# Auto Text Summarization and Categorization

## Richa Tanna[1], Akanksha More[2], Ankita Thakur[3], Prof. D. J. Dhangar[4]

*[1,2,3] Dept. of Computer Engineering, Rajiv Gandhi Institute of Technology, Mumbai, Maharashtra, India*
*[4]Professor, Dept. of Computer Engineering, Rajiv Gandhi Institute of Technology, Mumbai, Maharashtra, India*

---------------------------------------------------------------------***---------------------------------------------------------------------

**Abstract -** *Automatic text summarization is undergoing wide research and gaining importance as the availability of online information is increasing. Automatic text summarization is to compress the larger original text into shorter text called as summary. Abstraction and Extraction are the two main methods to carry out text summarization. Our approach uses extractive method. Summarization by extraction involves identifying important features and extracting sentences based on their scores. 30 documents from news based URLs are used as input. After pre-processing the input document, five features are used to calculate their score for each sentence. In this paper fuzzy logic method is proposed for improvement in the extraction of summary sentences.*

**Key Words:** Categorization, Feature matrix, Fuzzy Logic, Deep Learning Algorithm, Text Summarization.

## 1. INTRODUCTION

What is a Summary?
Summary is basically the shorter version of the larger documents and it is just showing the most important part of the documents or giving us the brief idea of what the document has to convey.

Goal of Auto Text Summarization:

With the super growth in the quantity and the complex versions of the documents on the internet on the every single topic, it has become very important to provide more developed mechanism for the  user to find the exact and particular important information that a document want to convey. Text summarization has become an important and handy tool for the users by helping users to interpret the large volumes of text data that is available in the documents.

Summarization Methods:

There are two main important methods of summarization. Firstly, there is abstractive summarization and other one being the extractive summarization. Here our approach uses extractive method. Summarization is carried out by extracting the features from each sentences and based on the higher score of feature matrix, we generate the summary.

## 1.1 Purpose

Business associates, academic researchers, students, and almost every person go through many documents every day to keep him/her updated, and a large amount of their time is spent just to figure out what exactly the documents are about and what they want to convey. By extracting important sentences and creating summaries, it is possible to quickly assess whether or not a document is worth reading. Auto text summarization is most useful for students and authors. Imagine that you are able to automatically generate an abstract based for your research paper or for any chapter in a book in a clear and easy way that is faithful to the original source material. With the growing amount of data in the world, interest in the field of auto summary generation has been widely increasing as it is reducing the manual effort of a person working on large amount of documents.

## 1.2 Scope

Auto summarization involves reducing a text file into a passage or paragraph that can convey the important points of the text file. The searching of important information or to understand what the large text file wants to convey is very difficult job as well as it is very time consuming for the users thus to automatically extract the important information or summary of the text file comes in as a handy tool. This auto generated summary helps the users to reduce time instead of reading the whole text file and it provides useful information quickly from the large documents. In today's world to extract information from the World Wide Web is very easy, but extracted information is a huge text file. With the rapid growth of the World Wide Web, information overload is becoming a problem for an increasing large number of people. Auto summarization can be a very important solution to reduce the information overload problem on the web.

## 2. SYSTEM ANALYSIS

## 2.1 Existing System

Yan Liu et al have proposed a document summarization framework via deep learning model, which has demonstrated distinguished extraction ability in document summarization. The framework consists of concepts extraction, summary generation and reconstruction validation. A query-oriented extraction technique has been concentrated information distributed in multiple documents to hidden units layer by layer. Then, the whole deep architecture was fine-turned by minimizing the information loss in reconstruction validation part. According to the concepts extracted from deep architecture, dynamic programming was used to seek most informative set of

sentences as the summary. Experiments on three benchmark dataset demonstrate the effectiveness of the framework and algorithm.

## 2.2 Disadvantages of Existing System

The existing system is not user friendly and it is difficult to compute summaries using the existing methods and summaries suffer from inconsistencies and lack of balance results in lengthy summary. Detailed Information is not present. Existing methods does not give quick overview. It provides authors view not the user friendly view. Existing systems cannot summarize multiple documents of the same type. Existing systems cannot handle different languages.

## 3. PROPOSED SYSTEM

Proposed System consists of:

**1) Preprocessing:** Initially, the input to the proposed approach is a set of document from DUC 2002.Dataset that has to be summarized. The document utilized for text summarization is organized by a set of pre-processing steps like, sentence segmentation, stop words removal and stemming.

**2) Segmentation:** It is performed by identifying the delimiter commonly denoted by "." called as full stop. This step is used to separate the sentences in the document. It is mainly useful for the user to understand each individual sentence which is there in the document.

**3) Stop Words Removal:** Stop words are removed mainly to reduce the insignificant and noisy words. These are predefined words such as a, an, in, by, etc., are called stop words which are filtered out before the pre-processing phase from the documents.

**4) Stemming:** Stemming is process of bringing the word to its base or root form for example using words singular form instead of using the plural. It basically removes the prefix and suffix of the concerned word to get the base form. There are many more number of algorithms, which are called as stemmers used to perform the stemming process.

**5) Training Phase:** On behalf of the training phase, the proposed approach defines five features sets. The feature sets are listed as follows,

(i) Title Similarity Feature: The ratio of the number of words in the sentence that occur in title to the total number of words in the title helps to calculate the score of a sentence for this feature and it is calculated by the formula given below

$$Title\ Feature\ (f1) = \frac{S \cap t}{t}$$

(ii) POS Tagger Feature: To calculate the positional score of sentence, the proposed approach considers the following conditions. If the sentence given is in the starting of the sentence or the last in the sentence of the paragraph then the feature value $f2$ is assigned as 1. Else if the sentence is in the middle of the paragraph then the feature value of $f2$ is assigned as 0.

(iii) Term Weight Feature: The Term Frequency of a word will be given by TF (f, d) where f is the frequency of the given word and d is text present the document. The Total Term Weight is calculated by Term Frequency and IDF for a document .Here IDF denotes the inverse document frequency which just implies that the term is common or rare across all documents.

$$IDF(t, D) = \log\left(\frac{D}{d \in D : t \in d}\right)$$

(iv)Concept Feature: The concept feature from the text document is retrieved using the mutual information and windowing process. In windowing process a virtual window of size „k" is moved over document from left to right. Here we have to find out the co-occurrence of words in same window and it can be calculated by following formula,

$$f_4 \Rightarrow MI(w_i, w_j) = \log 2 \frac{P(w_i, w_j)}{P(w_i) \times P(w_j)}$$

(v) POS Tagger Feature: Part of speech tagging is the process of categorizing the words of text on the basis of part of speech category such as noun, verbs, adverb, adjectives, they belong to. Algorithms such as hidden Markova models, using dynamic programming are used to perform this task. The POS Tags on each document is feature five of the given documents.

**6) Association of deep learning with fuzzy logic:** The sentence matrix S= (s$_1$, s$_2$,.....s$_n$) which is the feature vector set having element as s$_i$ which is set contains the all the five features extracted for the sentence s$_i$. Here this set of feature vectors S will be given as input to deep architecture of RBM as visible layer. Some random values is selected as bias 3_ where i = 1, 2 since a RBM can have at least two hidden layer. The whole process can be given by following equation: S= (s$_1$,s$_2$,....,s$_n$).where s$_i$= (f$_1$,f$_2$,....f$_n$), i<=n where n is the number of sentences in the document.

**7) Feature Matrix:** Here sentence matrix where S=(s$_1$,s$_2$,....s$_n$) where s$_i$= (f$_1$,f$_2$,....f$_n$), i<=n is the feature vector. The five features are the main attributes of the proposed text summarization algorithm. The whole documents under consideration are subjected for the feature extraction and a set of features are extracted accordingly. Now based on the

collected features a feature matrix is formed by mapping the features values. The feature matrix is constructed according to the sentences extracted from the multiple documents. In addition to the five features, an additional attribute also associated with the feature matrix. The addition feature associated with the feature matrix is the class labels for each sentence. The Figure below represents the feature matrix of the set of documents under consideration.

## 3.1 Characteristics of Proposed System

The proposed system is User Friendly. It easily generates summary. There is no problem with the size of the document. It boosts the summary generation time.

## 4. SYSTEM ARCHITECTURE



**Fig -1**: Architecture

This method of summarization take in to account all the characteristics of any text such as length of the sentence, similarity of sentences with the title, similarity of the sentences with the key words, etc. and then it is given as an input to the fuzzy classifier system. After the fuzzy classifier system, all the rules that are needed for summary generation are entered, in the knowledge base of the system. After that a value from zero to two is obtained for each sentence. The obtained scores of each sentence in the output determine the degree of the importance of the sentence during the final summary generation. The input membership function for each feature is divided into three membership functions which are composed of insignificant values. The important sentences are then extracted using IF-THEN rules according to the feature criteria.



**Fig -2**: Overall block diagram of text summarization

## 5. SOFTWARE ENVIRONMENT

### 5.1 Front End

Windows XP, Windows 7, 8

Visual Studio 2010

Windows Operating System

### 5.2 Back End

Windows XP, Windows 7,8

Visual Studio 2010

MS SQL Server 2008

Windows Operating System

## 6. RESULT



**Fig -3**: Login



**Fig -4**: Main GUI Screen



**Fig -5**: Uploading a file



**Fig -6**: View an uploaded file



**Fig -7**: Pre Processing



**Fig -8**: Stop Word Clearing



**Fig -9**: Stemming



**Fig -10**: POS Tagger

**Fig -11**: Summary Generated



**Fig -12**: Feature Matrix Calculation

## 7. CONCLUSIONS

In this paper we have implemented the auto text summarization which involves extraction of features of all sentences using fuzzy logic method. We have extracted five different features of all the sentences. Later using the sentence score we have generated the feature matrix. Then based on the scores of individual sentences and final feature matrix value we have generated the summary. We have tested the system with multiple numbers of documents. The results obtained show that the use of fuzzy logic in summary generation improves the quality of summary generated compared to other method of summary generation. We have applied our method of summary generation for single document which could be extended for multi-document summarization. In the input we used multiple documents; these documents can be further categorized like sports, politics, weather, music etc. hence the system will have the categorization of the text document feature which will help to know the text belongs to which category.

## REFERENCES

[1] Allan Borra, Almira Mae Diola, Joan Tiffany T. Ong Lopez, Phoebus Ferdiel Torralba,Sherwin So, "Using Rhetorical Structure Theory in Automatic Text Summarization for Marcu-Authored Documents", In titaniaaddueduph, 2010.

[2] B.Arman kiani, M.R. Akbarzadeh, "Autmatic Text Summarization using: Hybrid Fuzzy GA-GP" IEEE International Conference on Fuzzy Systems,2006.

[3] C.Gordon and R.Debus, "Developing deep learning approaches and personal teaching efficacy within a preservice teacher education context," Web Intelligence and Intelligent Agent Technology, British Journal of Educational Psychology,Vol. 72, No. 4, 2002, PP. 483-511.

[4] F. kyoomarsi, H. khosravi, E. eslami and M. davoudi, "Extraction-based text summarization using fuzzy analysis," Iranian Journal of Fuzzy Systems, Vol. 7, No. 3, 2010, PP. 15-32.