

Grammatical Error Detection Model for Assamese Sentences

Hirakjyoti Sarma¹, Debasish Das², Kishore Kashyap³

Student, Department of Information Technology, Gauhati University, Guwahati, India ^{1,2}

Assistant Professor, Department of Information Technology, Gauhati University, Guwahati, India ³

Abstract: Assamese is an Eastern Indo-Aryan language used mainly in the state of Assam. This paper presents an introduction to the Grammatical analysis model of the structures of Assamese sentences which has evolved from extensive computational, linguistic, and psycholinguistic research, provides a simple set of rules for describing the common properties of all natural languages and the particular properties of individual languages. Parsing is important in Linguistics and Natural Language Processing to understand the syntax and semantics of a natural language grammar. Parsing natural language text is challenging because of the problems like free word order, ambiguity and inefficiency. We proposed one model which is based on Top down parsing method and a context free grammar (CFG) for Assamese language with some limited domain.

Keywords: Assamese, Grammar Checker, CFG, Parsing.

I. INTRODUCTION

Natural language can be describe by formal language, in 1950 Noam Chomsky first attempted to give a precise characterization of the structure of natural languages, he tries to describe the syntax of a natural language in terms of mathematical rules. Which is later known as Context free grammar (CFG).Where grammar rule are represented by production rules and those grammar or production rules allow us to write a parser to check whether an input string is grammatically correct or not.

We have consider some simple Assamese sentence and design production rule for them .Using left factoring we remove left recursion problem of our proposed production rules. We have taken the top down parsing scheme non recursive predictive parsing algorithm to parse our proposed grammar for Assamese language. We have taken tagged sentences as input of our designed parser.

II. A BRIEF OVERVIEW OF ASSAMESE LANGUAGE

The Assamese language is cordially associated with the most important Indo-European Language Group. We have to study the Indo-European Language group to find the origin of Assamese language though it seems that it is made up of the Proto-Astrolied and Sino-Tibetan Language Group. Ascoli divided the Indo-European Language group into two main group viz. Satam and Centum. Indo-Aryan Languages are derived from the Indi-Iranian group which is one of the four sub-division of Satam. Assamese Language has also come through the three stages [(1) Old Indo-Aryan 1500BC-600BC→ (2) Middle Indo-Aryan 600BC-1000AD→ (3) New Indo-Aryan 1000AD-till now] of Indo-Aryan Language as the other Modern Indian Languages. The Indo-Aryan Languages, viz. Assamese, Bangla, Oriya etc., are derived from Avahattha through MagadhiApravhransa. The earliest evidence of Assamese dates back to the literature of the Charyapadas, written by a few Buddhist scholars.

The Assamese language present in the Charyapadas, reflects the initial stages of the evaluation of the Assamese language. There are eight vowel phonemes and twenty-one consonant phonemes including two semi-vowels in Standard Colloquial Assamese.

III.METHODOLOGY

DESIGNING CONTEXT FREE GRAMMAR FOR ASSAMESE LANGUAGE:

$G = (\Sigma, V, S, P)$ Where Σ is a finite nonempty set called the terminals, which will be parts of speech for Assamese language. For our proposed parser we consider six types of parts of speech. They are Noun, Pronoun, Article, Adjective, Conjunction, and Verb.

$S \rightarrow NPVP NP$
$NP \rightarrow$ noun conjunction noun nounART Noun pronoun conjunction pronoun noun pronoun noun nounpronoun noun noun noun conjunction pronoun nounnoun noun pronoun pronoun nounadjective noun pronoun adjective noun noun conjunction noun noun noun adjective pronoun conjunction noun Pronoun ART pronoun pronoun conjunction pronoun pronoun pronoun noun pronounpronoun pronoun pronoun noun conjunction pronoun pronounnoun pronoun pronoun pronoun pronounadjective pronoun pronoun adjective pronoun noun conjunction noun pronoun noun adjectiveVP \rightarrow noun verb noun noun verb verbnoun noun verb verbnoun noun verb verb noun verb noun pronounverb pronoun noun verb verbnoun prooun verb verbnoun pronoun verb verb pronoun verb noun verb verbadjective verbnoun adjective noun verb pronoun adjective noun verb

Fig 1: Proposed production rules for context free grammar for Assamese language

S→NP S1 S1→VP/ ε NP→NP5 NP1/adjectiveNP3 NP1→conjunction noun/ART/ ε /AP/noun NP4/pronoun NP2 NP2→conjunction pronoun/NP5/ ε /adjective NP3→ ε /adjective/NP5/conjunction adjective NP4→conjunction NP5/ ε NP5→noun/pronoun	AP→adjectiveAP1 AP1→ ε /conjunction adjective VP→NP5 VP1/verbVP3/adjective noun verb VP4 VP1→verbVP2/noun verb verb noun VP2→verb VP5/VP5 VP3→ adjective/VP5 VP4→pronoun/ ε VP5→noun/ ε
---	---

Fig 2: Production Rule after Left Factoring

First(S)= { noun,pronoun,adjective } First(S1)= { verb,adjective,noun,pronoun , ε } First(NP1)= { conjunction ,ART, ε,adjective,noun,pronoun } First(NP2)= { conjunction , ε,adjective,noun,pronoun } First(NP3)= { conjunction , ε,adjective,noun,pronoun } First(NP4)= { conjunction, ε } First(NP5)= { noun,pronoun } First(NP)= { noun,pronoun,adjective }	First(VP)= { verb,adjective,noun } First(VP1)= { verb,noun } First(VP2)= { verb, ε,noun } First(VP3)= { adjective,noun, ε } First(VP4)= { pronoun, ε } First(VP5)= { noun, ε } First(AP)= { adjective } First(AP1)= { ε,conjunction }
Follow(S1)= Follow(S)={ \$ } Follow(NP)= Follow(NP1))=Follow(NP2))= Follow(NP3)=Follow(NP4) = Follow(NP5)= { verb,adjective,noun,pronoun,conjunction ,ART , \$ }	Follow(AP)=Follow(AP1) Follow(VP)= Follow(VP1)= Follow(VP2)= Follow(VP3)= Follow(VP4)= Follow(VP5)={ \$ }

Fig 3: first and follow of our proposed grammar

Stack	Input	Production rule
\$\$	Pronoun noun verb\$	S-> NPS1
\$\$1 NP	Pronoun noun verb\$	NP->NP5NP1
\$\$1 NP1 NP5	Pronoun noun verb\$	NP5->Pronoun
\$\$1 NP1 Pronoun	Pronoun noun verb\$	
\$\$1 NP1	Noun verb\$	NP1->noun NP4
\$\$1 NP4 noun	Noun verb \$	
\$\$1 NP4	Verb \$	
\$\$1	Verb \$	NP4-> ε
\$\$VP	Verb \$	S1->VP
\$\$VP3 verb	Verb \$	VP->verb VP3
\$\$VP3	\$	
\$	\$	VP3-> ε

Fig 4: Parsing of “মইভাতখালোঁ”

	Noun	Pronoun	Verb	Adjective	Art	Conjunction	\$
S	NP S1	NP S1		NP S1			
S1	VP	VP	VP	VP			ε
AP	adjectiveAP1						
AP1	ε	ε	ε	ε		conjunction adjective	ε
NP	NP5 NP1	NP5 NP1		adjectiveNP3			ε
NP1	noun NP4	pronoun NP2	ε	AP	Art	conjunction noun	ε
NP2	NP5	NP5		Adjective		conjunction pronoun	ε
NP3	NP5	NP5	ε	Adjective		conjunction adjective	ε
NP4	ε	ε	ε	ε		conjunction NP5	ε
NP5	Noun	pronoun	ε				
VP	NP5 VP1	NP5 VP1	verbVP3	adjective noun verbVP4			
VP1	noun verb verb noun		verbVP2				ε
VP2	VP5		VerbVP5				ε
VP3	VP5		Verb	Adjective			ε
VP4		pronoun					ε
VP5	Noun						ε

Fig 5: predictive Parsing Table for our proposed Assamese grammar

For developing production rule for Assamese language we consider 31 types of subject and 16 types of predicate.

Using left factoring we eliminate left recursion of our developed grammar. After eliminating left recursion we rewrite our production rules as shown in Fig 2. After that we compute FIRST and FOLLOW for our proposed grammar using standard rules for computing first and follow, which results as shown in Fig 3. With the help of FIRST and FOLLOW we create parsing table using standard rules as shown in Fig 5.

Non recursive Predictive parsing algorithm

Let a be the first symbol of w and X be the top stack symbol;

```
While (  $X \neq \$$  ) { if (  $X = a$  ) pop the stack and let  $a$ 
be the next symbol of  $w$ ;
else if (  $X$  is a terminal ) error();
elseif (  $M[X,a]$  is an error entry ) error();
elseif (  $M[X,a] = X \rightarrow Y_1 Y_2 \dots Y_k$  ) { output the
production  $X \rightarrow Y_1 Y_2 \dots Y_k$  ;
pop the stack;
push  $Y_k Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$ 
onto the stack, with  $Y_1$  on top; }
Let  $X$  be the top stack symbol ; }
if  $X = \$$  , Sentence is Accepted.
```

IV. EXPERIMENTAL RESULT

We test our parser for subject part, between 31 inputs our parser will be able to parse 27 types of input sentence which is 87.09% of total. So based on subject part we can say that implemented parser is 87.09% efficient.

Successful parsing of an Assamese sentence "মই ভাত খালো" which is in the form of "Pronoun-Noun-Verb" is shown in Fig 3

V. CONCLUSION

This paper presents a top down predictive parsing approach to parse simple Assamese sentences. We only consider only six types of parts of speech while developing our proposed grammar for Assamese sentences. In future work we have to consider as many structures of Assamese sentences as we can and some more parts of speech.

REFERENCES

- [1] Deka and Kalita, "Adhunik rosona bisitra," 9th edition, Assam Book Dipo, 2007
- [2] Satyanath Bora, "Bohol biyakoron," 1st edition, Bina Library, 2012
- [3] Utpal Sharma, "Unsupervised Learning of Morphology of a Highly Inflectional Language," Ph.D. thesis, Department of Computer Science and Information Technology, Tezpur University
- [4] Navanath Saharia, Utpal Sharma, Jugal Kalita, "A First Step Towards Parsing of Assamese Text," 10 Dec.
- [5] K. M. Azharul Hasan, Al-Mahmud, Amit Mondal, Amit Saha, "RECOGNIZING BANGLA GRAMMAR USING PREDICTIVE PARSER," 10 Dec 2013
- [6] Monica S. Lam, Alfred V. Aho, Ravi Sethi, Jeffrey D. Ullman, "Compilers: Principles, Techniques and Tools," 2nd edition, Pearson (2008)
- [7] Tao Jiang, Ming Li, Bala Ravikumar, Kenneth W. Regan, "Formal Grammars and Languages," 10 Dec 2013

- [8] Kishore Kashyap, Hirakjyoti Sarma, Shikhar Kumar Sarma, "Luitspell: Development of an Assamese Language Spell Checker for Open Office Writer", vol. 2 Issue. 5, EJAET, 2015.
- [9] Mirzanur Rahman, Sufal Das and Utpal Sharma, "Parsing of part-of-speech tagged Assamese Texts," 10 Dec 2013.