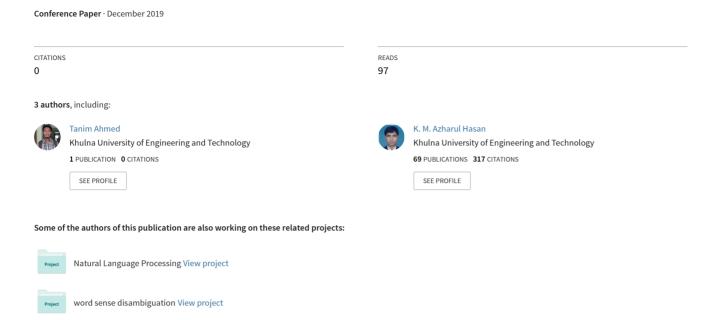
A formal method for designing a Bangla Stemmer using rule based approach



A formal method for designing a Bangla Stemmer using rule based approach

MD Shahidul Salim (Shakib)
Computer Science and Engineering
department
Khulna University of Engineering &
Technology
Khulna,Bangladesh
salim1507034@stud.kuet.ac.bd

Tanim Ahmed
Computer Science and Engineering
department
Khulna University of Engineering &
Technology
Khulna,Bangladesh
ahmed1507113@stud.kuet.ac.bd

Prof. Dr. K. M. Azharul Hasan Computer Science and Engineering department Khulna University of Engineering & Technology Khulna,Bangladesh az@cse.kuet.ac.bd

Abstract—In ML-based NLP makes your training data denser. It reduces the very large dictionary (number of words used in the corpus) simple. For Natural Language Processing stemming is necessary. Although there have many techniques for stemming different language. But in Bengali, there has no expected algorithm for stemming the Bengali Language. We have proposed a system for stemming Bangla word. It is possible to build up a search engine using this corpus since search engine works with keywords. Our system can help to find those root word. It also has other applications such as sentiment analysis, spam filter etc.

Keywords—stemming, search engine, sentiment analysis, corpus, NLP

I. INTRODUCTION

Stemming is the process for removing the commoner morphological and inflectional endings from words. Although there are different kinds of the algorithm for different languages for this task. Like as Porter Stemmer is a famous algorithm for stemming English words. But for Bengali almost there has no efficient algorithm. In the present, there have some works for Bangla Stemming but they are almost they follow brute force system. And the accuracy of this system is very poor. For the development of Bengali in Natural Language Language Understanding(NLU) and Natural Language Processing(NLP), stemming of Bengali Language is necessary. But it has some difficult issue to solve this problem. Because if you follow the brute force approach to solve this problem you accuracy and complexity will be very poor. But we have proposed a better solution for solving this problem. We have created some important rules to solve this problem. We have divided the words some categories(NOUN, VERB, NUMBER, etc.). And proposed an algorithm that will perfectly stem the Bengali words from dataset. And our algorithm provides better efficiency. The efficiency of our algorithm almost 94.35%. Table 0 shows some example how the words are stemmed with our techniques.

বিভক্তি(verb	(বিচন)Number	Others
inflection) Example	Example	Example
ায় → ε # করায়	টি> ε # মেয়েটি	তে $\rightarrow \varepsilon$ # ছাড়তে
াও → ε #	(ব্যতিক্রমঃ খিটিমিটি, বৃষ্টি)	(ব্যতিক্রমঃ হাতে,ভাতে)
করাও,পড়াও		জের → জ # তেজের
াইস $ ightarrow$ ϵ #	টি> ৪ # ছেলেটি	ার →া # বাঁধবার
করাইস	টা> ε # কলমটা	ের → ε # যুগের
াচ্ছ → ε # করাচ্ছ	খানা>ε #	তেই → ε # ভক়তেই
াচ্ছি $ ightarrow$ ϵ #	খাতাখানা	
করাচ্ছি	খানি>ε # বইখানি,	খি → খ # দেখি
াছে $ ightarrow$ ϵ #	রা> ε # ছাত্ররা	তে → ε # শেয়াকুলিতে,
করাইয়াছে,	রা> ε # পাখিরা,	
	রা>ε# পিপীলিকারা,	নের →ন # আয়তনের
করিয়েছ	,	ের →ε # গৌরবের
াক $ ightarrow$ ϵ # করাক	গুলো> ε # আমগুলো.	তে $ ightarrow$ ϵ $\#$ দিঘিতে
াও $ ightarrow$ ϵ # করাও	গুলো>ε #	টাই $ ightarrow$ ϵ # নামটাই
াইস → ε #	টাকাণ্ডলো,	তার → তা # দেবতার
করাইস	গুলো>ε #	নের $→$ ন $\#$ বিসর্জনের
ালেন $ ightarrow$ ϵ #	ময়ুরগুলো,	
করালেন ালি $ ightarrow$ ϵ #	গণ> ε # দেবগণ,	তার → তা # রাস্তার
করালি	,	পে → প # মাপে
ালাম → ε # করালাম	সমূহ> ε # বৃক্ষসমূহ	ার → া # প্রতিমার
াইতে → ε #	বলি> ε # পুস্তকাবলি	
করাইতে	গুচ্ছ> ε # কবিতাগুচ্ছ	
াতে → ε # করাতে	দাম> ε # কুসুমদাম	
(ব্যতিক্রমঃ উৎপাতে)		

Table 1

II. RELATED WORKS

In 1980 "Porter Stemmer" was developed by Martin Porter[1]. Porter stemmer defines some rules these rules applied in a word for removing suffix. Five steps are followed in porter stemmer[2] to stem any word. Several stemming algorithms exist for English such as Lovins Stemming[3], Praice/Husk[4], Dawson[5] and so on. Porter

Stemmer is most recognized one that is applied English and others languages that are Indonesian[6],Malay[7],Dutch[8],Slovene[9],Turkish[10] and Latin[13].For Bengali, there remain some works in this field. Bengali is a language that is highly inflected. It is commonly inflected with noun, verb, adjective. A Light Weight Stemmer for Bengali and its use in spell checker proposed in [15].A design a rule-based stemmer for natural language text proposed in [14].They have another one rule-based stemmer developed for Hindi and bengali[16].

III. METHODOLOGY

The stemming approaches follows some rules. These rules show how the stemming will be performed. And further an algorithm is developed. This algorithm works based on the following rules.

Number (বচন) CFG:

Rule 1: A # Remove A if it appears at the end of word (where A = টি , টা,খানা,খানি (table 1))

Example:

ছেলেটি -> ছেলে

বইখানা -> বই

Example with sentence:

কলমটি সন্দর

কলম সুন্দর

Bivokti(বিভক্তি) CFG:

Rule 1: A # Remove it if A appears at the end of word (where A =াই, াও,াইস table 1)

Example:

করাই -> কর

পড়াও -> পড়

Example with sentence:

কাজটি করাও

কাজ কর

শিক্ষক পডাই

শিক্ষক পড়

Others(বিবিধ) CFG:

Rule 1: A # Remove it if A appears at the end of word (where A = তেই,তে(table4))

Example:

শুরুতেই -> শুরু

ভাতে -> ভাত

Example with sentence:

মাছে ভাতে বাঙালি

মাছ ভাত বাঙালি

Rule 2: B -> C # Replace B with C if B appears at the end of a word. (B = 7%, 7%, C = %, %(table 4))

Example:

মাপে -> মাপ

দেখি -> দেখ

Example With Sentence:

মাপে ঠিক আছে

মাপ ঠিক আছে

Don't stemmed:

IF the data of the table 4 is contained in the suffix then there is no need no change the word. Because if we strip these words then they will lose their real meaning These word will be directly used as root word.

Example: খোটা,যারা,তারা (table 5)

ক্রিয়া-বিভক্তি

েছ | িয়েছ | ায় | াছি | ায়েছ | াছে | াতিস | ালাম | ালেন | াইলে | াইবি | াবেন | াইবে | াতেন | াছ | াইলি | াতাম | াইবে | াইব | ালি | ালে | াউক | াবো | াইস | ায়ো | াবে | াইও | াইয় | াবি | াতে | াছে | াইয়েছিলাম | াইয়েছিলাম | াইতেছিস | াইয়েছিলাম | াইতেছিস | াইয়েছিলেন | াইতেছিস | াইয়েছিলেন | াইয়েছিলে | াইয়াছিলে | াইয়েছিল | াইয়াছিলে | াইয়েছিল | াইয়েছেল | াইতেছেন | াইয়েছি | াইয়েছে | াইয়েছে | াইতেছেন | াইয়েছে | াইবেন | াইতাম | াইতিস | াইবেন | াইলাম | াইছে | াইয়েছ | াইয়েছ | াইলেন | াইলাম | াইরেছ | াইলেন | াইলাম | াইরেছ | াইলেন | াইলাম | াইরেছ | াইলেন | াইলেন | াইলাম | াইরেছ | াইলেন | াইলেন | াইলাম | াইরেছ | াইলেন | াইলেন | াইলাম | াইরেছি | াইলি | াইলি | াইলি | াইলি | াইল → ε

Table 2

বচন(বিশেষ্য-সর্বনাম)

বৃন্দ | মন্ডলী | কুজ | পুন্জ | গুচ্ছ | বৃন্দ | সমুদয় | সমূহ | বৰ্গ | রাশি | আবলি | খানা | গুলি | রাজি | নিকর | খানি | নিচয় | গুলো | মালা | যূথ | সকল | বলি | দের | এরা | দিগ | পাল | দাম | কূল | টা | রা | সব | গণ | টি | ান → ε

Table 3

	বিবিধ	
তে → ε	তার> তা	লেই> লা
ঠল → ঠ	বার> ε	নে> ন
টাই> ε	ার> া	াচ্ছিল> ε
ইতেছি> દ	পে> প	রিতে> রি
নোর> ε	মার> মা	ললে> ল
টার> ε	াতে> া	্ৰেছি> ε
সছে> স	তে> তা	লেম> ৪
লের> ল	দ্দর> দ্দ	খে> খ
3 < లయి	ঠল> ল	ৰ্বে> ৰ্ব
3 < &	মল> ম	নীর> নী
রলে> র	দের> ε	লের> ল
য়ে> ε	(**> **)	নের> ন
রে> র	ড়ে> ড়	লায়> লা
ড়াতে> ড়	েছে> દ	ঞে> প্
তেই> ε	েৱ ই > ε	লায়> লা
খি> খ	র ছে > র	ময়ে> ময়
ঠে> ঠ	টিয়ে> টা	জের> জ
েধেছে> াধ	াতেও> દ	রের> র
3 <	জে> জ	য়ে> য়
টা> ε	লতে> ল	ের> হ
কুল> ε	ারে> র	েছে> દ

Table 4

Do not stem

ক্লাস ,দিয়ে, শাসন , পরিচয়, আমাদের, নতুন ,তাদের ,থেকে, এদের, তারা ,যারা, সন্ধান, মার ,পারে ,দিতে, দরকার ,খিটিমিটি, নিয়ে ,খোঁটা

Table 5

Algorithm 1:

- 1: procedure: BanglaStemmer(string1, string2, string3)

 /*An algorithm to remove all possible suffixes from a word*/
- 2: Input: st: String of Bangla sentence having n words , string1:Set of all longest Bivokti(বিভক্তি) suffix , string2: Set of all smallest Bivokti(বিভক্তি) suffix, string3:Set of all bochon(বচন) suffix, string4: Set of all Extra suffix, do_not_stem: Dataset of do not stem these word
- 3: Output: A sentence with stemming word
- 4: data[i] \leftarrow Tokenize (st)/* i=1 to n */
- 5: data2[i] ← Remove stop word from the data[i]
- 6: for each word ← data2[i] do

7: if (wo	rd doesn't match with do_not_stem dataset) then	
8:	stemmed_word←Take Suffix of this word and match with the string1, return stemmed word	
9:	stemmed_word Take Suffix of this word and match with the string2, return stemmed word	
10:	stemmed_word Take suffix of this word and match with the string3, return stemmed word	
11:	stemmed_word Take suffix of this word and match with the string4, return stemmed word	
12: end if		
13: end for		

IV. RESULTS

For checking accuracy we have taken dataset from "যোগাযোগ(রবীশ্রনাথঠাকুর)"[18]. Firstly we have found the accuracy close to 85%. Then we have modified our rules. After some iteration we have found accuracy 94.35%. Table 5 shown accuracy test.

Accuracy test		
Total: 172 stemming word(target)		
stemmed :162 (True positive)		
false stemmed:5 (True negative)		
Should be stemmed but not stemmed : 5(False positive)		
Accuracy: 94.35% ((TP+TN)/(T+N)) (Using confusion		
matrix) [17]		

Table 5

FUTURE WORK

We will add if any rules that we have not yet added such the accuracy will be close to 100%. And try to propose our modified algorithm such that this algorithm gives best time complexity and space complexity and we use this to build Bangla Word2Vec.

CONCLUSION

Finally, we have perfectly build up an efficient algorithm and technique for Bengali stemming. That is more perfect than the exists technique in the perspective of accuracy and complexity. Although we have faced some problem during the rules generation after adding some exceptional rules we have overcome it successfully. This techniques algorithm can be used in many aspects for the development of the system(Sentiment analysis, Search engine, Spam filter, etc.) that is based on the Bengali Language.

REFERENCE

- [1] M.F. Porter, "An algorithm for suffix stripping", Program, 14(3) 1980, pp.130–137.
- [2] C.J. van Rijsbergen, S.E. Robertson and M.F.Porter, "New models in probabilistic information retrieval", *British Library Research and Development Report*, no. 5587, 1980.
- [3] J.B.Lovins, "Development of a stemming algorithm", *Mechanical Translation and Computational Linguistics* 11, 1968, pp. 22-31.
- [4] C.D. Paice, "An evaluation method for stemming algorithms", In the *Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval*, 1990, pp. 42 50.
- [5] J. Dawson, "Suffix removal and word conflation", *ALLCbulletin*, 2(3), 1974, pp. 33–46.
- [6] V. Berlian, S.N. Vega, and S. Bressan, "Indexing the Indonesian web: Language identification and miscellaneous issues", In the *Tenth International World Wide Web Conference*, Hong Kong. 2001.
- [7] S.Y. Tai, C.S. Ong, and N.A. Abdullah, "On designing an automated Malaysian stemmer for the Malay language", (Poster) In the *Proceedings of the fifth international workshop on information retrieval with Asian languages*, Hong Kong, 2000, pp. 207-208.
- [8] W. Kraaij and R. Pohlmann, "Viewing stemming as recall enhancement", In the *Proceedings of ACM SIGIR96*, 1996, pp. 40-48.
- [9] M. Popovic and P.Willett, "The effectiveness of stemming for natural language access to Slovene textual data", *JASIS*, 43 (5), 1992, pp. 384-390.

- [10] F.C. Ekmekcioglu, M.F. Lynch and P.Willett, "Stemming and n-gram matching for term conflation in Turkish texts", *Information Research News*, 7 (1), 1996, pp. 2-6.
- [11] I. Moulinier, A. McCulloh and E. Lund, "Non-English monolingual retrieval in Cross language information retrieval and evaluation", In the Proceedings of the *CLEF 2000 workshop*, C. Peters, Ed.: Springer Verlag, 2001, pp. 176-187.
- [12] C. Monz and M.de Rijke, "Shallow morphological analysis in monolingual information retrieval for German and Italian in Cross-language information retrieval and evaluation", In the *Proceedings of the CLEF 2001 workshop*, C. Peters, Ed., Springer Verlag, 2001.
- [13] M. Greengrass, A.M. Robertson, S. Robyn, and Willett, "Processing morphological variants in searches of Latin text", *Information research news*, 6 (4), 1996, pp. 2-5.
- [14] S. Sarkar and S. Bandyopadhyay, "Design of a rule-based stemmer for natural language text in bengali," in Proceedings of the IJCNLP-08 workshop on NLP for Less Privileged Languages, 2008.
- [15] M. Islam, M. Uddin, M. Khan, et al., "A light weight stemmer for bengali and its use in spelling checker," 2007.
- [16] D. Ganguly, J. Leveling, and G. J. Jones, "Dcu@ fire-2012: rule-based stemmers for bengali and hindi," 2012.
- $[17]\ https://github.com/shahidul034/stemmer-accuracy-test$
- [18]https://bn.wikisource.org/wiki/%E0%A6%AF%E0%A7%8B%E0%A6%97%E0%A6%BE%E0%A6%AF%E0%A7%8B%E0%A6%97/%E0%A7%A7
- [19]বাংলা ব্যাকরণ (নবম ও দশম শ্রেণি) Bangla Bakaron Class 9-10.