian distance. gives a real prediction. Enhanced FCM is one of the procedures that is used to extract data for the assembly.

Improved FCM algorithm by using two steps, first through a decision tree approach that modifies data in a wrong and sequential way and second by creating noise-free data. Bulk power functionality is often used to evaluate grouping performance in a different index, and even two different assembly methods. Among the criteria there are inevitable types of FCM in which the sample is set on the basis of the fuzzy partition. The basic idea of a valid function based on a mysterious section is that the lowest partitioning divisions is the best performance. An enhanced FCM algorithm is developed and applied to manipulate the data set at the top of the traditional FCM algorithm. The proposed algorithm has improved the classic FCM algorithm by approving a different strategy for selecting the first cluster center to solve the FCM problem. This algorithm is modified FCM called the enhanced fuzzy C-mean  algorithm based on the first cluster selection and initial membership value. If a good primary cluster center is found, it is close to the actual final block center. The pyramid distance is an extension method, and this technique is expected to perform better than the other methods.

Enhanced FCM is used to form the cluster group by satisfying some dataset. The quality of the cluster is calculated by detecting the distance scale. The suggested FCM is used to form the cluster with a smaller number of iterations. Data sets are taken and the cluster is set up, and a variable is usually randomly chosen but here is a tension to avoid the choice of the variable arbitrarily. Essentially, a fuzzy assembly is appropriate to remove a complex, multidimensional dataset. Where the member is incomplete or unclear. Among the changing variable technology, the FCM algorithm is the most general where the data contains an incomplete part of the membership with each of the pre-defined clustercenters.

## 3.2 Algorithm For Proposed work

An improved fuzzy C-mean algorithm is developed and applied to manipulate the data set on top of a traditional fuzzy C-mean algorithm.

1. First, the initial blurred partitioning matrix is created and the ambiguous Initial Block Center (C) is calculated.

2. At each step of the iteration the cluster center and the membership degree point are updated and the objective function is minimized to find the best mode for the group.

3. Improved FCM technology is expected. It is used to resolve the minimum distance by using the Hermitian distance method formula:

$$X=[(x_1 - y_1), (x_2 - y_2), (x_3 - y_3) \ldots \ldots (x_n - y_n)]$$

$$D = (x * x^t)^2$$

Where **X** is a matrix and **X$^t$** is a transpose of matrix .

4. The process stops when the maximum number of iterations is reached or when the target function improvement between two consecutive iterations is less than the minimum data specified.



**Fig 3.1: Flow Diagram of Improved Fuzzy C-Means algorithm**

The result is more detailed, accurate, better improved and focused on improvement. It would actually be applicable to technical similarity. The update is performed in the iteration by using the membership degree plus the center of the cluster that represents the change of the parameters where the steps are duplicated until a specified point is reached, called the maximum or the action stops at the maximum number of iterations being accessed or when the target function is improved between Two consecutive iterations less than the least amount of development specified.

## 4. Result Analysis

## 4.1 Evaluation parameters

DONE Index: The evaluation and parameter of the quality index is measuring the unification of clusters. The goal is to identify the dense and well-closed mass. It is defined as the ratio between the minimum distances between groups to the maximum distance within the cluster for each partition. Specifies the cluster consolidation.

DBIndex:It is an estimation parameter whose value belongs to 0-1. It is also known as an surprising index coefficient. It summarizes the similarity of the data. The algorithm that produces the block will be the distance within the low block and the high distance between the groups will contain the Davis-two low-born index.

## 4.2 sample dataset



## 4.3 Comparative Results

| Method | DB Index | DONE Index |
|---|---|---|
| K-Means | 0.25397 | 0.5 |
| Fuzzy C-Means | 0.24282 | 0.5 |
| Improved Fuzzy C-Means | 0.20853 | 0.5 |

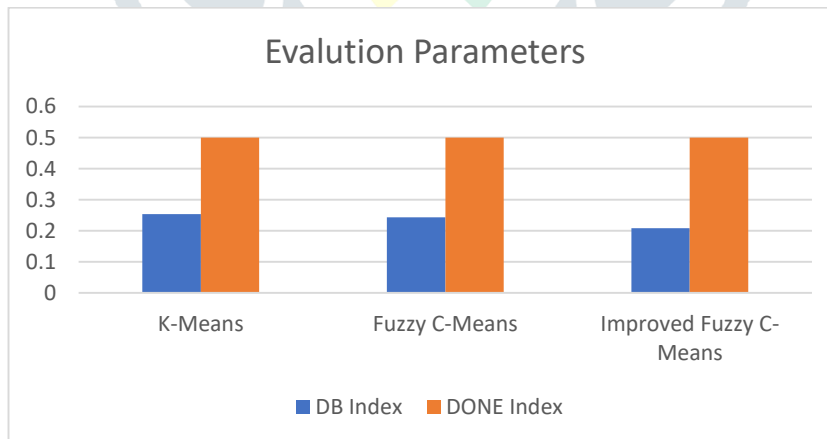**Table 1 : Comparative results of clustering techniques**



**Figure4.1:Simulated study ofproposedworkusing chartoutline**

# 5.  Conclusion

Data collection plays a key role in grouping the corresponding data type into an exact cluster. The objective of the cluster analysis is to classify groups of similar organisms, thus supporting the detection of dispersion in patterns and stimulating links in large data sets. A promiscuous assembly is an extension of the cluster analysis, which represents the correlation of data points to groups by membership. In this search, a virtual analysis was performed for K-means, fuzzy C-means, and an improved fuzzy C-means assembly using the Fuzzy relational database to detect clusters on the data set of the conscious URL. Supports in the database are broken in such a way that similar IP addresses are in the same group. It was found in the table that the improved fuzzy C-means method yields better results in terms of the DB index and the DONE index. MATLAB was used to describe the function of membership, the threshold equation, and the discovery of many groups. Algorithms are created in MATLAB for investigation and comparison. The result produced a higher neutral resolution and required less calculation. The objective of the algorithm is to provide a classification of some well-known aggregation algorithms. In conclusion, it turns out prove that the improved FCM algorithm technique performance better-quality than other methods.

## REFERENCES

[1] Thomas, Tiji K., and K. S. Sudeep. "An improved Page Rank algorithm to handle polysemous queries" In the proceeding of International Conference on In Computing, Analytics and Security Trends (CAST), pp. 106-111, 2016.

[2] Antoniou, Dimitris, Yannis Plegas, Athanasios Tsakalidis, Giannis Tzimas, and EmmanouilViennas. "Dynamic refinement of search engines results utilizing the user intervention " In the proceeding of Journal of Systems and Software 85, no. 7 , 1577-1587, 2012.

 [3] Lu, Yan, Yuanyi Li, Meng Xu, and Weihui Hu. "A user model based ranking method of query results of meta-search engines " In the proceeding of International Conference on In Network and Information Systems for Computers (ICNISC), pp. 426-430 , 2015.

[4] Vispute, Sushma R., Shrikant Patil, Sagar Sangale, Akshay Padwal, and AvinashUkarde. "Parallel Processing System for Marathi Content Generation " ,In the proceeding of 2015 International Conference on In Computing Communication Control and Automation (ICCUBEA), pp. 575-579, 2015.

[5] Vesanen, Jari, and Mika Raulas. "Building bridges for personalization: a process model for marketing." In the proceeding of Journal of Interactive Marketing 20, no. 1, pp. 5-20, 2006.

 [6] Singh, Ranveer, and Dilip Kumar Sharma. "Enhanced-RatioRank: Enhancing impact of inlinks and outlinks ", In the proceeding of Conference on In Information & Communication Technologies (ICT), pp. 287-291, 2013.

 [7] Hassan, Osama Al-Haj, Thamer Al-Rousan, Anas Abu Taleb, and Adi Maaita. "An efficient and scalable ranking technique for mashups involving RSS data sources ", In the proceeding of Journal of Network and Computer Applications 39 , 179-190, 2014.

[8] Chahal, Poonam, Manjeet Singh, and Suresh Kumar. "Ranking of web documents using semantic similarity ", In the proceeding of International Conference on In Information Systems and Computer Networks (ISCON) ,  pp. 145-150 , 2013.

[9] Dhivya, G., K. Deepika, J. Kavitha, and V. Nithya Kumari. "Enriched content mining for web applications ", In the proceeding of International Conference on In Innovations in Information, Embedded and Communication Systems (ICIIECS),  pp. 1-5 , 2015.

[10] Eirinaki, Magdalini, and Michalis Vazirgiannis. "Web mining for web personalization." ACM Transactions on Internet Technology (TOIT) 3, no.1, pp. 1-27, 2003.

[11] Khanchana, R., and M. Punithavalli. "Web page prediction for web personalization: a review", In the proceeding of global Journal of Computer Science and Technology, 2011.

[12] Shukla, Rajesh K., Sanjay Silakari, and P. K. Chande. "Existing Trends and Techniques for Web Personalization ", In the proceeding of International Journal of Computer Science Issues 9, no. 4, 2012.

[13] Hawalah, Ahmad, and Maria Fasli. "Dynamic user profiles for web personalisation." Expert Systems with Applications 42, no. 5,  pp. 2547-2569, 2015.

[14] Chirita, Paul-Alexandru, Daniel Olmedilla, and Wolfgang Nejdl. "Pros: A personalized ranking platform for web search", In Adaptive Hypermedia and Adaptive Web-Based Systems, pp. 431-461, 2004.

[15] Kenneth Shum, "Notes on PageRank Algorithm", In Proceedings of IEEE Conference on EDU, pp. 1-9, 2013.

[16] Bhawiyuga, Adhitya, and AnnisaPuspaKirana. "Implementation of page rank algorithm in Hadoop MapReduce framework", In the proceeding of International Seminar on In Intelligent Technology and Its Applications (ISITIA),  pp. 231-236 , 2016.

[17] No'aman, M., Ahmed M. Gadallah, and Hesham A. Hefny. "A hybrid recommendation model for web navigation ", In the proceeding of Seventh International Conference on In Intelligent Computing and Information Systems (ICICIS), pp. 552-560 , 2015.

[18] Bansal, Divya, Joy Bose, and Ankit Kumar, "EEG based detection of area of interest in a web page", In the proceeding of International Conference on In Advances in Computing, Communications and Informatics (ICACCI),  pp. 1320-1325 , 2015.

[19] Blanco, Roi, Michael Matthews, and Peter Mika, "Ranking of daily deals with concept expansion", In the proceeding of Information Processing & Management 51, no. 4, 359-372, 2015.

[20] Ševa, Jurica, Markus Schatten, and Petra Grd. "Open Directory Project based universal taxonomy for Personalization of Online (Re) sources." In the proceeding of Expert Systems with Applications 42, no. 17, pp. 6306-6314, 2015.

[21] Selvan, Mercy Paul, A. Chandra Shekar, Deepak R. Babu, and A. Krishna Teja. "Efficient ranking based on web page importance and personalized search ", In the proceeding International Conference on In Communications and Signal Processing (ICCSP), pp. 1093-1097, 2015.

[22] Ganeshiya, Deepak Kumar, and Dilip Kumar Sharma. "A survey: hyperlink analysis in webpage ranking algorithms." In the proceeding of International Conference on In Soft Computing Techniques for Engineering and Technology (ICSCTET) ,  pp. 1-8., 2014.