

# A SURVEY ON CLUSTERING SENTENCE-LEVEL TEXT USING A NOVEL FUZZY RELATIONAL CLUSTERING ALGORITHM

**Mr. Bhosale A.T.**

Computer Engineering  
Dattkala Group of Institutes,  
Bhigwan, India.  
bhosaleamit3671@gmail.com

**Prof. S.S. Bere**

Computer Engineering  
Dattkala Group of Institutes,  
Bhigwan, India.  
sachinbere@gmail.com

## ABSTRACT

In the world of pattern recognition and information retrieval clustering algorithm is used find the location of information. Clustering algorithm plays an important role when there is little or no any knowledge of the pattern can be derived from the large data. There are different fuzzy technique like fuzzy c-means can be used for pattern recognition. Clustering text at the document level is well established in the Information Retrieval (IR) literature, where documents are typically represented as data points in a high dimensional vector space in which each dimension corresponds to a unique keyword, leading to a rectangular representation in which rows represent documents and columns represent attributes of those documents. In the literature there exist several approaches to form the clustering. The clustering tool, simfinder, organizes small parts of information from multiple documents in to single cluster. Simfinder place the same type of data into single cluster is usefull for the subsequent content selection or generation component to reduce each cluster to a single sentence, either by extraction or by reformulation.

## Keywords

Frecca algorithm, fuzzy clustering, pattern recognition.

## 1. INTRODUCTION

There are many text processing activities to find the desired text from large sentences using sentence clustering. There are many types of clustering algorithms, when data object is in the exactly in only one group in a set of data objects is called partitional algorithm and for clustering analysis hierarchical clustering is uses. Hierarchical clustering algorithms were developed for to overcome the disadvantages of partitional clustering algorithms.

### 1.1 Page Ranking algorithm

Page ranking algorithm when we apply to the cluster and interpreting the Page-Rank score of an object within some cluster as a likelihood, then use the

Expectation-Maximization (EM) framework to determine the model parameters (i.e., cluster membership values and mixing coefficients). The result is a fuzzy relational clustering algorithm which is generic in nature, and can be applied to any domain in which the relationship between objects is expressed in terms of pairwise similarities.

## 1.2 Fuzzy Relational Clustering Algorithms (FRECCA)

A fuzzy relational clustering concept is used to produce clustered sentences, in which same content is present in same cluster. The output of clustering represent that there is tight connection between the data element. This algorithm that is a novel fuzzy relational clustering algorithm (FRECCA) is proposed by Andrew Skabar and Khaled Abdalgar. This algorithm is divided into three steps: Initialization, Expectation and Maximization.

## 2. LITERATURE SURVEY

### 2.1 History

Clustering text at the document level is well established in the Information Retrieval (IR) literature, where documents are typically represented as data points in a high dimensional vector space in which each dimension corresponds to a unique keyword. The conventional fuzzy clustering approaches based on prototypes or mixtures of Gaussians are generally not applicable to sentence clustering. The project is to develop an application which is used to clustering sentence level text. In this project implement a new fuzzy clustering algorithm that operates on relational input data. The algorithm uses a graph representation of the data and operates in an Expectation-Maximization framework. The algorithm HFRECCA is capable of identifying overlapping clusters of semantically related sentences

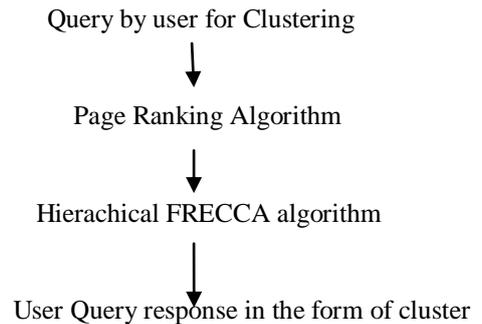
## 2.2 Existing System

In existing system high dimensional vector space is used for information retrieval in which documents are represented as data points. Each data point is similar to unique keyword in the form of rectangular. In this rectangular form rows of the rectangle represents the documents and column represents the attribute of the document. In existing system vector space technique is very useful. In this semantic measure technique is used. In this technique similarity measures such as cosine similarity, then also apply relational clustering algorithms such as Spectral Clustering and Affinity Propagation, which take input data in the form of a square matrix where is the relationship between the data object. To distinguish it from attribute data, refer to such data as relational data. A broad range of hierarchical clustering algorithms can also be applied. The vector space model has been successful in IR because it is able to adequately capture much of the semantic content of document-level text. This is because documents that are semantically related are likely to contain many words in common, and thus are found to be similar according to popular vector space measures such as cosine similarity, which are based on word co-occurrence. What is important to note is that these measures do not represent sentences in a common metric space, and this means that prototype-based clustering algorithms such as those described above are generally not applicable. The topic of interest, therefore, is fuzzy relational clustering, i.e., fuzzy clustering based on (pairwise) relational input data.

## 2.3 Proposed System

In the proposed system the FRECCA algorithm is replaced by the hierarchical FRECCA called as HFRECCA. The HFRECCA is an improvement over the FRECCA algorithm. HFRECCA algorithms give the best solution for clustering on the data objects. HFRECCA divides the data objects into the number of clusters. HFRECCA uses the page ranking algorithm to rank the retrieved pages. HFRECCA is capable of identifying softer clusters than ARCA, without sacrificing performance as evaluated by external measures. It is important to have a data mining system that can mine multiple kinds of patterns to accommodate different user expectations or applications. Data mining systems should be able to discover patterns at various granularities i.e. different levels of abstraction.

Following figure represents the flow of HFRECCA clustering process.



## 3. CONCLUSION

This paper gives a survey on different clustering techniques for clustering sentences. Performance of clustering techniques is totally dependent on the quality of input we are taking and similarity measures that are considered.

## 4. ACKNOWLEDGMENTS

I am very thankful to Prof. Bere S.S. who gives me very important guidelines during this survey.

## 5. REFERENCES

- [1] R. Kosala and H. Blockeel, "Web Mining Research: A Survey," ACM SIGKDD Explorations Newsletter, vol. 2, no. 1, pp. 1-15, 2000.
- [2] V. Hatzivassiloglou, J.L. Klavans, M.L. Holcombe, R. Barzilay, M. Kan, and K.R. McKeown, "SIMFINDER: A Flexible Clustering Tool for Summarization," Proc. NAACL Workshop Automatic Summarization, pp. 41-49, 2001.
- [3] Pabitra Kumar Dey, Gangotri Chakraborty, and Suvobrata Sarkar "Cluster Detection Analysis Using Fuzzy Relational Database" . International Journal of Information and Electronics Engineering, Vol. 3, No. 2, March 2013
- [4] Rui Xu, Student Member, IEEE and Donald Wunsch II, Fellow, IEEE, "Survey of Clustering Algorithms", IEEE TRANSACTIONS ON NEURAL NETWORKS, VOL. 16, NO. 3, MAY 2005
- [5] J.C. Dunn, "A Fuzzy Relative of the ISODATA Process and its Use in Detecting Compact Well-Separated Clusters," J. Cybernetics, vol. 3, no. 3, pp. 32-57, 1973
- [6] Tapas Kanungo, Senior Member, IEEE, David M. Mount, Member, IEEE, Nathan S. Netanyahu, Member, IEEE, Christine D. Piatko, Ruth Silverman, and Angela Y. Wu, Senior Member, IEEE, "An Efficient k-Means Clustering Algorithm: Analysis and Implementation", IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL. 24, NO. 7, JULY 2002