# Discovering Latent Semantics in Web Document Using Fuzzy Clustering

[1]Nanthini.B, [2]Sabari Ramachandran.M, [3]Balasubramanian.N,[4]Mohamed Rafi.M

[1]Final year MCA, Mohamed Sathak Engineering College, Kilakarai.
[2]AssistantProfessor, Dept of MCA, Mohamed Sathak Engineering College,Kilakarai.
[3]Associate Professor, Dept of  MCA, Mohamed Sathak Engineering College,Kilakarai.
[4]Professor, Dept of MCA, Mohamed Sathak Engineering College,Kilakarai.

**ABSTRACT:**Web is the massive collection of data source in that finding the correct information is not easy. Searching and web mining techniques are used to retrieve the information from the web. Web document clustering is the most useful technique to improve the efficiency of information searching problem. The traditional web mining techniques has various difficulties in handling the data which are not clear. Web usage mining is a data mining technology to mining the data of the web server. It can find the searching patterns of the user and some kind of associations between the web pages. Web usage mining gives the support for the website design, providing personalization server and other business making decision, etc.Text clustering plays an important role in providing automatic navigation and browsing mechanisms by organizing large sets of documents into a small number of significant clusters.Web is the huge collection of data source in that finding the correct information is not easy. Examining and web mining techniques are used to retrieve the information from the web. Web document clustering is the most useful technique to improve the efficiency of information searching problem.

*Keywords:*Document, K-means algorithm, Natural Language Processing Tool Kit (NLTK), Principal Component Algorithm (PCA), Text Mining.

## 1.INTRODUCTION

The modern world is becoming more and more digital with the development of new smart technologies. Web search is very common practice among the people all over the world. A huge quantity of information is published on the internet in every second. So, the web is enriched with massive number of literal documents every day. With this information overloaded, looking up for the precise and relevant information resource and extracting the key concept from the resource has become challenging within a very short timModern days, Web document retrieval systems like search engines play vital role of people who search pursuing related web document in the over the Internet. In the sense lot of new algorithms are proposed for due to the complexity recovery web documents in clusters. In the past several years, the medical data have been growing explosively. For example, the number of papers published in PubMed was increased from 112,177 in 1960 to 2,019,238 in 2013 and the annual average number of discharges between 2007 and 2010 is around 35 million1. Recently, various text mining techniques have been introduced into the medical domain. One fundamental objective of those techniques is to process the unstructured medical data into a proper format for better utilization to recognize explicit facts.Text clustering plays an important role in providing in-built navigation and browsing mechanisms by organizing large sets of documents into a small number of meaningful clusters. Many fuzzy clustering algorithms, such as K-means, deal with documents as bag of words. The bag of words representation method used for these clustering is often unsatisfactory because it ignores the semantic of wordsThe data mining techniques are used for analyzing data, in order to find significant information among available raw data. These techniques are developed through the computational algorithms that help to automate processes required to analyze the data. According to the nature of data and their application requirements the different kinds of algorithms can be applied on the data.Clustering is an unverified machine erudition technique. The unsupervised feature makes it more proper for clustering search result as it is not possible to determine as to how many categories are there in search result. Clustering of web search involves four basic steps: a) search result acquisition, b) result preprocessing, c) cluster formation and d) labeling of clusters. Some clustering.Search engines are indispensable tools to find, filter, and extract the desired information, which attempt to aid users in gathering relevant contents from web. Surveyed and compared search engines, such as AltaVista, Excite, HotBot, Infoseek, Lycos, Northern Light,

and Google. Both studies have shown the weaknesses of these search engines and proposed some solutions to overcome some of the flaws. A refined mechanisms to recognize the latent semantics in the returned search results is crucial to enhance completeness of a search engine.

## 2.PROPOSED SYSTEM

The proposed algorithm is Fuzzy Latent Semantic Clustering (FLSC) that covers the latent semantics of web documents that can applicable in text domains, it can be extended to the applications such as Data mining Bio informatics, Content based or collaborative information filtering. Latent Semantic Clustering (LSC) is a technique in overlapping the cluster processing.The proposed approach utilizes the semantic relationship between words to create concepts in the semantic-based model such as semantic vector space model (SVSM). It exploits the WordNet ontology in turn to create low dimensional feature vector which allows us to develop an efficient clustering algorithm. A new semantic-based model that analyzes documents based on their meaning is proposed. The proposed model analyzes terms and their corresponding synonyms and/or hypernyms in the documents. The proposed agent aims at increasing the performance of IR process by enhancing the document clustering. The accuracy of clustering has been computed before and after combining ontology with vector space model (VSM). Experimental results determine that the newly developed semantic-based model enhances the clustering quality of sets of documents substantially. Fuzzy logic is based on the theory of fuzzy sets, a theory which relates to classes of objects with unsharp boundaries in which membership is a matter of degree. Documents, queries and their characteristics could easily be viewed as fuzzy granular classes of objects with un-sharp boundaries and fuzzy memberships in many concept areas .Since the concept of fuzzy logic is quite intuitive, the fuzzy logic model provides a framework that is easy to understand for a common user of IR system.The proposed System extracts features from the web documents using semi-supervised learning schemes called named entitiesand builds a fuzzy linguistic topological space based on the associations of features.

## 3.RELATED WORK

There are two major research areas in mining medical documents. The first one tracks concepts by looking for frequency of words. The second area categorizes the concepts to find latent variables in medical documents. The first approach leads to high sparse dimensionality data therefore, researchers have been motivated to use the second approach such as topic modeling. Among topic models,

LDA is a popular and effective unsupervised topic model. In the medical domain, LDA has been leveraged in a wide range of applications. Used LDA for comparing the topics of patient notes, and used LDA in FDA drug side effects labels to cluster drugs. One of the methods that has not been fully considered in medical text mining is fuzzy set theory.Drakshayani and Prasad proposed a new model for text document representation. The proposed model follows parsing, preprocessing and assignment of semantic weights to document phrases to reflect the semantic similarity between phrases and k-means clustering algorithm. They evaluated the proposed model using 5 different datasets in terms of FMeasure, Entropy, and Purity for K-Means clustering algorithm. The results demonstrated a performance improvement compared to the traditional vector space model and latent semantic indexing model. More NLP techniques may be included to enhance the performance of the text document clustering.This approach authors proposed two-dimensional scatter plots of the projected data which helps in easier visual assessment. The scatter plots for the IRIS data from each algorithm as can be easily seen, Sammon"s algorithm works effectively for Irisas the size of the data set is relatively small. Structure preservation is achieved as there is slight overlap between the two classes and the other is distinctly separate. For the fuzzy rule based model the results are similar to that of Sammon"s algorithm and the projection is good.
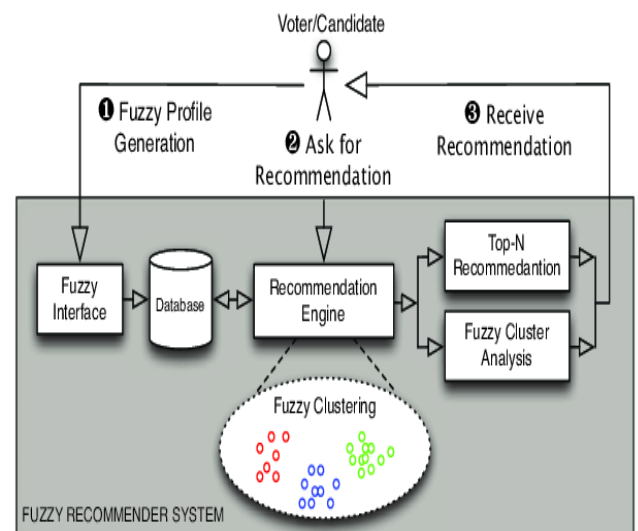
## 4.SYSTEMARCHITECTURE



Fig. 1. This figure illustrates the structure of fuzzy linguistic topological space. Documents are composite of meaningful features that can be categorized into different topics with different possibilities. Herein, the fuzzy linguistic function

denotes the possibility of a set of features belonging to a semantic topic.

The proposed approach utilizes the semantic relationship between words to create concepts in the semantic-based model such as semantic vector space model. It exploits the WordNet ontology in turn to create low dimensional feature vector which allows us to develop an efficient clustering algorithm. A new semantic-based model that analyzes documents based on their meaning is proposed. The proposed model analyzes terms and their corresponding synonyms and/or hypernyms in the documents.

# 5. SYSTEM COMPONENTS
## 5.1 Agent System

Agent technology is a new algorithm model, which is highly intelligent, easy to construct distributed system and having strong reusability. The concept of agent and technology has appeared in the development of distributed applied system and shown its remarkable effectiveness. From some research about agent and developing work in the aspect of distributed application, the meaning of the concept and technology.

## 5.2 Fuzzy K-Means:

Clustering of objects is as ancient as the human need for describing the salient characteristics of men and objects and identifying them with a type. Therefore, it embraces various scientific disciplines: from mathematics and statistics to biology and genetics, each of which uses different terms.

## 5.3 Clustering words:

Clustering is the grouping of a particular set of objects based on their characteristics, aggregating them according to their similarities. Regarding data mining, this methodology partitions the data implementing a specific join algorithm, most suitable for the desired information analysis. For document clustering, it is the process of grouping of contents such as words, word phrases or documents based on their content to extract knowledge. Word clustering is the similar process of clustering that groups the words according to its semantic values. The word belonging to one cluster is similar and the word belonging to another cluster is dissimilar. Grouping of words into individual cluster can be described by its own concept. 2pologies formed using this analysis.

## 5.4 Stemming:

Another important way to reduce the number of words in the representation is to use stemming. This is based on the observation that words in documents often have many morphological variants. For example, we may use the words computing, computer, computation, computational, computes, computable, computability etc. in the same document.

# 6. IMPLEMENTATION

The Implementation model is involved during the implementation phases are described in detail in Section.

## 6.1 Document Preprocessing

After getting input in a textual format, the document is preprocessed with document preprocessing technique. The detail concept is described in section. This step represents the text document with a set of index terms. Pre-processing significantly reduces the size of the input text documents by tokenization, stop word removal and stemming.

## 6.2 Document Representation

Document representation is the most crucial and challenging part of our implementation model. The remaining words of the document are to be represented by vector value. Each vector value of a word also represents the semantic relationship among them.

## 6.3 Text Extraction

The first step of in preprocessing process is extracting textual data from the web pages. Then convert each page into individual text document to apply text preprocessing techniques on it. This step is applied on input Web documents dataset by scanning the web pages and categorizing the HTML tags in each page. Then exclude the tags that contain no textual information like formatting tags and imaging tags.

---

# 7. RESULTS AND DISCUSSIONS
## 7.1 RESULT

Here this approach is taken three fuzzy based cluster algorithms to retrieve the web documents and compare with FLSC algorithm. Our FLSC algorithm used to consider new clusters based on the similar web documents found in the different clusters. Once Data admin upload the web document into the data server its created a new cluster and storing into the data server whether user given query search all the cluster and retrieve the web documents but it take time to retrieve the web documents and also it is not retrieve effectively to the method.

## 7.2 DISCUSSION

FLSC is an iterative algorithm. The aim of FLSC is to find clusters which very similar web documents found in the different clusters that minimize a dissimilarity function. In Fuzzy clustering each member is associated some membership value, that indicate the strength of association between a data element and a particular cluster. FLSC find clusters centers that minimize a dissimilarity function. FLSC iteratively moves the cluster centers to the "right" location within a dataset. Fuzzy set allows for degree of membership A single point can have partial membership in more than one class. There can be no empty classes and no class that contains no data points.

# 8. CONCLUSION

Analysing a large volume of medical data is important to advance state-of-the-art healthcare. Due to the unstructured

nature of free-text format for the medical data, text mining techniques such as topic modeling are widely adopted to extract latent semantic properties of a medical corpus. Despite the usefulness of topic models for medical data analysis, existing topic models such as LDA still suffer from several critical issues, such as extremely high computational complexity and unsatisfactory performance for data analytical tasks.Text mining is the most powerful and efficient technique for organizing the text corpus. The thesis based on this technique to discover an idea of a document. The thesis was accomplished by finding the similarities among the words to extract a general concept of the document. Word clustering technique is easier, less time consuming and will give a better concept of an ambiguous text document.

The proposed semantic-based model framework provides improved performance and makes a clustering be efficient. It also overcomes the problems existing in the VSM commonly used for clustering. The clustering result based on semanticbased model has higher efficiency values and faster than those based on the traditional VSM.

## 9.REFERENCES:

[1] Yang Yan, Lihui Chen, William-Chandra Tjhi, "Fuzzy semi-supervised co-clustering for text documents", Fuzzy Sets and Systems 215 (2013)

[2] Chien-Liang Liu, Tao-Hsing Chang, Hsuan-Hsun Li, "Clustering documents with labeled and unlabeled documents using fuzzy semi-Kmeans", Fuzzy Sets and Systems 221 (2013) 48–64, 2013 Elsevier B.V. All rights reserved.

[3] Vishal Gupta, Gurpreet S. Lehal (2009), "A Survey of Text Mining Techniques and Applications", Journal of Emerging Technologies in Web Intelligence, Vol.1, No.1.

[4].Magne Setnes "Supervised Fuzzy Clustering for Rule Extraction" Fuzzy Systems Conference Proceedings, 1999.

[5] Yanjun Li Congnan Luo (2008), "Text Clustering with Feature Selection by Using Statistical Data" IEEE.

[6] Navathe, Shamkant B. and Elmasri Ramez, (2000), "Data Warehousing and Data Mining", in "Fundamentals of Database Systems", Pearson Education pvt Inc, Singapore, 841-872.

[7]. Ajith Abraham "Natural Computation for Business Intelligence from Web Usage Mining" 3 IEEE Congress on Evolutionary Computation (CEC2003), Australia, IEEE Press, ISBN 0780378040, pp. 1384-1391,2005.

[8] Fellbaum, C. (2010). WordNet. Theory and Applications of Ontology: Computer Applications, 231, PP: 231-243, Springer Science+Business Media B.V.

[9] Drakshayani, B. & Prasad, E. (2012). Text Document Clustering based on Semantics. International Journal of Computer Applications (0975 – 8887). Vol. 45– No.4.

[10] Luo, C., Li, Y. & Chung, S. (2009). Text document clustering based on neighbors. Data & Knowledge Engineering 68 (2009) 1271–1288. Elsevier B.V.