

An Adaptive Hierarchical Clustering Algorithm for Segmenting Sentence level Text

¹Gandikota Gopi, ²Mrs.T.Suneetha Rani

¹M.Tech Student, QIS College of Engineering & Technology, Ongole.

²Associate Professor, Dept. Of CSE, QIS College of Engineering & Technology, Ongole.

ABSTRACT:

Text segmentation is designed to group documents with high levels of similarity. It has found applications in several fields of text mining and data retrieval. The digital data accessible nowadays has steadily grown in tremendous volume and retrieving useful information from that is quite challenge. Text clustering has discovered an important usage to organize information and to extract useful information from the available corpus. In this proposed system, we have method for clustering the text documents. In the initial phase characteristics are chosen using a preprocessing based method. In the next phase the extracted keywords are clustered by means of hybrid algorithm. In this proposed work, graph based Fuzzy EM framework is implemented to cluster the sentences with the corpus. Experimental results show proposed approach has better performance in terms of cluster rate and time are concern.

I. INTRODUCTION

Text mining is a crucial process with respect to information retrieval [1]. Text mining comprises of a wide array of processes like text clustering, classification, text summarization and automatic organization of text documents. Documents that may on the net are increasing each day which is the most of those are loosely structured. Clustering has turned out to be an important and used widely text mining tool to structure these documents to make sure that similar documents are clustered into your same group and dissimilar documents are separated into different groups. Text clustering is undoubtedly an unsupervised learning method where similar documents are grouped into clusters. It has been considered as method of finding sets of similar objects inside the data. The similarity between objects is calculated using various similarity functions. Clustering might be very beneficial numerous text domains, wherein the objects to remain clustered are of assorted types which can include paragraphs, sentences, documents or terms. Clustering makes it possible to organize the fax which can further helped to improve information retrieval and support browsing. The nature of any text mining methods namely classification and clustering is extremely rely on the noisiness of a given features whose purpose is regarding the process. Therefore, the

basic features ought to be selected effectively to further improve the clustering quality. A few of the frequently used feature selection methods are document frequency based selection method, term strength and entropy based ranking. Following the features have also been selected, any text mining tasks which can include classification, clustering, summarization can easily be applied.

The method of text clustering usually contains two phases. The very first phase is keyword extraction leading to clustering those keywords. When it comes to the first phase, the text documents are usually represented by using the vector space model, where each row can be seen as documents, and of course the column relates to the assorted attributes of the document. Like the document size increases, overall performance the VSM decreases [5]. To beat this, within this paper, we propose a brand new keyword extraction method along with a clustering algorithm.

The primary requirements associated with a good clustering algorithm are:

- 1) This relationship between words ought to be displayed prominently inside the document model.
- 2) Clusters ought to be identified with a meaningful label.
- 3) The high dimensionality files ought to be reduced efficiently.

The intention of text summarization will be to present the most essential information inside a shorter variety of the first text while keeping its main content and assists the user to quickly understand large volumes of data. Text summarization addresses both the concern about picking the biggest sections of text plus the issue of generating coherent summaries. This system is significantly differ from that of human based text summarization since human can capture and relate deep meanings and themes of text documents while automation of that is a method is amazingly problematic to implement. Automatic text summarization researchers since Luhn work [6], they're aiming to solve or at least relieve that problem by proposing techniques for generating summaries. Regardless of the precise task (e.g., summarization, text mining, etc.), most documents will contain interrelated topics or themes, and several sentences will probably be linked to some extent to your wide range of these. The hard work described in this paper is motivated through belief that successfully with the ability to capture such

fuzzy relationships will help you see a rise within the breadth and coverage problems that sentence clustering can possibly be applied. However, clustering text with the sentence level poses specific challenges not present when clustering larger segments of text, namely documents. We now highlight some important differences between clustering at both of these levels, and examine some existing approaches to fuzzy clustering.

II BACKGROUND AND RELATED WORK

Contradiction Analysis is probably one of the popular text-mining operations in which a document whose content is contradictory towards the theme of an arrangement documents is identified [4]. This is actually a means to identifying Outlier documents that don't enlighten the whole sense conveyed by other documents. The majority of the existing techniques perform document-level comparisons, ignoring the sentence-level semantics, often resulting in lack of vital information. Applications in domains like Defence and Healthcare require large amounts of accuracy and identification of micro-level contradictions are vital. In this particular paper, we propose an algorithm for identifying contradictory documents using sentence-level clustering technique in association with an optimization feature. A new visualization scheme is additionally suggested to present the outcome to an end-user.

Sentence clustering plays a pivotal aspect in theme-based summarization, which discovers topic themes known as the clusters of highly related sentences to refrain from redundancy and canopy more diverse information [5]. Clearly as the period of sentences is small and of course the content it contains is restricted, the bag-of-words cosine similarity traditionally utilized for document clustering no longer is suitable. Special treatment for measuring sentence similarity is required. In this post, we go about it in the sentence-level clustering problem. After exploiting concept- and context-enriched sentence vector representations, we develop two co-clustering frameworks to reinforce sentence-level clustering for theme-based summarization—integrated clustering and interactive clustering—both allowing word and document to be an explicit function in sentence clustering as independent text objects as an alternative to using word or concept as elements of a sentence inside a document set. In all framework, we experience two-level co-clustering (i.e., sentence-word co-clustering or sentence-document co-clustering) and three-level co-clustering (i.e., document-sentence-word co-clustering). Compared against concept- and context-oriented sentence-representation reformation, co-clustering shows a

clear advantage in each of intrinsic clustering quality evaluation and extrinsic summarization evaluation conducted upon the Document Understanding Conferences (DUC) datasets.

Multi-document summarization aims to generate a concise summary which has salient information given by a multitude of source documents. With this field, sentence ranking has hitherto been the subject of most concern [4-7]. Since documents often recognise a large number of topic themes with each theme represented using a cluster of highly related sentences, sentence clustering was recently explored inside the literature so that you can provide more informative summaries. Existing cluster-based ranking approaches applied clustering and ranking as an isolated condition. Subsequently, the ranking performance will surely be inevitably influenced by the clustering result. With this paper, we propose a reinforcement approach that tightly integrates ranking and clustering by mutually and simultaneously updating another to ensure the performance

of both might be improved. Experimental results toward the DUC datasets demonstrate its effectiveness and robustness.

Clustering algorithms can be used in plenty of Natural Language Processing (NLP) tasks. They are indeed to become often used effective tools to employ to discover categories of similar linguistic items [10]. In this particular exploratory paper, we propose a brand new clustering algorithm to automatically cluster together similar sentences driven by sentences part-of-speech syntax. The algorithm generates and merges together the clusters making use of a syntactic similarity metric based upon a hierarchical organization of a given parts-of-speech. We demonstrate benefits features of the algorithm by implementing it in an issue type classification system, as a way to find the positive or negative impact of various changes to the algorithm.

Clustering text along at the document range is more successful among the Information Retrieval (IR) literature, where documents are typically represented as points of information inside a high dimensional vector space by which each dimension relates to an exceptional keyword [6], resulting in an oblong representation through which rows represent documents and columns represent attributes of those documents (e.g., tf-idf values of this very keywords). The sort of data, which we check with as "attribute data," is amenable to clustering by the large range of algorithms. Since points of information lie in a metric space, we are able to readily apply prototype based algorithms for instance k-Means [7], Isodata [2], Fuzzy c-Means

(FCM) [3], [4] and of course the closely related mixture model approach [1], all of these represent clusters in relation to parameters which can include means and covariances, and hence assume a typical metric input space. Since pairwise similarities or dissimilarities between points of information can readily be calculated that are caused by the attribute data using similarity measures namely cosine similarity, we are able to also apply relational clustering algorithms which can include Spectral Clustering [12] and Affinity Propagation [13], which take input data comprising of a square matrix $W = w_{ij}$, where w_{ij} happens to be the (pairwise) relationship connecting the i th and j th data object.

III. PROPOSED FRAMEWORK

The proposed system is based on Hierarchical fuzzy relational clustering algorithm. We first describe the use of distance as a general graph centrality measure, and review the objective function, Optimizing Memberships, Optimizing Weights and Hierarchical Fuzzification Degree clustering approach. We then describe how fuzzification can be used within the hierarchical framework to construct a complete relational fuzzy clustering algorithm.

A. Fuzzy Objective function

The objective function of Fuzzy is to classify a data point, cluster centroid has to be closest to the data point of membership for estimating the centroids, and typicality is used for alleviating the undesirable effect of outliers. The function is composed of two expressions:

- The first is the fuzzy function and uses a distance exponent,
- The second is possibilistic function and uses a typical fuzziness weighting exponent; but the two coefficients in the objective function are only used as exhibitor of membership and typicality. The objective function is to discover nonlinear relationships among data, kernel methods use embedding mappings that map features of the data to new feature spaces.

Input : Documents D,

Output: Clusters.

Algorithm:

1. Read Documents D.
2. For each d in D
 - Do
 - for each line l in d
 - do
 - if(l.endswith("."))
 - List(d,tokens):= StringTokenizer(l);
 - done
 - done
3. Find normalization to each tokens list in
4. the corresponding phrase of document d.
5. **Tokenization**
 - -Divide the given text into words/phrases
 - by following the below constraints
 - 1) do not split on hyphens,
 - 2) do not split on single quotation marks,
 - 3) do not split on commas, and
 - 4) do not split on parentheses and brackets.

5. Stemming : It is a process of constructing a root word from the given word.

6. Stop word removal : Take a pre defined stop word list from wordnet and remove those words from given Corpus.

7. Duplicates removal : Apply Duplicate Removal Algorithm to remove the duplicate strings.

$$\begin{aligned}
 L(t_1, t_2, \dots, t_n) &= \prod_{i=1}^n f(t_i) \\
 &= \prod_{i=1}^n \lambda e^{-\lambda t_i} \\
 &= \lambda^n \cdot e^{-\lambda \sum_{i=1}^n t_i}
 \end{aligned}$$

Fuzzy Hierarchical Clustering:

Start

Assign sqrt(n) objects as initial for next step

Get initial clusters using HAC algorithm

Make these clusters as initial k document clustering

Assign Cluster centers randomly

For each object in document collection

If new object

then

check the nearest mean,

join the cluster with nearest mean,

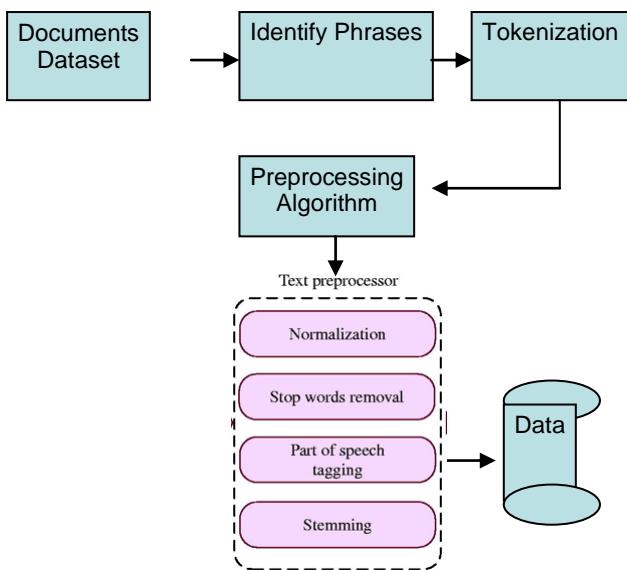
calculate the new mean for the cluster

end if

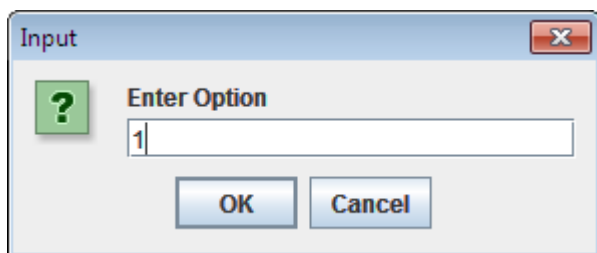
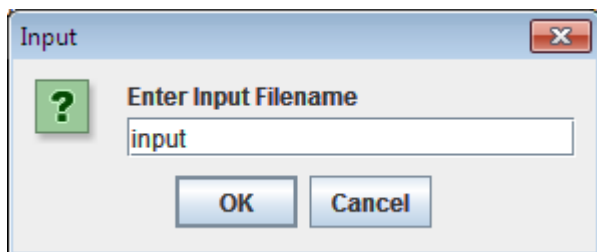
end for

```

Given:
A set X of objects  $\{x_1, \dots, x_n\}$ 
A distance function  $dist(c_i, c_j)$ 
for  $i = 1$  to  $n$ 
     $c_i = \{x_i\}$ 
end for
 $C = \{c_1, \dots, c_n\}$ 
 $l = n+1$ 
while  $C.size > 1$  do
    -  $(c_{min1}, c_{min2}) = \text{minimum } dist(c_i, c_j)$  for all  $c_i, c_j$  in  $C$ 
    - remove  $c_{min1}$  and  $c_{min2}$  from  $C$ 
    - add  $\{c_{min1}, c_{min2}\}$  to  $C$ 
    -  $l = l + 1$ 
end while
    
```



Experimental Results:



String one : " java is a nice programming environment".
 String Two: " The Java Runtime Environment (

JRE) is what you get when you download Java software. The JRE consists of the Java Virtual Machine (JVM), Java platform core classes, and supporting Java platform libraries. The JRE is the runtime portion of Java software, which is all you need to run it in your Web browser."

Score :0.5801954953637369

String one : " java is a nice programming environment".

String Two: " Java is a computer programming language. It enables programmers to write computer instructions using English based commands, instead of having to write in numeric codes. Once a program has been written, the high-level instructions are translated into numeric codes that computers can understand and execute."

Score :0.4425710648397979

String one : " java is a nice programming environment".

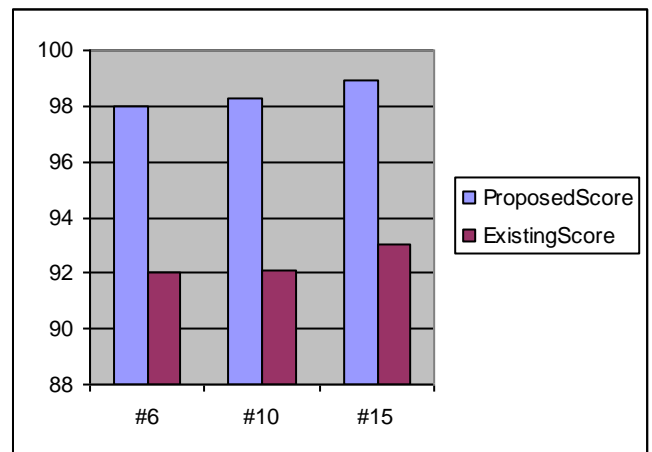
String Two: " A high-level programming language developed by Sun Microsystems. Java was originally called OAK, and was designed for handheld devices and set-top boxes. Oak was unsuccessful so in 1995 Sun changed the name to Java and modified the language to take advantage of the burgeoning World Wide Web."

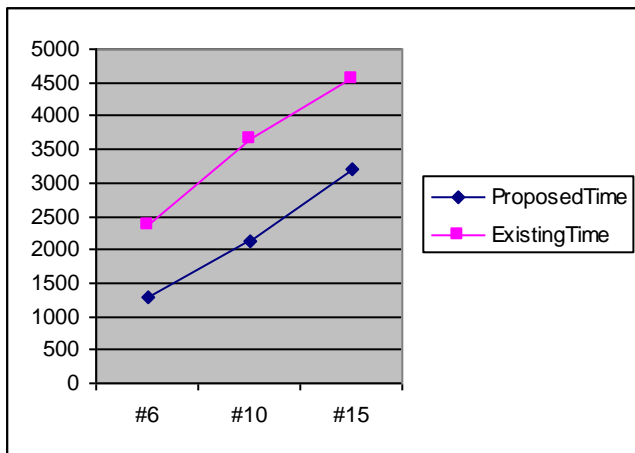
Score :0.6572729784684463

String one : " java is a nice programming environment".

String Two: " jAvA is one of the robust programming language".

Score :0.9999263993686861





IV. Conclusion

In this particular paper we present a sentence clustering based multidocument summarization system whose performance is comparable to the peak performing multi-document summarization systems participated on task2 on DUC 2004. We also investigate alternatively variants of this product. Our work specializes the style of causing successful

clustering based summarization and of course the related issues namely learn how to cluster sentences, how you can order clusters as well as how to select representative sentences seen from the clusters. Our experiment has verified that behavior a clustering based multi-document summarization is possible competitive with one of the best top performing multidocument clustering systems. To create the our system portable to new domain and new language, we have not apply stemming upon the input therefore we didn't incorporate features namely length, sentence position, cue phrase in this particular work though these features are tested to effective among the news domain. So, concerning domain independency and language independency our approach is likewise better. In this proosed work, graph based Fuzzy EM framework is implemented to cluster the sentences with the corpus. Experimental results show proposed approach has better performance in terms of cluster rate and time are concern.

REFERENCES:

- [1] 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.
- [2] H. Zha, "Generic Summarization and Keyphrase Extraction Using Mutual Reinforcement Principle and Sentence Clustering," Proc. 25th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 113-120, 2002.
- [3] D.R. Radev, H. Jing, M. Stys, and D. Tam, "Centroid-

Based Summarization of Multiple Documents," Information Processing and Management: An Int'l J., vol. 40, pp. 919-938, 2004.

- [4] R.M. Aliguyev, "A New Sentence Similarity Measure and Sentence Based Extractive Technique for Automatic Text Summarization," Expert Systems with Applications, vol. 36, pp. 7764-7772, 2009.
- [5] R. Kosala and H. Blockeel, "Web Mining Research: A Survey," ACM SIGKDD Explorations Newsletter, vol. 2, no. 1, pp. 1-15, 2000.
- [6] Ms. Seema V. Wazarkar, Ms. Amrita A. Manjrekar, "Text Clustering Using HFRECCA and Rough K-Means Clustering Algorithm", International Conference on Advances in Computer Engineering & Applications (ICACEA-2014) at IMSEC, GZB.
- [7] K.Sathishkumar, E.Balamurugan, and D.Kavin , "Sentence Level Clustering Approaches and its Issues in Various Applications", International Journal of Applied Research "