

Annotated Bangla News Corpus and Lexicon Development with POS Tagging and Stemming

Abdul Matin¹, Tasnim Haider Chaudhury² and M.S. Hossain³

1

Received: 15 December 2016 Accepted: 2 January 2017 Published: 15 January 2017

Abstract

In this paper, we have developed a mono-linguistic Bengali news corpus using knowledge based AI (Artificial Intelligence) technique from some widely read Bengali newspapers which will be used as a reference corpus and will be very useful for lexicon development, morphological analysis, and automatic parts of speech detection. The corpus contains 74,698 word forms. The words in the lexicon are annotated with a combination of manual tags addressing Parts-of-Speech, Stemming, Morphemes, and other grammatical features are very important for almost all Natural Language Processing (NLP) applications. The lexicon contains around 14 thousand entries. In this paper we present some statistical analysis on some Bengali newspapers Prothom-Alo, Daily Janakantha, Daily Kalerkantho and Amardesh online from 1st January, 2012 to 31st January, 2012 those are the most popular Bengali newspapers in Bangladesh. We proposed a user friendly software interface to the user to annotate a large existing Bengali word set for the lexicon build up process.

Index terms— corpus, POS, tagging, stemming, lexicon.

1 Introduction

The significance of large annotated corpus is a widely known fact. It is an important tool for researchers in Machine Translation (MT), Information Retrieval (IR), Speech Processing, Knowledge Based Computer System and Natural Language Processing (NLP). But in Bengali language we do not have large annotated corpus. The development of corpus creation and distribution of language resources and its availability is must for enhancing Language processing capabilities and research in this field [1]. A corpus is also an essential language resource for creating automatic dictionary from a huge collection of language text [2]. It is the central repository of data for all language processing applications. Researchers are taking this field as a huge sector of researching. In [3, 4, and 5] they focus on automatic Bangla corpus creation by combination of Bangla font. It contains information for human consumption as well as computer programs. The book of "Corpus linguistics and Language Technology" [6] is a warehouse for corpus related studies with special attention to Bangla, where discussed almost every linguistic features of this language. In this paper, we are trying to present a reference corpus and a new approach into language investigation to understand how a text corpus database is utilized to obtain new result on a language or its properties. Bangla Corpus can be used for several purposes including spell checkers and morphological analysis for Bangla language.

A Bangla corpus can be extracted systematically from a Bangla corpus since it is considered a source of all words which is important for verification of Bangla sentence structure [7,8]. This paper proposes another process to manually build up a corpus, which is essentially a list of all words in the language and tag the words sufficiently with features such as word meaning, Parts-Of-Speech (POS) and all other grammatical features. All these information need to be stored in a database and properly formatted before display to end users. The aim of the project is to formalize a procedure for a collaborative effort by different individuals or groups towards producing a tagged Bangla corpus. This requires a POS tagging interface, both web based and standalone

95 ?????,????????? ? ??????(Here quotation mark has eliminated) ????? '????????????? ??????????????????
 96 ?????????????????????????????? (Here Apostrophe has eliminated) F. We replace all sequences of single
 97 or more spaces by a single space. str = str.replaceAll(" "+, " "); G. Then output file will save as
 98 "D:\PreProcessed_File.txt" and this pre-processed text will be used for POS tagging, Stemming and Mor-
 99 phological Analysis. H. This pre-processed text file will be referred to as our Raw Corpus. This corpus will
 100 be useful to evaluate the performance of annotated text corpus and also will be used for input text of machine
 101 translation system.

102 We proposed a user-friendly software interface to the user to annotate a large existing Bangla word set for the
 103 lexicon build up process. The effort will be a significant progress towards development of a properly annotated
 104 lexicon. This user interface has two distinct parts -one for building corpus and show text information such
 105 as source, source type, category, date, title and news content. Another parts for issuing manual addressing
 106 Parts-of-Speech, Stemming, Morphemes, issue clitic, ambiguous condition and other grammatical features.

107 A supervised machine learning method has been used for lexicon development from the Bengali news corpus.
 108 No extensive knowledge about the language is required except the knowledge of the different inflections that can
 109 appear with the different words in Bengali. To make proper annotation of word form, we accomplished each
 110 word form with POS tag, stem form, suffix, prefix, ambiguous condition (if exists), statistical counting. Initially,
 111 all the words (inflected and uninflected) are extracted from the pre-processed text and added to a database with
 112 proper POS, stem, prefix, and suffix. The system retrieves the words from the preprocessed text and creates a
 113 database of distinct word forms with fully annotation. Here is given the structure of lexicon development process
 114 and some sample word forms shown in Fig. ??.

115 12 Fig.2: Structure of Lexicon development

116 The Part-of-Speech (POS) tagging is the process of assigning each word of a text with an appropriate parts of
 117 speech tag. POS tags often signify the morphological [9], phonological and contextual properties of a word, and
 118 also provide information about neighboring words. In Bengali, there are five different POS namely, noun, pronoun,
 119 verb, adjective, and indeclinable (prepositions, con-junctions, and interjections). Noun, verb and adjective belong
 120 to the open class of POS in Bengali. In this lexicon analysis, we use seven parts-of-speech by extraction main five
 121 parts-of-speech in Bangla language. As, we know, there are a lot of word with proper noun are used in Bangla
 122 language. So we keep the proper noun distinct from other noun to get us more detail of the word form which is
 123 in the range of noun. To handle clitic which is one of the most common ambiguous situation in natural language
 124 Processing (NLP), we define a new POS form named clitic. If we add a word form in our lexicon database without
 125 setting any POS to that word form, our corpus creator software automatically set its POS as UNKNOWN. Noun
 126 and verb words are tagged by looking at their infections. Some infections may be common to some word form.
 127 In these cases, more than one POS may be generated for few words form. But here we set the mechanism for
 128 only one POS of a word form.

129 We only suggest a procedure to handle this ambiguity for initial level where POS ambiguity is resolved by
 130 checking he number of occurrences of these possible root words along with the POS tags as derived from same
 131 word forms. Pronoun and indeclinable are basically closed class of POS [10] in Bengali and these are added to
 132 the lexicon manually. It has been observed that adjectives in Bengali generally occur in four different forms based
 133 on the suffixes attached. For simplicity of counting or detecting sentence, we propose a user define POS named
 134 EOL (end of line). Here we detect POS for Purnacched, Question Mark, and Exclamation mark as EOL. The
 135 short description of POS categories is given in Table ??.

136 Stemming is an operation that splits a word into the constituent root part and affix without doing complete
 137 morphological analysis. It is used to improve the performance of spelling checkers and information retrieval
 138 applications, where morphological analysis would be too computationally expensive. Terms with common stems
 139 tend to have similar meaning. So it can drastically reduce the dictionary sized used in various NLP applications,
 140 especially for highly inflected languages. We handle this stemming process manually like previous POS tagging
 141 process. After tokenizing preprocessed text into individual word form we manually set root of that word form
 142 by removing prefix and suffix of it. Then this stem form is stored in STEM field of lexicon database.

143 In our process, we first stripped off the suffix part from Bengali words depending upon the type of suffixes.
 144 Then we checked for the validity of the suffix stripped word as root word, using a Bengali dictionary. If it is not
 145 sufficient we strip the affix part of the remaining part of the word form. It can bring a set of word with same
 146 root form in a series to learn about them easily. We can get almost similar word by retrieving word which has
 147 same root/stem form.

148 A smallest meaningful linguistic unit is consisting of a word or a word element that can't be divided into
 149 smaller meaningful parts. At the time we set stem word for a word from, we also store the stripped part of the
 150 word form as morphemes. First of all we split the suffix part of a word and store it in out lexicon database. Then
 151 we check out remaining part of our word form whether there is any prefix part of that word. If exist, we strip it
 152 from the remaining word part and sore it to database by fully manually.

153 A Corpus from linguistic point of view is defined as a collection of transcribed speech or written text compiled
 154 mainly to enhance linguistic research. The key resource to any linguistic research is a trained, annotated corpus
 155 which can elevate language processing capability such as automatic part-of-speech tagging, machine translation,
 156 questionanswering, stemming etc.

157 We design and develop a view of annotated corpus which is mainly based on knowledge based representation
158 (Knowledge based AI technique). Here we used our Lexicon as knowledge reference for our corpus. Our lexicon
159 is the collection of word forms with fully annotated where each word form is accomplished with parts-of-speech,
160 stem, Morphemes (suffix, prefix) and statistical counting. When we add a word form to corpus, we bring all the
161 morphological and grammatical information from lexicon and add this information with that word form. Corpus
162 procedure and flow are shown in Fig. 3 First of all we take pre-processed text as an input of our corpus creation.
163 2. Then tokenize this text into word forms. Then all these word forms is stored in an iterative list. 3. This
164 iterative list is looped and gets each word forms as a sequence they were in pre-processed text. Then for each
165 word, POS, Stem is brought from lexicon database and adds this information following the word form separating
166 with for slash (/). Then we make a small change in lexicon, we increase the count value of lexicon by 1 of a word
167 from each time we find this word form. This helps us to find the number of occurrence of a word form. 4. For
168 defining end of a sentence, we use EOL as word form and EOL as POS.

169 ? If sentence is an Assertive sentence, we use as stem. Example: ???????/UNK/UNK 6. Before adding
170 tokenized word form of preprocessed text to corpus, we add the entire information associate with this new text
171 with some predefine TAG to raw corpus. Tag format of our news corpus has given in Table ??I.

172 13 a) Statistical Analysis

173 Regardless of the size of the corpus, it may subjected to both qualitative as well as quantitative analysis using
174 various methods of statistics . Both these types of corpus analysis have different perspectives. Quantitative
175 analysis focuses classifying different linguistic properties where qualitative analysis aims to give some complete
176 and detailed description of the observed phenomena. We wish to focus on some simple quantitative analysis using
177 U-Gram model.

178 We develop our corpus development program in such efficient away where researcher can easily get a lot of
179 common and most focused perspective statistical output without any further processing. Here also some user
180 define output generator where user can get output with is desire requirement.

181 Here we divide our statistical output generator procedure in two distinct parts:

182 ? One for automated query based information.

183 ? Another for user defines query based information.

184 As result of automated query base perspective statistical output, we provide twelve statistical counting results.
185 This type of statistical counting will be very helpful for linguistic analysis, machine translation, Morphological
186 analysis, spelling variations, morphological structure, and word sense analysis. These statistical counting are,

187 ? Number of source from where this corpus data collected and there list.

188 ? Number of source type and their list of this source of data.

189 Table ??I Corpus is considered as basic resource for language analysis and research for many foreign languages.
190 This reflects both ideological and technological change in the area of language research. The effort will be a
191 significant progress towards development of a properly annotated lexicon. The outcome of the research will
192 significantly be helpful for future analyzer in the processes of Morphological Analysis, Automatic grammar
193 Extraction and Machine Translation for Bangla.

194 14 Global

195 1 2

¹© 2017 Global Journals Inc. (US)

²Year 2017 J

Structure:

Example:

WORD/SETM/POS

????/???/ADJ/??????/?????/NN/????/? Individual Word form ??/NN/????/????/PRO/???? /???? /ADJ

????/???/ADV/?????/?????/NN/????? ,????

?/ADJ/???-??/? ?/?/NN/?????/?????/NN/?????????/?/?????/Write "WORD/"

N/?????/?????/ADJ/??????/?????/NN/????/

????/PRO/?? ?/? ?/NN/??/?/VRB/????/?/VRB/

???? / ??? /ADJ/?/?/VRB/EOL/AS/EOL Increase

5. count by1

Tag Description Noun(Except Proper Noun) Proper Noun Adjective Adverb

Figure 1: ?

Annotated Bangla News Corpus and Lexicon Development with POS Tagging and Stemming

Year

2017

10

XVII VII. Experiment

Issue and Data

I Ver- Analysis

sion

I

()

Vol-

ume

J

Global

Jour-

nal

of Re-

searches

in

Engi-

neer-

ing

Tag Name

Tag Description / Purpose

<ENTRY>

<SOURCE></SOURCE>

<TYPE></TYPE>

<DATE></DATE>

<CATAGORY></CATAGORY>

<TI-

TLE></TITLE>

<CON-

TENT></CONTENT>

</ENTRY>

Statistical counting of our annotated Bangla text provide some qualitative analysis aims to give some complete and detailed description of the observed corpus is shown in Table III. Our Corpus program also To define start of a new news information/data. Source of data. (www.prothom-alo.com) Source type of data (news, blog) Date of collection of data (11-01-12) Genres of that data (sports, crime) Title of news/data Main content of the news. To define end of this news information/data.

phenomena which include word level frequency analysis, behavior of bangle word, use of non-Bangla word etc. These type of information can be get by using user defines query based annotated text corpus program interface.

b) Word frequency Analysis
Study of frequency calculation can provide important information about the usage of words in a

© 2017 Global Journals Inc. (US)

[Note: text]

Figure 2: :

III

Figure 3: Table III :

IV

Figure 4: Table IV :

V

Annotated Bangla News Corpus and Lexicon Development with POS Tagging and Stemming

Serial No	Information
1	Number of source
2	Number of source type
3	Number of fields/genres
4	Number of Raw word/Number
5 6 7	Number of Unique word Number of Unique Stem word Total Number of Sentence
8	Number of Assertive Sentence
9 10 11 12	Word Number of Interrogative Sentence. Number of Exclamatory Sentence. Number of Clitic N
à ? ? ?	

???	??	???	???	1.21	1.15	0.95	0.52	0.47	Percentage	0.4	???	????	??	???	?	????	Word Percentage	0.30	0.
??	Word	???	??	0.25	0.23	0.20													
???	????																		

????	0.18	??
????	0.16	???
???	0.16	??
??	0.16	???

[Note: © 2017 Global Journals Inc. (US)]

Figure 5: Table V :

VI

Figure 6: Table VI :

VII

Figure 7: Table VII :

		POS Name	Percentage		
		NN	56.43		
		VRB	20.53		
		ADJ	16.41		
		PN	13.71		
		ADV	5.94		
		PRO	3.39		
		IND	1.98		
		CLK	0.104		
		UNK	1.35		
Year 2017	Prefix	Percentage	Suffix	Percentage	
12	?	9.30	?	15.70	
I	??	4.07	??	15.43	
J () Volume XVII Issue	Conclusion	2.23 2.23 1.74	? ? ?	8.53 4.72 4.63	
I Version	?? ??? ?				
Journal of Researches in Engineering					

Figure 8:

-
- 196 [Toutanova and Cherry] ‘A global model for joint lemmatization and part-of-speech prediction’. Kristina
197 Toutanova , Colin Cherry . *Proceeding on ACL ’09 Proceedings of the Joint Conference of the 47th Annual*
198 *Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the*
199 *AFNLP*, (eeding on ACL ’09 eedings of the Joint Conference of the 47th Annual Meeting of the ACL and
200 the 4th International Joint Conference on Natural Language essing of the AFNLP) 1 p. .
- 201 [Md et al.] *Analysis of and Observations from a Bangla News Corpus*, Khair Md , Yeasir Arafat , Md Majumder
202 , Naushad Islam , Mumit Uz Zaman , Khan . Dhaka, Bangladesh. Center for Research on Bangla Language
203 Processing, BRAC University
- 204 [Asif Iqbal Sarkar et al.] *Automatic Bangla Corpus Creation*, Dewan Asif Iqbal Sarkar , Mumit Shahriar Hossain
205 Pavel , Khan . Dhaka, Bangladesh. BRAC University
- 206 [Dash ()] *Corpus Linguistics and Language Technology*, N S Dash . 2005. New Delhi. (Mittal)
- 207 [Cieri and Liberman] *Issues in Corpus Creation and Distribution: The Evolution of the Linguistic Data*
208 *Consortium, University of Pennsylvania and Linguistic Data Consortium Philadelphia*, C Cieri , M Liberman
209 . Pennsylvania, USA.
- 210 [Hasan ()] *Master’s thesis, School of Computer Science and Information Technology*, J Hasan . 2001. RMIT
211 University (Automatic dictionary construction from large collections of text)
- 212 [M Asaduzzaman and Ali ()] ‘Morphological Analysis of Bangla Words for Automatic Machine Translation’. M
213 M Asaduzzaman , Muhammad Masroor Ali . *th International Conference on Computer and Information*
214 *Technology (ICCIT) 2003. Jahangirnagar University*, (Dhaka, Bangladesh) 2003. p. .
- 215 [Md Hanif Seddiqui et al.] ‘Parts of speech tagging using morphological analysis in bangla’. Md Hanif Seddiqui
216 , Abdullah Al Rana , Taufique Mahmud , Sayeed . *Proceeding of the 6th International Conference on*
217 *Computer and Information Technology (ICCIT)*, (eeding of the 6th International Conference on Computer
218 and Information Technology (ICCIT)Bangladesh)
- 219 [Bharati et al. (1998)] ‘Some Observations Regarding Corpora of Some Indian Languages’. A Bharati , R Sangal
220 , S M Bendre . *Proc. Intl. Conf. Knowledge Based Computer Systems (KBCS98)*, (Intl. Conf. Knowledge
221 Based Computer Systems (KBCS98)NCST, Mumbai) 19 Dec. 1998. p. 17.
- 222 [Nur Hossain Khan et al. ()] ‘Verification of Bangla Sentence Structure using N-Gram’. Md Nur Hossain Khan ,
223 Md Khan , Md Islam , Habibur Rahman , Bappa Sarker . *Global Journal of Computer Science and Technology*
224 2014. 14. (Issue 1 Version 1.0 Year)