# Document Summarization Using Clustering and Text Analysis

**Mrs.ShahanaBano[1], B.Divyanjali[2], A.K M L R V Virajitha[3], M.Tejaswi[4]**

[1,2,3,4]*Department of CSE, K L E F, Guntur, India*
*\*Corresponding author E-mail: shahanabano_cse@kluniversity.in*

## Abstract

Document summarization is a procedure of shortening the content report with a product, so as to make the outline with the significant parts of unique record.Now a days ,users are very much tired about their works and they don't have much time to spend reading a lot of information .they just want the maximum and accurate information which describes everything and occupies minimum space.This paper discusses an important approach for document summarization by using clustering and text analysis. In this paper, we are performing the clustering and text analytic techniques for reducing the data redundancy and for identifying similarity sentences in text of documents and grouping them in cluster based on their term frequency value of the words. Mainly these techniques help to reduce the data and documents are generated with high efficiency.

*Keywords*:*clustering, word count, term frequency, sentence score, document summarization.*

## 1. Introduction

Text summarization is a outline of rundown of the source form of unique content while keeping its principle substance and help the client to rapidly comprehend the huge volumes of data. It is the process of extracting information from a given text .The main idea of summarization is toward removing unusual data from a given content. The synopsis is to discover subset of information which contains data of whole content. It finds principally the most critical striking and featured words.

Text summarization can be delegated into two sorts. They are abstract and extract. In extract,The synopsis comprising of various vital content units chosen from the input. It is adaptable and expands less measure of time when contrasted with the abstract. In abstract, it speaks to the submit paper of the article with the content units, which are created by reformulating the imperative units chose from unit.As per single to multi-document summarization[6] it describes in 3 stages for redundancy of data.It works on stigma words.

In this paper, we present the document summarization system, which is used as removing and reducing the data redundancy using clustering and it produces the effective summary at the end. We also implement by using text analysis which uses natural processing tools like dictionaries, pos tagger and word net. Term frequency is a factual measure utilized as a part of ascertaining importance of document. it enlightens something regarding the document all in all concerning a client inquiry. We utilize a term frequency model to numerically describe the significance of terms. This model is utilized to archive the featured data sentences from the reports.According to this k-means algorithm [7] which can be calculate the similarity between sentences. Here in this paper we used r program which reduces the data.Where there are techniques there for summarization. Like sentence ordering[9] and by utilizing grouping calculations here we done by compressing infor-

mation in view of r program.According to[15]this multi-document summarization it has different types of documents.
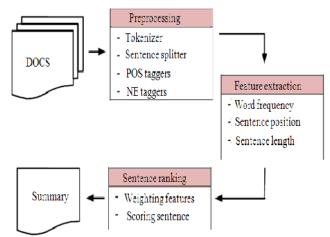
## 2. Literature Survey

Multi-document summarization[1] using symmetric non negative matrix factorization. It mainly describes about recognizing and adopting the data and also identifies the contrasts among report and it mainly focus on enlightening the substances however it could be reasonable.

The clustering algorithm feature profile[2] is used to extract the most important sentences from multiple documents. In clustering based documentation performance based on factors like i) clustering sentence ii)clustering ordering iii) selection of representative sentences.

Improving Multi document using text classification[3] mainly describes on novel summarizations and it proposes DUC structure for classifying plentiful of data in categories of documents.

Multi document Summarization using A[\*] search[4]. It proposes a model of algorithm for getting the best scoring summary. It also decomposes the model parameters for evaluation purpose. Automatic text summarization[8] describes about shorter form of text. Text can be divided into phases and eliminates the words which are not required. It proposes normalization where duplicates of data are not allowed. According to this keywords can be identified by using term frequency and inverse sentence frequency. Machine learning approach to sentence ordering[9] it describes about the summarization can be done by using one of the ordering technique that is sentence ordering where it states that ordering can be on position of sentences in a given document. Sentence Ordering based on Cluster Adjacency[10] which describes about the sentence ordering which is not easy in multi document summarization. It proposes sentence ordering method called cluster-adjacency ordering based. Text summarization using clustering[11] defines that it compares similarity between documents and finds the loca-

tion of the sentence for eliminating the repeated sentences. Multi-document summarization using tf-idf[12] allow the users to search for files and for upload files. Where term frequency helps for summarizing data very easily.Sentence fusion[13] for multi-document news summarization describes that clusters that means divides data in the document into subparts where the input document divided into themes. Where sentence fusion is a technique where the data selection can also be done from multiple documents. Where Towardscoherent of multi-document[14] describes that extractive summarizer is based on a pipeline for selection of sentence from documents. Where selecting a sentence is not an eassy thing where some of the data may be incoherent. From [15]Multi-document summarization by sentence extraction. It describes about what are the techniques used in single document that are can be used for multi-document summarization. And also it describes about different types of documents basically those are about sections.Where From[16]Improving Chronological Sentence Ordering by Precedence Relation states that the proper arrangement has to be done to get a good summary of data from multiple documents. It describes about text structure which is also important.

## 3. Proposed Work



Here, we are proposing the document summarization based on the clustering and text analysis for reducing the data and for extracting the information queried by the user.

Clustering is a process of partitioning a set of data into a set of meaningful sub-classes, called clusters. it helps user to understand the natural grouping into structure in a data set. As we know that clustering has some techniques which are very important for summarization. There are as follows.

1.Agglomerative hierarchical clustering.
2. K means clustering.
3. DBSCAN (Density Based Spatial clustering application with noise).

These are the known clustering techniques which are used for summarizing data. We have done summarization by using k means clustering which is very easy for summarizing the document or a given input data. Basically we all know that k means is a unsupervised learning method. Where we have done k means by using r program. Where the dataset in r contains the data where summarization can be done. Where the data is defined in r program as iris in r which contains length and width. Here we have constructed an algorithm the total number of words can be reduced. That means unwanted words for the sentence or a para of given information can eliminated. Where text analysis or we can call it as a natural language process which is also very important in summarizing data and speech recognition and also it is used for

other mechanisms. We use text analysis for removing stop words and other repeated words can be removed. That means it will also useful in making a document into a shorter version which makes the user easy way to understand. Where in this analysis they derived term frequency and inverse document frequency. Term frequency eliminates that any other document contain the same text it will removed or eliminated. As per this we have under gone through these analysis summarization is a method to make a shorter version of particular document that makes easy to understand.

### Advantages of k means

1. K means will compute fast when compared with other clustering techniques.
2. Simple unsupervised algorithm.
3. It can run many times to reduce difficulty in data grouping.

### Advantages of text analysis

1. Easily determined.
2. Eliminating the sentences in term frequency plays an important role.
3. Smart way for summarizing any kind of document.
4. And easy method for manipulating data.

### 3.1 Word Tokenization

The following is the R code for algorithm

By R basics with tabular data we have done o reduce the data for user purpose. R programming language is statistical analysis and it will analyze data quickly.

Firstly, we have to read the data into R console.
Text <- paste (" This is text summarization which analyses the data in a short way and is extracted in form of user query ").
The first step in processing text is tokenize_words function from the tokenizespackage to split the text.
Words<-tokenize_words(text).
We use the length function directly on words object
Length(words).
To view the words inside the document , we use
Length(words[[1]]).
We apply table function for using the result
Tab<- table(words[[1]]).
Tab<- data_frame(word=names(tab), count = as. Numeric(tab)).
Tab.

### 3.2 Sentimental Analysis

Breaking down the feeling or tone of what individuals are saying in regards to your organization on media. There are three methods for approaching this sort of wistful investigation. The first is extremity examination, where you just recognize if the tone of correspondences' is sure or negative. The second level is classification, where apparatus get all the more fine grained and recognize if nearly ones befuddled or furious.

### 3.3 Topic Modeling

Topic modeling is a helpful technique for recognizing overwhelming subjects in an immense range of records and for managing substantial corpus of content. For instance, may need to experience a large number of archives utilized as a part of enormous cases. This is the place subject demonstrating can prove to be useful.

# 4. Experimental Results and Analysis

The summarization system is measured with the no of inputs words in the source document and no of words in output summary file and its reduced words percentage is given below.
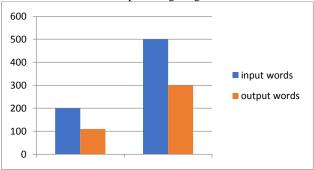


**Fig1**: reducing of word count.

# 5. Conclusion

Research has been done in field of text summarization. we have developed an approach which reduces the redundancy and it proves sentence simplification which means important data can be shown to the users .we proved and neglected the data redundancy and also generated an efficient number of documents.

# References

[1] Multi document based summarization "International journal of advanced research in electrical" vol3, issue 4,April 2004.

[2] Text features weighting for summarization of documents "International journal of computer science", Ahmed ridha,2012.

[3] .Improving Multi document summarization via text classification "springer", Trever kohn.

[4] .Multi document summarization using A* algorithm "Internation-al journal of computer science", Robert zaisagukas.

[5] .Klaus zechner" A literature survey on information extraction and text summarization", computational linguistics program, carneige Mellon university,1997.

[6] Chin-yiewlin and Eduard Hovy," from single to multi docu-ment summarization: A prototype System and its evaluation", proceedings of the ACL conference, Philadelphia, PA.2002.

[7] .Rene Arnulfo Garcia- Herandez and Yulia Ledeneva," word sequence models for single text summarization",IEEE,2009.

[8] .Jimmy Lin, "summarization", Encyclopedia of database system-sHeidelberg, germany : springer- verlag, 2009.

[9] A Machine Learning Approach to sentence ordering for multi-document summarization and its evaluation, university of Tokyo, japan.

[10] Sentence ordering based on cluster adjacency in multi-document summarization, institute for infocomm research Singa-pore, 119613.

[11] Text summarization using clustering technique by "Interna-tional journal of engineering trends and technology".

[12] Multi-document summarization using TF-IDF algorithm by "International journal of engineering and computer science".

[13] Sentence fusion for multi-document news summarization by Regina Barzilay* Massachusetts Institute of Technology.

[14] Towards Coherent Multi-Document SummarizationComputer Science & Engineering.

[15] Multi-Document Summarization By Sentence Extraction *Language Technologies Institute Carnegie Mellon University.

[16] Improving Chronological Sentence Ordering by Precedence Relation Naoaki OKAZAKI The University of Tokyo.