

An Innovative Lemmatization Technique for Bangla Nouns by using Longest Suffix Stripping Methodology in Decreasing Order

Alok Ranjan Pal

Detp. of Computer Sc. and Engg.
College of Engg. & Mgmt, Kolaghat
Kolaghat, India
chhaandasik@gmail.com

Niladri Sekhar Dash

Linguistic Research Unit
Indian Statistical Institute
Kolkata, India
ns_dash@yahoo.com

Diganta Saha

Dept. of Computer Sc. and Engg.
Jadavpur University
Kolkata, India
neruda0101@yahoo.com

Abstract—In this proposed work, an attempt is made to find out the root part from inflected Bangla nouns by applying an innovative technique by using longest suffix stripping methodology in decreasing order. The test data is generated from a Bangla text corpus developed in the TDIL Project of the Govt. of India. The exhaustive suffix list obtained from the research work carried out at *Linguistic Research Unit of Indian Statistical Institute, Kolkata* while the Bangla non-inflected noun list used in this work is obtained from a wordlist generated by *Pashchimbanga Bangla Akademi, Kolkata* and available in the net. The algorithm is applied on randomly selected 1273 noun instances and accuracy is achieved around 94%.

Keywords—*lemmatization; suffix stripping; part-of-speech tagging; morphology*

I. INTRODUCTION

Nouns in most of the natural languages are available either as single bases in their non-inflected forms or in affixed bases in inflected forms. In case of the second option, nouns are inflected in two ways: inflectional morphology (e.g., cows, baby-babies, ox-oxen, etc.) and derivational morphology (e.g., organization-organizational, computer-computerized, etc.). In most cases the range of inflection is restricted within a few types. In Bangla - a language of the Indian sub-continent, the range of inflection types is quite large. Here a noun can have several types of inflection. For instance, the words *māthā* “head” can have several inflected forms, such as, *māthāte* “in head”, *māthāā* “the head”, *māthāy* “on head”, *māthār* “of head”, *māthāgulo* “heads”, *māthārā* “heads”, *māthāi* “the head”, *māthāder* “of heads”, *māthāri* “of head itself”, etc. This happens in case of single inflection only, which can be doubled if the processes of double inflection addition is invoked, such as, *māthātei* “in head itself”, *māthāāi* “the head itself”, *māthāyo* “on the head also”, *māthāri* “of head itself”, *māthāgulate* “on the heads”, *māthārāi* “heads themselves”, *māthāio* “the head also”, *māthāderi* “of heads themselves”, etc. There are few determiners like, -ā, -āi, -khānā, -khāni, -khānek as well as a few emphatic markers, like -i and -o, etc. which are often tagged with inflected nouns for double inflection. So, finding base forms of Bangla nouns is a challenging task. The situation becomes far more complicated in case of nouns made

with double suffixes like *māthāgulate* “in heads” (two suffixes: -gulo and -te) and multiple suffixes, such as, *māthāgulatei* “in heads themselves” (three suffixes: -gulo, -te and -i).

In the proposed approach, an algorithm has been designed to find the root part from inflected nouns by stripping the longest suffix in reserve order. To increase the efficiency of the algorithm, an alphabetically sorted list of non-inflected words has been used in this work. The exhaustive suffix list obtained from the research work carried out at *Linguistic Research Unit of Indian Statistical Institute, Kolkata* while the Bangla non-inflected noun list used in this work is obtained from a wordlist generated by *Pashchimbanga Bangla Akademi, Kolkata* and available in the net.

The organization of the paper is as follows: Section 2 describes the survey on lemmatisation and stemming in Indian languages; nature of Bangla nominal morphology is discussed in Section 3; Section 4 depicts the overall approach in detail; results and corresponding evaluations are presented in Section 5; Section 6 concludes the discussion.

II. SURVEY ON LEMMATISATION AND STEMMING IN INDIAN LANGUAGES

Due to limited number of morphological variants many works on lemmatization and stemming have been done for English. For Bangla, however, not much work is done due to its morphological complexities involved in word formation. Since the language is very rich in morphological diversity, less amount of work on lemmatisation or stemming has been done in it. This does not stand true for Bangla alone, this is equally true for most of the Indian languages as the following references show:

Ramanathan and Rao [1] designed a lightweight stemmer for Hindi language in 2003. The approach was tested on 35977 unique words taken from different sources like politics, sports, business, health, etc. The error in this method was 4.68% for under-stemming and 13.84% for over-stemming task. Majumder et al. [2] proposes a statistical model for stemming named as YASS (Yet Another Suffix Stripper). It is a clustering-based approach that works on string distance measures sans linguistic knowledge. The model clusters the

lexicons generated from a text corpus into some homogenous groups and each group represents an equivalent class containing morphological variants of a single root word. Dasgupta and Nag [3] designs a strategy for unsupervised morphological processing of Bangla words. This method is tested on a set of 4110 human-segmented Bangla words with 83% success rate. Pandey and Siddiqui [4] proposes an unsupervised stemming method for Hindi words based on Goldsmith approach. This method is evaluated on 1000 words randomly selected from the Hindi WordNet and it records 89.9% accuracy. Majgaonker and Siddiqui [5] proposes an unsupervised stemming approach for Marathi words. When the algorithm is tested on 1500 manually stemmed words, it comes out with 82.5% accuracy. Suba et al. [6] develop two kinds of stemmers for Gujarati words. The first one is a light-weight inflectional stemmer based on a hybrid approach, while the second one is a heavy-weight derivational stemmer based on a rule-based approach. While the inflectional stemmer produces around 90.7% accuracy, the derivational stemmer gives around 70.7% accuracy. Zahurul Islam, et al. [7] proposes a light-weight stemmer for Bangla which can be used in a spell checker with some satisfactory outputs. Dash [13] has recently proposed a simple linguistic strategy for lemmatizing Bangla inflected nouns, which may be further converted into a rule-based approach studded with algorithms for stemming for lemmatizing not only inflected nouns but also inflected verbs and adjectives of the language. Along this line some works may be credited to Chaudhuri [10], Sarkar and Bandyopadhyay [11], and Dasgupta and Khan [12]. Most of these works are based on rule-based approach for stemming and lemmatization.

III. NATURE OF BANGLA NOMINAL MORPHOLOGY

In Bangla, the scale of morphological variety is so large that no defined rule can cover all types of inflections for nouns. For example, the following instances can be solved by single-suffix stripping approach:

- (1) baiguli = bai + (-guli)
- (2) baite = bai + (-te).

But, this becomes more challenging when more than one suffix remains appended with the word. For example, the following instances can be solved by multiple-suffix stripping approach:

- (3) baigulite = bai + (-guli) + (-te)
- (4) baitite = bai + (t-i) + (-te)

And, the following instances are the obstacles in the lemmatization process, as the identification of accurate suffix part is quite complicated:

- (5) maayer = maa + (-yer) is right
- (6) chabiTi = chabi + (-Ti) is right
- But,
- (7) samayer = sama + (-yer) is wrong.
- (8) maaTi = maa + (-Ti) is wrong.

The proposed approach takes the tagged text as input and generates the lemmatized form of the individual noun.

IV. THE PROPOSED APPROACH

In the proposed approach, the tagged nouns are given as input and the base forms of the inflected nouns are generated as output.

First, the Bangla text corpus, developed in the TDIL project of the Govt. of India is taken as an input. But, the texts in the corpus are non-normalized in nature. So, these texts are passed through a series of manual text normalization steps (refer section V.A).

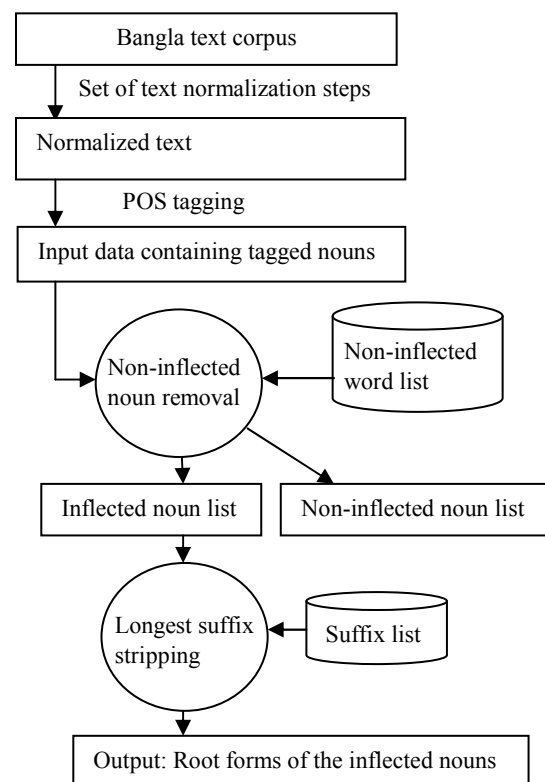
After that, the nouns in the normalized texts are tagged (refer section V.D).

Next, the non-inflected nouns are detected and removed from the test data by the help of the non-inflected word list developed by the *Bangla Academy*. This procedure was applied to increase the accuracy of the result, as the last characters of several non-inflected nouns, match with different suffixes.

Then, the inflected nouns are passed through the longest-match suffix stripping algorithm. The suffix list is collected from the *Linguistic Research Unit of Indian Statistical Institute, Kolkata* [13]. To match the longest suffix with the end part of the inflected nouns, the suffixes are rearranged according to their lengths in decreasing order (refer section V.C).

Thus the suffix stripped root forms of the inflected nouns are obtained as output. The overall approach is represented by the following diagram (refer Fig. 1).

Fig.1. Pictorial representation of the overall approach



A. Non-normalized texts

Fig.2. A sample non-normalized text from the Bengali corpus



Fig. 3. A sample normalized text



length (refer Fig. 4). According to the morphological feature of the Bengali nouns, they end with no-suffix or single suffix or multiple suffixes. This algorithm strips the suffix of longest match from a noun to find the root form of it.

Fig. 4. The suffixes in decreasing length

দিগেরও, খানাতেই, খানাতেও, খানিকেই, খানিকেও, খানিতেই, খানিতেও, টুকুকেই, টুকুকেও, টুকুতেই, টুকুতেও, গুলিকেই, গুলিকেও, গুলিতেই, গুলিতেও, গুলোকেই, গুলোকেও, গুলোতেই, গুলোতেও, গুলাকেই, গুলাকেও, গুলাতেই, গুলাতেও, দিগেরও, খানাকে, খানাতে, খানিকে, খানিতে, টুকুকে, টুকুতে, গুলিকে, গুলিতে, গুলোকে, গুলোতে, গুলাকে, গুলাতে, দিগের, টাকেরই, টাকেরও, খানারও, খানায়ও, খানিরই, খানিরও, টুকুরই, টুকুরও, গুলিরই, গুলিরও, গুলোরই, গুলোরও, গুলোয়ই, গুলোয়ও, গুলারই, গুলারও, গুলায়ও, দিগেরও, দিগেরই, দিগেরই, খানারই, খানায়ই, গুলায়ই, দেরকেই, দেরকেও, দিগকেই, দিগকেও, খানাই, খানাও, খানিই, খানিও, টুকুই, টুকুও, গুলিই, গুলিও, গুলোই, গুলোও, গুলাই, গুলোও, খানার, খানায়, খানির, টুকুর, গুলির, গুলোর, গুলোয়, গুলার, গুলায়, দেরকে, দিগকে, দিগের, টাকেই, টাকেও, টাতেই, টাতেও, টিকেই, টিকেও, টিতেই, টিতেও, খানা, খানি, টুকু, গুলো, গুলি, গুলো। রেবাই, রেবায়, দেবই, দেবও, দিগই, দিগও, টাকে, টাতে, টিকে, টিতে, তেই, তেও, টাইই, টাইও, টায়ই, টায়ও, টিরই, টিরও, য়েরই, য়েরও, রে, দিগ, তে, য়ের, টিই, টিও, টাই, টাও, রাই, রাও, টায়, টার, টির, কেই, কেও, তেই, তেও, রেই, রেও, ায়ই, ায়ও, ারই, ারও, ায়ই, ায়ও, টা, টি, ারে, কে, তে, কোর, ায়, হেই, হেও, রই, রও ই, ও, রা, ায়ে

D. Partial view of the input data

The nouns in the input data were tagged in the following way. The input text contained total 1273 tagged nouns (refer Fig. 5).

Fig. 5. The tagged input data

আধুনিক হিসাবশাস্ত্র/noun ক্রমশঃই জটিল হইতে জটিলতর আকার/noun ধারণ/noun করিতেছে,শতাব্দীর/noun পর শতাব্দী/noun ধর্ম্মা ব্যবসায়িগণ ও রাষ্ট্রের/noun রাজস্ববিভাগ কর্তৃক ব্যবহারিক ক্ষেত্রে হিসাবের/noun সর্বাধিক প্রয়োগের/noun ফলে বর্তমানে হিসাবশাস্ত্র ব্যবহারিক বিভাগে/noun মর্যাদা/noun লাভ করিয়াছে,হিসাবশাস্ত্রের/noun ফলে হিসাবশাস্ত্রের/noun পুঙ্খ নুপুঙ্খ রাষ্ট্র/noun ও ব্যবসায়িগণের/noun ভূমিকা/noun অনস্বীকার্য,পরিবর্তিত পরিহিসিতে/noun ন্যূন অর্থাৎ সময়েঃ/noun পরিবর্তনের/noun সঙ্গে সঙ্গে মূলতঃ ইহাদের ব্যবহারিক সুবিধা/noun এবং উপযোগের/noun দিকে লক্ষ্য/noun রাখিয়াই নূতন নূতন নীতি/noun হিসাবশাস্ত্রের/noun সমঝিতি হইতেছে,এই সকল নীতি ও সূত্রগুলি/noun শিক্ষণীয় এবং শিক্ষাপ্রাপ্ত ব্যক্তিগণের অভিজ্ঞতালব্ধ মূল্যবান ফসল/noun,হিসাবশাস্ত্র/noun সম্পর্কে পূর্ণাঙ্গ জ্ঞান/noun লাভ করিতে হইলে ইহার ক্রমবিকাশের/noun ইতিহাস ইহার প্রকৃতি কার্যাবলী/noun ইত্যাদি সম্পর্কে পরিষ্কার ধারণা/noun থাকা প্রয়োজন/noun,যেইজন্য এই অধ্যায়ে/noun উক্ত বিষয়াদি সম্পর্কে বিস্তারিত/noun আলোচনা/noun করা হইল,হিসাবশাস্ত্রের/noun ভিত্তি মানবসভ্যতার/noun ক্রমবিকাশের/noun ইতিহাসের/noun মতই হিসাবশাস্ত্রের/noun ক্রমবিকাশের/noun ইতিহাস/noun অতি প্রাচীন ও ভৌতিকগণ,প্রাচীন যুগ/noun

E. Partial view of the output file

The partial view of the output file is given below (refer Fig. 6).

Fig. 6. The generated output data

আদুনিচ হিসাবশাস্ত্র /NOUN/হিসাবশাস্ত্র ক্রমশঃই জটিল হইতে জটিলতর আকার /NOUN/আকার ধারণ /NOUN/ধারণ করিতেছে। শতাব্দীর /NOUN/শতাব্দী পূর শতাব্দী /NOUN/শতাব্দী ধায়ায় ব্যবসায়িগণ ও রাষ্ট্রের /NOUN/রাষ্ট্র রাজস্ববিভাগ কর্তৃক ব্যবহারিক ক্ষেত্রে হিসাবের /NOUN/হিসাব সম্বন্ধিক প্রয়োজের /NOUN/প্রয়োজ বহনোরে হিসাবশাস্ত্র ব্যবহারিক বিজ্ঞানের /NOUN/বিজ্ঞান মর্যাদা /NOUN/মর্যাদা লাভ করিয়াছে। হিসাবশাস্ত্রের /NOUN/হিসাবশাস্ত্র ক্রমবিকাশের /NOUN/ক্রমবিকাশ পিঠের /NOUN/পিঠের /NOUN/রাষ্ট্র ও ব্যবসায়িগণের /NOUN/ব্যবসায়িগণ ভূমিক /NOUN/ভূমিকা অনবধিক। পরিবর্তিত পরিস্থিতিতে /NOUN/পরিস্থিতি অর্থাৎ সময়ের /NOUN/সময় পরিবর্তনের /NOUN/পরিবর্তন সঙ্গে সঙ্গে মূলতঃ ইহাদের ব্যবহারিক সুবিধা /NOUN/সুবিধা এবং উপযোগের /NOUN/উপযোগ দিকে লক্ষ্য /NOUN/লক্ষ্য রাখিয়াই নতুন নতুন নীতি /NOUN/নীতি হিসাবশাস্ত্রের /NOUN/হিসাবশাস্ত্র সংযোজিত হইতেছে। এই সকল নীতি ও সূত্রগুলি /NOUN/সূত্রগুলি শুল্ক শিক্ষণার্থ এবং শিক্ষাপ্রাপ্ত ব্যক্তিগণের অভিজ্ঞতাগুলি অনুমান ফসল /NOUN/ফসল হিসাবশাস্ত্র /NOUN/হিসাবশাস্ত্র সম্পর্কে পূর্ণাঙ্গ জ্ঞান /NOUN/জ্ঞান লাভ করিতে হইলে ইহার ক্রমবিকাশের /NOUN/ক্রমবিকাশ ইতিহাস ইহার প্রকৃতি কার্যাবলী /NOUN/কার্যাবলী ইত্যাদি সম্পর্কে /NOUN/সম্পর্কে পরিকার ধারণ /NOUN/ধারণ থাকা প্রয়োজন /NOUN/প্রয়োজন। সেইজন্য এই অধ্যায়ে /NOUN/অধ্যায়ে উক্ত বিষয়াদি সম্পর্কে বিস্তারিত /NOUN/বিস্তারিত আলোচনা /NOUN/আলোচনা করা হইল।

VI. RESULT AND DISCUSSION

The model was evaluated on 1273 noun instances and 94% accuracy was achieved.

The efficiency of the algorithm is measured based on the three parameters “Precision”, “Recall”, and “F-Measure” as well. Precision (P) is the ratio of the “matched instances according to the human decision” and “number of instances responded by the system”. Recall value (R) is the ratio of “number of instances matched with the human decision” and “total number of instances in the dataset” and F-Measure is evaluated as “ $(2 * P * R / (P + R))$ ”.

As the system responded to all the target instances present in the dataset (either correctly or wrongly), the Precision and the Recall values are same. $P=R=1197/1273=0.94$ and F-Measure = $1.7672/1.88=0.94$.

According to the morphological feature of the Bengali nouns, they end with no-suffix or single suffix or multiple suffixes. As the algorithm stripped the longest suffix from the nouns, the root of the nouns were derived in most of the cases, except a few, which are discussed in section VI B.

Generally, the Bengali nouns are not deformed with a prefix part, like western languages. So, prefix stripping is not considered for Bengali noun lemmatization.

A. Complexity of the algorithm

The complexity of the algorithm is $O(n^2)$, as for every noun in the test set, the whole list of the suffixes is traversed individually.

B. Few close observations

A few obstacles in this approach have been discussed in Section 3. According to the nature of the Bangla nouns, few words end with the characters which are used as suffix as well. For example (refer Fig. 7)

Fig. 7. Few obstacles in Bengali morphology

হিসাবশীর্ষক/noun/হিসাবশীর্ষক	কার্যকরিত্ব/noun/কার্যকরিত্ব	মাপকাঠি/noun/মাপকাঠি
সময়েতি/noun/সময়েতি	ব্যবহার/noun/ব্যবহার	কপালে/noun/কপালে
মানবসত্তা/noun/মানবসত্তা	ব্যবহার/noun/ব্যবহার	রচনা/noun/রচনা
জন্ম/noun/জন্ম	পরিচালনা/noun/পরিচালনা	লেখা/noun/লেখা
সূচনা/noun/সূচনা	পরিচালনা/noun/পরিচালনা	এলাকা/noun/এলাকা
প্রস্তাব/noun/প্রস্তাব	কার্যকরিত্ব/noun/কার্যকরিত্ব	পালনে/noun/পালনে
সত্যতা/noun/সত্যতা	ব্যবহার/noun/ব্যবহার	হাতে/noun/হাতে
আশ্রয়/noun/আশ্রয়	পল্লি/noun/পল্লি	বৃক্ষ/noun/বৃক্ষ
জীবন/noun/জীবন	সংকেত/noun/সংকেত	চোখে/noun/চোখে
জীবিকা/noun/জীবিকা	বাংলা/noun/বাংলা	
চেনা/noun/চেনা	কাগজ/noun/কাগজ	
সমস্যা/noun/সমস্যা	পল্লি/noun/পল্লি	
বিনিময়ে/noun/বিনিময়ে	বুকে/noun/বুকে	

VII. CONCLUSION AND FUTURE WORK

According to the structure of the Bengali nouns, the degradation in accuracy of the performance was obvious. We wish to deal with these erroneous situations by forming new learning rules.

Acknowledgement

The authors acknowledge Mr. Saptarshi Datta, Miss. Trisha Saha and Miss. Sudipta Bose for their spontaneous help to this work.

Reference

- [1] Ramanathan and D.D. Rao. A Lightweight Stemmer for Hindi. Workshop on Computational Linguistics for South-Asian Languages, EACL, 2003.
- [2] P. Majumder, M. Mitra, S.K. Parui, G. Koley, P. Mitra, and K. Datta, “YASS: Yet Another Suffix Stripper, Association for Computing Machinery Transactions on Information Systems,” 25(4), 2007, pp.18-38.
- [3] S. Dasgupta and V. Nag, “Unsupervised Morphological Parsing of Bengali. Language Resources and Evaluation,” 40(3-4), 2006, pp. 311-330.
- [4] A.K. Pandey and T.J. Siddiqui, “An Unsupervised Hindi Stemmer with Heuristic Improvements. Proceedings of the Second Workshop on Analytics for Noisy Unstructured Text Data,” 2008, pp. 99-105.
- [5] M.M. Majgaonker and T.J. Siddiqui, “Discovering Suffixes: A Case Study for Marathi Language,” International Journal on Computer Science and Engineering, 2(8), 2010, pp. 2716-2720.
- [6] K. Suba, D. Jiandani and P. Bhattacharyya, “Hybrid Inflectional Stemmer and Rule-based Derivational Stemmer for Gujarati,” Proceedings of the 2nd Workshop on South and Southeast Asian Natural Language Processing (WSSANLP), IJCNLP 2011, Chiang Mai, Thailand, 2011, pp. 1-8.
- [7] M. Z. Islam, M.N. Uddin and M. Khan, “A Light Weight Stemmer for Bengali and its Use in Spelling Checker,” Proceedings of the 1st Intl. Conf. on Digital Comm. and Computer Applications (DCCA07), Irbid, Jordan, March, 2007, pp. 19-23.
- [8] M.F. Porter, “An algorithm for suffix stripping. Program,” 14(3), 1980, pp. 130-137.
- [9] P. Kundu and B.B. Chaudhuri, “Error Pattern in Bengali Text,” International Journal of Dravidian Linguistics, 1999.
- [10] B.B. Chaudhuri, “Reversed word dictionary and phonetically similar word grouping based spell-checker to Bengali text,” Proceedings of LESAL Workshop, 2001.
- [11] S. Sarkar and S. Bandyopadhyay, “Study on Rule-Based Stemming Patterns and Issues in a Bengali Short Story-Based Corpus,” ICON-2009 (Poster).
- [12] S. Dasgupta and M. Khan. Morphological parsing of Bangla words using PCKIMMO. In: ICCIT 2004. (2004)
- [13] N.S. Dash, “Back to Basics: A Road to Return to Nominal Base through Lemmatization,” Proceedings of Abstracts of the 36th International Conference of the Linguistic Society of India (ICOLSI-36), 1-4 December 2014, Dept. of Linguistics, University of Kerala, Trivandrum, Kerala, India. Pp. 94-104.