

A Survey of Clustering Sentence-Level Text Using A Novel Fuzzy Relational Clustering Algorithm

Mr. Amit T. Bhosale¹, Mr. Bere Sachin S.²

^{1,2} Department of Computer Engineering

^{1,2} Dattakala group of Institutes, Bhigwan, India

Abstract- A Cluster is a group of object which is forms on the basis of similarities of cluster and dissimilarities of other clusters. Clustering can be applied in many research areas. Clustering algorithm is used to arrange the large amount of data in to group, which are the clusters. Sentence clustering helps to avoid content overlapping problem. A simifinder which is clustering tool that arrange small parts of information from multiple document into single cluster, is useful for the subsequent content selection or generation of component to reduce each cluster to a single sentence by extraction. Clustering text at the sentence level is well established in the Information Retrieval literature, in which large documents are considered as data objects in a multi-dimensional vector space in which every dimension considered to be to a unique keyword, propose to a rectangular representation in which rows represent documents and columns represent attributes of those documents. In proposed work advanced hierarchical fuzzy relational clustering algorithm at the sentence level is done. This algorithm is better than previous algorithm and produce desired output.

Keywords- Pattern Recognition; Page Rank Algorithm; Hierarchical Clustering; EM Algorithm; Fuzzy Relational Clustering

I. INTRODUCTION

In many text processing activities sentence clustering plays very important role. For example, various authors have argued that incorporating sentence clustering into extractive multi-document summarization helps avoid problems of content overlap, leading to better coverage. Data mining used in very research area where clustering plays very important role. However, sentence clustering can also be used within more general text mining tasks. Irrespective of the specific tasks like summarization most documents will contain interrelated topics or themes, and many sentences will be related to some degree to a number of these.

A. Clustering Algorithms

Clustering algorithm plays very important role in data mining. Each algorithm possesses the distinct feature from another clustering algorithm. Each clustering algorithm forms

the cluster. Each cluster contains the data objects which are very close to each other or similar to each other. Clustering plays an important role in different area like, Image Processing, bioinformatics, etc. The main objective of clustering algorithm is to form a cluster in a such way that each cluster must contain data object which are different from other formed cluster.

B. Finding Similarities

Each cluster contains the similar data objects. Data objects of one cluster are different from other cluster. These data objects are retrieved from one or more documents. The main aspect here is that to find the similar data objects from multiple documents to form a cluster. To find the similarity between the different data objects tools like simifinder is used. Highly similar data objects forms the tight cluster.

II. MOTIVATIONAL SURVEY

A. FRECCA Algorithm

A. Skabar[1] presents data is to be clustered in the form of graph representation. It uses page rank algorithm. FRECCA means fuzzy relational clustering algorithm uses the graph centrality algorithm. This algorithm find how node are important in graph. FRECCA create the pair of data objects to find the similarity between data objects.

B. SIMIFINDER

V. Hatzivassiloglou[2] presents clustering tool to find the similarity between the data objects to form the cluster. Simifinder refers one or multiple documents and extracts the text from it, then it manage extracted text into the hard cluster. Simifinder puts the text which is very close to each other in the same cluster. Here we puts the suggestion for improvement of components used by simifinder.

C. Pattern Recognition using Fuzzy c-means technique

Samrjit Das[3] This paper state that Fuzzy c-means technique is primary used for analysis of large geostatistical

data. This technique forms the partitions called as fuzzy partitions and prototypes for any set of numerical data. These fuzzy partitions are useful for suggesting substructure in unexplored data. Primly due to use of human perception and lack of standard mathematics in patter recognition system forms a complex system. Pattern recognition mainly follows the three steps preprocessing, feature extraction and selection. Euclidean distance, Hamming distance are used in fuzzy c-means algorithm to find the membership values of objects in clusters.

D. Centroid-based summarization of multiple documents.

Dragomir R.Radev [4] present the summarization tool called as MEAD. MEAD has three basic components these are feature extractor, sentence scorer and sentence reranker.

III. EXISTING SYSTEM

Existing system uses the high dimensional vector space model 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 which documents corresponds to rows and attributes of documents corresponds to column of the rectangular. Vector space technique used in the existing system semantic measure technique such as cosine similarity is used. After this technique relational clustering algorithm such as spectral clustering and affinity propagation, which takes input data in the form of a square matrix where is the relationship between the data objects.

IV. PROPOSED SYSTEM

In proposed system hierarchical fuzzy relational clustering algorithm (HFRECCA) is used which is extension of FRECCA algorithm in the existing system. HFRECCA has advanced feature over the FRECCA algorithm. HFRECCA uses the page rank algorithm for ranking of pages which are retrieved. HFRECCA is better solution over the FRECCA algorithm. HFRECCA algorithm finds the softer clusters than ARCA, with a good performance as evaluated by external measures. Data mining system should be able to discover patterns at various levels of abstraction.

Following figure shows the steps of HFRECCA Clustering process.

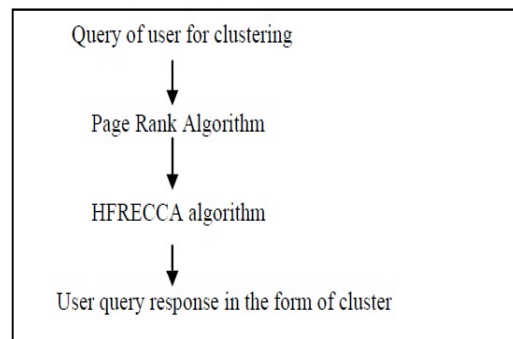


Fig. 1. Steps of Hierarchical FRECCA clustering process

A. Steps of algorithm

- 1) Let each sentence be a cluster and let the membership ship degree of each clustering belonging to its another sentence level text equal to 1.
- 2) Find the proximity matrix of distances.
- 3) Assign sequence numbers to clustering and then perform fuzzy relational clustering algorithm for sentence clustering.
- 4) Check whether the membership value is smaller than perticular threshold value or not using value obtaine from FRC if yes then compute the new membership value and update value and update the mixing coefficient. Otherwise go to next step
- 5) Combine the first sentence cluster with second cluster in to a new sentence cluster.
- 6) Then again find the degree of similarity between any two sentence cluster s is smaller than the threshold value, then stop.
- 7) Otherwis we go to step 5 and repeat the process.

V. SYSTEM ARCHITECTURE

Text data in the form of xml file i.e. hierarchical form is clustered using the Hierarchical FRECCA algorithm. The output of the HFRECCA will be the cluster, which are grouped from text data present in the given document.

Following figure shows the system architecture of HFRECCA algorithm.

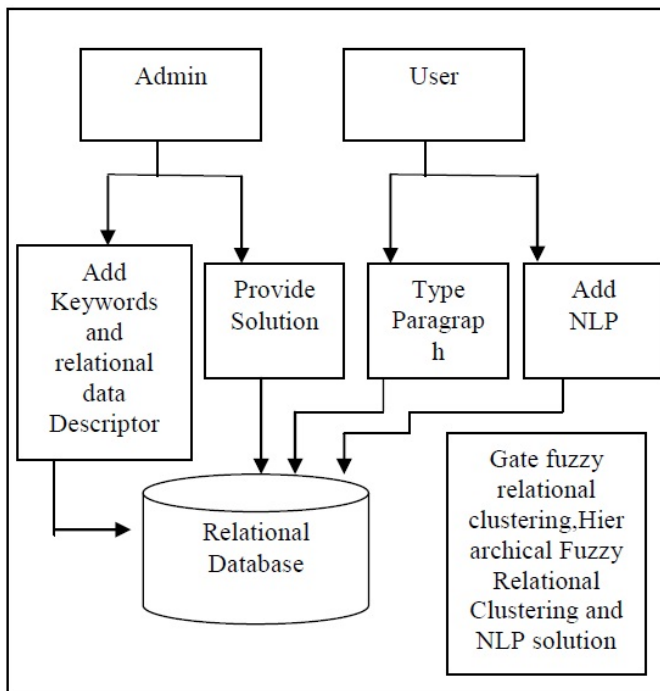


Fig. 2. System architecture of Hierarchical FRECCA

VI. CONCLUSION

This paper gives the use of Hierarchical FRECCA algorithm which finds the sentences which are semantically close to each other. It also useful for finding overlapping clusters. It finds the cluster of several clusters of same meaning.

ACKNOWLEDGMENT

I am very thankful to my guide Prof. Bere S.S. for his measure support and important advice for this paper.

REFERENCES

- [1] A. Skabar, "Clustering Sentence –Level Text Using A Novel Fuzzy Relational Clustering Algorithm", IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING ,VOL.25,NO.1,JANUARY 2013
- [2] V. Hatzivassiloglou, J. L. Klavans, M. L. Holcombe, R.Barzilay,M.Kan and K. R. McKnown, " SIMIFINDER; A Flexible Clustering Tool for Summarization", Proc.NAACL workshp automatic summarization,pp.41-49,2001.
- [3] Samarjit Das, "Pattern recognition using the fuzzy c-means techniu.e." International Journal of Energy, Information and Communication Vol.4,Issue 1,February,2013

- [4] Dragomir R. Radev, Hongyan Jing, Magorzata stys, Daniel Tam, "Centroid-based summarization of multiple documents", Information Processing and Management 40 (2004)918-938