

# Sentence Similarity Estimation for Text Summarization Using Deep Learning.

Sheikh Abujar<sup>1</sup>, Mahmudul Hasan <sup>2</sup>, Syed AkhterHossain<sup>3</sup>

<sup>1</sup>Daffodil International University, Dhanomondi, Dhaka, Bangladesh  
sheikh.cse@diu.edu.bd

<sup>2</sup>Comilla University, Comilla, Bangladesh  
mhasanraju@gmail.com

<sup>3</sup> Daffodil International University, Dhanomondi, Dhaka, Bangladesh  
aktarhossain@daffodilvarsity.edu.bd

**Abstract.** One of the key challenges of Natural Language Processing (NLP) is to identify the meaning of any text. Text Summarization is one of the most challenging applications in the field of NLP where appropriate analysis is needed of given input text. Identifying the degree of relationship among input sentences will help to reduce the inclusion of insignificant sentences in summarized text. Result of summarized text always may not identify by optimal functions, rather a better summarized result could be found by measuring sentence similarities. The current sentence similarity measuring methods only find out the similarity between words and sentences. These methods state only syntactic information of every sentence. There are two major problems to identify similarities between sentences. These problems were never addressed by previous strategies: provide the ultimate meaning of the sentence and added the word order, approximately. In this paper, the main objective was tried to measure sentence similarities, which will help to summarize text of any Language, but we considered English and Bengali here. Our proposed methods were extensively tested by using several English and Bengali Texts, collected from several online news portals, blogs, etc. In all cases, the proposed sentence similarity measures mentioned here was proven effective and satisfactory.

**Keywords:** Sentence Similarity, Lexical Analysis, Semantic Analysis, Text Summarization, Bengali Summarization, Deep Learning.

## 1 Introduction

Text Summarization is a tool that attempts to provide a gist or summary of any given text automatically. It helps to understand any large document in a very short time, by getting the main idea and/or information of entire text from a summarized

text. To produce the proper summarization, there are several steps to follow, i.e. Lexical Analysis, Semantic analysis and Syntactic analysis. Possible methods and research findings regarding sentence similarity is stated in this paper. Bengali language has very different sentence structure and analyzing those Bengali alphabets may found difficult in various programming platforms. The best way of preprocessing both Bengali and English sentences before deep analysis, is using Unicode [2]. Sentence could be identified in a standard form, it will help to identify sentence or words structure as needed. The degree of measuring sentence similarity is being measured by method of identifying sentence similarity as well as large and short text similarity. Sentence similarity measures should state information like: if two or more sentences are either fully matched in lexical form or in semantic form, sentence could be matched partially or we could found any leading sentence. Identifying centroid sentence is one of the major tasks to accomplish [1]. Few sentences can contain some major or important words which may not be identified by words frequency. So, only depending on word frequency may not always provide the expected output, though several times most frequent words may relate with the topic models. Meaningfully same but structurally different sentences have to avoid while preparing a better text summarizer [3]. But related or supporting sentences may add a value to the leading sentences [4]. Finally most leading sentence and relationship between sentences could be determined.

In this paper, we have discussed several important factors regarding assessing Sentence and text similarity. Major findings are mentioned in details and more importantly a potential deep learning methods and model were stated here. Several experimental results were stated and explained with necessary measures.

## 2 Literature Review

The basic feature of text summarization would be either abstractive or extractive, approach. Extractive method applies several manipulation rules over word, sentence or paragraph. Based on weighted values or other measures, extractive approach choose appropriate sentence. Abstractive summarization requires several weights like, sentence fusion, constriction and basic reformulation (Mani &Maybury, 1999; Wan, 2008)[5].

Oliva et al. (2011) introduced a model SyMSS[6], which measure sentence similarity by assessing, how two different sentences syntactic structure influence each other. Syntactic dependence tree help to identify the rooted sentence, as well as the similar sentence. This methods state that, every word in a sentence has some syntactic connections and this will create a meaning of every sentence. The combination of LSA (Deerwester et al., 1990)[7] and WordNet (Miller, 1995)[9] to access the sentence similarity in between every words were proposed in Han et al.(2013)[8]. They have proposed two different methods to measure sentence similarity. First one makes a group of words – known as the align-and-penalize approach and the Second one is known as SVM approach, where the method applies

different similarity measures using n-gram and by using Support Vector Regression (SVR), they use LIBSVM (chang and Lin, 2011)[10], as another similarity measure.

A threshold based model always returns the similarity value between 0 and 1. Mihalcea et al. (Mihalcea et al., 2006)[11] represents all sentences as a list of bag of words vector and they consider first sentence as a main sentence. To identify word-to-word similarity measure, they have used highest semantic similarity measures in between main sentence and next sentence. The process will continue repeated times until the second main sentence could be found, during this process period. Das and Smith (Das and Smith, 2009)[12] introduced a probabilistic model which states Syntax and semantic based analysis. Heilman and Smith (Heilman and Smith, 2010)[13] introduces as new method of editing tree, which will contain syntactic relations between input sentences. It will identify paraphrases also. To identify sentence based dissimilarity, a supervised two phase framework has been represented using semantic triples (Qiu et al., 2006)[14]. Support Vector Machine (SVM) can combine distributional, shallow textual [15]-[17] and knowledge based models using support vector regression model.

### 3 Proposed Method

This Section represents a new proposed sentence similarity measuring model for English and Bengali language. The assessing methods, sentence representation and degree of sentence similarity have been explained in detail. The necessary steps required specially for Bangla language, has been considered while developing the proposed model. This model will work for measuring English and Bengali sentence similarity. The sentence structure and lexical form are very different for Bangla language. The semantic and syntactic measures also can add more values in this regards. The concept of working with all those necessary steps will help to produce better output, in every aspect. In this research - lexical methods have been applied, and untimely a perfect expected result has been found.

#### A. Lexical Layer Analysis:

The lexical layer has few major functions to perform, such as: Lexical representation and Lexical similarity. Both of these layers have several other states to perform. The Fig. 1 is the proposed model for lexical layer.

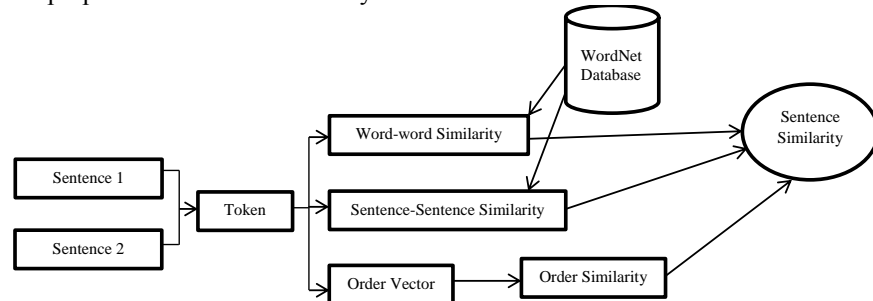


Fig. 1. Lexical Layer analysis model.

Figure 1 introduces the sentence similarity measures for lexical analysis. Different sentences will be added into a token. A word-to-word and sentence-to-sentence analyzer will perform together. An order vector will add all those word and/or sentence order in a sequence based on similarity measures. With the reference of weighted sum, the order of words and sentence will be privileged. A WordNet database will send lexical resources to word-to-word and sentence-to-sentence processes. Ultimately based on the order preference, the values from three different states (Word-word Similarity, Sentence-Sentence Similarity and Order Similarity) will generate the similar sentence output. The methods was followed by one of the popular deep learning algorithm – Text Rank.

1. *Lexical Analysis*: This Sate splits sentence and words into different tokens for further processing.
2. *Stop Words Removal*: Several value holds representative information. Such as article, pronoun, etc. These types of words could be removed while considering text analysis.
3. *Lemmatization*: This is a step to convert and/or translates each and every token into a basic form, exactly from where it belongs to. The very same verb form in the initial form.
4. *Stemming*: Stemming is the state of word analysis. Word-word and sentence-to-sentence both methods need all their contents (text/word) in a unique form. Here every word will be treated as a rooted word. Such as : play, player – both words are different as word, though in deep meaning those words could be considered as a branch words of the word “Play”. By using a stemmer, we could have found all those text in a unique form before further processing. The confusion of getting different words in structure but same in inner meaning will reduce. So, it is a very basic part of text preprocessing modules.

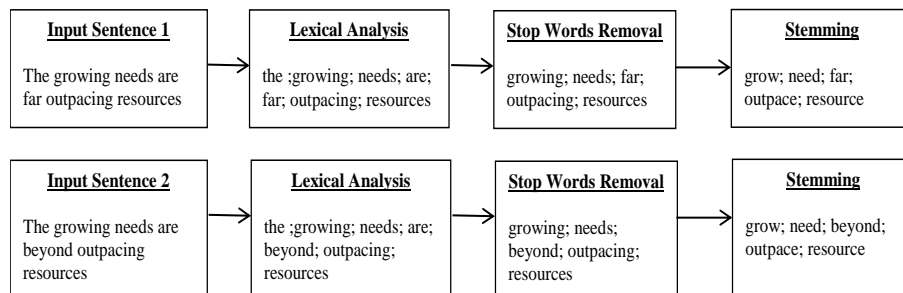


Fig.2. Lexical Layer processing of input sentences.

Figure 2. states, how lexical steps had been processed, with appropriate example. All the necessary processes as: Lexical analysis, stop words removal and stemming had been done as per the mentioned process. Those sentences will be used for further experiments in this paper.

B. *Sentence similarity:*

**Path measure** helps to sense the relatedness of words from the hierarchies of WordNet. It calculates and replies the path distance between two words. Path measure will be used to identify similarity scores between two words. Path measure could be calculated through Eq. (1).

$$Path\_measure(token1, token2) = 1 / Path\_length(token1, token2) \quad (1)$$

*Path\_measure* will send two different tokens as: token1 and token2. Both tokens are assigned the value of a single sentence after splitting. *Path\_length* will return the distance of two different concepts from WordNet.

**Levenshtein Distance** (Lev.) algorithm has been used to identify the similarity matrix in between of words. To identify the sentence similarity, measuring words similarity pays more importance. Lev. counts the minimum number of similarity required for the operation of insertion, deletion and modification of every character which may require transforming from a sentence to another sentence. Here it was used to identify distance and/or similarity measure between words. LCS (Long common subsequences) has also implemented though expected output was found using Lev. Here LCS does not allow substitutions. The distance of sentences followed by Lev. will be calculated based on the Eq. (2).

$$LevSim = 1.0 - (Lev.Distance(W1, W2) / maxLength(W1, W2)) \quad (2)$$

The degree of relationship helps to produce a better text summarizer by analyzing text similarity. The degree of measurement could be word-word, word – sentence, sentence – word and sentence – sentence. In this research, we had discussed the similarity between two different words. Such as there are a set of Words (after splitting every individual sentence):  $W = \{W1, W2, W3, W4, \dots, Wn\}$ . Lev.Distance calculate the distance between two words: W1 and W2, and maxLength will reply the score of maximum character found in between W1 and W2. Only similarity will be checked between two different words. The similarity between words could be measured by algorithm 1.

**Algorithm 1.** Similarity between Words

```

1:      W1= Sentence1.Split(" ")
2:      W2= Sentence2.Split(" ")
3:      if Path_measure(W1,W2) < 0.1 then
4:          W_similarity= LevSim(W1,W2)
5:      else
6:          W_similarity = Path_measure(W1,W2)
7:      end if

```

In Algorithm 1. the value of path will be dependent of distance values and LevSimairty (LevSim) value could be found from Eq. 1. The words similarity

score less than 0.1 will be calculated through the LevSim method else the score will be accepted from the path measure algorithm.  $W\_similarity$  will receive similarity score of between two words. The range of maximum and minimum score is in between  $\{0.00 -1.00\}$ . Table 1. represents the similarity value of words from sentence 1.

**Wu and Palmer measure (WP)** use the WordNet taxonomy to identify the global depth measures (relatedness) of two similar or different concepts or words by measuring edge distance as well as will calculate the depth of LCS (Least-Common-Subsumer) value of those two inputs. Based on Eq. (3), WP will return a relatedness score if any relation and/or path exist in between on those two words, else if no path exist – it will return a negative number. If the two inputs are similar then the output from synset will only be 1.

$$WP\_Score = 2 * Depth(LCS) / (depth(t1) + depth(t2)) \quad (3)$$

In Equation 3, t1 and t2 are token of Sentence 1 and sentence 2. Table. (2) States the WP similarity values of given input (as mentioned in Fig. 2).

**Lin measure (Lin.)** will calculate the relativeness of words or concepts based on information content. Only due to lack of information or data, output could become zero. Ideally the value of Lin would be zero when the synset value is the rooted node. But if the frequency of the synset is zero then the result will also be zero but the reason will be considered as, lack of information or data. The Eq. (4) will be used to measure the Lin. Value and table 3 will state the output values after implementing the input sentences on Lin. Measures.

$$Lin\_Score = 2 * IC(LCS) / (IC(t1) + IC(t2)) \quad (4)$$

In equation 4, IC is the information content.

A new similarity measure algorithm was experimented where all those mentioned algorithm and/or methods will be used. Equation (5) states the new similarity measure process.

$$total\_Sim(t1,t2) = (Lev\_Sim(t1,t2) + WP\_Score(t1,t2) + Lin\_Score(t1,t2)) / 3 \quad (5)$$

In Equation 5, a new total similarity values will be generated based on all mentioned lexical and semantic analysis. Edge distance, global depth measure and analysis of information content is very much essential. In that purpose, this method has applied and experimented out is shown in table (4).

**Algorithm 2.** A proposed similarity algorithm

```

1:      matrix = newmatrix(size(X)*size(Y))
2:      total_sim = 0
3:      i=0
4:      j = 0

```

```

5:         for i ∈ A do
6:         for j ∈ B do
7:             matrix(i, j) = similarity_token(t1,t2)
8:         end for
9:     end for
10:    for has_line(matrix) and has_column(matrix) do
11:        total_Sim = (Lev_Sim(matrix) + WP_Score(matrix) + Lin_Score(matrix)) / 3
12:    end for
13:    return total_Sim

```

The Algorithm-2 receives the token on two different X,Y as input text. Then it will create a matrix representation of m\*n dimensions. Variable total\_sim (total similarity) and i, j (which are the values for iteration purpose) will initially become 0. Initially, matrix(i,j) will generate the token matrix, where values will be added. The variable total\_sim will record and update calculate the similarity of pair of sentences based on token matrix – matrix(i,j).

## 4 Experimental results and discussion

Several English and Bengali texts were tested though the proposed lexical layer to find out the sentence similarity measure. Texts are being collected from online resource, for example: [www.prothom-alo.com](http://www.prothom-alo.com), [bdnews24.com](http://bdnews24.com), etc. Our python web crawler initially saved all those web (html content) data into notepad file. We have used Python – Natural Language Tool Kit (NLTK: Version– 3) and WS4J (a java API, specially developed for WordNet use). All the experimented results are stated in this section, below.

Table.1: Similarity score between words using path measure and LevSim.

	grow	need	far	outpace	resource
grow	1.00	0.25	0.00	0.14	0.00
need	0.25	1.00	0.11	0.17	0.14
far	0.00	0.11	1.00	0.00	0.09
outpace	0.14	0.17	0.00	1.00	0.00
resource	0.00	0.14	0.09	0.00	1.00

Table.2: Similarity score between words using Wu and Palmer measure (WP)

	grow	need	far	outpace	resource
grow	1.00	0.40	0.00	0.25	0.00
need	0.40	1.00	0.43	0.29	0.57
far	0.00	0.43	1.00	0.00	0.38
outpace	0.25	0.29	0.00	1.00	0.00
resource	0.00	0.57	0.38	0.00	1.00

Table.3: Similarity score between words using Lin measure (Lin.)

	grow	need	far	outpace	resource
grow	1.00	0.40	0.00	0.25	0.00
need	0.40	1.00	0.43	0.29	0.57
far	0.00	0.43	1.00	0.00	0.38
outpace	0.25	0.29	0.00	1.00	0.00
resource	0.00	0.57	0.38	0.00	1.00

Table 1, table 2 and Table 3, states the experimented result of similarity measure by using path measure and LevSim, Wu and Palmer measure (WP) and Lin measure (Lin.) consecutively. All those methods are either applied in lexical analysis or semantic analysis. In this research article, the proposed method of identifying sentence similarity using a hybrid model is being stated in Table 4.

Table.4: New Similarity score

	grow	need	far	outpace	resource
grow	1.00	0.21	0.00	0.13	0.00
need	0.21	1.00	0.18	0.15	0.34
far	0.00	0.18	1.00	0.00	0.15
outpace	0.13	0.15	0.00	1.00	0.00
resource	0.00	0.34	0.15	0.00	1.00

This method was also applied in Bengali language using Bengali WordNet. Experimented results are shown in Table (5).

Table.5: New Similarity score (Applied in Bengali Sentence).

	দৈন	সিট	ভাড়া	গন্তব্য
দৈন	1	0.78	0.88	0.16
সিট	0.78	1	0.31	0.24
ভাড়া	0.88	0.31	1	0.23
গন্তব্য	0.16	0.24	0.23	1

## 5 Conclusion and lines for further work

This paper has presented sentence similarity measure using lexical and semantic similarity. Degree of similarity were mentioned and implemented in the proposed method. There are few resources available for Bengali language. More development on Bengali language is just more than essential. Bengali WordNet is not stable as



like other WordNet available for English language. This research found suitable output in the unsupervised approach though a huge dataset will be required to implement the supervised learning methods. There are other sentence similarity measures, could be done by more semantic analysis and syntactic analysis. Both of these analysis if could be done together including lexical similarities, a better result could be found. More importantly, for a better text summarizer, we need to identify the leading sentences. Centroid sentences could optimize the analysis of post processing of text summarization. Evaluating system developed summarizer before publishing as final form is more important. Backtracking methods could possibly be a good solution in his regards.

## 6 Acknowledgment

We would like to thanks, Department of Science and Engineering of two universities: Daffodil International University and Comilla University, Bangladesh for facilitating such joint research.

## References

- [1] Rafael Ferreira et al. "Assessing Sentence Scoring Techniques for Extractive Text Summarization", Elsevier Ltd., Expert Systems with Applications 40 (2013) 5755-5764.
- [2] Sheikh Abujar, Mahmudul Hasan, "A Comprehensive Text Analysis for Bengali TTS using Unicode" 5th IEEE International Conference on Informatics, Electronics and Vision (ICIEV), 13-14 May 2016, Dhaka, Bangladesh.
- [3] Sheikh Abujar, Mahmudul Hasan, M.S.I Shahin, Sayed Akter Hossain "A Heuristic Approach of Text Summarization for Bengali Documentation" 8th IEEE ICCCN 2017, July 3 -5, 2017, IIT Delhi, Delhi, India.
- [4] Lee, Ming Che. "A novel sentence similarity measure for semantic-based expert systems." Expert Systems with Applications 38.5 (2011): 6392-6399.
- [5] Mani, Inderjeet, and Mark T. Maybury, eds. Advances in automatic text summarization. Vol. 293. Cambridge, MA: MIT press, 1999.
- [6] Oliva, Jesús, et al. "SyMSS: A syntax-based measure for short-text semantic similarity." Data & Knowledge Engineering 70.4 (2011): 390-405.
- [7] Deerwester, S., Dumais, S.T., Furnas, G.W., Landauer, T.K., Harshman, R., 1990. Indexing by latent semantic analysis. J. Am. Soc. Inf. Sci. 41 (6), 391–407.
- [8] Han, L., Kashyap, A.L., Finin, T., Mayfield, J., Weese, J., 2013. UMBC EBIQUITY-CORE: semantic textual similarity systems. Volume 1, Semantic Textual Similarity, Association for Computational Linguistics, Atlanta, Georgia, USA, June, pp. 44–52.
- [9] Miller, G.A., 1995. Wordnet: a lexical database for English. Commun. ACM 38, 39–41.

- [10] Chang, C.-C., Lin, C.-J., 2011. LIBSVM: a library for support vector machines. *ACM Trans. Intell. Syst. Technol.* 2 (May (3)), 27, 1–27.
- [11] Mihalcea, R., Corley, C., Strapparava, C., 2006. Corpus-based and knowledge-based measures of text semantic similarity, *National Conference on Artificial Intelligence - Volume 1*. AAAI Press, Boston, Massachusetts, pp. 775–780.
- [12] Heilman, M., Smith, N.A., 2010. Tree edits models for recognizing textual entailments, paraphrases, and answers to questions, *Annual Conference of the North American Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, Stroudsburg, PA, USA, pp. 1011–1019, 2010.
- [13] Heilman, M., Smith, N.A., 2010. Tree edit models for recognizing textual entailments, paraphrases, and answers to questions. *Human Language Technologies*, Stroudsburg, PA, USA, pp. 1011–1019, 2010.
- [14] Qiu, L., Kan, M.-Y., Chua, T.-S., 2006. Paraphrase recognition via dissimilarity significance classification, *EMNLP*. Association for Computational Linguistics, Stroudsburg, PA, USA, pp. 18–26.
- [15] Dzikovska, Myroslava O., et al. “Intelligent tutoring with natural language support in the Beetle II system.” *Sustaining TEL: From Innovation to Learning and Practice*. Springer Berlin Heidelberg, 2010.620-625.
- [16] Jurgens, David, Mohammad TaherPilehvar, and Roberto Navigli. ”SemEval-2014 Task 3: Cross-level semantic similarity.”*SemEval 2014* (2014): 17.
- [17] Mikolov, Tomas, et al. ”Extensions of recurrent neural network language model.” *Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2011.