# DISCOVERING WEB DOCUMENT CLUSTERING USING WEIGHTED SCORE MATRIX AND FUZZY LOGIC

**\*Pranali R. Raut** [1] | **Prof. Nilesh R. Khochare** [2]

[1] Student, Department Of Computer, JSPM NTC, India. \*Corresponding Author

[2] Professor, Department Of Computer, JSPM NTC, India

## ABSTRACT

In computer analysis, many files are usually analysis much of the data but in those files consists of unstructured data, and that data examined by computer examiners are more difficult to be performed. For this purpose it uses clustering documents, it can gives new and useful knowledge from the documents under analysis. Recently used clustering algorithms have some disadvantage like data preparation, outliers. The main theme is web documents is converted into clustering documents with the help of data preprocessing, features extraction, and weighted scores matrix techniques.

**KEYWORDS:** Web Documents, Web Crawler, Fuzzy Logic, Weighted score Matrix, Feature extraction, Clustered Document.

**INTRODUCTION:**

The basic idea of web clustering is hundreds of thousands of files are usually examined to come to a conclusion. So there is a need to discover the fast method that can group the required documents. Web documents are complex and heterogeneous. Web documents contain many data which is in not structured format so that data is analysis by any machine is not possible. So, automated methods of document examined are of great interest.

Applying clustering techniques on web document contains huge number of the web pages and that make the complexities in the user to cluster them semantically. This need can be for many different purposes like detection of mood, forensic analysis of the web users and many other events. Main approach is that applies document clustering algorithms to forensic analysis of computers seized in police investigations. In this paper, approach by carrying out extensive experimentation with six well-known clustering algorithms K-medoids, Complete Link, Single Link, Average Link, K-means and CSPA.

Representation of text is useful for the selecting features to represent text that will be clustered. Feature selection is a process of identifying the most effective subset of the original features to be used in clustering. Extraction of features is the process of using linear or non-linear transformations on original features to generate projected features to be used in clustering [3].

Web mining has fuzzy characteristics, so it is good suitable in comparison with conventional clustering. Fuzzy clustering contains two methods, first it is on fuzzy c-partitions, is called a Fuzzy C-Means clustering and another based on the fuzzy equivalence relations, is called a Fuzzy Equivalence Clustering. Data mining technique called association examined, it is more important for discovering interesting relationship hidden in large data set also useful for clustering. There are two broad principles use for association analysis. One is Apriori and another is Frequent Pattern growth principle. FP-growth is a divide and conquers strategy that mines a complete set of frequent item sets without candidate generation. FP-growth outperformance Apriori because Apriori incurs considerable I/O overhead since it requires making several passes over the transaction data set. In this paper a method is being proposed of web document clustering based on FP-growth and FCM that helps the search engine to retrieve relevant web documents needed for any user. Documents in the FCM are strongly corre-lated; however traditional FCM clusters are sensitive to the initialization of membership matrix and centre. It also needs the number of clusters to be formed as initial parameter. Our approach handles all this by using FP-growth approach which initializes this for FCM.

**Specified many key requirements for Web document clustering methods:**

1) **Relevance:** Produce clusters that group documents relevant to the user's query separately from irrelevant ones.

2) **Browsable Summaries:** The user examined whether cluster documents contents are useful or not useful. Therefore the method has to gives concise and accurate explanations of the clusters.

3) **Overlap:** Documents have many topics, so it is important to avoid overlapping of documents that is one document is present in only one cluster.

4) **Speed:** The clustering method able to cluster up to one thousand snippets in a few seconds.

5) **Incrementality:** For the saving time purpose, the method should start to process each snippet as soon as it is received over the Web.

This paper can be classified as follows: Section I dedicated for Introduction. Section II reserved for Related Work, Section III is allocated for System Description and finally section IV is done with conclusion.

**RELATED WORK:**

To put forward the idea of "Discovering Web Document Clustering Using Weighted Score Matrix And Fuzzy Logic". This paper analyzes many concepts of different authors as mentioned below:

N. L. Beebe introduces this paper introduces an approach to overcome trace- ability issues in digital forensic investigation process. Digital crime inflicts immense damage to users and systems and now it has reached a level of sophistication that makes it difficult to track its sources or origins especially with the advancements in modern computers, networks and the availability of diverse digital devices. Forensic has an important role to facilitate investigations of illegal activities and inappropriate behaviors using scientific

methodologies, techniques and investigation frameworks. Digital forensic is developed to investigate any digital devices in the detection of crime. This paper emphasized on the research of trace- ability aspects in digital forensic investigation process. This includes discovering of complex and huge volume of evi- dence and connecting meaningful relationships between them. The aim of this paper is to derive a traceability index as a useful indicator in measuring the accuracy and completeness of discovering the evidence. This index is demonstrated through a model (TraceMap) to facilitate the investigator in tracing and mapping the evidence in order to identify the origin of the crime or incident. In this paper, tracing rate, mapping rate and offender identification rate are used to present the level of tracing ability, mapping ability and identifying the offender ability respectively. This research has a high potential of being expanded into other research areas such as in digital evidence presentation. [1].

S. Decherchi had analyzed Authorship verification can be checked using stylometric techniques through the analysis of linguistic styles and writing characteristics of the authors. Stylometry is a behavioral feature that a person exhibits during writing and can be extracted and used potentially to check the identity of the author of online documents. Although stylometric techniques can achieve high accuracy rates for long documents, it is still challenging to identify an author for short documents in particular when dealing with large authors populations. These hurdles must be addressed for stylometry to be usable in checking authorship of online messages such as emails, text messages, or twitter feeds. In this paper, we pose some steps toward achieving that goal by proposing a supervised learning technique combined with n-gram analysis for authorship verification in short texts. Experimental evaluation based on the Enron email dataset involving 87 authors yields very promising results consisting of an Equal Error Rate (EER) of 14.35% for message blocks of 500 characters [2].

Dr.T.Nalini explains clustering is the process of grouping the documents and the grouping is performed by finding similarities between data based on their properties. These groups are termed as Clusters. Clustering algorithms across two different data items is performed. The performance of the many different clustering algorithms is compared based on the time taken to form the estimated clusters. Thus it can be resulted as the time taken to form the clusters increases as the number of cluster increases. The first clustering algorithm takes very less time to cluster the data items whereas the simple K-Means takes the more time to perform clustering [3].

George Forman describes most research in speed up the text mining involves algorithm effectiveness to induction algorithms. For one or more large scale applications, means in classifying as well as indexing big document repositories. The time spent reduced word features from texts can it greatly increased the initial training time. This paper explains very fast techniques for text feature extraction that get together Unicode conversion, forced lower casing, word boundary detection, and string hash computation. It show empirically that integer hash features result in classifiers with equivalent statistical performance to those built using string word features, but require far least computation and less memory[4].
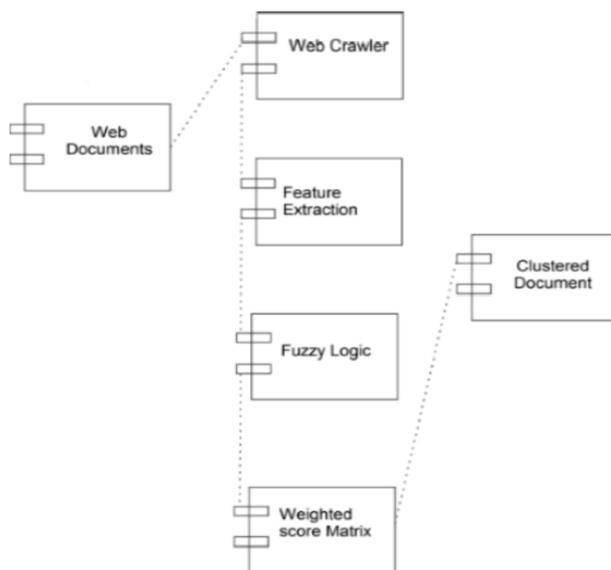
Giridhar N S explains Information Retrieval (IR) is essentially a matter of deciding which documents in a collection should be retrieved to satisfy a user's need for information. The user needs for information is described by an any query, and contains many search terms, add few additional information contains weight of the words. So, the retrieval decision is done by comparing the each terms of the query with the index terms appearing in the each documents. The decision may be binary such that retrieve or rejected, or it may consist estimating the degree of relevance that the document has to query. The words that apply in documents and in queries often have many morphological variants. After that the infor-

mation retrieval from the documents the stemming techniques are applied on the target data set. So it makes reduce the size of the data set which will increase the effectiveness of IR System. In this paper, surveys of stemming techniques have been presented [5].

**SYSTEM DESCRIPTION:**
Volume of data in digital world is growing increasingly, which has badly impact on forensic analysis. So there is a need to find the quick method that can group the required documents. Many clustering algorithms are used. So, system is pre-process unstructured format to structured format, and then extracts 4 important features of each document like numeric words, proper nouns title sentences and term weights. This method is very easy than others methods. Then system neglecting unwanted extension's considering only extensions which are rich in text. In the last step of clustering, system makes a score matrix of each document by comparing with one another to yield a score matrix which contains aggregate feature score. The group of these scored values. Then that scored values represents the most correct clustered documents.

This system first creates an interactive web crawler which eventually parses the web pages and collects the data and saves in .txt file format. Then the folder in which these web data is stored is given as the input to the system which then preprocesses this data to extract the features and then fuzzy logic is applied to get the feature scores classification pattern and then this is feed to the weighted matrix method to create semantic clusters for the web page documents.



**Figure1: Overall System Diagram**

Main aim is to convert many web documents to clustered documents. In this web document cluster contains web crawler, data preprocessing, feature extraction and weighted score matrix. Web crawler contains many web pages that will be converted into clustered information. In data preprocessing contains special symbol removing, stop word removing, stemming. Feature extraction starts from an initial set of measured data and builds derived values (features) intended to be more information and the subsequent learning, in few cases it is to good human interaction. Feature extraction it is to be dimensionality reduced. If input data to an algorithm is too big to be processed and it is suspected to be redundant, after that it is transformed into an extract set of features. The extracted features are assumed to contain the more relevant information from the input data, so that the task performed by using this reduced representation instead of the complete initial data.

Fuzzy logic can be used as an interpretation model for the properties of neural networks; it gives a more easy description of their performance. We will show that fuzzy operators can be conceived as generalized output functions of computing units. Fuzzy logic used to specify networks without having to apply a learning algorithm. An expert in a certain field can sometimes produce a simple set of control rules for a dynamical system with less effort than the work involved in training a neural network.

Weighted score matrix used to define the level of criteria. Assign meaning to weighting factors is subjective. For this reasons, keep the number of weighting factors small.

Cluster analysis or clustering is the task of grouping a set of objects in such a way that objects in the same group are more similar to each other than to those in other groups. It is a more important task of exploratory data mining, and technique for statistical data analysis. Examined of clustering is not one particular algorithm, but that task to be solved. It can be achieved by different algorithms that differ significantly in their notion of what constitutes a cluster and how to effectively find them. Popular notions of clusters include groups with small distances among the cluster participant. Clustering can be formulated as a multi-objective optimization problem.

**CONCLUSION:**
This paper successfully accumulates most of the techniques of many authors as described in section II of related work. So, by analyzing all methods it seem to be like number of method is perfect in providing solution for "Discovering Web Document Clustering Using Weighted Score Matrix And Fuzzy Logic".

As an effort to this, this paper tries to improve the concept of "Discovering Web Document Clustering Using Weighted Score Matrix And Fuzzy Logic" by introducing clustering based techniques is to extract features from the web documents using conditional random field methods and build a fuzzy linguistic topological space based on the associations of features.

**REFERENCES:**

1. "Digital forensic text string searching: Improving information retrieval effectiveness by thematically clustering search results," N. L. Beebe and J. G. Clark, Digital Investigation, Elsevier, vol. 4, no. 1, pp. 49–54,2007.Year of publication: 2007.

2. "A Cluster-based Approach to Browsing Large Document Collections", Cutting, D. R., Karger, D. R., Pedersen, J. O., and Tukey.W Proc. Of SIGIR'92 (pp. 318–329).

3. "Web document clustering: a feasibility demonstration," O. Zamir and O. Etzioni, in Proceedings of 19th international ACM SIGIR conference on research and development in information retrieval (SIGIR 98), 1998, pp. 4654.

4. "Searching the world wide web," S. Lawrence and C. L. Giles, Science, vol. 280, no. 5360, pp. 98100, 1998.

5. "Search technologies for the internet," M. Henzinger, Science, vol. 317, no. 5837, pp.468471, 2007.

6. "Semantic web content analysis: A study in proximity-based collaborative clustering," V. Loia, W. Pedrycz, and S. Senatore, IEEE T. Fuzzy Systems, vol. 15, no. 6, pp. 12941312, 2007.

7. "An approach to flexible information access systems using soft computing," H. L. Larsen,  in Proc. of the 32nd Annual Hawaii International Conference on System Sciences, Hawaii, 1999, p. 231.

8. "Document clustering with cluster refinement and non-negative matrix factorization," S. Park, D. U. An, B. R. Cha, and C. W. Kim, in Proceedings of the 16th International Conference on Neural Information Processing, Bangkok, Thailand, 2009,pp. 281288.

9. "A sentence-to-sentence clustering procedure for pattern analysis," S. Lu and K. Fu, IEEE Transactions on Systems, Man and Cybernetics, vol. 8, pp. 381389, 1978.

10. "A Text Mining Technique Using Association Rules Extraction," Mahgoub Hany, Nabil Ismail, Torkey, Fawzy, v-4, pp.21-28, 2008.