# Automatic text summarization based on frequency count for Marathi e-Newspaper

Mr. Shubham Bhosale[1], Ms. Diksha Joshi[2], Ms. Vrushali Bhise[3], Rushali A. Deshmukh[4]

[1,2,3,]BE Student, Computer Engineering, JSPM's Rajashri Shahu College of Engineering, Pune
[4]Professor, Computer Engineering, JSPM's Rajashri Shahu College of Engineering, Pune

**Abstract-** In today's world, Time is of essence and nobody desires to waste their time not even on knowing what's happening around him. Most used and reliable thanks to forward news is newspaper till they was fictional. Nobody desires to browse a page article describing whole event of any incident happened around them. Everybody desires that everything around them ought to be as quick as potential. Therefore however will we have a tendency to create newspaper articles as little as potential and simple to browse and perceive. Account is that the method of minimizing the text content from a bigger text document and displaying solely the knowledge that is most vital. The intention of text account is to specific the content of a document during a condensed type that meets the wants of the user. Most techniques extract sentences that have bundles of keywords than the rest. Keyword extraction usually is finished by extracting proper words having a more robust frequency than others.

## I. INTRODUCTION

Growing quality of internet newspapers is accepted by user. There are several well-liked Marathi e-newspapers offered freely within the net, like geographical region Times, Lokmat, Sakal, Prabhat, Pudhari etc. These newspapers ar extracting all the mandatory data from newspapers. Extracting data from newspapers may be a troublesome job for each person.Necessity for a tool that extracts solely required data from these knowledge sources. Automatic keyword extraction is that the method of choosing words and phrases from an editorial that may at the best project the core sentiment of the article with none human intervention reckoning on the model . There are chiefly 2 sorts of account
1.Extractive summarization[4]
2.Abstractive summarization.[4]
In Extractive account, necessary sentences are known and solely those sentences ar enclosed in outline, desired outline length is obtained by use of compression quantitative relation. just in case of theoretical  account, in step with grading criteria acknowledge relevant sentences and method these sentences thus because the sentences are often enclosed in outline. Theoretical  account includes deep understanding of linguistic communication and it's additionally supported compression. Largely automatic account deals with extractive form of account.
The system works by checking the word frequency count generated from the document to be summarized. After selecting words with greater frequency , weselect lines containing those specific words.The lines will be displayed by the order they appear within the document.and after combining those selected lines with their respective order the summarized document will be created.

## II. RELATED WORKS

In this section we have a tendency to square measure discussing totally different approaches for Keyword extraction,Keyword ranking/classification,Text summarisation and algorithms to try to to therefore.
[1] Kumar Mohapatra proposes a technology which might work on telugunews papers and summarize their content. In their approach they have human intervention to coach the system to seek out the probable keywords. These keywords square measure then passed to the POS tagger to additional analyse. In learning stage system uses newspaper cut-outs, this cut-out is taken into account as target document and helpful statistics like nouns, verbs, adverbs, etc. square measure calculated. By victimisation keyword extraction algorithmic program key phrases square measure extracted by victimisation chance distribution, and by victimisation those key phrases new summarized document.
[2] For doing effective text summarisation authors, before doing extraction text preprocessing is important means that the information that is in unstructured type ought to be 1st regenerate to structured type. They additional discuss concerning the four stages that square measure required for txt preprocessing, that square measure choosing candidate terms, filtering, vector house model, and ranking. For choosing candidate terms text is tokenized and candidate terms square measure designated by victimisation n-gram approach. In filtering vocabulary pruning is
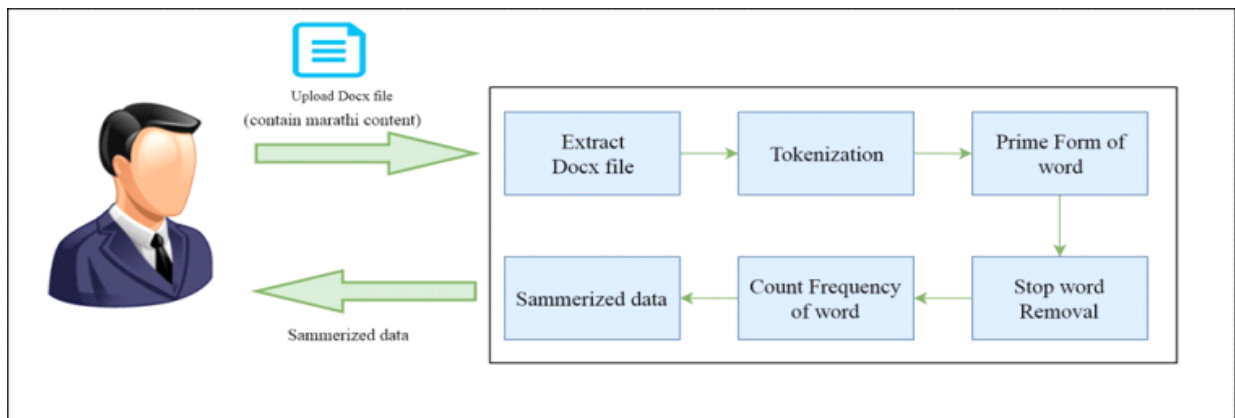
completed to additional minimize the candidate terms and predefined common terms square measure far away from the list. In vector house model the text document is portrayed in an exceedingly style of vector. The most reason to use VSM is to spot the particular which means of associate degree term in line with the context. And ranking deals with the statistics creation method.

[3] Steven J. Simske, Marcelo Riss mentioned concerning associate degree grammatically correct English summery creation system that was developed by Lins. This technique works on the news text extracted from the CNN web site. The system principally produces 3-5 lines news highlights instead of manufacturing a full summarized article. This is often done by adding all news information into associate degree XML document and enumeration every sentence and paragraph. Tis XML document is then passed to the feature extractor wherever a feature vector is generated by victimisation it. And this vector is employed for hard standards like mixture similarity, sentence length, correct nouns, TF/IDF , Uppercases, Word Frequencies. And by victimisation these standards 3-5 sentences square measure fashioned and presented the user.

[4] In authors Sheetal Shimpikar, Sharvari Govilkar discuss concerning technologies victimisation that we are able to develop text summarizer for Indian regional languages. During this paper they discuss concerning 2 main summarisation techniques that square measure theoretic and Extractive summarisation and make a case for thoroughly however they work. Theoretic text summarisation contains totally different strategies like Structure based mostly approach, Tree based mostly methodology, model based mostly methodology, metaphysics based mostly methodology, Rule based mostly methodology, and linguistics based mostly approach. And extractive text summarisation techniques like Term frequency Inverse Document Frequency, Cluster based mostly methodology like k-means algorithmic program, Maximum, connexion multi document(MMR-MD), and graph speculative. And every technology works otherwise relying upon the content provided thereto.

### III. PROPOSED SYSTEM

In this section we'll fathom however the task of keyword extraction, summarisation is completed within the system. Tasks area unit any divided for simplification of development method. beginning with exceptive the article the user desires to summarize, it are often done by accessing article directly from newspaper web site, by uploading the doc file, or by copy pasting it onto the location.



System Architecture
Once exceptive knowledge, it's analysed for locating often occurring words and classifying them. Victimisation these classified words we tend to choose the lines containing those words to form a summarized article, that is main objective of this method.
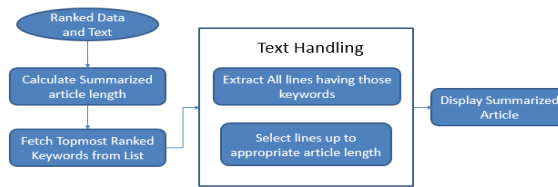
*3.1 Module*
System is divided into two main modules.
Word Extraction module mainly deals with text provided by the user. First of all text provided by the user is tokenized for ease of further word processing. After getting tokenized data it is filtered to remove punctuations and numbers. Filtered data is then ranked and components such as stopword list, frequencies, etc. are created.
Summarization Module, As name suggests this module deals with text summarization. First it calculates average article length which can be calculated and it is average of 30% to 40 % size of original article. Then according to size calculated in previous step topmost ranked keywords are selected which can be in range of 5-15 keywords. Then all the lines containing those keywords are fetched and a new summarized article is displayed to the user.

## Summarization Module



*3.2 Algorithms*

Algorithm 1: Extract Keywords
1. Input document (as copied text or docx file upload)
2. In the first step we take the text and tokenize it, for further operations.
3. Then the tokenized data will be provided to the stopword checking and removing function
4. After removing stopword we take process the data to check primeword and their frequency.
5. After we get a list of primeword and their frequency,we sort that list to take words with maximum frequency count.
6. Then we extract the lines containing selected words from the original text as it was uploaded by the user prior to processing.

Algorithm 2: Stop Words Removal Approach

Step 1: Text in the uploaded document is tokenized by spaces and words are stored within structures.

Step 2: Then system fetches single stop word from the list for further processing.

Step 3: Then this stop word is matched with the whole text by using sequential search technique.

Step 4: If it matches, the word in array is removed, and the comparison is continued till length of array.

Step 5: After removing one stopword completely, next stopword is taken from stopword list and step 2 is repeated until all the stopwords are matched with text and removed.

Step 6: Resultant text after operation is displayed, also required statistics like stopword removed, no. of stopwords removed from target text, total count of words in target text, count of words in resultant text, individual stop word count found in target text is displayed.

## IV. RESULT



Home page



About Project page

Upload Text


Result of System

## V. ACKNOWLEDGEMENTS

It gives us great pleasure in presenting the paper on "Automatic text summarization based on frequency count for Marathi e-Newspaper". We would like to take this opportunity to thank my internal guide Prof. Rushali A. Deshmukh, Professor, Computer Engineering, Rajashri Shahu Maharaj College Of Engineering, Pune for giving me all the help and guidance We needed. We am really grateful to them for their kind support. Their valuable suggestions were very helpful. such as laboratory with all needed software platforms, continuous Internet connection for Our Project.

## VI. CONCLUSION

In this paper, the planned work deals with e newspaper articles for account. The keyword extraction algorithmic program works to seek out the highest scored words expeditiously, and by victimization this knowledge the account module produces summarized article that principally rely upon the scale of original article. For choosing words applied mathematics approach is employed as a result of its higher performance and fewer complexness.

## VII. REFERENCES

[1] Naidu, Reddy, et al. "Text Summarization with Automatic Keyword Extraction in Telugu e-Newspapers." Smart Computing and Informatics. Springer, Singapore, 2018. 555-564.
[2] Hanumanthappa, M., M. Narayana Swamy, and N. M. Jyothi. "Automatic Keyword Extraction from Dravidian Language." Dept of Computer Science, Bangalore University, IJISET-International journal of Innovative science and technology 1.8 (2014).
[3] Silva, Gabriel, et al. "Automatic text document summarization based on machine learning." Proceedings of the 2015 ACM Symposium on Document Engineering. ACM, 2015.
[4] Shimpikar, Sheetal, and Sharvari Govilkar. "A Survey of Text Summarization Techniques for Indian Regional Languages." International Journal of Computer Applications 165.11 (2017).
[5] Litvak, Marina, and Mark Last. "Graph-based keyword extraction for single-document summarization." Proceedings of the workshop on Multi-source Multilingual Information Extraction and Summarization. Association for Computational Linguistics, 2008.
[6] Rahma, Abdul Monem S., Suhad M. Kadhem, and Alaa Kadhim Farhan. "Finding the Relevance Degree between an English Text and its Title." Eng. & Tech. Journal 30.9 (2012): 1625-1640.
[7] Matsuo Y, Ishizuka M. Keyword extraction from a document using word co-occurrence statistical information. Transactions of the Japanese Society for Artificial Intelligence. 2002;17:217-23.

VIII. Biographies

Prof. Rushali A. Deshmukh completed M.E. & B.E. in Computer Science & Engineering. She is pursuing Ph.D. from Dr. Babasaheb Technological University Lonere, Maharashtra. She is currently working as Associate Professor in Rajarshi Shahu college of Engg. Pune India , Department of Post Graduate Computer engineering with the total Experience of about 18  years.
She has published 09 papers in National Conference/ Journal and 21 papers in International Conference/ Journals. She has written 3 books. She has also filed one patent.  Her key research interest includes Natural language Processing, Machine Learning, Data Mining, Compilers

Ms. Diksha D. Joshi completed diploma in computer science.She is pursuing engineering from Rajarshi Shahu college of Engg. Pune India .
She has published 01 paper in journal paper.Her major field is computer science.

Ms. Vrushali P. Bhise is pursuing engineering from Rajarshi Shahu college of Engg. Pune India .
She has published 01 paper in journal paper.Her major field is computer science.

Mr. shubham s. bhosale completed diploma in computer science.He is pursuing engineering from Rajarshi Shahu college of Engg. Pune India .
He has published 01 paper in journal paper.His major field is computer science.