

# Parikh matrix on the Context-Free Grammar for Natural Languages

Amrita Bhattacharjee<sup>1</sup>, Bipul Syam Purkayastha<sup>2</sup>

<sup>1</sup>Department of Computer Science,  
Assam University Silchar 788011, India,  
[amritabhattacharjee10@gmail.com](mailto:amritabhattacharjee10@gmail.com)

<sup>2</sup>Department of Computer Science,  
Assam University Silchar 788011, India,  
[bipulsyam@gmail.com](mailto:bipulsyam@gmail.com)

**Abstract:** In this paper Parikh Matrices over context-free languages are investigated. Context-free grammars for Natural languages are a developing area of investigation. Parikh matrix is a significant tool of Formal languages. Context-free language is a kind of formal language. Parikh matrix can be used in context-free language. A context-free grammar for Bengali language is also a developing area of investigation. As a case study in this paper Bengali letters, Bengali words and Bengali sentences are studied by using Parikh matrix.

**Keywords:** Parikh matrix, Subword, M- ambiguous words, Formal grammar, Context-free grammar.

## 1. Introduction

Natural languages are those languages which are used by human beings either vocally or in written form in their day to day life for communication. Natural language and formal language are different to each other with respect to their configuration and utility. Much work has been done to establish the interrelation between natural language and formal language. A formal language is often defined by means of a formal grammar such as a regular grammar or context-free grammar. In this present work Parikh matrix, a tool of formal language is used in a Natural language. Parikh matrix is defined on formal language. The Parikh mapping or Parikh vector is an old and important tool in the theory of formal languages introduced by R.J.Parikh [1]. This notion is an important tool in the theory of formal languages. With the help of this tool properties of words can be expressed numerically. For example, for the word  $w = abbaccac$  the Parikh vector is  $(3, 2, 3)$ . In 2001 Mateescu et al. [2] introduced the notion of Parikh matrix. With every word over an ordered alphabet, a Parikh Matrix can be associated and it is a triangular matrix. In recent decades many techniques have been developed to solve complex problems of words using Parikh Matrix. The notion of Parikh matrix is an extension of Parikh Mapping. We cite a few examples [3, 4, 5 ...17] which have used Parikh matrix for solving the problems of word.

An ordered alphabet is a set of symbols  $\Sigma = \{a_1, a_2, \dots, a_n\}$  where the symbols are arranged maintaining a relation of order (" $<$ ") on it. For example if  $\{a_1 < a_2 < \dots < a_n\}$ , then we use notation  $\Sigma = \{a_1, a_2, \dots, a_n\}$ . With every word over an ordered

alphabet, a Parikh matrix can be associated and it is a triangular matrix. All the entries of the main diagonal of this matrix is 1 and every entry below the main diagonal has the value 0 but the entries above the main diagonal provide information on the number of certain sub-words in  $w$ . As for example the tertiary word

$$\xi_1 = abc \underbrace{\dots ca}_{25} \underbrace{\dots ad}_{10} \underbrace{\dots db}_{15} \underbrace{\dots b}_{10} abcd$$

has the Parikh matrix

$$\Psi_{M_4}(\xi_1) = \begin{pmatrix} 1 & 12 & 123 & 148 & 523 \\ 0 & 1 & 12 & 37 & 412 \\ 0 & 0 & 1 & 26 & 401 \\ 0 & 0 & 0 & 1 & 17 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

Natural Language Processing (NLP) is a theoretically inspired variety of computational techniques. It is used for analysing and representing naturally occurring texts at one or more levels of linguistic analysis. Bengali is a natural language. Bengali language is an emerging area of investigation of NLP. Many research works are going on in the field of Bengali language. In this paper study is done on the same field. The fact used in this paper is that a context-free grammar for Bengali language can be generated.

The paper is organized as follows. The following section 2 reviews the related works on Parikh Matrix and Bengali Language Processing. Section 3 goes toward reviewing the basic preliminaries of Parikh Matrix and computational linguistics. Section 4 gives representation of Bengali letters by Parikh matrices; In Section 5, representation of Bengali words by Parikh matrices is introduced; in section 6, representation of Bengali sentences by Parikh matrices are presented. We conclude the paper in Section 7 by summarizing the observations.

## 2. Preliminaries

Throughout this paper  $Z$  will denote the set of natural numbers including zero. First we recall some definitions.

**Ordered alphabet:** An ordered alphabet is a set of symbols  $\Sigma = \{a_1, a_2, a_3, \dots, a_n\}$  where the symbols are arranged maintaining a relation of order (" $<$ ") on it. For example if  $a_1 < a_2 < a_3 < \dots < a_n$ , then we use notation:

$$\Sigma = \{a_1, a_2, a_3, \dots, a_n\}$$

**Word:** A word is a finite or infinite sequence of symbols taken from a finite set called an alphabet. Let  $\Sigma = \{a_1, a_2, a_3, \dots, a_n\}$  be the alphabet. The set of all words over  $\Sigma$  is  $\Sigma^*$ . The empty word is denoted by  $\lambda$ .

$|w|_{a_i}$ : Let  $a_i \in \Sigma = \{a_1, a_2, a_3, \dots, a_n\}$  be a letter. The number of occurrences of  $a_i$  in a word  $w \in \Sigma^*$  is denoted by  $|w|_{a_i}$ .

**Sub-word:** A word  $u$  is a sub-word of a word  $w$ , if there exists words  $x_1 \dots x_n$  and  $y_0 \dots y_n$ , (some of them possibly empty), such that  $u = x_1 \dots x_n$  and  $w = y_0 x_1 y_1 \dots x_n y_n$ . For example if  $w = abaabcac$  is a word over the alphabet  $\Sigma = \{a, b, c\}$  then  $baca$  is a sub-word of  $w$ . Two occurrences of a sub-word are considered different if they differed by at least one position of some letter. In the word  $w = abaabcac$ , the number of occurrences of the word  $baca$  as a sub-word of  $w$  is  $|w|_{baca} = 2$ .

**Parikh vector:** The Parikh vector is a mapping  $\Psi: \Sigma^* \rightarrow Z \times Z \times \dots \times Z$  where  $\Sigma = \{a_1, a_2, a_3, \dots, a_n\}$  and  $Z$  is the set of natural numbers including 0, such that for a word  $w$  in  $\Sigma^*$ ,  $\Psi(w) = (|w|_{a_1}, |w|_{a_2}, |w|_{a_3}, \dots, |w|_{a_n})$  with  $|w|_{a_i}$  denoting the number of occurrences of the letter  $a_i \in w$ . For example, for the word  $w = abaabcac$  the Parikh vector is  $(4, 2, 2)$ .

**Triangle matrix:** A triangle matrix is a square matrix  $m = (m_{ij})_{1 \leq i, j \leq n}$  such that:

1.  $m_{i,j} \in Z$  ( $1 \leq i, j \leq n$ ),
2.  $m_{i,j} = 0$  for all  $1 \leq j < i \leq n$ ,
3.  $m_{i,i} = 1$  ( $1 \leq i \leq n$ ).

**Parikh matrix:** Let  $\Sigma = \{a_1 < a_2 < a_3 < \dots < a_n\}$  be an ordered alphabet, where  $n \geq 1$ . The Parikh matrix mapping, denoted  $\Psi_{M_n}$ , is the homomorphism  $\Psi_{M_n}: \Sigma^* \rightarrow M_{n+1}$  defined as follows:

if  $\Psi_{M_n}(a_q) = (m_{i,j})_{1 \leq i, j \leq n+1}$  then  $m_{i,i} = 1$ ,  $m_{q,q+1} = 1$  and all other elements are zero.

**M-ambiguous or Amiable words:** Two words  $\alpha, \beta \in \Sigma^*$  ( $\alpha \neq \beta$ ) over the same alphabet  $\Sigma$  may have the same Parikh matrix. Then the words are called amiable or M-ambiguous.

The words  $baaabaa$  and  $ababaaa$  has the same Parikh Matrix  $\begin{pmatrix} 1 & 5 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix}$ . So these two words are amiable.

**M-unambiguous words:** A word  $w$  is said to be M-unambiguous if there is no word  $w'$  for

which  $\Psi_{M_n}(w) = \Psi_{M_n}(w')$ .

**Formal grammars:** A formal grammar of this type consists of:

- a finite set of terminal symbols,
- a finite set of nonterminal symbols,
- a finite set of production rules (left-hand side  $\rightarrow$  right-hand side) where each side consists of a sequence of these symbols,
- a start symbol.

**Context-free grammar:** A context-free grammar  $G$  is defined by the 4-tuple:  $G = \{V, \Sigma, R, S\}$  where

1.  $V$  is a finite set; each element  $v \in V$  is called a non-terminal character or a variable. Each variable represents a different type of phrase or clause in the sentence. Variables are also sometimes called syntactic categories. Each variable defines a sub-language of the language defined by  $G$ .
2.  $\Sigma$  is a finite set of terminals, disjoint from  $V$ , which make up the actual content of the sentence. The set of terminals is the alphabet of the language defined by the grammar  $G$ .
3.  $R$  is a finite relation from  $V$  to  $(V \cup \Sigma)^*$ . The members of  $R$  are called the rules or productions of the grammar.
4.  $S$  is the start variable, used to represent the whole sentence (or program). It must be an element of  $V$ .

The asterisk represents the Kleene star operation.

### 3. Related works

Since the introduction of the notion of Parikh vector in 1966 [1] continuous research works are going on in this field. In this introductory paper [1] certain properties of context-free or type 2 grammars are investigated. In particular, questions regarding structure, possible ambiguity and relationship to finite automata are considered. Some important results are also presented. A sharpening of Parikh mapping namely Parikh matrix is introduced in [2] and this matrix representation gives more information than Parikh vector does. With the extension of Parikh matrix an interesting interconnection between mirror images of words and inverses of matrices are investigated in [3]. Researchers [4] have presented ratio property of two words. This property is a sufficient condition for the words  $uv$  &  $vu$  to have the same Parikh matrix. In this paper [5] universal languages for Parikh matrices is introduced and studied. In [6] M-unambiguity is studied both in terms of words and matrices and several criteria for M-unambiguity are provided in both cases. In this paper [7] palindromic amicable words are studied in the context of binary words. Researchers [8] have introduced subword condition. Various characterization and decidability results for languages subword conditions are discussed. In this paper [9] Parikh Matrices over tertiary alphabet are investigated. Algorithm is developed to display Parikh Matrices of words over tertiary alphabet. A set of equations for finding tertiary words from the respective Parikh matrix is introduced. In this paper [10] the notion of a subword history closely related to Parikh matrices is introduced and obtained a sequence of general results. A general inequality





every Bengali letter is a 51x51 matrix. All the entries of the main diagonal of this matrix is 1 and every entry below the main diagonal has the value 0 but the entries above the main diagonal provide information on the concerning Bengali letter . Every word is a matrix product of these matrices. The entries above the main diagonal provide information on the concerning Bengali word and Bengali sentence. It is an effort to use the tool Parikh matrix to Bengali language. With the advance of development of context-free grammar for Bengali language this effort will result more fruitful. Many tools of context-free grammar can be used to Bengali language. Various results of Parikh matrix can also be applied to Bengali language.

## References

- [1] R.J.Parikh, (1966) "On the context-free languages" *Journal of the Association for Computing Machinery*, Vol No.13 pp 570-581.
- [2] A. Mateescu, A. Salomaa, K. Salomaa, S.Yu, (2001) "A sharpening of the Parikh Mapping" *Theoret. Informatics Appl.*, Vol No. 35 pp 551-564.
- [3] A. Mateescu, A. Salomaa, K. Salomaa, S.Yu: On an extension of the Parikh mapping, *T.U.C.S Technical Report No 364*.
- [4] K. G. Subramanian, A. M. Huey, A. K. Nagar, (2009) "On Parikh matrices" *Int. J. Found. Comput. Sci.* Vol.20 No.2 pp211-219.
- [5] C. Ding, A. Salomaa, (2006) "On some problems of Mateescu concerning sub word occurrences" *Fundamenta Informaticae* Vol No. 72 pp1-15.
- [6] V.N. Serb̃anut̃a, (2006) "Injectivity of the Parikh matrix mappings revisited" *Fundamenta Informaticae* Vol No. XX pp1-19, IOS Press.
- [7] A. Atanasiu, C. M.Vide, A. Mateescu, (2001) "On the injectivity of the Parikh matrix mapping" *Fundam. Informa.* Vol No.46 pp 1-11.
- [8] A. Salomaa et al, (2006) "Subword conditions and subword histories" *Information and Computation* Vol No.204 pp1741-1755.
- [9] Amrita Bhattacharjee, Bipul Syam Purkayastha, (2014) "Parikh Matrices and Words over Tertiary Ordered Alphabet" *Int. J. of Computer Applications*, Vol 85 No.4 pp10-15.
- [10] A. Mateescu, A. Salomaa, S. Yu, (2004) "Subword histories and Parikh matrices" *J. Comput. Syst. Sci.*, Vol No.68 pp 1-21.
- [11] T.-F. S, erb̃anut̃a, (2004) "Extending Parikh matrices" *Theoretical Computer Science*, Vol No.310 pp 233-246.
- [12] S. Fosse, G. Richmomme, (2004) "Some characterisations of Parikh matrix equivalent binary words" *Information Processing Letters*. Vol.92 No.2 pp77-82.
- [13] A. Salomaa, (2005) "Connections between subwords and certain matrix mappings" *Theoretical Computer Science*, Vol No.340 pp 188-203.
- [14] A. Salomaa, (2006) "Independence of certain quantities indicating subword occurrences" *Theoretical Computer Science*. Vol.362 No.1 pp222-231.
- [15] Adrian Atanasiu, Radu Atanasiu, Ion Petre, (2008) "Parikh matrices and amiable words" *Theoretical Computer Science Vol No.390* pp 102-109.
- [16] A. Atanasiu, (2007) "Binary amiable words" *Int. J. Found. Comput. Sci* Vol.18 No.2 pp387-400.
- [17] Amrita Bhattacharjee, Bipul Syam Purkayastha, Application of ratio property in searching of M-ambiguous words and its generalization, 3rd International Conference on Soft Computing for Problem Solving ,AISC (SocProS 2013) 258, pp857-865. December 27-29, Greater Noida Extension Centre of IIT Roorkee, India.
- [18] J. E. Hopcroft, Rajeev Motwani and J. D. Ullman, "Introduction to Automata Theory Languages and Computation," Pearson Education Publishing Company, Second Edition (book style)
- [19] Apurbalal Senapati, Utpal Garain "Bangla Morphological Analyzer using Finite Automata" ISI @FIRE MET 2012
- [20] K. M. Azharul Hasan, Al-Mahmud, Amit Mondal, Amit Saha "RECOGNIZING BANGLA GRAMMAR USING PREDICTIVE PARSER" IJCSIT Vol 3 December 2011.
- [21] Muhammad Nasimul Haque and Mumit Khan "Parsing Bangla using LFG: An Introduction" (IJCSIT) Vol 3, No 6, Dec 2011
- [22] M. A. Karim (Old Dominion University, USA), M. Kaykobad (Bangladesh University of Engineering & Technology, Bangladesh) and M. Murshed (Monash University, Australia) chapter 7 "Parsing Bangla Grammar Using Context-free Grammar".
- [23] M.S. Hasan ,Amit Mondal, Amit Saha "A context-free grammar and its predicative parser for Bangla grammar recognition" IEEE Xplore , computer and information technology,(ICCIT), 2010, 13<sup>th</sup> International conference.
- [24] Amrita Bhattacharjee, Bipul Syam Purkayastha, "Some Alternative Ways To Find M-Ambiguous Binary Words Corresponding To A Parikh Matrix", International Journal of Computational Science, AIRCC, 2014, 4(1), 53-64.
- [25] Amrita Bhattacharjee, Bipul Syam Purkayastha, Parikh Matrices and words over ternary alphabet, 4th International Conference on Soft Computing for Problem Solving, AISC (SocProS 2014) 335,135-145. Springer India, December 26-28, National Institute of Technology, Silchar, India.
- [26] **Author Profile**  
Amrita Bhattacharjee received M. Sc. Degree in Pure Mathematics from Gauhati University and at present a Ph.D. Research scholar of Assam University.