# NOVEL MODIFIED FPCM FOR WEB LOG MINING BY REMOVING GLOBAL NOISE AND WEB ROBOTS

**[1]P.NITHYA, [2] DR.P.SUMATHI**

[1]Doctoral Student in Computer Science, Manonmanaiam Sundaranar University, Tirunelveli
[2]Assistant Professor, PG & Research Department of Computer Science, Govt.Arts College, Coimbatore.
E-mail:  nithi.selva@gmail.com

## ABSTRACT

Nowadays, internet is a useful source of information in everyone's daily activity. Hence, this made a huge development of World Wide Web in its quantity of interchange and its size and difficulty of websites. Web Usage Mining (WUM) is one of the main applications of data mining, artificial intelligence and so on to the web data and forecast the user's visiting behaviors and obtains their interests by investigating the samples. Since WUM directly involves in large range of applications, such as, e-commerce, e-learning, Web analytics, information retrieval etc. Web log data is one of the major sources which contain all the information regarding the users visited links, browsing patterns, time spent on a particular page or link and this information can be used in several applications like adaptive web sites, modified services, customer summary, pre-fetching, generate attractive web sites etc. There are varieties of problems related with the existing web usage mining approaches. Existing web usage mining algorithms suffer from difficulty of practical applicability. So, a novel research is very much necessary for the accurate prediction of future performance of web users with rapid execution time. The main aim of this paper to remove the noise and web robots by novel approach and provide faster and easier data processing and it also helps in saving time and it resource. In this paper, a novel pre-processing technique is proposed by removing local and global noise and web robots. Anonymous Microsoft Web Dataset and MSNBC.com Anonymous Web Dataset are used for evaluating the proposed preprocessing technique. An Effective Web User Analysis and Clustering are analyzed using Modified FPCM. Then results are evaluated using Hit Rate and Execution time.

**Keywords:** *Preprocessing, Data Cleaning. Modified Fuzzy Possibilistic C Means, Fuzzy C Means, Hit Rate, Execution Time.*

## 1. INTRODUCTION

Millions of electronic data are included on hundreds of millions data that are previously on-line today. With this significant increase of existing data on the Internet and because of its fast and disordered growth, the World Wide Web has evolved into a network of data with no proper organizational structure. In addition, survival of plentiful data in the network and the varying and heterogeneous nature of the web, web searching has become a tricky procedure for the majority of the users. This makes the users feel confused and at times lost in overloaded data that persist to enlarge. Moreover, e-business and web marketing are quickly developing and significance of anticipate the requirement of their customers is obvious particularly. As a result, guessing the users' interests for improving the usability of web or so called personalization has turn out to be very

essential. Web personalization can be depicted as some action that builds the web experience of a user personalized according to the user's interest.

Generally, three kinds of information have to be handled in a web site: content, structure and log data. Content data contains of anything is in a web page, structure data is nothing but the organization of the content and usage data is nothing but the usage patterns of web sites. The usage of the data mining process to these dissimilar data sets is based on the three different research directions in the area of web mining: web content mining, web structure mining and web usage mining  (Jalali M., *et al.*, 2008) (Maratea A and Petrosino A, 2009).

Web usage mining (Jian Chen *et al.*, 2004) (Wu, K.L., Yu, P. S. and Ball man, A , 1998) consists of three main steps. Web log file is usually given as input.

Preprocessing is an important step because of the complex nature of the Web architecture which

takes 80% in mining process. The raw data is pretreated to get reliable sessions for efficient mining. It includes the domain dependent tasks of data cleaning, user identification, session identification, and path completion and construction of transactions. Data cleaning is the task of removing irrelevant records that are not necessary for mining. Data cleaning includes

- Elimination of Local and Global Noise,
- Removal of records of graphics, videos and the format information
- Removal of records with the failed HTTP status code
- Robots cleaning

User identification is the process of associating page references with same IP address with different users. Session identification is breaking of a user's page references into user sessions. Path completion is used to fill missing page references in a session. Classifications of transactions are used to know the users interest ad navigational behavior. The second step in web usage mining(Labroche N *et al.*, 2007) is knowledge extraction in which data mining algorithms like association rule mining techniques, clustering, classification etc. are applied in preprocessed data. The third step is pattern analysis in which tools are provided to facilitate the transformation of information into knowledge. Knowledge query mechanism such as SQL is the most common method of pattern analysis. This paper focuses on path completion process which is used to append lost pages and construction of transactions in preprocessing stage. In this study a referrer-based method is proposed to efficiently construct the reliable transactions in data preprocessing.

## 2. RELATED WORKS

The discovery of the users' navigational patterns using SOM is proposed by (Labroche N *et al.*, 2007 ). In (Etminani K. *et al.*, 2009) presented a Web usage mining technique based on fuzzy clustering in Identifying Target Group. In (Jianxi Zhang et al., 2009) suggests a complete idea for the pattern discovery of Web usage mining. In (Nina S.P et al., 2009) given a Web Usage Mining technique based on the sequences of clicking patterns in a grid computing environment. The author discovers the usage of MSCP in a distributed grid computing surroundings and expresses its effectiveness by empirical cases. In (Chih-Hung Wu et al., 2010)

proposed the usage of incremental fuzzy clustering to Web Usage Mining. Rough set based feature selection for web usage mining is proposed by (Aghabozorgi, S.R.and Wah, T.Y ,2009). In (Inbarani H.H , 2007)put forth a web usage mining technique based on LCS algorithm for online predicting recommendation systems. For providing the online prediction effectively, in (Shinde, S.K. and Kulkarni U.V ,2008) provides a architecture for online recommendation for predicting in Web Usage Mining System.

In (Zhang Huiying and Liang Wei, 2004) given an intelligent algorithm of data pre-processing in Web usage mining. In (Nasraoui O et al., 2008) provides a whole framework and findings in mining Web usage navigation from Web log files of a genuine Web site which has every challenging characteristics of real-life Web usage mining, together with evolving user profiles and external data describing an ontology of the Web content. In (Hogo et al., 2003) proposed the temporal Web usage mining of Web users on single educational Web site with the help of the adapted Kohonen SOM based on rough set properties. A development of data preprocessing technique for Web usage mining and the information's of algorithm for path completion are provided by(Yan Li et al., 2008).

(Baraglia, R. and Palmerini, P , 2002) Proposed a Web usage mining (WUM) system, called SUGGEST which continuously creates the suggested connections to Web pages of probable importance for a user. In (Chu-Hui Lee and Yu-Hsiang Fu, 2008) put forth a Web Usage Mining technique based on clustering of browsing characteristics.

## 3. METHODOLOGY

Web log data preprocessing is a complex process and takes 80% of total mining process. Log data is pretreated to get reliable data. The aim of data preprocessing is to select essential features clean data by removing irrelevant records and finally transform raw data into sessions.

### 3.1. Data Cleaning

The process of data cleaning is removal of outliers or irrelevant data. Analyzing the huge amounts of records in server logs is a cumbersome activity. So initial cleaning is necessary. If a user requests a specific page from server entries like gif, JPEG, etc., are also downloaded which are not useful for further analysis are eliminated. The

records with failed status code are also eliminated from logs. Automated programs like web robots, spiders and crawlers are also to be removed from log files. Thus removal process in the experiment includes

### 3.1.1. Elimination Of Local And Global Noise

Web noise can be normally categorized into two groups depending on their granularities:

Global Noise: It corresponds to the unnecessary objects with huge granularities, which are no smaller than individual pages. This noise includes mirror sites, duplicated Web pages and previous versioned Web pages, etc.

Local (Intra-Page) Noise: It corresponds to the irrelevant items inside a Web page. Local noise is typically incoherent with the major content of the page. This noise includes banner ads, navigational guides, decoration pictures, etc. These noises have to remove for better results.

### 3.1.2. The Records Of Graphics, Videos And The Format Information

The records have filename extension of GIF, JPEG, CSS, and so on, which can be found in the URI field of the every record, can be removed. This extension files are not actually the user interested web page, rather it is just the documents embedded in the web page. So it is not necessary to include in identifying the user interested web pages. This cleaning process helps in discarding unnecessary evaluation and also helps in fast identification of user interested patterns.

### 3.1.3. The Records With The Failed HTTP Status Code

The HTTP status code is then considered in the next process for cleaning. By examining the status field of every record in the web access log, the records with status codes over 299 or under 200 are removed. This cleaning process will further reduce the evaluation time for determining the used interested patterns.

### 3.1.4. Method Field

It should be pointed out that different from most other researches, records having value of POST or HEAD in Method field are reserved in present study for acquiring more accurate referrer information.

### 3.1.5. Robots Cleaning

Web robot (WR) (also called spider or bot) is a software tool that periodically scans a web site to extract its content. Web robots automatically follow all the hyperlinks from a web page. Search engines, such as Google, periodically use WRs to gather all the pages from a web site in order to update their search indexes. The number of requests from one WR may be equal to the number of the web site's URIs. If the web site does not attract many visitors, the number of requests coming from all the WRs that have visited the site might exceed that of human-generated requests.

Eliminating WR-generated log entries not only simplifies the mining task that will follow, but it also removes uninteresting sessions from the log file. Usually, a WR has a breadth (or depth) first search strategy and follows all the links from a web page. Therefore, a WR will generate a huge number of requests on a web site. Moreover, the requests of a WR are out of the analysis scope, as the analyst is interested in discovering knowledge about users' behavior.

Most of the Web robots identify themselves by using the user agent field from the log file. Several databases referencing the known robots are maintained [Kos, ABC]. However, these databases are not exhaustive as each day new WRs appear or are being renamed, making the WR identification task more difficult.

To identify web robots' requests, the data cleaning module implements two different techniques.

- In the first technique, all records containing the name "robots.txt" in the requested resource name (URL) are identified and straightly removed.
- The next technique is based on the fact that the crawlers retrieve pages in an automatic and exhaustive manner, so they are distinguished by a very high browsing speed. Therefore, for each different IP address, the browsing speed is calculated and all requests with this value more than a threshold are regarded as made by robots and are consequently removed. The value of the threshold is set up by analyzing the browser behavior arising from the considered log files.

This helps in accurate detection of user interested patterns by providing only the relevant web logs. Only the patterns that are much interested

by the user will be resulted in the final phase of identification if this cleaning process is performed before start identifying the user interested patterns.

## 3.2. The Fuzzy Clustering Algorithm

The basic theory about Fuzzy Clustering analysis based on Fuzzy Equivalent Relation is that Fuzzy Corresponding R(fuzzy relation) is a subset of U×U (universe of fuzzy set), when we are using λ-threshold, the subset of U×U will be just an ordinary equivalent relation and the objects in the U will be classified. When λ is reduced from 1 to 0, gradually the classification will merge forming a dynamic clustering dendritic diagram consequently. Thus, the establishment of fuzzy relation R is a key step of fuzzy analysis method.

Fuzzy Clustering is an iterative algorithm. The aim of Fuzzy Cluster is to find cluster centers (centroid) that minimize a dissimilarity function. This algorithm works by assigning membership to each data point corresponding to each cluster center on the basis of distance between the cluster center and the data point (J. Bezdek , 1981).

### 3.2.1. The Session Identification

User session identification is the process of analyzing usage data in order to extract useful information concerning user navigational behavior by structuring the requests contained into the Web log files. Access log files of a Web site consist in text files where the server stores all the accesses made by the users in chronological order. According to the Common Log Format, each log entry includes: the user's IP address, the request's date and time, the request method, the URL of the accessed page, the data transmission protocol, the return code indicating the status of the request, the size of the visited page in terms of number of bytes transmitted. Based on such information, we can determine the user sessions, i.e. the sequence of URLs that each user has accessed during his/her visit.

One record in web access log is written as: 219.144.222.253 - - [16/Aug/2004:15:36:11 +0800] "GET /images/1 r3 c2.jpg HTTP/1.1" 200 418 "http://202.117.16.119:8089/index.html" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1)". The meaning of every field is explained in Table 1

Here, a user session is defined as the finite set of URLs accessed by a user within a predefined time period (in our work, 25 minutes). Since the

information about the user login is not available, user sessions are identified by grouping the requests originating from the same IP address during the established time period. Finally, data are filtered in order to retain only the most relevant pages and user sessions. At the end of preprocessing, we obtain collection of $n_s$ sessions denoted by the set $S = \{s_1, s_2, \ldots \ldots s_{ns}\}$. Each session contains information about accesses to pages during the session time. Precisely, a user session is formally described as a triple $S_i = (u_i, t_i, p_i)$ where $u_i$ represents the user identifier, $t_i$ is the access time of the whole session, $p_i$ is the set of all pages (with corresponding access information) requested during the i-th session. Namely:
$$P^i = \{(p_{i1}, t_{i1}, N_{i1}), (p_{i2}, t_{i2}, N_{i2}), \ldots \ldots (p_{i,n_i}, t_{1n_i}, N_{1n_1})\}$$
With $P_{(ij)} \in P$ , where $N_{ij}$ is the number of accesses to page $P_{ij}$ during the i-th session and $t_{ij}$ is the total time spent by the user on that page during the i-th session.

## 3.3. Fuzzy Possibilistic C Means Algorithm

It is a method of clustering which allows one piece of data to belong to two or more clusters. This method was developed by Dunn in 1973(J.C. Dunn, 1973) and Modified by Bezdek in 1981 (J.C. Dunn, 1973)( J. Bezdek , 1981) and this is frequently used in pattern recognition. FPCM produces memberships and possibilities simultaneously, along with the usual point prototypes or cluster centers for each cluster. FPCM is a hybridization of possibilistic c-means (PCM) and fuzzy c-means (FCM) that often avoids various problems of PCM and FCM.

FPCM solves the noise sensitivity defect of FCM, overcomes the coincident clusters problem of PCM.The choice of an appropriate objective function is the key to the success of the cluster analysis and to obtain better quality clustering results; so the clustering optimization is based on objective function. To meet a suitable objective function, we started from the following set of requirements: FPCM algorithm merges the advantages of both fuzzy and Possibilistic c-means techniques. Memberships and typicalities are essential for the accurate characteristic of data substructure in clustering technique.

$$J_{(FPCM)}(U,T,C) = \sum_{(i=1)}^{c} \sum_{(j=1)}^{n} \left( (u)_{ij}^m + (t_{(ij)})^{\eta}(2) \right) (X_{(j)} - )^2_{ij}$$
*(1)*

With the following constraints:

$$\sum_{i=1}^{c} u_{ij} = 1, \forall j \in \{1, \dots, n\} \qquad (2)$$

$$\sum_{j=1}^{n} t_{ij} = 1, \forall i \in \{1, \dots, c\} \qquad (3)$$

A solution of the objective function can be obtained through an iterative process where the degrees of membership, typicality and the cluster centers are updated with the equations as follows.

$$u_{ij} = \left[ \sum_{k=1}^{c} \left( \frac{d^2(x_j, v_i)}{d^2(x_j, v_k)} \right)^{2/(m-1)} \right]^{-1}, 1 \le i ; \quad (4)$$

$$t_{ij} = \left[ \sum_{k=1}^{n} \left( \frac{d^2(x_j, v_i)}{d^2(x_j, v_k)} \right)^{2/(\eta-1)} \right]^{-1}, 1 \le i \le \quad (5)$$

$$v_i = \frac{\sum_{k=1}^{n} [(u]_{ik}^m + t_{ik}^\eta) x_K}{\sum_{k=1}^{n} [(u]_{ik}^m + t_{ik}^\eta)}, 1 \le i \le c. \quad (6)$$

FPCM constructs memberships and possibilities simultaneously, together with the normal point prototypes or cluster centers for every cluster. Hybridization of Possibilistic C-Means (PCM) AND Fuzzy C-Means (FCM) is the FPCM that frequently rejects several drawbacks of PCM, FCM. The noise sensitivity fault of FCM is solved by FPCM, which conquers the concurrent clusters drawbacks of PCM.

Wen-Liang Hung presented a new approach called Modified Suppressed Fuzzy c-means (MS-FCM), which drastically improves the performance of FCM owing to the use of prototype-driven learning of parameter α (Ming Yang and Hong Li, 2009). Exponential separation strength between clusters is the base for the learning process of α and is updated at each of the iteration. The parameter α can be calculated as

$$-\min_{i \neq k} \left[ \frac{[\|v]_i - v_k \| \square]^2}{\beta} \right] \quad (7)$$

In the above equation β is a normalized term hence β is chosen as a sample variance.

It is to be noted that the common value used for this parameter by all nodes at each iteration may induce an error. A new parameter is integrated which suppresses this common value of α and substitutes it by a new parameter like a weight to each vector. Or every point of the data set possesses a weight in association with every cluster. Accordingly this weight allows having a better selection. The following equation is used to calculate the weight.

$$w_{ji} = exp \left[ - \frac{\|x_j - v_i\|^2}{[\sum_{j=1}^{n} \|x_j - \bar{v}\|^2] * c/n} \right] \quad (8)$$

In the previous equation $W_{ji}$ denotes the weight of the point $j$ in relation to the class $i$. In order to modify the fuzzy and typical partition, this weight is exploited. The objective function is composed of two expressions: the first is fuzzy function which uses fuzziness weighting exponent and the second is possibilistic function which uses a typical weighting exponent; however the two coefficients in the objective function are only used as exhibitor of membership and typicality. A new relation, lightly different, enabling a faster decrease in the function and enhancement in the membership and the typicality when they tend toward 1 and decrease this degree when they tend toward 0. This relation is to add weighting exponent as exhibitor of distance in the two under objective functions.

The objective function of the MFPCM can be given as follows:

$$J_{MFPCM} = \sum_{i=1}^{c} \sum_{j=1}^{n} [(\mu]_{ij}^m w_{ij}^m d^{2m}(x_j, v) + t_{ij}^\eta w_{ij}^\eta d^{2\eta} [(x]_j, v_i)) \quad (9)$$

$U = \{\mu_{ij}\}$ represents a fuzzy partition matrix, is defined as:

$$\mu_{ij} = \left[ \sum_{k=1}^{c} \left( \frac{d? X_j, v_i}{d? X_j, v_k} \right)^{2m/(m-1)} \right]^{-1} \quad (10)$$

$T = \{t_{ij}\}$ represents a typical partition matrix, is defined as:

$$t_{ij} = \left[ \sum_{k=1}^{n} \left( \frac{d? X_j, v_i}{d? X_j, v_k} \right)^{2\eta/(\eta-1)} \right]^{-1} \quad (11)$$

$V = \{v_i\}$ represents c centers of the clusters, is defined as:

$$v_i = \frac{\sum_{j=1}^{n}\left[(\mu)_{ij}^{m}\, w_{ji}^{m} + t_{ij}^{\eta} w_{ji}^{\eta}\right] * X_j}{\sum_{j=1}^{n}\left[(\mu)_{ik}^{m}\, w_{ji}^{m} + t_{ik}^{\eta} w_{ij}^{\eta}\right)} \qquad (12)$$

## 4. EXPERIMENTAL RESULTS

In order to evaluate the proposed preprocessing phase with robots cleaning, experiments were carried out using UCI Machine Learning Repository (University of California, Irvine). This repository contains 211 datasets.

Two standard datasets from the UCI Machine Learning Repository datasets were selected for the evaluation purpose. Following are the data sets used for evaluating the proposed preprocessing phase with robots cleaning.

- Anonymous Microsoft Web Dataset [24],
- MSNBC.com Anonymous Web Dataset [25]

The web user analysis are evaluated using

- Hit Rate
- Execution time

### 4.1. Preprocessing

Evaluate the preprocessing phase by using Anonymous Microsoft Web Dataset and MSNBC.com Anonymous Web Dataset. In this Anonymous Microsoft Web dataset consists of 37711 records in the log file. Whereas MSNBC.com Anonymous Web dataset contains 989818 records in the log file. Then the data cleaning process is carried out. Initially, after removing records with local and global noise, graphics and videos format such gif, JPEG, etc., 29862 and 865412 records are obtained in this two type of datasets. Preprocessing phase is briefly described in (Nithya, P and Sumathi P, 2012) [19]

### 4.2. Hit Rate:

The proposed web user clustering uses the Modified Fuzzy Possibilistic C Means (MFPCM) algorithm. The accuracy of the proposed web user clustering system is compared with the previous web user clustering systems which uses the Modified Fuzzy C Means (MFCM) for clustering. Figure 2 shows the clustering accuracy comparison of the proposed and the existing approaches. The Hit Rate of FCM cluster view is from 75.3% on Top-1 relevance to 85 % on Top-5 relevance. Then the Hit Rate of MFPCM cluster view is from 79%

on Top-1 relevance to 89.2% on Top-5 relevance which is shown in Figure 2

Table 2 gives the comparison table of Hit Rate for proposed approaches against existing approaches. The Hit Rate of FCM cluster view is from 62% on Top-1 relevance to 72% on Top-5 relevance. Then the Hit Rate of MFCM cluster view is from 65% on Top- 1 relevance to 75 on Top- 5 relevance. From the four methods MFPCM, FPCM, MFCM, FCM, proposed approaches of MFPCM are proved better results against FPCM, MFCM, and FCM.

Figure 2 shows the comparison of Hit Rate for proposed approaches against existing approaches. From the Figure it is clearly noticed that the proposed method of MFPCM gives better result and shows higher performance than the FPCM, MFCM, and FCM.

### 4.3. Execution Time

Execution time shows the time taken by the standard methods to execute the process in all the datasets. The method that takes less time to execute is considered as the best one and it helps in increasing the growth of clustering results.

Table 3 gives the comparison table for the execution time for proposed method against existing method. That the Proposed method of MFPCM have 0.25 seconds is very low compared to other approaches followed by FPCM have 0.28, MFCM, 0.32 seconds, and FCM 0.37 seconds.

In the Figure 3, it shows graphical representation of the comparison of the time taken for the clustering classification proposed method against existing methods. It is observed from the graph that the time taken by the MFPCM is very less when compared to the FPCM, MFCM, FCM.

## 5. CONCLUSION

This paper mainly focuses on web usage mining with the help of clustering process. The break out the complexity in clustering the web data is initially preprocessed. It reduces the noises and web robots effectively. Here used Modified Fuzzy Possibilistic algorithm for the purpose of clustering. FPCM algorithm has advantages of both fuzzy and Possibilistic c-means techniques. The experimental results of time and hit rate of Prediction Model express these clustering effective suits for web usage mining. Data preprocessing treatment system for web usage mining has been analyzed and implemented for log data. Data cleaning phase includes the removal of records of graphics, videos

and the format information, the records with the failed HTTP status code and finally robots cleaning. Different from other implementations records are cleaned effectively by removing local and global noise and robot entries. This preprocessing step is used to give a reliable input for data mining tasks. Accurate input can be found if the byte rate of each and every record is found. The data cleaning phase implemented in this paper will helps in determining only the relevant logs that the user is interested in. Anonymous Microsoft Web Dataset and MSNBC.com Anonymous Web Dataset are used for evaluating the proposed preprocessing technique and it reveals that number of records.

The problem of web users clustering is to use web access log files to partition a set of users into clusters such that the users within a cluster are more similar to each other than users from different clusters are solved by using Modified Fuzzy Possibilistic C Means algorithm MFPCM and it is compared with FCM algorithm. And the experimental result shows that the proposed technique results in higher hit rate and less execution time. Thus, the proposed MFPCM approach is best suited for the web users clustering applications effectively.

## REFERENCES

[1] Aghabozorgi, S.R.(2009) Using Incremental Fuzzy Clustering to Web Usage Mining. In: Wah, T.Y(eds) *Proceedings of International Conference of Soft Computing and Pattern Recognition,* pp. 653-658.

[2] Baraglia, R. (2002) SUGGEST: a Web usage mining system. In: Palmerini, P (eds) *Proceedings of International Conference on Information Technology: Coding and Computing,* pp. 282-287.

[3] J. Bezdek (1981) Pattern Recognition with Fuzzy Objective Function Algorithms, New York: Plenum Press.

[4] Chih-Hung Wu (2010) Web Usage Mining on the Sequences of Clicking Patterns in a Grid Computing Environment. In: Yen-Liang Wu, Yuan-Ming Chang and Ming-Hung Hung(eds) *Proceedings of International Conference on Machine Learning and Cybernetics (ICMLC)*, Vol. 6, pp. 2909-2914.

[5] Chu-Hui Lee(2008) Web Usage Mining Based on Clustering of Browsing Features. In: Yu-Hsiang Fu(eds) *Proceedings of Eighth International Conference on Intelligent Systems Design and Applications*, Vol. 1, pp. 281-286.

[6] J.C. Dunn(1973)A Fuzzy Relative of the ISODATA Process and Its Use in Detecting Compact Well-Separated Clusters, *Journal of Cybernetics* Vol.3,pp.32-57.

[7] Etminani K. (2009) Web Usage Mining: Discovery of the Users' Navigational Patterns Using SOM. In: Delui, A.R., Yanehsari, N.R. and Rouhani, M(eds) *Proceedings of First International Conference on Networked Digital Technologies*, pp.224-249.

[8] Havens Timothy C (2011) Speedup of fuzzy and possibilistic kernel c-means for large-scale clustering. In:Chitta Radha,Jain Anil Kand Jin Rong(eds) *Proceedings of IEEE International Conference on Fuzzy Systems (FUZZ)*, pp. 463–470.

[9] Hogo (2003) Temporal Web usage mining. In: M., Snorek, M. and Lingras, P(eds) *Proceedings of International Conference on Web Intelligence*, pp. 450-453.

[10] Inbarani H.H (2007) Rough Set Based Feature Selection for Web Usage Mining. In: Thangavel, K. and Pethalakshmi A (eds) *Proceedings of International Conference on Conference on Computational Intelligence and Multimedia Applications*, Vol. 1, pp. 33-38.

[11] Jalali M.(2008) A Web Usage Mining Approach Based on LCS Algorithm in Online Predicting Recommendation Systems. In: Mustapha, N., Sulaiman, N.B. and Mamat, A (eds) *Proceedings of 12th International Conference Information Visualisation,* pp. 302-307.

[12] Jian Chen (2004) Discovering Web usage patterns by mining cross-transaction association rules. In: Jian Yin, Tung, A.K.H. and Bin Liu (eds) *Proceedings of International Conference on Machine Learning and Cybernetics,* Vol. 5, pp. 2655-2660.

[13] Jianxi Zhang (2009) Web Usage Mining Based On Fuzzy Clustering in Identifying Target Group. In: Peiying Zhao, Lin Shang and Lunsheng Wang(eds) *International Colloquium on Computing, Communication, Control, and Management*, Vol. 4, pp. 209-212.

[14] Labroche N (2007) A New Web Usage Mining and Visualization Tool. In: Lesot M.J. and Yaffi L (eds) *Proceedings of 19th IEEE International Conference on Tools with Artificial Intelligence*, Vol. 1, pp. 321-328.

[15] Maratea A (2009) An Heuristic Approach to Page Recommendation in Web Usage Mining.

In: Petrosino A (eds) *Proceedings of Ninth International Conference on Intelligent Systems Design and Applications*, pp. 1043-1048.

[16] Ming Yang (2009) User Analysis Based on Fuzzy Clustering.In: Hong Li (eds) *Proceedings of 2009 International Conference on Business Intelligence and Financial Engineering.*

[17] Nasraoui O (2008) A Web Usage Mining Framework for Mining Evolving User Profiles in Dynamic Web Sites. In: Soliman M, Saka E, Badia A and Germain R *Proceedings of IEEE Transactions on Knowledge and Data Engineering,* Vol. 20, No. 2, pp. 202-215.

[18] Nina S.P(2009) Pattern Discovery of Web Usage Mining.In: Rahman M, Bhuiyan K.I. and Ahmed K(eds) *Proceedings of International Conference on Computer Technology and Development,* Vol. 1, pp.499-503.

[19] Nithya, P(2012) Novel pre-processing technique for web log mining by removing global noise and web robots. In: Sumathi, P(eds) *Proceedings of National Conference on Computing and Communication Systems (NCCCS),* vol.,1, pp.21-22 .

[20] Shinde, S.K.(2008) A New Approach for on Line Recommender System in Web Usage Mining. In: Kulkarni U.V(eds) *Proceedings of International Conference on Advanced Computer Theory and Engineering,* pp. 973-977.

[21] Wu, K.L., Yu, P. S. and Ballman, A.,(1998) SpeedTracer: A Web usage mining and analysis tool. *IBM Systems Journal*, Vol. 37, No. 1, pp. 89-105.

[22] Yan Li(2008)Research on Path Completion Technique in Web Usage Mining. In: Boqin Feng and Qinjiao Mao b(eds) *Proceedings of International Symposium on Computer Science and Computational Technology,* Vol. 1, pp. 554-559.

[23] Zhang Huiying and Liang Wei,(2004) An intelligent algorithm of data pre-processing in Web usage mining. *Fifth World Congress on Intelligent Control and Automation*, Vol. 4, 3119- 3123.
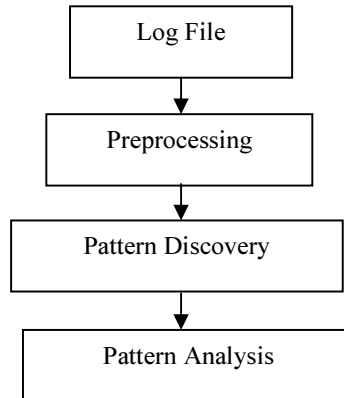
[24] http://archive.ics.uci.edu/ml/datasets/Anonymous+Microsoft+Web+Data

[25] http://archive.ics.uci.edu/ml/datasets/MSNBC.com+Anonymous+Web+Data

*Figure 1: Steps in Web Usage Mining*

*Table 1: Format of Web Access Log.*

| Field | Meaning |
|---|---|
| 219.144.222.253 | User's IP address (UIP) |
| [16/Aug/2004:15:36:11+0800] | The date and time of the request |
| GET | The method of the request |
| /images/1….r3… | The URL of the current request (URI) |
| Http/1.1 | The version of the transport protocol (Version) |
| 200 | The HTTP status code returned to the client (Byte) |
| 418 | The content length of the page transferred (Bytes) |
| http://202.11... | The URL requested just before (Refer URI) |
| Mozilla/4.0(… | Browser & Os (BrowserOS) |

*Table 2: Comparison of Hit Rate for Proposed Method with Existing Method*

| Top r value | FCM | MFCM | FPCM | MFPCM |
|---|---|---|---|---|
| 5 | 62 | 65 | 75 | 79 |
| 10 | 65 | 68 | 79 | 82 |
| 15 | 68 | 72 | 81 | 86 |
| 20 | 72 | 75 | 85 | 89 |



*Figure 2: Comparsion of Hit rate*

www.jatit.org

*Table 3: Comparison of Execution Time*

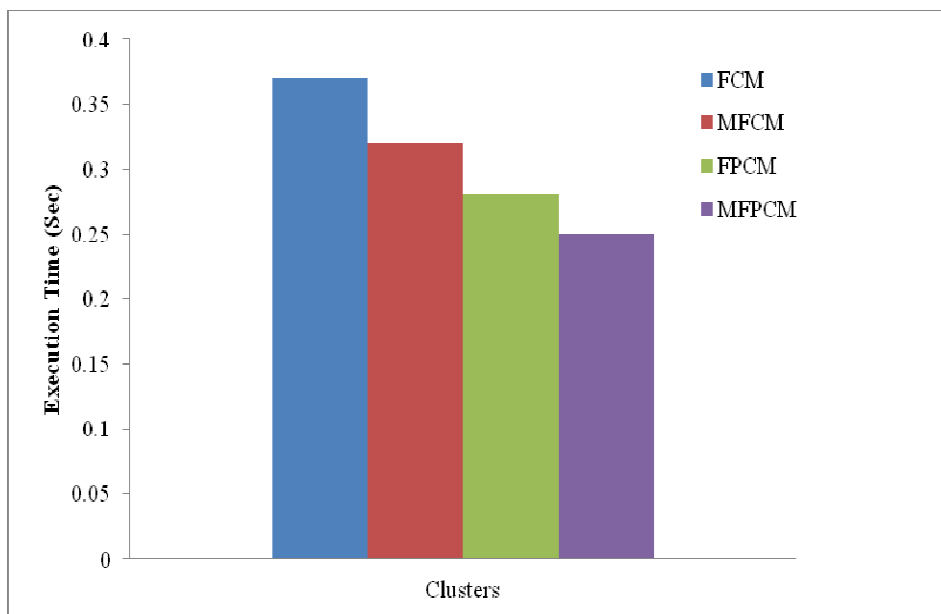| Methods | Execution Time |
|---------|----------------|
| FCM | 0.37 |
| MFCM | 0.32 |
| FPCM | 0.28 |
| MFPCM | 0.25 |



*Figure 3: Comparsion of Execution time*