

# Automatic Summarization of Document using Machine Learning

Vaishnav Bhardwaj, Shwetansh Sharma, Ajay Katiyan  
Guided by: Dr. Waqar Ahmed

**Abstract:-** Internet domain is flooded with text information/ documents and it is difficult to get what kind of information exactly we are looking. To get the information we are looking for is head tic job as to find right or wrong we have go through whole information or text. This system Automatic text summarization summarizes the whole text or paragraph and give result in form of natural language using machine learning. This paper aim to present a process of summarization by using Machine Learning algorithms based on extraction of text based on their features the data to be summaries. There are two types of features the algorithms looks for one the frequency of element in the text and second is the linguistic extracted structure of the text. We also give some computational results achieved by applying our summarizer to certain dataset, which is compare to some baseline summary processes.

## I. INTRODUCTION

Text Summarization being a part of “Natural Language Processing” (NLP) technology that will surely affect positively our lives. Using the advancement of digital media and the constant development of distribution, reading entire articles or books to assess if they will be useful for a certain task becomes easier with Text Summarization tactics.

The process of providing a brief summary while maintaining critical information from a large given text is known as automatic text summarization. Text summarizing is classified into two types: extractive summarization and abstractive summarization. Abstractive Summarizing employs advanced NLP algorithms to build an altogether new summary, whereas Extractive Summarization depends on locating the correct sentences for summarization. For the objectives of this study, we used an extractive method by implementing the Text Rank Algorithm.

The text rank algorithm provide a text summary of several sources. We will just use probability methodologies to apply scores to phrases and extract the most essential sentences based on the highest score in order to conduct extractive summarization. Our major goal is to summarise the many categories that comprise the most relevant sample sentences and illustrate the findings based on the summaries.

## II. LITERATURE REVIEW

When we go to the history we see that single-document summarization was more focused on technical documents. When we go deep into that we come to know that most cited paper on summarization is that of LUHN in 1958s that briefly describes the research done at IBM. In his work, Luhn proposed that the frequency of a particular word in an article provides a useful measure of its significance . If sentence is studied in a paragraph the significant factor shows importance in any sentence as the significance and frequency of word in a sentence is studied. not only Significance of word but linear distance among words is also important for summarization of paragraph. To Summaries the paragraph sentence are ranked according to their significance of word, This hyposthis was introduced first in 1969 by Baxendale at IBM. in 1969 Edmondson experiment with extraction of document by implementing a typical structure for summarization of paragraph.

## III. PROBLEM DEFINITION

When we come to the new period, the period of 21<sup>st</sup> Century where tremendous amount of data is available on the Web. It is extremely intense for every individual to physically pick the synopsis of expansive archives of content. So there is an issue of scanning for vital reports from the accessible archives and discovering essential data. Along with these programmed content rundown is the need of great importance. Content rundown is the way that recognise the most vital important data in a record or set of related archives.

With the increasing amount of data day by day in the world, interest in the field of automatic summarization generation has been widely increasing so as to reducing the personal effort put by a person working on it. This thesis focuses on the comparison of various algorithms which are already present for the summarization of text passages.

## IV. PROPOSED WORK AND METHOD

Text Rank is a ranking of text in a paragraph which uses graph-based ranking system with NLP. This make pillar of frequency of word in a paragraph. If graph of word is aligned in a horizontal line we can visually see the level of frequency of significance of word. Higher graph represents the high importance of graph. Text Rank determines how similar each sentence is to the rest of the text.

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

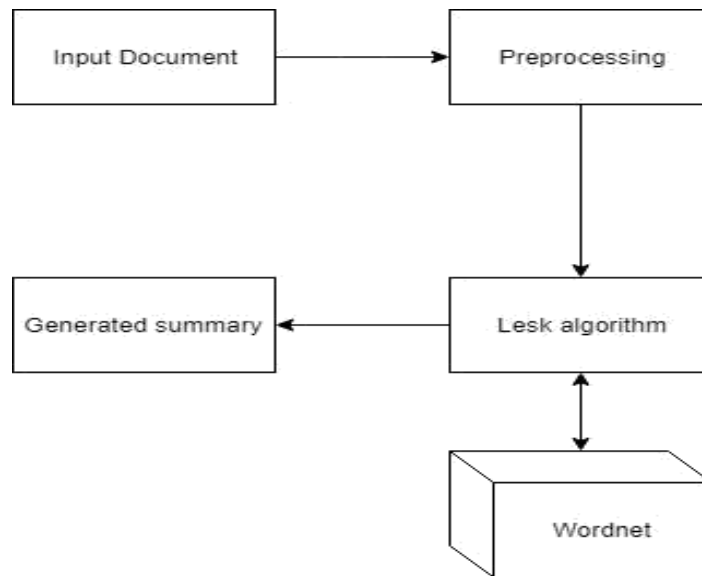


Fig. 1

**V. SYSTEM ARCHITECTURE FOR EXTRACTIVE APPROACH:-**

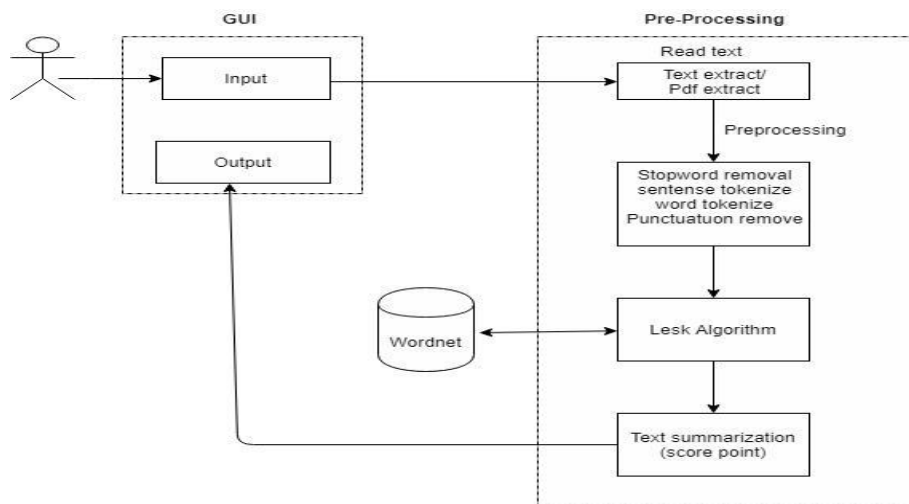


Fig. 2

In an example below 4 sentences are shown. This matrix is made for computation of similarity value of cosine. all values in diagonal are assigned with 0 other are assigned 1. After similarity is made in a matrix the values are converted into graphs on a horizontal line according to their frequency and significance. Here  $WS(V_i)$  is text rank of sentence  $V_j$  are sentences that have inline factor

$M =$

	w1	w2	w3	w4
w1				
w2				
w3				
w4				

Fig. 3

## VI. RESULT

In extraction based summarization the key word is extracted with selection of individual word which is tag in a document inside a paragraph. and in summation of document bases on paragraph summary. By implementing clustering algorithm as k-mean many experiments have been performed for document of large data to be summarized. But summarization based on this technique is used by removing full stop etc. In this technique large data of document is represented in the form of clusters of similar text. But this method leads to loss meaning of content. In our proposed method frequencies based on significance of word is implemented and formed in rank based on a single line then NLP is used for meaning of text.

## VII. CONCLUSION AND FUTURE SCOPE

The clustering algorithm is used by as plays a important role for finding the most unique ideas in the text . The outcome of the method shows that employing of multiple factors in the summarization can help to find the diversity or we can say that uniqueness in the text because of isolation of all similar sentences in one group can solve apart of the redundancy problem among the document sentences and the different part of that problem is solved by the diversity based method.

In future work abstractive method scan be implemented. In abstractive method build an internal semantic representation and then use natural language which we have and use of generation techniques to create a meaningful summary.

## REFERENCES

- [1.] Mehdi Kalahari, Seyedamin Pouriyeh, Mehdi Assefi, Saeid Safaei, Elizabeth D. Trippe, Juan B. Gutierrez, Krys Kochut, “ Text Summarization Techniques: A Brief Survey”, (IJACSA) International Journal of Advanced Computer Science and Applications
- [2.] Pankaj Gupta, Ritu Tiwari and Nirmal Robert, ”Sentiment Analysis and Text Summarization of Online Reviews: A Survey” International Conzatiferece on Communication and Signal Processing, August 2013
- [3.] Vishal gupta, Gurpreet Singh Lehal, ”A Survey of Text Summarization Extractive Techniques.” JOURNAL OF EMERGING TECHNOLOGIES IN WEB INTELLIGENCE, VOL. 2, NO. 3, AUGUST 2010
- [4.] Jiwei Tanai, Xiaojn Wan, Jiaguo Xiao Institute of Computer Science and Technology, Peking University “Abstractive document summarization with a Graph- Based attentional neural model.
- [5.] Seoggi Rang, Graduate from school of Information science and technology, University of East Asia Takeshi Abekawa, National institute of informatics “Framework of automatic text summarization using Reinforcement learning”
- [6.] Tian shi, Yaser Keneshloo, Naren ramakrishnan, Chandan K. Reddy, Senior member, IEEE “ Neural Abstractive text summarization with sequence-to -sequence models”
- [7.] Jianpeng Cheng, ILCC, school of informatics, University of Edinburgh Mirella Lapata, 10 crichton street, Edinburgh “Neural Summarization by extracting sentences and words”
- [8.] Alexander M. rush, Facebook AI research/ Harvard SEAS Sumit Chopra, Facebook AI research Jason Weston,

Facebook AI research “ A Neural Attention Model for Abstractive Sentence summarization”