

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/261212570>

Automated Bangla text summarization by sentence scoring and ranking

Conference Paper · May 2013

DOI: 10.1109/ICIEV.2013.6572686

CITATIONS

19

READS

1,010

3 authors:



Md. Iftekharul Alam Efat

Noakhali Science & Technology University

10 PUBLICATIONS 34 CITATIONS

[SEE PROFILE](#)



Mohammad Rahimee Ibrahim

Asia e University AeU

7 PUBLICATIONS 38 CITATIONS

[SEE PROFILE](#)



Humayun Kayesh

Griffith University

8 PUBLICATIONS 30 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Detecting Event Causality from Social Media Short Text [View project](#)



Financial Analytics [View project](#)

Automated Bangla Text Summarization by Sentence Scoring and Ranking

Md. Iftekharul Alam Efat, Mohammad Ibrahim, Humayun Kayesh

Institute of Information Technology (IIT)

University of Dhaka, Dhaka-1000, Bangladesh

Email: iftekhar.efat@gmail.com, ibrahim_iit@yahoo.com, hkayesh@gmail.com

Abstract - In Natural Language Processing (NLP) the document summarization is an area that is getting interest of modern researchers. Though there are many techniques that have been proposed for English language but a few notable works have been done for Bangla text summarization. This paper deals with the development of an extraction based summarization technique which works on Bangla text documents. The system summarizes a single document at a time. Before creating the summary of a document, it is pre-processed by *tokenization*, *removal of stop words* and *stemming*. In the document summarization process, the countable features like *word frequency* and *sentence positional value* are used to make the summary more precise and concrete. Attributes like *cue words* and *skeleton of the document* are included in the process, which help to make the summary more relevant to the content of the document. The proposed technique has been compared with a human generated summary of documents and performance is 83.57%.

Keywords- Document Summarization, Text extraction, Word Tokenization, Word Steeming, Sentence scoring, Sentence Ranking

I. INTRODUCTION

Bangla is one of the most spoken languages in the World and the national language in Bangladesh. Over time the number of Bangla documents is increasing in a large amount. Reviewing these documents and evaluating them from any specific perspective is a gigantic task for a reviewer. It would take a lot time and efforts. Therefore, a system that automates the manual and tiresome process of summarizing documents is necessary. It helps saving a lot of time by reviewing the documents.

English document summarization systems are already there and serving with satisfactory accuracy. But there is no complete system for Bangla document summarization. This can help someone evaluating a large amount of Bangla documents or writings and giving necessary information about the content of the documents. For example, a blogger posted interesting topic in a Bangla blog and got a large number of responses from the readers with thousands of comments. But he does not have enough time to review all those comments. An automated summarization system can help him in this case to get a digest of the responses from the readers. The work presented in this paper is intended to generate an extraction based summary from a Bangla document.

The rest of the paper is organized as follows: In section II, we discuss the previous research works in this area. Next in

section III, we describe our proposed method for text summarization technique. Sentence scoring and summarization with pre-processing has been described in this section. Section IV illustrates the experimental results and discussion. Section V concludes the paper and provides direction for future work.

II. RELATED WORK

The earliest work on single-document summarization proposed the *frequency* of a particular word in a document to be a useful measure of significance described by Luhn in [1]. Though Luhn's methodology was a preliminary step towards the summarization, but many of his ideas are still found to be effective for text. In the first step, all the stop words were removed and rest of the words were stemmed to their root forms. A list of content words then compiled and sorted by decreasing frequency, the index providing an important measure of the word. From every sentence a *significance factor* was extracted that reflects the number of occurrences of significant words within a sentence, and correlation between them is measured due to the intervention of non-significant words. All the sentences are positioned in order to their significance factor, and the top positioned sentences are finally selected to form the automatically generated abstract.

Baxendale et al. marked the position of a sentence in a paragraph as an important feature in finding a prominent part of documents [2]. They inspected about 200 paragraphs and noticed that in 85% of the time the sentence containing main idea of the paragraph came as the first one and in 7% of the time it was the last sentence. This positional feature has been used in many complex machine learning based systems.

Edmundson et al. proposed a typical structure of text summarization methodology in [3]. They incorporated both *word frequency* and *positional value* ideas generated by two of his previous works. The first of two other features used was the presence of cue words (occurrence of words like *significant*, or *hardly*), and the second one was the skeleton of the document (the sentence is a title or heading). The sentences were scores based on these features to extract sentences for summarization.

All these research works are conducted for English, however the same procedure can be followed for other languages (e.g. Bangla) as proposed by Kamal Sarkar [4]. This approach used *word frequency* and *sentence position* in the

document as significant features to rank the sentences in the document.

Another work on Bangla text summarization proposed by Amitava and Sivaji used features such as Part of Speech (POS), Title Words, First Paragraph Words, Words from Last Two Sentences, etc. [5]. They used theme clustering to create a reasonable set of clusters for a given set of documents and as a way of extracting sentences based on their importance which regulates the quality of the output summary.

III. PROPOSED METHOD

The Bangla document summarizer is a Natural Language Processing (NLP) application which is proposed to extract the most important information of the document(s). In automatic summarization, there are two distinct techniques either text extraction or text abstraction. Extraction is a summary consisting of a number of sentences selected from the input document(s). An abstraction based summary is generated where some text units are not present into the input document(s). With extraction based summary technique, some more features are added based on Information Retrieval. However, the total system is alienated into three segments: pre-processing the test document, sentence scoring based on text extraction and summarization based on sentence ranking.

Input to a summarization process can be one or more text documents. When only one document is the input, it is called single document text summarization but in multi-document summarization the input is a group of related text documents. The text summarization can also be classified based on the types of users the summary is wished-for: User focused (query focused) summaries are adapted to the requirements of a particular user or group of users and generic summaries are aimed at a broad community of readers [6].

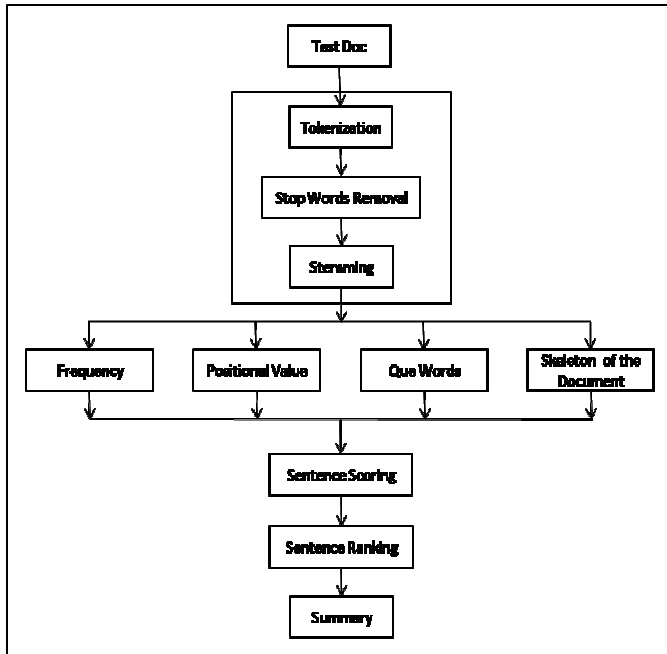


Figure 1. Steps of the proposed text summarization technique

A. Pre-processing

In Bangla document summarization process, some pre-processing is needed before executing the sentence scoring algorithm. By the pre-processing, the documents are prepared for ranking and summary generation. The pre-processing done on the documents are as follows:

Tokenization – A document is the combination of sentences and a sentence consists of some words. Here every word is considered as a token. A document is treated as a chain of tokens (marks).

Stop words removal – In Bangla words like এবং (And), অথবা (Or), কিন্তু (But), etc. are used frequently in sentences which have little significance in the implication of a document. These words can simply be removed for classification process.

Stemming – A word can be found in different forms in the same document. These words have to be converted to their original form for simplicity. The stemming algorithm is used to transform words to their canonical forms, like বাংলাদেশ, বাংলা, বাংলাদেশকে, বাংলাদেশেরদেশেও, etc. should be converted to their original form বাংলাদেশ. In this work, we use a lightweight stemmer that splits a word into its root form using a predefined suffix list [9].

B. Sentence Ranking & Summarization

After an input document is tokenized and stemmed, it is split into a collection of sentences. The sentences are ranked based on four important features: Frequency, Position value, Cue words and Skeleton of the document.

Frequency – Frequency is the number of times a word occurs in a document. If a word's frequency in a document is high, then it can be said that this word has a significant effect on the content of the document. The total frequency value of a sentence is calculated by sum up the frequency of every word in the document. The equation used to estimate the total frequency value of a sentence k is:

$$STF_k = \sum_{i=1}^n W_F \quad (1)$$

Where W_F (Word Frequency) is the total frequency of a word in the document, n is the number of words in a sentence and STF stands for Sentence Total Frequency.

Positional Value – The position of a sentence in a document has a considerable influence over the content of the document. The positional value of a sentence is computed by assigning the highest value to the first sentence and the lowest value to the last sentence of the document. The position value PV is calculated using the formula:

$$PV_k = \frac{1}{\sqrt{k}} \quad (2)$$

Where, k is the actual positional value of a sentence in the document.

Cue Words – Cue words are connective expressions (such as *therefore*, *hence*, *lastly*, *finally*, *meanwhile* or *on the other hand*) that links spans of communication and signals semantic

relations in a text. This is one of the summarization strategies which involve the use of “Cue Words” to select important sentences. The examples of “Cue Words” in Bangla are মোটকথা, অবশেষে, ইতিমধ্যে, যেহেতু, অতএব etc.

Skeleton of the Document – The skeleton of the document consists of the words in titles and headers. These words are considered having some extra weights in sentence scoring for summarization.

Sentence Scoring – The final score is a Linear Combination of frequency, positional value, weights of Cue Words and Skeleton of the document. The formula used to produce the final score of a sentence k is as follows:

$$S_k = (\alpha \times STF_k) + (\beta \times PV_k) + \gamma + \lambda, 0 \leq \alpha, \beta, \gamma, \lambda \leq 1 \quad (3)$$

Where, α and β are two co-factors of Sentence Total Frequency and Positional Value respectively. On the other hand, γ and λ symbolizes the weights of Cue Words and Skeleton of the documents correspondingly. The values of α , β , γ and λ are 0 to 1.

Summary Making – After ranking the sentences based on their total score the summary is produced selecting X number of top ranked sentences where the value of X is provided by the user. For the readers’ convenience, the selected sentences in the summary are reordered according to their original positions in the document.

IV. EXPERIMENTAL RESULTS & DISCUSSION

To test our summarization system, we collected 45 Bangla documents from the Bangla daily newspapers, Daily Prothom Alo. The documents are typed and saved in the text files using UTF-8 format. For each document we consider only one reference summary for evaluation. Evaluation of a system generated summary is done by comparing it to the reference summary.

It is difficult to summarize a document automatically compared to human summarization technique. In document summarization technique, total frequency of a sentence is more important than its positional value. If any cue word exists in the sentence then we need to consider it with high priority as a summary sentence.

For the most excellent results, it needs a fine tune of appropriate threshold value of the coefficient ($\alpha, \beta, \gamma, \lambda$) factors. We have chosen 10 random documents from the 45 test documents those we had used to tune these parameters. Initially, we set the value of α to 0.1 as it only multiply with the total frequency of a sentence, so we give this value a little weight whereas the cue word feature γ is primarily given a weight of 0.7. Beside this, we also consider that positional co-factor β to 0.2 as the important sentences are naturally kept in the early portion of the document. To include the scoring more efficiently the last co-factor λ is weighted 0.4 as it compare to the documents with the headlines structure which generally moves to the most important key words of the document.

We compiled our program on the test data sets to tune the parameters by summarizing those documents and compared the accuracy with human summarization process. Figure 03

describes the accuracy on changing the values of those parameters where the accuracy percentages are taken from Table 1.

In Fig 2 randomly the values of the co-factors are changed in the range of the primary assigned weight to find the accuracy. From the graph it is seen that the α curve moves like a sine curve where the amplitude is decreasing.

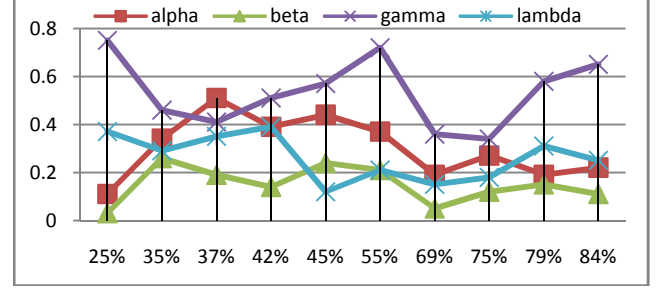


Figure 2. Performance vs Co-factors graph

It means that the Sentence Total Frequency has only 22% impact of its total frequency in Bangla text summarization technique. In the other hand, β factor has not changed much from the range and almost saturated in the more accurate points. It has also 11% impact of its original positional value corresponding to this document. To contrast, the λ factor was initially changed due to varieties of documents types but finally in the end point it acts like an inundated point. The most frequently changed value of co-factor is γ because weighting the cue words in a small dictionary is harsh to determine. The value of γ is comparatively high (0.65) considering with other co-factors because line which contains titles and headers’ words has more probability be selected as a summary sentence. Finally, we have set the weight of the co-factors are given in Table 1.

TABLE I. WEIGHT OF CO-FACTORS OF SENTENCE SCORING

Co-factors	Values
α (alpha)	0.22
β (beta)	0.11
γ (gamma)	0.65
λ (lambda)	0.25

We compared our algorithms’ summaries with the human summaries, computing the following scores. For each document we let k_h be the length of the human summary, k_m the length of the machine generated summary and the r the number of sentences they share in common. We defined precision (P), recall (R) and F_1 as metrics to compare the two summaries by:

$$P = 100 \frac{r}{k_h} \quad (4)$$

$$R = 100 \frac{r}{k_m} \quad (5)$$

$$F_1 = 100 \frac{2PR}{P+R} = 100 \frac{2r}{k_h+k_m} \quad (6)$$

To measure the accuracy of our algorithm we tested on 10 documents with the above equation of F_1 . For this we give the human summarize sentence line as input to compare with the sentences generated by the proposed system to find the F_1 value. Finally the average accuracy of Bangla text summarization is 83.57% corresponding with human generated summarization. The accuracy on basis of F_1 is given Table 2 where N is the number of lines in the document:

TABLE II. ACCURACY MEASUREMENT WITH F_1 VALUE

Doc No.	N	k_h	k_m	r	$F_1 = 100 \frac{2r}{k_h + k_m}$
1	172	31	35	25	75.76%
2	157	29	32	26	85.25%
3	166	35	34	30	93.75%
4	184	34	37	29	81.69%
5	145	32	29	23	75.41%
6	191	42	39	34	83.95%
7	178	32	36	31	91.18%
8	169	39	34	29	79.45%
9	188	35	38	33	90.42%
10	183	34	37	28	78.87%

V. CONCLUSION & FUTURE WORK

In this paper, we discussed extraction based Bangla text summarization method in a single document. The performance of this proposed system is 83.57% in generating summaries that agree well with human generated summaries, despite using minimal natural language processing (NLP) information.

In the future, this work can easily be extended as an abstraction based summarization. The performance of the proposed system may further be improved by improving learning model and adding more features in sentence scoring and ranking.

ACKNOWLEDGMENT

This research is completed with the support of the teachers of Institute of Information Technology, University of Dhaka.

Deepest thanks should be expressed to these person, who had kept their immense contribution in the Bangla Text Summarization field by their knowledge and effort that helps us to get some idea and help those helps in the research completion.

It would be a pleasure to thank **B. M. Mainul Hossain** because of his valuable suggestion and direction contribution that helps us in many ways to complete this research. Again, the commencement of the research on Artificial Intelligence course coordinated by **Ahmedul Kabir** is also necessary in this acknowledgement. Without that initiative of him it will not possible for us to research in this interesting and unique topic.

Another special thanks to **Shah Mostafa Khaled**, who inspired us with guidance, time and also supervises our working progress. It is completely impossible to development of this thesis without his motivation and help.

Finally the gratitude goes to the great almighty **ALLAH** who gives us the patient to do such a work like *Automated Bangla Text Summarization by Sentence Scoring and Ranking*.

REFERENCES

- [1] Luhn, H. P. (1958). The automatic creation of literature abstracts. IBM Journal of Research Development, 2(2):159-165.
- [2] Baxendale, P. (1958). Machine-made index for technical literature - an experiment. IBM Journal of Research Development, 2(4):354-361.
- [3] Edmundson, H. P. (1969). New methods in automatic extracting. Journal of the ACM, 16(2):264.
- [4] *Bangla Text Summarization by Sentence Extraction*, Kamal Sarkar, Department of Computer Science and Engineering Jadavpur University.
- [5] *Topic-Based Bangla Opinion Summarization*, Amitava Das and Sivaji Bandyopadhyay, Department of Computer Science and Engineering Jadavpur University, 2010.
- [6] Mani, I. (2001). *Automatic summarization*, Volume 3 of Natural language processing, Amsterdam/Philadelphia: John Benjamins Publishing Company.
- [7] A Light Weight Stemmer for Bangla and Its Use in Spelling Checker, Md. Zahurul Islam, Md. Nizam Uddin and Mumit Khan.