# A Suffix-based Noun and Verb Classifier for an Inflectional Language

Navanath Saharia Department of CSE Tezpur University Assam, India 784028 nava.nath@yahoo.in Utpal Sharma Department of CSE Tezpur University Assam, India 784028 utpal@tezu.ernet.in Jugal Kalita Department of CS University of Colorado Colorado Springs, USA 80918 kalita@eas.uccs.edu

Abstract—Nouns and verbs pose the major challenge in partof-speech tagging exercises. In this paper we present a suffix based noun and verb classifier for Assamese, an inflectional, relatively free word order Indic language. We used a tiny dictionary of frequent words to increase the accuracy. We obtained F-score of around 85%.

Keywords-Assamese, Noun, Verb, POS tagging.

## I. INTRODUCTION

In any language, nouns and verbs are the most crucial parts of a sentence. The noun class is always an open lexical category; its members can occur as the head word in the subject of a clause, the object of a verb, or the object of a preposition. In this paper we discuss our research towards part-of-speech (POS) tagging, an important step in any natural language processing task, where the goal is to automatically assign lexical category to each lexical object occurring in a given text. For any thorough study in computational linguistics, we need to POS tag each and every word of a sentence. Two factors determine the syntactic category of a word. The first is lexical information that is directly related to the category of the word, and other is contextual information related to the environment of the word. In this paper we mainly focus on the properties of a lexicon that may help in POS tagging. We choose Assamese, a morphologically rich inflectional Indic language, for the experiments. Spoken by about 30 million people, Assamese is the lingua-franca of north-eastern region of India, and the official language of the state of Assam.

This paper is organized is as follows. In section II we give a brief survey of literature related to this work, the state-of-theart and some linguistic characteristics of Assamese. In section III we present our methodology and a brief description of the corpora used. In section IV we report our experimental results and a discussion thereof. Section V concludes the paper with hints of future work.

## II. LITERATURE SURVEY

## A. Existing work

In [1] all POS tagging algorithms are categorised into three basic categories- rule based, stochastic, and hybrid. Grammatical and morphological rules are defined for rule based POS tagger. Most taggers, either rule based, stochastic or hybrid, are initially developed for English, and afterwards adapted to other languages. Brill's tagger [2], is a widely discussed linguistically motivated rule based POS tagger for English. In the two stage architecture of Brill's tagger, in the first stage the input tokens are initially tagged with their most likely tags; after that lexical rules are employed to assign tags to unknown tokens. On the other hand, TnT [3], a widely discussed statistical POS tagger based on a second order Markov model, was developed for English and German. It calculates the lexical probabilities of unknown words based on their suffixes. Comparison between statistical and linguistic rule based taggers shows that for the same amount of remaining ambiguity, the error rate of a statistical tagger is one order of magnitude greater than that of the rule based one [4]. The taggers described above are specifically designed for relatively fixed word order languages, where position of the word plays an important role. For relatively free word order languages, Dincer et. al [5] described a suffix based POS tagging approach for Turkish. They use the well-known Hidden Markov Model with a closed lexicon that consists of a fixed number of letters from word endings, and obtained accuracy 90.2%.

Indian languages are highly inflectional, morphologically rich and relatively free word order. Morphological richness and free word order nature make morphological analysis a crucial task in tagging of Indian language texts. While in some Indic languages such as Assamese, most case markers occur as suffixes, in others such as Hindi they occur as separate words, leading to local word grouping. Beyond that Indic languages have similar degree of free word order. Table I shows some of the reported POS taggers for Indian languages.

#### B. Assamese noun and verb morphology

As Assamese nouns and verbs are open lexical categories, if we can tag words in these classes correctly, tagging the remaining words in a text will be facilitated. In this work, we consider only the morpho-syntactic properties of Assamese words. Assamese words can be categorized into inflected classes (noun, pronoun, adjective and verb) and un-inflected classes (adverb and particle). Among inflected classes two

The Department of Computer science & Engineering, Tezpur University is funded by university grant commission (UGC)'s departmental research support (DRS) Phase-I under special assistance programme (SAP).

		TABLE I	I			
REPORTED POS	TAGGING	RESULTS	WITH	INDIAN	LANGUAGH	ES

Author	Approach	Language	Accuracy
[6]	HMM	Bengali	84.37%
		Hindi	70.67%
[7]	CRF	Bengali	65.47%
		Telegu	65.85%
		Hindi	69.98%
[8]	HMM	Bengali	67.52%
		Telegu	68.32%
		Hindi	71.65%
[9]	CRF	Bengali	80.63%
		Telegu	53.15%
[10]	HMM	Assamese	85.64%
[11]	CRF	Moninuri	72.04%
	HMM	Manipun	74.38%
[12]	Rule Based	Manipuri	69%

main types of inflection are noun inflection and verb inflection.

## C. Noun inflection

The inflection model of the noun in Assamese is depicted in Figure 1. Noun inflection represents gender, number and case in Assamese. For example, nouns ending with -जान (*jan*) and -जाने (*janI*) are identified as masculine and feminine nouns, respectively. All rules applied to noun inflection can be applied to pronouns, adjectives and even numerals. Table II and Table III show formation of compound and derivational words, respectively. Most compounds in Assamese are noun; although other forms are also not rare. Derivation takes place for suffixes, prefixes, or a combinations of both. The base for derivation can be a simple word or a compound word. Only suffixes can change the word category, prefixes do not change the category of a word [cf. Table III].



Fig. 1. Assamese noun inflection model

## D. Verb inflection

Assamese verbs are inflected with tense, aspect and modality (TAM). Traditionally, Assamese verbs are categorised as either *finite* or *non-finite*. Verb roots are in non-finite form for which tense, person or grammatical markers are added. In comparison to nouns, Assamese verb inflection is complex. [13], [14] reported 520 inflectional forms for root verb  $\exists \forall (bH : to sit)$ . Table IV shows some inflectional forms of verb  $\exists \forall (kr : to do)$ . An Assamese verb conjugator

TABLE II Formation of compound word in Assamese

Stem(Category)	Stem(Category)	New word(Category)
কথা (NN)	ছবি (NN)	কথাছবি (NN)
(talk : kathA)	(picture : Cbi)	(cinema : kathACbi)
কৃষ্ণ (NN)	পক্ষ (NN)	কুষ্ণুপক্ষ (NN)
(dark : krishna)	(fortnight : pakhya)	(dark fortnight :
		krishnapakhva)

TABLE III
Formation of derivational noun and verb in Assamese

Prefix	Root	Category	Suffix	New word	New
					category
-	আঙুলি	NN	আ	আঙুলিয়া	VB
-	কৰ	VB	আ	কৰা	VB
-	চল	VB	অন	চলন	NN
-	ডাঙৰ	ADJ	জনী	ডাঙৰজনী	NN
ন	হয়	VB	-	নহয়	VB
ন	কৰ	NN	এ	নকৰে	VB

is available in *http://www.tezu.ernet.in/~nlp/res.htm*. Table IV summarizes of some inflectional form of verb  $\overline{\Phi a}$  (*kr : to do*) and Table V shows some suffixes and their categories with example.

#### III. OUR APPROACH

We use a part of the EMILLE Assamese text corpus<sup>1</sup> (5300 sentences), jointly developed by Lancaster University and CIIL-Mysore. We tokenized our corpus as far as possible considering white space as word separator and punctuations ( $|,?,1\rangle$ ) as sentence terminator. Some examples of the difficulties in this task are given below.

- 1) Quite frequently the same place name is written in two ways, such as নতুনপাৰা and নতুন পাৰা. পাৰা (*pArA*) and বাৰী (bArI) are among most popular Assamese suffixes placed after village or neighbourhood names. In নতুন পাৰা, নতুন is considered separate adjective word, which qualifies পাৰা whereas actually নতুনপাৰা when used as a single word is a noun. Thus irregularities in placing white space and hyphen make tokenizing process a complex job.
- 2) Foreign words are very commonly used in Assamese, especially in news reports and scientific or technical writing. Foreign words take Assamese suffix. Based on the category of suffix we can identify the foreign words. For Example- U.G.C.' or U.G.C.- (of U.G.C). This makes the number of OOV words high.

In our method, we follow the following three basic steps to tagged tokenized text.

1) **Brute-force determination of suffix sequences:** *In this step, we obtain all possible sequences of noun suffixes following our model shown in Figure 1.* Assamese nouns and pronouns take more than one suffix in a sequence, though not all suffix sequences are grammatically

<sup>1</sup>http://www.emille.lancs.ac.uk/

TABLE IV Some Inflectional form of kr verb with respect to tense and person.

কিৰ (kr : to do)	1 <sup>st</sup> Person	2 <sup>nd</sup> Person (Familiar)	2 <sup>nd</sup> Person (Respect)	3 <sup>rd</sup> Person
Present	কৰোঁ (karo)	কৰ (kar)	কৰক (karaka)	কৰা (karA)
Past	কৰিলোঁ (karilo)	কৰিলি (karili)	কৰিলে (karile)	কৰিলা (karilA)
Future	কৰিম (karim)	কৰিবি (karibi)	কৰিব (kariba)	কৰিবা (karibA)
Present Perfect	কৰিছোঁ (karicho)	কৰিছ (karicha)	কৰিছে (kariche)	কৰিছা (karichA)
Past Perfect	কৰিছিলোঁ (karichilo)	কৰিছিলি (karichili)	কৰিছিল (karichil)	কৰিছিলা (karichilA)
Causative	-	কৰাবা (karAbA)	কৰোঁৱাওক (karowAok)	কৰোৱা (karowA)
Future Conditional	কৰিমচোন (karimson)	কৰিবিচোন (karibison)	কৰিবচোন (kkaribason)	কৰিবাচোন (karibAson)

TABLE V EXAMPLE OF SUFFIXES WITH CATEGORIES IN ASSAMESE

Case Marker	-ক, -্ৰ, -ত, -ই, এ,	মানুহৰ
	-ৰে, -লৈ, -ৰপৰা etc.	
Plural suffix	-বোৰ, -হত, -মথা, -সোপা	মানুহবোৰ
	-সমুহ, -গুন etc.	
Classifiers	-জন, -জনী, -কণ, -থন	মানুহজন
	-ডान, -পাত -টো, -টা etc.	
Verbal suffix	-ইছিলোঁ, -ইছিলি, -ইবি,	কৰিছিলোঁ
	-ইম, -আ etc.	

correct. For example

নাতিনীয়েককেইজনীমানেহে (nAtinIyekkeijanImAnehe) → নাতিনী + য়েক + কেই + জনী + মান + এ + হে

Noun+ inflected form of kinship  $noun^2$  + infix + feminine marker + plural marker + nominative case marker + particle.

We can obtain all possible sequences of noun suffixes from Figure 1. Some suffixes are always used with words from a small class of roots. For example the suffix -আখি (aAkhi) is always placed after केल (kal : banana) that is kalaAkhi. So in the next step, i.e., in sequence pruning we try to minimize the search space.

- 2) Suffix sequence pruning: In this step we filter out the non-valid suffix sequences from among all the sequences obtained in Step 1. Though a number of suffix sequences can possibly occur after a root word, we usually find only 3 suffix sequence in the corpus we studied, though there are exceptions. All suffix sequences are not valid. So if we list most of the valid suffix sequences beforehand using our linguistic knowledge, we need not go through all possible combinations of the suffixsequence. Figure 1 depicts a model to order legal noun suffixes to form sequences. From this model we obtain rules of suffix sequences attached to a noun.
- 3) Suffix stripping: In this step, we identify the noun

and verb roots based on the single suffix that occurs immediately after the root. For example, if we found the word মানুহবোৰৰ, it will first identify the suffix sequence বোৰৰ (বোৰ+ৰ: plural marker + genitive marker), which is a noun suffix and hence মানুহবোৰৰ is tagged as genitive plural noun and মানুহ as noun, in suffix free form. We find that the suffix 'a' (Nominative case marker for noun and the endings of  $2^{nd}$  person (familiar) past and present perfect tense marker) is the most ambiguous as it applies to both nouns and verbs.

All modules mentioned here are developed in Java. Table VII shows the statistics of test corpus.

## IV. RESULTS AND DISCUSSION

We are not able to categories some most frequent verbs with the method describe above. For example the inflectional form of  $\operatorname{AII}(jA : to go)$  verb is not the same as the inflectional form of  $\operatorname{AII}(kha : to eat)$  or  $\operatorname{AII}(gA : to sing)$  verb. Therefore to increase the accuracy we add 300 most frequently used verb root words in the form of a tiny dictionary. Table VIII shows obtained results. We mentioned above that the principle applied to noun can also be applied to pronoun and adjective. The inflected pronouns and adjectives in the corpus are tagged as nouns. As a result the tagging accuracy of nouns comes down. If we can embed the contextual information in our method, we hope, it will help increase POS tagging accuracy.

Ours is one of the earliest work on Assamese. The performance values of the few earlier works are not available.

#### V. CONCLUSION

We have implemented a suffix based noun and verb tagging approach for Assamese. We find that the performance of this method is better than stochastic approaches, such as HMM technique should be useful where required linguistic knowledge is available, but resources to prepare a large tagged corpus for training are not available. It will be interesting to compare results of this approach with those of stochastic approaches for other inflectional languages.

As the word order of Assamese is relatively free, we can not use positional information like in fixed word order languages. So a morpho-syntactic approach gives batter results in comparison to [6], [10]. Another important observation from this experiment is that though Assamese is relatively free word order, some parts of speech do not occur in the initial or final positions of the sentence. As a future work we will try to

Kinship Noun	1 <sup>st</sup> Person	$2^{nd}$ Person (Familiar)	2 <sup>nd</sup> Person (Respect)	3 <sup>rd</sup> Person
দেউতা (deutA : father)	দেউতা (deutA)	দেউতাৰ (deutAr)	দেউতাৰা (deutArA)	দেউতাক (deutAk)
তাই (bhAi : younger brother)	ভাই (bhAi)	ভায়েৰ (bhAyer)	ভায়েৰা (bhAyerA)	ভায়েক (bhAyk)
ককাই (kakAi : elder brother)	ককাই (kakAi)	ককায়েৰ (kakAyer)	ককায়েৰা (kakAyerA)	ককায়েক (kakyek)
শাহ (sAhu : mother in law)	শাহ (sAhu)	শাহিবেৰ (sAhuWer)	শাহুৱেৰা (sAhuwerA)	শাহিৱেক (sAhuwek)
$\overline{\Phi}\overline{\Phi}$ (kakA : grand father)	কিকা (kakA)	ককাৰ (kakAr)	ককাৰা (kakArA)	কিকাক (kakAk)

 TABLE VI

 Some inflected form of Assamese kinship noun with person.

#### TABLE VII Statistics of used corpus.

Total tokens

Sentences

5300 48027 28923 2629 TABLE VIII

Total nouns

Total verbs

TABLE VIII PRECISION, RECALL AND F-MEASURE OF OUR APPROACH.

Category	Precision	Recall	F-measure
Noun	0.87	0.79	0.82
Verb	0.91	0.87	0.88

embed the linguistics word agreement rules in tagging nouns and verbs in Assamese.

#### REFERENCES

- Daniel Jurafsky and James H. Martin, SPEECH and LANGUAGE PROCESSING, An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition, Pearson Education, 2000.
- [2] Eric Brill, "Transformation based error-driven learning and natural language processing: A case study in part of speech tagging," *Computational Linguistics*, 1995.
- [3] Thorsten Brants, "Tnt a statistical part-of-speech tagger," in Proceedings of the Sixth Applied Natural Language Processing (ANLP-2000), 2000.
- [4] Christer Samuelsson and Atro Voutilainen, "Comparing a linguistic and a stochastic tagger," in *Proceedings of the Thirty-Fifth Annual Meeting of the Association for Computational Linguistics and Eighth Conference of the European Chapter of the Association for Computational Linguistics*. 1997, pp. 246–253, ACL.
  [5] Taner DINCER, Bahar KARAOGLAN, and Tarik Kisla, "A suffix
- [5] Taner DINCER, Bahar KARAOGLAN, and Tarik Kisla, "A suffix based part-of-speech for Turkish," in *Proceedings of 5th International conference on Information Technology: New Generations*, 2008.
- [6] Sandipan Dandapat and Sudeshna Sarkar, "Part-of-speech tagging for Bengali with Hidden Markov Model," in *Proceedings of NLPAI ML* workshop on Part of speech tagging and Chunking for Indian language, 2006.
- [7] Dilip Rao and David Yarowsky, "Part of speech tagging and shallow parsing of Indian languages," in *Proceedings of IJCAI-07 workshop on Shallow Parsing for South Asian Languages (SPSAL)*, 2007.
- [8] G. M. Ravi Sastry, Sourish Chaudhuri, and P Nagender Reddy, "A HMM based part-of-speech and statistical chunker for 3 Indian languages," in Proceedings of IJCAI-07 workshop on Shallow Parsing for South Asian Languages (SPSAL), 2007.
- [9] Asif Ekbal, Samiran Mandal, and Sivaji Bandyopadhyay, "POS tagging using HMM and rule based chunking," in *Proceedings of Workshop on Shallow Parsing for South Asian Languages (SPSAL)*, 2007.
- [10] Navanath Saharia, Dhrubajyoti Das, Utpal Sharma, and Jugal Kalita, "Part-of-Speech Tagger for Assamese Text," in *Proceedings of the 47th* Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, 2009.
- of the Asian Federation of Natural Language Processing, 2009.
  [11] T. D. Singh and S. Bandyopadhyay, "Manipuri POS tagging using CRF and SVM: A language independent approach," in *Proceedings of International Conference on Natural Language Processing(ICON)*, 2008.

- [12] T. D. Singh and S. Bandyopadhyay, "Morphology driven Manipuri POS tagger," in *Proceedings of IJCNLP-08 workshop on NLP for Less Privilege Languages*, 2008.
- [13] Utpal Sharma, Unsupervised Learning of Morphology of A Highly Inflectional Language, Ph.D. thesis, Tezpur University, 2007.
- [14] Utpal Sharma, Jugal K. Kalita, and Rajib K. Das, "Acquisition of morphology of an Indic language from text corpus," ACM Transactions on Asian Language Information Processing (TALIP), vol. 7, 2008.
- [15] Lilabati Saikia Bora, Asamiya Bhasar Ruptattva, M/s Banalata, Guwahati, 2006.