



# Name Entity Recognition and Natural Language Processing for Improvised Fuzzy clustering in Web Documents

**Kalyani Ramesh Pole<sup>1</sup>, Vishakha R. Mote<sup>2</sup>**

<sup>1</sup>Computer Science (CSE), ME (Appear), (PES College of Engineering Aurangabad), (India)

<sup>2</sup>Information Technology (IT), ME, PhD (Appear), (PES College of Engineering  
Aurangabad), (India)

## ABSTRACT

Web documents are heterogeneous and complex. There are complicated associations within a single Web document and there can be complex relations with other documents as well. The high interactions between the terms of the documents show merely vague and thinly ambiguous meanings. Efficient and efficient grouping methods are required to discover latent and consistent meanings in context. This article presents a diffuse linguistic topology space with a diffuse cluster algorithm to identify and discover the basic contextual understatement and meaning in Web documents. The proposed algorithm extracts the functionality of Web documents using random conditional field methods and creates a diffuse linguistic topological space according to the associations of characteristics. The intrinsic associations of words that have the attribute to occur again in hierarchy of chained semantic compound terms called as CONCEPTS, where a diffuse linguistic measure is applied to each complex to evaluate 1) the relevance of a document belonging to a subject and 2) the difference between the other subjects. Web content can be grouped into subjects in the hierarchy according to their diffuse linguistic measures; Internet users can further explore the CONCEPTS of web content accordingly. In addition to the applicability of the algorithm in Web text fields, it can be extended to other applications, such as data mining, bioinformatics, content or collaborative information filtering, etc.

The internet or as we call it World Wide Web is termed as the most important information store of recent years. The growth of the Web is greatly expanded with new technologies. In case of search engines they are termed as inefficient when the number of documents on the web has been propagated. In more or less similar way, query recovery, most of which there is no relation to what the user was looking for. The documented varied and multifaceted Web, there are difficult relationships with a web document and a link to others. This research focused on the clustering algorithm to discover and identify latent association of semantics in the text based corpus from a diffuse linguistic point of view. In addition, applicability in text fields can be extended to applications such as data mining, bioinformatics, content-based or collaborative information screening, and so on. Second, the recovery document belongs to a research topic that should differ from other issues the difference between other topics. Web content can be grouped into subjects in the hierarchy based on their diffuse linguistic measures.

**Keywords:** collaborative information filtering, linguistic topological space, Web content

## **I. INTRODUCTION**

Grouping is a technique where objects in a group are similar to each other and different from objects in other groups. Grouping is an usually a unsupervised machine evaluation technique. The unmonitored feature makes the cluster search result more appropriate because it is not possible to determine the number of categories available in the search result. Combining web search involves four basic steps: a) acquiring research results, b) pre-processing results, c) clustering, and d) labelling clusters. Some cluster engines acquire search results from one or more search engines, and then merge them into a single unified result set.

In pre-processing, each document on the search result page is transformed into a stream of words or phrases or phrases based on the attributes of the clustering method. Other tasks performed during pre-processing are word removal, deletion, filtering, and so on. Many methods that include k-means, nearest neighbour group and hierarchical cluster are used to select a set of key terms or phrases to organize function vectors based on differences between documents to capture semantics to adapt user Intentions.

Suffix-tree clustering is a sentence-based approach, which performs clustering of documents based on similarities between documents. Fuzzy c means and fuzzy hierarchical clustering needs prior knowledge about the "cluster number" and the "initial cluster centroid", which are considered to be threat full and complex limitations of these approaches.

Based on the actual vector based space model, the basic similarity between the two documents is measured by using vectorial distance, such as the Euclidean distance, the Manhattan distance, and so on. A fuzzy hierarchical clustering approach to discover a set of frequent high-fuzzy object sets represented to represent the candidate clusters. Decomposition technique by Breaking the network organized by keywords co-occurrence with the centroid, i.e. the nodes with the maximum degree and a cut-off threshold in several clusters. Data clustering algorithm based on the unique hidden Markov model, which identifies an appropriate amount of clusters under provided dataset by excluding usage of prior knowledge about the number of clusters.

## **II. RELATED WORK**

This section provides a study of the various contributions recently made and the research that is helping to design an effective technique. In this article, Yang Yan et al [1] has suggested a new semi-supervised heuristic algorithm for fuzzy co-clustering (SSHFCR) for the selection of large Web documents. In this approach, the clustering process is accomplished by integrating prior knowledge in the form of user-supplied paired constraints in the fuzzy co-cluster framework. Each constraint specifies whether a pair of documents "should" or "cannot" be grouped together.

Deep Experimental readings and observations on a number of large reference datasets showcases the power and actual potentials in real terms of accuracy with stability and efficiency, compared with few recent yet popular techniques of semi-supervised clustering approaches. With focus on specific document based clustering, work represents another semi-supervised method of clustering algorithm called "semi-fuzzy K-means". The semi-K-means fuzzy is an extension of the K-means clustering model, and is inspired by an EM algorithm and a Gaussian mixture model. In addition, the fuzzy half-K makes it possible to use different fuzzy membership functions to measure the distance between the data.



Chien-Liang Liu et al [2] use a Gaussian weighting function to carry out experiments, but the similarity function of the cosine can also be used. This work performs experiments on three sets of data and compares the fuzzy semi-K means with several methods. Experimental results indicate that fuzzy semi-K-means can generally surpass other methods. There are two main areas of research in mining medical records. The first follows concepts by looking for the frequency of words (Poulin et al., 2014). The second area classifies concepts to find latent variables in medical documents (Lin, Karakos, Demner-Fushman and Khudanpur, 2006). The first approach leads to data of large dimensions and large dimensions (Aggarwal & Zhai, 2012); Therefore, researchers were motivated to use the second approach, such as modeling subjects. Among the thematic models, LDA (Latent Dirichlet Allocation) is a model of unsupervised popular and effective subject (Halpern, Hornig, Nathanson, Shapiro and Sontag, 2012).

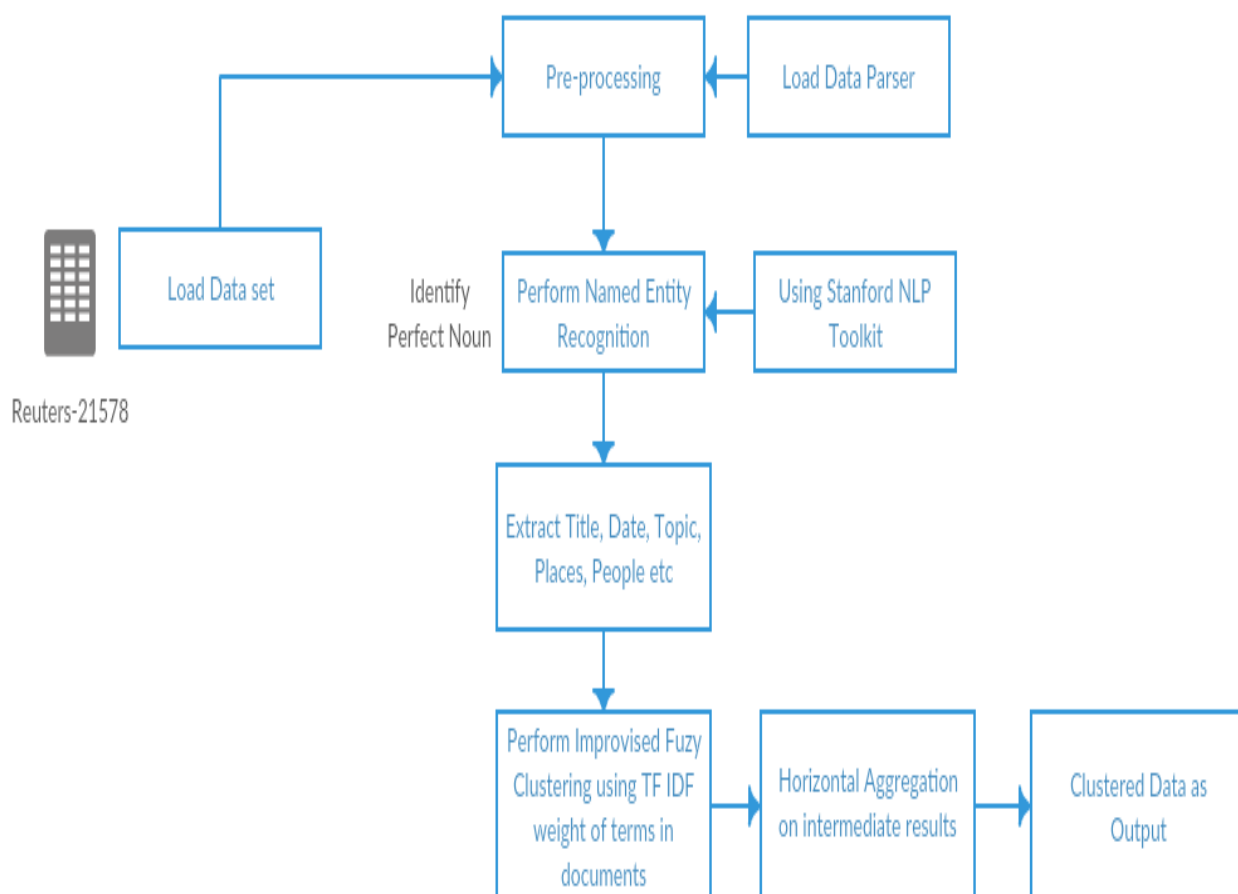
In the medical field, LDA has been exploited in a wide range of applications. For example, Arnold et al. (2010) used LDA to compare the subjects of patient scores (Arnold, El-Saden, Bui and Taira, 2010) and Bisgin et al. (2011) used LDA in the FDA drug side effects labels for cluster drugs (Bisgin et al., 2012). One of the methods that has not been fully taken into account in the exploration of medical text is the theory of fuzzy sets. Fuzzy set theory has been used to model systems that are difficult to define precisely. It incorporates the inaccuracy and subjectivity of human decision-making in the formulation and solution process of the model (Karami & Guo, 2012). Since Bellman and Zadeh (Bellman & Zadeh, 1970) have described the method of decision making in fuzzy environments, an increasing number of studies have dealt with uncertain fuzzy problems by applying the theory of fuzzy sets (Karami & Guo, 2012 Karami, Yazdani, Beiryae, & Hosseinzadeh, 2010). Some work was done in the exploration of medical text using fuzzy clusters to group or classify the documents without any knowledge of the latent semantics of the documents (Ben-Arieh & Gullipalli, 2012, Fenza, Furno and Loia, 2012). In addition, we recently used the fuzzy cluster as a feature transformation approach (dimension reduction) for medical text data (Karami & Gangopadhyay, 2014). Among the fuzzy grouping methods, Fuzzy C-means (FCM) (Bezdek, 1981) is the most popular (Bataneh, Naji, & Saqer, 2011). In this research, we propose a new method using a fuzzy cluster to extract latent semantic features from medical documents.

In computer forensic analysis, hundreds of thousands of files are usually examined. Much of the data in these files consists of unstructured text, which is difficult to analyze by computer examiners. In this context, automatic analytical methods are of great interest. In particular, algorithms for grouping documents can facilitate the discovery of new and useful knowledge from the documents analyzed. Luís Filipe da Cruz Nassif et al [3] present an approach that applies algorithms to group documents to the forensic analysis of computers seized in police investigations. The authors illustrate the proposed approach by conducting in-depth experimentation with six well-known clustering algorithms (K-means, K medoids, Single Link, Complete Link, Average Link and CSPA) applied to five real datasets obtained from Computers seized in real world research. The experiments were carried out with different combinations of parameters, which gave rise to different instantiations of algorithms. In addition, two relative validity indices were used to automatically estimate the number of clusters. Related studies in the literature are much more limited than our study. Based on experience, the Average Link and Complete Link algorithms provide the best results for the application domain. If properly initialized, the partial algorithms (K-means and K-medoids) can also give very good results. Finally, the authors also present and discuss several

practical results that can be useful to forensic computing researchers and practitioners. D.Renukadevi et al [4] studied the clustering technique and discussed their observations because advances in information technology and the increasing ease of use of the Internet radically alter all areas of activity in the modern day. As a result, a very large number of people would be required to interact more frequently with computer systems. To make human-machine interaction more effective in such situations, it is desirable to have systems capable of managing inputs in various forms, such as printed paper / manuscript documents. The computer must process scanned images of printed documents effectively, the techniques must be more sophisticated. Text documents are pre-processed; Term Frequency and Reverse Document Frequency (TF-IDF) are used to classify the document. Finally, similar information is grouped using the Fuzzy C-Means Clustering algorithm.

**III. PROPOSED ARCHITECTURE**

The proposed work is motivated by a research paper [6]. Based on the findings and the literature review, the following key issues and challenges are addressed to improve the traditional text-gathering technique.



**Fig 1: Proposed Architecture**

1. The length of textual records is not similar; therefore, evaluation of individual textual content requires a significant amount of computing resources.



2. Extracting functionality from different documents is different in nature and length, so that measuring the similarity of a data object to another object is a complex task.
  3. Cluster formation of documents must select some centroid for accurate group formation, but the random and fluctuating centroid selection in text documents can increase the process time and clustering accuracy.
  4. The approximation of similarity in text extraction must compare the text document with their important characteristics, but directional information on similarity is still calculated to optimize clustering performance.
- In order to solve the obtained issues and challenges for document clustering the following solution is proposed for further investigation and design.

1. Design of a strong pre-processing technique for refining the noisy contents form the documents learning set.
2. Design a new feature extraction and selection technique for optimizing the performance of document content analysis and their comparisons.
3. Enhance the traditional fuzzy c-means in order to achieve higher accuracy over the text content analysis and their clustering.
4. Implement the modifications on the fuzzy c-means clustering to demonstrate the hierarchical relationship among the documents.

Pseudo code for proposed module:

1. Load data set Reuters 21578
2. Perform data parsing and segregate data attributes.
3. Identify Perfect Noun, i.e. Name, Place, People, Country etc using Stanford NLP Toolkit.
4. Perform Named Entity Recognition to create a level 0 simplex with all named entities.
5. Calculate Term frequency / Inverse Document Frequency of document along with weight score.
6. Using TF-IDF score, perform clustering.

Require:  $V = \{x_1, x_2, \dots, x_n\}$  be the vertex set of all reserved named entities generated from  $W$  associated with

their categories  $_$  in a collection of documents.

Ensure:  $H$  is the hierarchy of connected components.

Let  $S = \{\{x_1\}, \{x_2\}, \dots, \{x_n\}\}$  be the set of all 0-

Simplexes initially.

TF - IDF Given two weights  $a$  and  $b$

Let  $k \leftarrow 0$ .

While  $k \leq n$  do

Let  $S_i$  and  $S_j$  be two  $n$ -simplexes in  $S$ .

While  $\delta(S_i, S_j) \geq a$  and  $\sigma(S_i, S_j) \leq b$  do

$S' \leftarrow S_i \cup S_j$ .

Add  $S'$  to  $S$

end while

$k \leftarrow (k + 1)$

end while

7. Perform horizontal aggregation on results.
8. Collect output and calculate Precision, Recall and F measure.
9. The input from the matrix is the generated score. We calculate the smallest and biggest scores. We calculate exactly five ranges starting from smallest the value and ending with the largest value. Now assign the score to score calculated in master matrix step and check the score in these five ranges. Once, the scores have been calculated a threshold of say '2' is set. The file having threshold more than two is added to the cluster and discards the file which fails to satisfy the condition.

#### **IV. CONCLUSION**

In general fuzzy clustering algorithm is use to identify the latent semantic in web documents and Natural language processing and name entity recognition technique is very efficient for cluster the documents.

After successfully implementation of the proposed technique of document clustering approach the following outcomes are expected.

1. An improved approach of fuzzy c-means clustering for making accurate document clustering using weighted technique.
2. A comparative performance study with fuzzy c-means clustering and strength evaluation of the proposed methodology.
3. A new technique for document domain identification with less resource consumption (running time) as compared to traditional document clustering approach.

#### **REFERENCES**

- [1] Yang Yan, Lihui Chen, William-Chandra Tjhi, "Fuzzy semi-supervised co-clustering for text documents", Fuzzy Sets and Systems 215 (2013) 74–89, 2012 Elsevier B.V. All rights reserved.
- [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] Luís Filipe da Cruz Nassif and Eduardo Raul Hruschka, "Document Clustering for Forensic Analysis: An Approach for Improving Computer Inspection", IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY, VOL. 8, NO. 1, JANUARY 2013.
- [4] D. Renukadevi , S. Sumathi, "TERM BASED SIMILARITY MEASURE FOR TEXT CLASSIFICATION AND CLUSTERING USING FUZZY C-MEANS ALGORITHM", International Journal of Science, Engineering and Technology Research (IJSETR), Volume 3, Issue 4, April 2014.



- [5] I-Jen Chiang, Charles Chih-Ho Liu, Yi-Hsin Tsai, and Ajit Kumar, "Discovering Latent Semantics in Web Documents Using Fuzzy Clustering", IEEE TRANSACTIONS ON FUZZY SYSTEMS, VOL. 00, NO. 0, 2015.
- [6] Athman Bouguettay, Qi Yu, Xumin Liu, Xiangmin Zhou, Andy Song, "Efficient agglomerative hierarchical clustering", Expert Systems with Applications 42 (2015) 2785–2797.
- [7] V. Loia, W. Pedrycz, and S. Senatore, "Semantic web content analysis: A study in proximity-based collaborative clustering," IEEE T. Fuzzy Systems, vol. 15, no. 6, pp. 1294–1312, 2007.
- [8] H. L. Larsen, "An approach to flexible information access systems using soft computing," in Proc. of the 32nd Annual Hawaii International Conference on System Sciences, Hawaii, 1999, p. 231.
- [9] W. B. Frakes and R. Baeza-Yates, Information Retrieval Data Structures and Algorithms. Englewood Cliffs, NJ: Prentice Hall, 1992.
- [10] S. Park, D. U. An, B. R. Cha, and C. W. Kim, "Document clustering with cluster refinement and non-negative matrix factorization," in Proceedings of the 16th International Conference on Neural Information Processing, Bangkok, Thailand, 2009, pp. 281–288.
- [11] T. Kohonen, Self-Organization Maps. Berlin Heidelberg: Springer- Verlag, 1995.
- [12] J. MacQueen, "Some methods for classification and analysis of multivariate observations," in Proceedings of the fifth Berkeley symposium on mathematical statistics and probability, vol. 1. Berkeley, CA: University of California Press, 1967, pp. 281–297.
- [13] A. K. Jain and R. C. Dubes, Algorithms for Clustering Data. Prentice Hall, 1988.
- [14] S. Lu and K. Fu, "A sentence-to-sentence clustering procedure for pattern analysis," IEEE Transactions on Systems, Man and Cybernetics, vol. 8, pp. 381–389, 1978.
- [15] O. Zamir and O. Etzioni, "Web document clustering: a feasibility demonstration," in Proceedings of 19th international ACM SIGIR conference on research and development in information retrieval (SIGIR 98), 1998, pp. 46–54.
- [16] P. Lingras, R. Yan, and C. West, "Fuzzy c-means clustering of web users for educational sites," Lecture Notes in Computer Science, vol. 2671, pp. 557–562, 2003.
- [17] R. R. Papalkar and G. Chandel, "Fuzzy clustering in web text mining and its application in iee abstract classification," International Journal of Computer Sciences and Management Research, vol. 2, no. 2, pp.1529–1533, 2013.
- [18] A. B. Raut and G. R. Bamnote, "Web document clustering using fuzzy equivalence relations," Journal of Emerging Trends in Computing and Information Sciences, vol. 2, pp 22-27, 2010.