# Automatic Text Summarization Using Natural Language Processing

## Pratibha Devihosur[1], Naseer R[2]

*[1] M.Tech. student, Dept. of Computer Science and Engineering, B.I.E.T College, Karnataka, India*

*[2] Assistant Professor, Dept. of Computer Science and Engineering, B.I.E.T College, Karnataka, India*

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

**Abstract -** *Automatic Text Summarization is the technique by which the huge parts of content are retrieved. In this paper The Automatic Text Summarization plays out the summarization task by unsupervised learning system. The significance of a sentence in info content is assessed by the assistance of Simplified Lesk calculation. As an online semantic lexicon WordNet is utilized. Word Sense Disambiguation (WSD) is a critical and testing system in the territory of characteristic dialect handling (NLP). A specific word may have distinctive significance in various setting. So the principle task of word sense disambiguation is to decide the right feeling of a word utilized as a part of a specific setting. To begin with, Automatic Text Summarization assesses the weights of the considerable number of sentences of a content independently utilizing the Simplified Lesk calculation and orchestrates them in diminishing request as indicated by their weights. Next, as indicated by the given level of rundown, a specific number of sentences are chosen from that requested rundown. The proposed approach gives best outcomes up to 50% summarization of the first content and gives attractive outcome even up to 25% outline of the first content.*

*Key Words***:** Automatic Text Summarization, wordnet, Streamlined lesk Calculation, Word Sense Disambiguation

## 1. INTRODUCTION

Automatic Text Summarization [1] H. Dalianis, [2] M. Hassel, is the plan to get an important data from a huge amount of information. The amount of data accessible on internet is increasing every day so it turns space and time expanding matter to deal with such huge amount of information. So, managing that large amount of data is makes a major problem in different and real data taking care of uses. The Automatic Text Summarization undertaking makes the users simpler for various Natural Language applications, like, Data Recovery, Question Answering or content decreasing etc. Automatic Text Summarization assumes an inescapable part by creating significant and particular data from a lot of information.

Filtering from heaps of reports can be troublesome and tedious. Without a summary or rundown, it can take minutes just to make sense of what the people will discuss in a paper or report. So the Automatic Text Summarization that concentrates a sentence from a content record, figures out which are the most imperative, and returns them in a readable and organized way. Automatic Text Summarization is a piece of the field natural language processing, which is

the manner by which the PCs can break down, and get importance from human dialect.

Automatic Text Summarization that uses the classifier structure and its rundown modules to look over huge amount of reports and returns the sentences that are helpful for producing a summary. Programmed outline of content works by taking the overlapping sentences and synonymous or sense from wordnet most overlapping sentences are considered as high score words [3] H. Seo, H. Chung, H. Rim, S. H., Myaeng, S. Kim, [4] A. J. Cañas, A. Valerio, J. Lalinde-Pulido, M. Carvalho, M. Arguedas. The higher recurrence words are considering most worth. And the top most worth words and are taking from the content and sorted according to its recurrence and generate a summary.

Lesk algorithm [5] S. Banerjee, T. Pedersen, [6]M. Lesk, is used for evaluating the waits for the input text using online semantic dictionary wordnet and it also uses the word sense disambiguation to identifying the most overlapping sentences in the input content that type of sentences are called equivocal words. Those types of words or sentences are having higher recurrences during the summarization.

In numerous normal dialects, a word can speaks to numerous implications/sense, and such type of word is called a homograph. WSD is the route toward making sense of which sentiment a homograph is used as a piece of given setting. WSD is a long-standing issue in computational linguistics, and has a come bonafide application including machine elucidation, information extraction, and information recuperation. Gener-accomplice, WSD use the setting of a word for its sense disambiguation, and setting information can begin from either clarified/unannotated content or other learning resources, for instance, responsive view point word expert, parallel corpora.

### 1.1 Natural Language Processing

Natural Language Processing technique using the nltk for building a main stage for python projects to work with human dialect information. This gives the easier to-utilize by giving the interfaces to one or more than 40 corpora and lexicon assets, for libraries for characterization, for splitting paragraphs sentences, to get its original form of words, labeling, parsing, and vocabulary thinking, and wrappers for modern thinking quality common dialect handling libraries, and for dynamic discourse discussion.

The NLTK is going to use an enormous tool compartment, and is going for make a favour for people with the entire

common dialect handling procedure. This will going to help people with all thing from part sentences from passages, to part up words, seeing the syntactic components of those words, marking the essential topics, doing this is helps to your machine b appreciating what really matters to the substance.

## 1.2 Streamlined Lesk Calculation

Calculation 1: This calculation compresses a single report content utilizing unsupervised learning approach. In This approach , the heaviness of each sentence in a content is determined utilizing Improved Lesk calculation and WordNet. The summarization procedure is performed as indicated by the given level of summarization [4]A. J. Cañas , A. Valerio, J. Lalinde-Pulido, M. Carvalho, M. Arguedas.

Info: Single-report input content.

Yield: Summarized content.

Step 1: The list of distinct sentences of the content is prepared.
Step 2: Repeat steps 3 to 7 for each of the sentences.
Step 3: A sentence is gotten from the list.
Step 4: Stop words are expelled from the sentence as they don't take an interest straightforwardly in sense assessment system.
Step 5: Glosses(dictionary definitions) of all the important words are extricated utilizing the WordNet.
Step 6: Intersection is performed between the sparkles and the information content itself.
Step 7: Summation of all the crossing point comes about speaks to the heaviness of the sentence.
Step 8: Weight appointed sentences are arranged in descending request concerning their weights.
Step 9: Desired number of sentences are chosen by the level of summarization.
Step 10: Selected sentences are re-orchestrated by their real sequence in the info content.
Step 11: Stop.

## 1.3 Advantages

• Reading the whole document, dismembering it and isolating the critical thoughts from the crude content require some serious energy and exertion. Perusing a document of 600 words can take no less than 10 minutes. Programmed outline programming condense writings of 500-5000 words in a brief instant. This enables the client to peruse less information yet get the most essential data and make strong conclusion.
• It reduces the human effort while creating a synopsis. A few vital products compress records as well as website pages.
• The persons quickly determine which points are imported for reading.

## 2. PROPOSED SYSTEM

In the Automatic Text summarization, we are using a solitary or single input content is going to outlined by the given rate of summarization utilizing unsupervised learning. In any case, the streamlined lesk's computation is associated with each of the sentences to find the guarantees of each sentence. After that, sentences with induced weights are composed in sliding solicitation concerning their weights. Presently as per a particular rate of summarization at a specific occurrence, certain quantities of sentences are chosen as an outline.

The proposed computations, abridges solitary or single report content utilizing unsupervised learning approach. Here, the heaviness of every sentence in a substance is resolved using streamlined Lesk's computation and wordnet. After that, summarization procedure is performed as indicated by the given rate of synopsis. In which, we are taking solitary info content and display summarization as yield. First info content is passed, to the lesk' computation and wordnet, where the weights of each sentences of the content are inferred utilizing and semantic investigation of the concentrates are performed. Next, weight doled out sentences is passed to derive the final summary according to the percentage of synopsis, where the last abridged outcome is assessed as and showed.
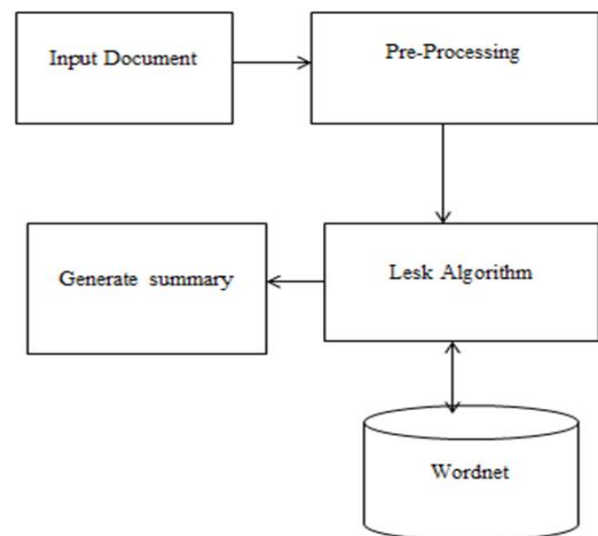


**Fig -1**: Overall Representation for Automatic Text Summarization Using Natural Language Processing.

## 1.2 System Architecture Of The Proposed System

The proposed system depicts the three stages for Automatic Text Summarization and they are listed below.

Stage 1: Data Pre-Processing

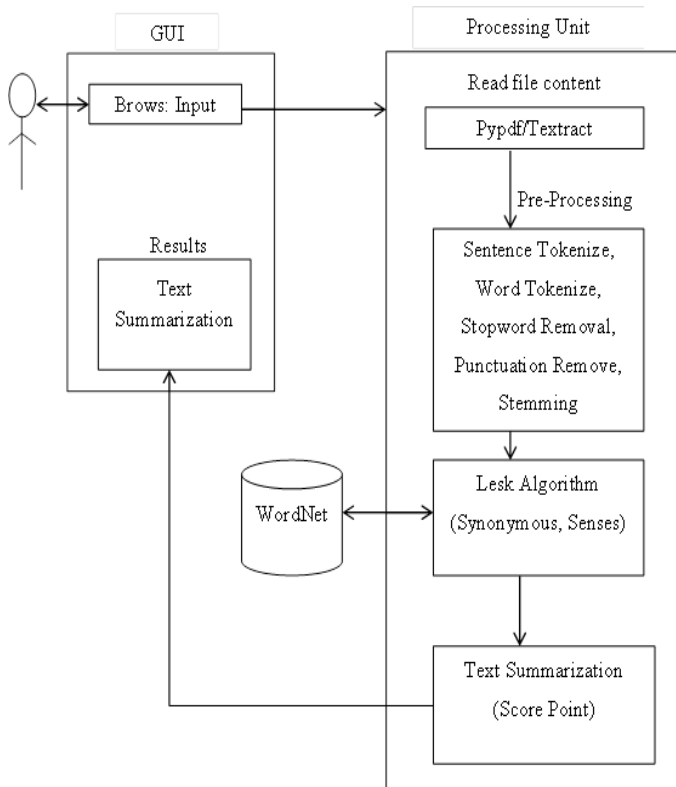Stage 2: Evaluation of weights

Stage 3: Summarization

**Fig -2**: System Architecture For Automatic Text Summarization Using Common Handling Dialect.

Stage 1: Data Pre-Processing

Programmed record outline generator is for clearing the undesirable things which exist in the substance. Henceforth it will additionally process it will performing sentence part, tokenisation, empty stopword, clear accentuation and perform stemming.

Stage 2: Evaluation of weights

This stage processes the repeat of the sentences of a substance utilizing lesk count and wordnet. In the first place finding the total number of spreads between a particular and the radiance this philosophy is performed for the all n number of sentences. By then once-over a particular sentence of the substance is set up for each of the sentences. A sentence is snatched from the once-over. Stopwords are removing from the sentence as they don't take an intrigue particularly in sense task method. Sparkles of each vital word removed using wordnet. Union is performed between the sparkles and the data content itself. Once-over of all the intersection guide comes to fruition talks toward the largeness of the sentence.

Stage 3: Summarization

This stage evaluates the last outline of a substance and the introductions the yield, which is surveyed at the period of arranging the sentences. In the first place it select the once-

over of weight named sentences are planned in jumping demand concerning their weights. Pined for number of sentences is picked by the rate of summary. Picked sentences are re-composed by their genuine gathering in the information content. The modified substance summary will gathers a substance without depending upon the association of the substance, rather than the semantic information lying in the sentence. Modified substance once-over is without vernacular. To remove the semantic information from a sentence, only a semantic word reference in the last vernacular is required.

## 3. OUTPUT AND DISCUSSION

Trial consequences of the venture for pre-preparing, assessment of the weights and showing the outline stage are executed. The results of following of these stages are represented in roar figure. In this approach we are using the word document and pdf document as input source.

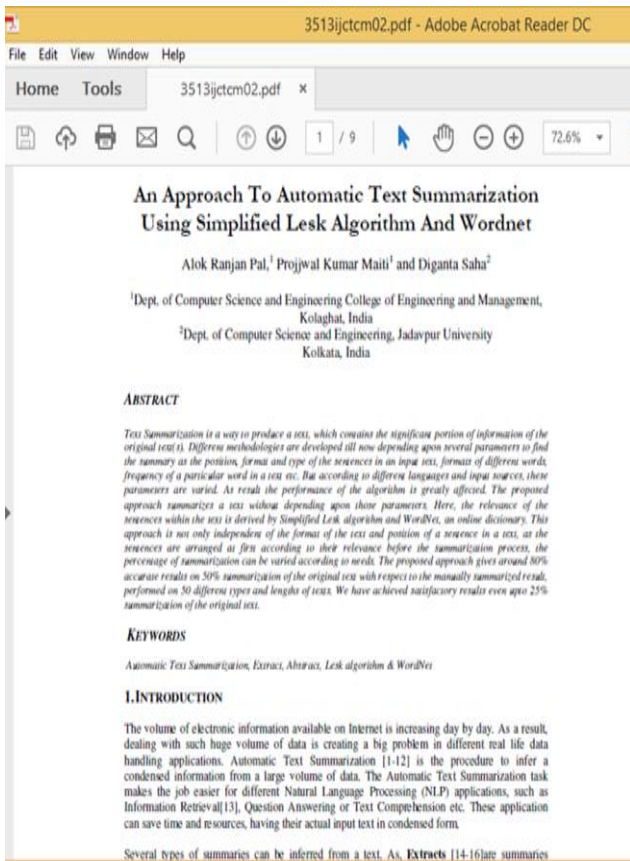



**Fig -3**: Input File for Word Document.
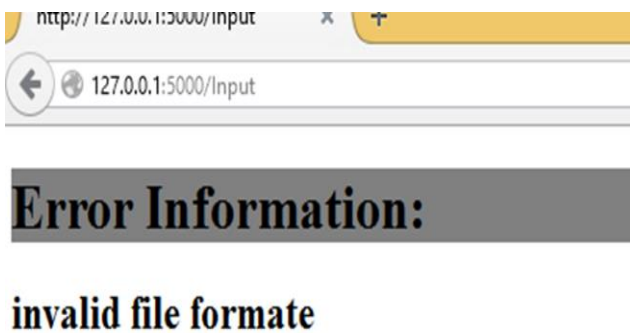
**Fig -4**: Input File For pdf Document.



**Fig -5**: Input File For Other than pdf or Word Document.

If info record is other than .pdf or .docx organize blunder will show like invalid data and invalid document design
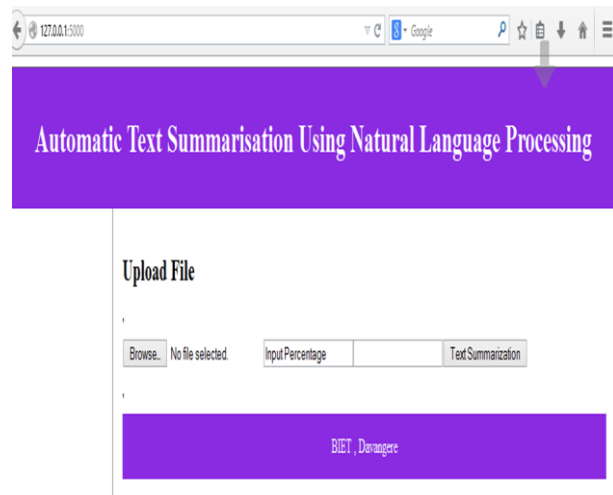


**Fig -6**: User Interface Form.

The User interface shape comprises of 2 catches, Browse and Text Summarization. The Brows catch will open a document to compress and Text Summarization is to begin procedure of the summarization.
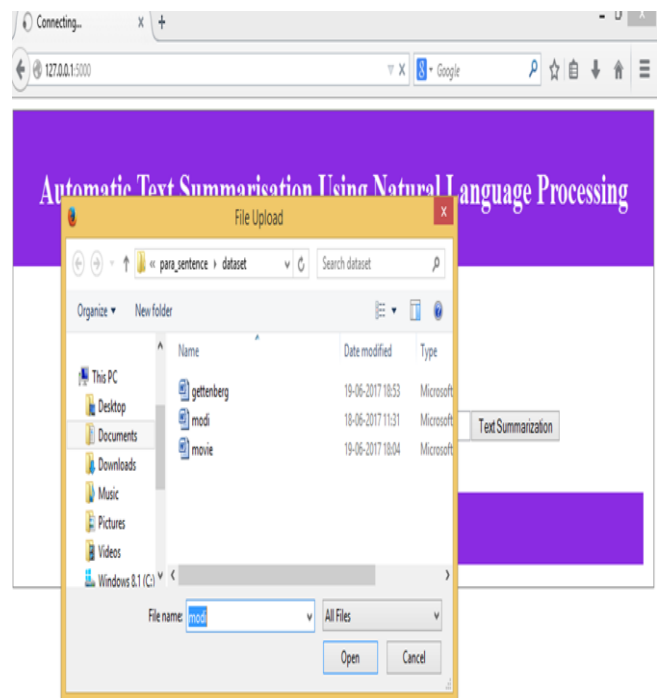


**Fig -7**: Brows Catch will Brows the file.

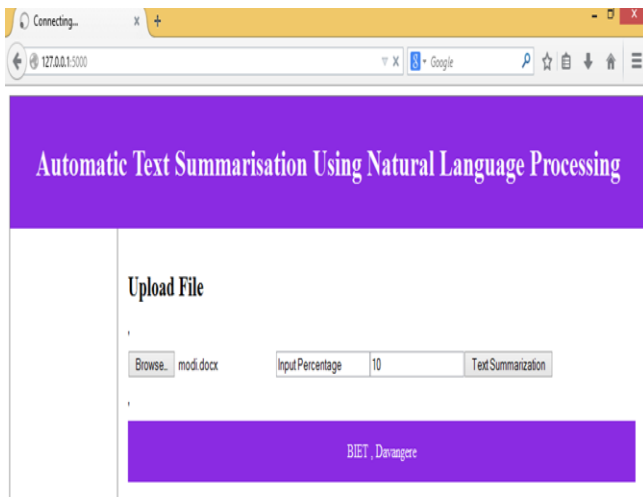The brows catch will select the input file to give summarization process

**Fig -8**: Input Percentage.

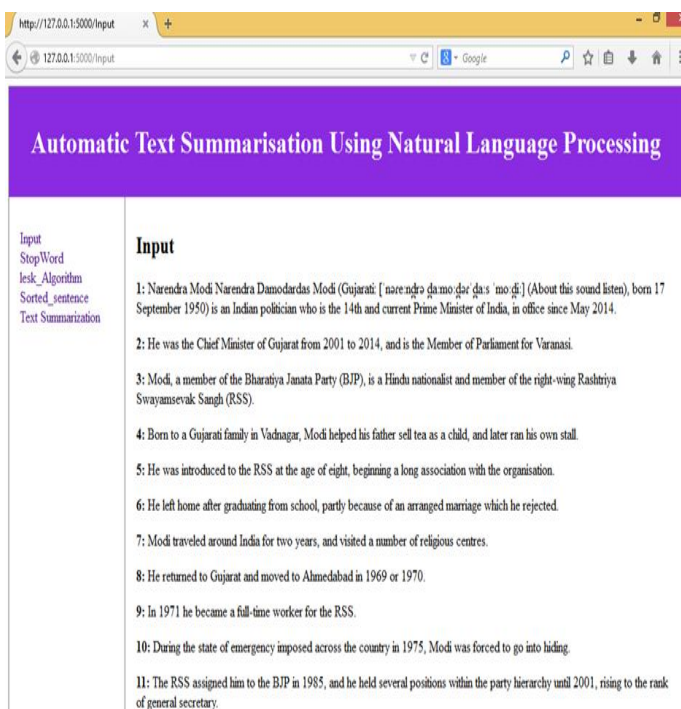After that client needs to give rate, how much summary need to show.



**Fig -9**: Brows Catch will Brows the file.

Therefore In Pre-handling the tokenization is parts the contribution as sentences or words.
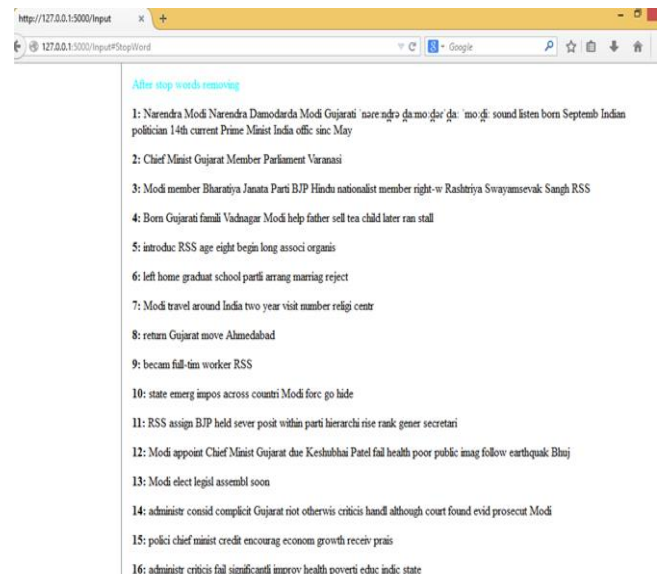


**Fig -9**: Brows Catch will Brows the file.

After it will list the sentences in the wake of evacuating the stopwords.
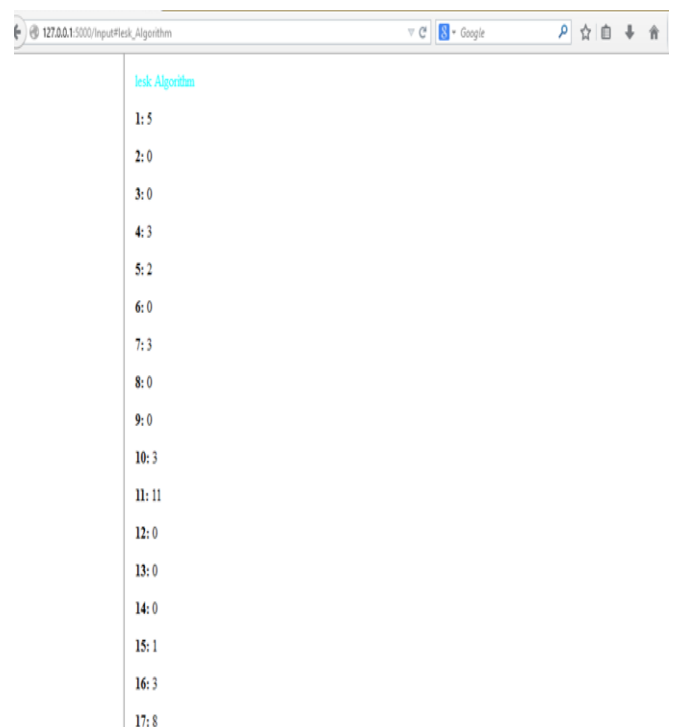


**Fig -10**: Lesk Calculation.

It will show weights for the input sentences according to its most important sentences
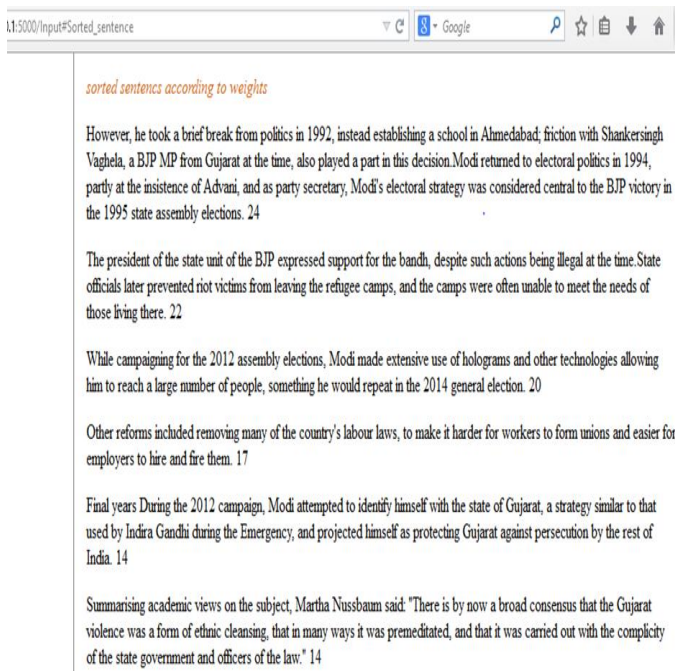
**Fig -11**: Brows Catch will Brows the file.

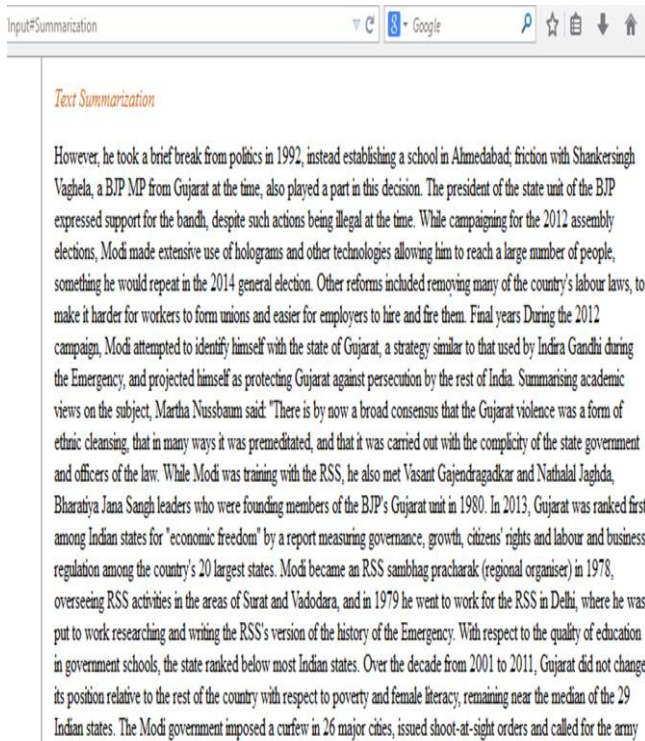After it demonstrates the arranged sentences According to weights.



**Fig -12**: Brows Catch will Brows the file.

Finally it will show the section of sentences constrained by rate.

## 4. CONCLUSION AND FUTURE SCOPE

Automatic Text Summarization approach depends on upon the semantic data of the concentration in a substance. So this way, gathered parameters like approaches, spots of different substances are not considered. In this recommendation, Lesk mean for word sense disambiguation by utilizing the vocabulary definitions to the electronic dictionary information base on utilizing wordnet. This goal is clear from covering sentence, couple of fusing words that give the setting of the word, in this not utilizing the late using the definitional shines of those words, other than those of words related to them through with the unmistakable relations portrayed in wordnet. So furthermore we are endeavoring to use other enlightening record away by wordnet for each word. For example, design sentences and identical words et cetera.

Among future work is the use of all the more balanced gathering to upgrade occurs additionally. Attempting diverse things with more tongue specific segments for instance, morphological parsers, printed entailment and anaphoric assurance is an open research for more updates later on. Programmed content summarisations should be possible for various archives. Client can be given an office to print the record from the interface specifically. A point of confinement to re-synopsis alternative perhaps included for record Shorter long. Additional line hole acquired in the outline can be evacuated. Spare as choice can be added to the application for the client to spare the synopsis in various arrangement.

## REFERENCES

[1] H. Dalianis, "SweSum – A Text Summarizer for Swedish," Technical report TRITA-NA-P0015,IPLab-174, NADA, KTH, October 2000.D.

[2] M. Hassel,"Resource Lean and Portable Automatic Text Summarization. PhD thesis, Department of Numerical Analysis and Computer Science," Royal Institute of Technology, Stockholm, Sweden 2007.

[3] H. Seo, H. Chung, H. Rim, S. H., Myaeng, S. Kim, "Unsupervised word sense disambiguation using WordNet relatives," Computer Speech and Language, Vol. 18, No. 3, pp. 253-273, 2004.

[4] A. J. Cañas , A. Valerio, J. Lalinde-Pulido, M. Carvalho, M. Arguedas, "Using WordNet for Word Sense Disambiguation to Support Concept Map Construction," String Processing and Information Retrieval, pp. 350-359, 2003.

[5] S. Banerjee, T. Pedersen,"An adapted Lesk algorithm for word sense disambiguation usingWordNet," In Proceedings of the Third International Conference on Intelligent Text Processing and Computational Linguistics, Mexico City, February, 2002.

[6]  M. Lesk,"Automatic Sense Disambiguation Using Machine Readable Dictionaries: How to Tell a Pine Cone from an Ice Cream Cone," Proceedings of SIGDOC, 1986.

## BIOGRAPHIES

Pratibha Devihosur (M.Tech). student, Dept. of Computer Science and Engineering, B.I.E.T College, Karnataka, India.

Naseer R Assistant Professor, Dept. of Computer Science and Engineering, B.I.E.T College, Karnataka, India.