

Verification of Bangla Sentence Structure using N-Gram

Nur Hossain Khan¹

¹ Islamic University, Kushtia, Bangladesh

Received: 10 December 2013 Accepted: 5 January 2014 Published: 15 January 2014

Abstract

Statistical N-gram language modeling is used in many domains like spelling and syntactic verification, speech recognition, machine translation, character recognition and like others. This paper describes a system for sentence structure verification based on Ngram modeling of Bangla. An experimental corpus containing one million word tokens was used to train the system. The corpus was a part of the BdNC01 corpus, created in the SIPL lab. of Islamic university. Collecting several sample text from different newspapers, the system was tested by 1000 correct and another 1000 incorrect sentences. The system has successfully identified the structural validity of test sentences at a rate of 93

Index terms— N-gram, sentence structure, corpus, witten-bell smoothing, word error.

1 Introduction

The goal of Statistical Language Modeling is to build a statistical language model that can estimate the distribution of natural language as accurate as possible. A statistical language model (SLM) is a probability distribution $P(s)$ over strings S that attempts to reflect how frequently a string S occurs as a sentence. By expressing various language phenomena in terms of simple parameters in a statistical model, SLMs provide an easy way to deal with complex natural language in computer. Therefore N-gram based modeling finds extensive acceptance to the researchers working with structural processing of natural language. An n-gram model is a type of probabilistic model for predicting the next item in such a sequence. More Microsoft Office Suite grammar checker, is also not Abstract-Statistical N-gram language modeling is used in many domains like spelling and syntactic verification, speech recognition, machine translation, character recognition and like others. This paper describes a system for sentence structure verification based on Ngram modeling of Bangla. An experimental corpus containing one million word tokens was used to train the system. The corpus was a part of the BdNC01 corpus, created in the SIPL lab. of Islamic university. Collecting several sample text from different newspapers, the system was tested by 1000 correct and another 1000 incorrect sentences. The system has successfully identified the structural validity of test sentences at a rate of 93%. This paper also describes the limitations of our system with possible solutions. gram". For a sequence of words, for example "the dog smelled like a skunk", the trigrams would be: "# the dog", "the dog smelled", "dog smelled like", "smelled like a", "like a skunk" and "a skunk #". N-Grams are typically constructed from statistics obtained from a large corpus of text using the co-occurrences of words in the corpus to determine word sequence probabilities. N-Grams have the advantage of being able to cover a much larger language than would normally be derived directly from a corpus. Open vocabulary applications are easily supported with N-Gram grammars [1]. Within the many application areas, an important application is to assess the probability of a given word sequence appearing in text of a language of interest in pattern recognition systems, speech recognition, OCR Intelligent Character Recognition (ICR), machine translation and similar applications [2]. By converting a sequence of items to a set of n-grams, it can be embedded in a vector space, thus allowing the sequence to be compared to other sequences in an efficient manner. The idea of n-gram based sentence structure verification has come from these opportunities provided by n-grams. Sentence structure verification is the task of testing the syntactical correctness of a sentence. It is mostly used in word processors and compilers. For applications like compiler, it is easier to implement because the vocabulary is finite for programming languages but for a natural language it is challenging because of infinite vocabulary. Three methods are widely used for

6 V. EXPERIMENTAL RESULTS AND DISCUSSION

46 grammar checking in a language; syntax-based checking, statistics-based checking and rule-based checking. In
47 syntax based grammar checking [3], each sentence is completely parsed to check the grammatical correctness of
48 it. The text is considered incorrect if the parsing does not succeed. In statistics-based approach [4], a corpus is
49 used to train a model. Some sequence will be very common others will probably not occur at all. Uncommon
50 sequences in the training corpus can be considered incorrect in this approach. In rule-based approach [5], a set of
51 hand crafted rules is matched against a text which has at least been POS tagged. This approach is very similar to
52 statistics-based approach, but the rules are developed manually. However, one of the most widely used grammar
53 checkers for English, above controversy [6]. It demonstrates that work onx x x x n i i i ? ? ? ? .., ,..... ,**3 2**
54 **1**

55 .
56 In Probability terms, this is nothing but ()x x x x n i i i P ? ? ? .., ,..... , | 2 1
57 . An n-gram of size 1 is referred to as a "unigram"; size 2 is a "bigram", size 3 concisely, an n-gram model
58 predicts x i based on is a "trigram"; and size 4 or more is simply called an "n-Year 2014 grammar checker in
59 real time is not very easy task; so starting the implementation for language like Bangla structural verification of
60 a sentence is a major feat. In our work, an effort has been made to develop system to verify Bangla sentence
61 structure using statistical or more specifically n-gram based method. This is because, this approach does not
62 need language resources like handcrafted grammatical rules, except for a corpus to train the language model
63 (LM). Given the scarcity of language resources for Bangla, proposed approach may be the only reasonable one
64 for the foreseeable future.

65 2 II. Techniques Adopted in the Proposed System

66 Similarly, the possible quad-grams (N-grams with N=4) are:

67 After training a model using above concept it was used to design a test system. For the purpose of testing
68 whether a sentence is correct or not, the number of N-grams (2, 3, or 4) in the sentence was counted first.
69 Using all the N-grams of the sentence, we have generated a score for the sentence. If the score is greater than
70 a predefined threshold, the sentence is syntactically correct. On the other hand, if the score is not greater than
71 the threshold, the sentence is syntactically incorrect.

72 3 III.

73 4 Training the N-gram Model

74 The first step to compute N-grams is counting unigrams. The unigram count and necessary software tools was
75 ready in the laboratory and the work was started from bigram count. After updating the existing software tools
76 bigrams, trigrams and quad-grams were identified, counted and stored in separated disk files. In all cases input
77 to the software was the sample corpus contained in file corpus. In statistical approach we can simply measure
78 the probability of a sentence using n-gram analysis. For example, using bigram probability of the sentence "????"
79 "?? ?? ??????" is, To estimate the structural correctness of a sentence, we calculate the probability of a sentence
80 using the formula above. If the value of the probability is above some threshold then we consider the sentence
81 to be structurally correct. Now if any of these three words are not in the corpus then the probability of the
82 sentence will become zero because of multiplication. To solve this problem, Witten-Bell smoothing [7] was used
83 to calculate the probability of a sentence in our work. A sample corpus was used in this work that is a part
84 of another corpus under construction in the speech and image processing lab of Islamic University, Bangladesh.
85 We have developed necessary programs to assemble sequences of N tokens into Ngrams. Typically N-grams are
86 formed of contiguous tokens that occur one after another in the input corpus. IV.

87 5 The Test System

88 For the purpose of testing whether a sentence is correct or not, at first, all the number of bigrams of the sentence
89 was counted. Getting probabilities from the respective models, Witten-Bell smoothing was applied to compute
90 a set of probabilities contained all nonzero values. Multiplying all the bigrams of the sentence, a score for the
91 sentence was generated. If the score is greater than a predefined threshold, the sentence is syntactically correct.
92 The functional block diagram of the system is shown in figure ???. For the trigram or quadgram models, the same
93 algorithm was followed by replacing only the bigrams with trigrams or quad-grams respectively.

94 6 V. Experimental Results and Discussion

95 In our experiment, 1000 sentences collected from the web edition of a daily newspaper to form a test set. The
96 test set was disjoint from the training corpus. All of these 1000 sentences were structurally correct. Taking these
97 correct sentences as input, the result generated by the test system is shown in table-1. For another experiment,
98 All of these 1000 sentences were modified to make structurally incorrect and presented again as input to the test
99 system. The result generated by second experiment is also shown in table-1.

100 **7 Discussion**

101 The word-error in Bangla can belong to one of the two distinct categories, namely, non-word error and real-word
102 error. A string of characters separated by spaces without a meaning is a non-word. By real-word error we mean a
103 valid but not the intended word in the sentence, thus making the sentence syntactically or semantically ill-formed
104 or incorrect. The developed system can identify both types of errors with an failure rate of 6.9% on average. The
105 major cause of this error is the volume of training corpus. As large as the volume of training corpus so will be
106 success rate.

107 **8 VII.**

108 **9 Conclusion**

109 We have developed a statistical Sentence structure verifier for Bangla, which has a reasonably good performance
110 as a rudiment Sentence verifier. By increasing the volume of training data the performance of the system can be
improved and a hybrid system combining both statistical and rule based system can be developed. ¹



Figure 1: T © 2014

111

¹© 2014 Global Journals Inc. (US)

Unigram	Frequency
আবার	২
জমজমাট	২
রাজনীতি	২
দীর্ঘদিনের	১
জড়তা	১
কাটিয়ে	৩
ফের	১
সরগরম	১
রাজপথ	১

Figure 2:

$$P(\text{“রহিম ফুটবল খেলে”}) = P(\text{রহিম} | \langle s \rangle) * P(\text{ফুটবল} | \text{রহিম}) \\ * P(\text{খেলে} | \text{ফুটবল})$$

Figure 3:

(রহিম, ফুটবল, খেলে)

Figure 4: a

If we consider a bangla sentence "আমরা যে দেশে বাস করি তার নাম বাংলাদেশ", the possible bigrams (N-grams with N=2) are: আমরা যে, যে দেশে, দেশে বাস, বাস করি, করি তার, তার নাম, নাম বাংলাদেশ

Bigram probability, $P(\text{আমরা} | \text{যে}) * P(\text{যে} | \text{আমরা}) * P(\text{দেশে} | \text{যে}) * P(\text{বাস} | \text{দেশে}) * P(\text{করি} | \text{বাস}) * P(\text{তার} | \text{করি}) * P(\text{নাম} | \text{তার}) * P(\text{বাংলাদেশ} | \text{নাম})$

and possible trigrams (Ngrams with N=3) are:

আমরা যে দেশে, যে দেশে বাস, দেশে বাস করি, বাস করি তার, করি তার নাম, তার নাম বাংলাদেশ

Trigram probability, $P(\text{আমরা} | \text{যে} | \text{দেশে}) * P(\text{যে} | \text{দেশে} | \text{আমরা}) * P(\text{দেশে} | \text{আমরা} | \text{যে}) * P(\text{বাস} | \text{যে} | \text{দেশে}) * P(\text{করি} | \text{দেশে} | \text{বাস}) * P(\text{তার} | \text{বাস} | \text{করি}) * P(\text{নাম} | \text{করি} | \text{তার}) * P(\text{বাংলাদেশ} | \text{তার} | \text{নাম})$

Figure 5: Figure 1 (Figure 2 :Figure 3 :

1

Results with correct sentences

Models	No. of Sentences	No. of success	Performance
Bigram	1000	900	90%
Trigram	1000	905	90.5%
Quadrigram	1000	907	90.7%
Results with incorrect sentences			
Bigram	1000	950	95%
Trigram	1000	961	96.1%
Quadrigram	1000	963	96.3%
		Average	93.1%

VI.

Figure 6: Table 1 :

112 [Krishnamurthy] *A Demonstration of the Futility of Using Microsoft Word's Spelling and Grammar Check*,
113 Sandeep Krishnamurthy . <http://faculty.washington.edu/sandeep/check>

114 [Naber ()] *A Rule-Based Style and Grammar Checker*, Daniel Naber . 2003. Computer Science -Applied,
115 University of Bielefeld (Diploma Thesis)

116 [Atwell and Elliott ()] *Dealing with illformed English text, The Computational Analysis of English*, Eric Atwell ,
117 Stephen Elliott . 1987. Longman.

118 [Chen ()] 'Linear Networks and Systems (Book style)'. *Natural Language Processing, the PLNLP approach*, W.-K
119 Chen (ed.) (Belmont, CA) 1993. 1993. Wadsworth. p. .

120 [Jurafsky and Martin (1999)] *Speech and Language ProcessingAn Introduction to Natural Language Processing: Computational Linguistics and Speech Recognition*, Daniel Jurafsky , James H Martin . September 28. 1999.
121 Englewood Cliffs, New Jersey 07632: Prentice Hall.

123 [Brown et al. (2010)] 'Stochastic Language Models (N-Gram) Specification'. Michael K Brown , Andreas Kellner
124 , Dave Raggett . <http://www.w3.org/TR/ngram-spec> W3C/Openwave, (Access date) 8th Dec. 2010.

125 [Wikipedia (2010)] Wikipedia . <http://en.wikipedia.org/wiki/N-gram> n-gram, (Access date) 17th Dec.
126 2010.