

## Sanskrit lemmatizer for improvisation of morphological analyzer

Jaideepsinh K. Raulji & Jatinderkumar R Saini

To cite this article: Jaideepsinh K. Raulji & Jatinderkumar R Saini (2019) Sanskrit lemmatizer for improvisation of morphological analyzer, Journal of Statistics and Management Systems, 22:4, 613-625, DOI: [10.1080/09720510.2019.1609186](https://doi.org/10.1080/09720510.2019.1609186)

To link to this article: <https://doi.org/10.1080/09720510.2019.1609186>



Published online: 25 Jun 2019.



Submit your article to this journal [↗](#)



Article views: 6



View related articles [↗](#)



View Crossmark data [↗](#)



## **Sanskrit lemmatizer for improvisation of morphological analyzer**

**Jaideepsinh K. Raulji \***

*School of Computer Studies*

*Ahmedabad University*

*Ahmedabad 380009*

*Gujarat*

*India*

*and*

*Dr. Babasaheb Ambedkar Open University*

*Ahmedabad 382481*

*Gujarat*

*India*

**Jatinderkumar R Saini**

*Department of Computer Science*

*Narmada College of Computer Application*

*Bharuch 392011*

*Gujarat*

*India*

---

### **Abstract**

The process of stripping off affixes from a word to arrive at root word or lemma is known as Lemmatization. The usefulness of lemmatizer in natural language operations cannot be overlooked especially if the language is rich in its morphology. A lexicon cum rule based lemmatizer is built for Sanskrit Language. The Lemmatizer has profound applications in NLP main stream tasks like Information Retrieval, Morphological Analyzer ,POS taggers, Question-Answering Systems, Machine Translation Systems etc. Here the Sanskrit rule based morphological analyzer augmented with lemmatizer proved, analyzing Sanskrit word more precise and accurate to its morphological characteristics.

---

**Keywords:** *Information Retrieval (IR), Natural Language Processing (NLP), Sanskrit, Inflection, Tokenization, Morphology, Lemma, Case Markers (Vibhaktis)*

---

\*E-mail: [jaideepraulji@gmail.com](mailto:jaideepraulji@gmail.com)



## 1. Introduction

Lemmatizers are costly in terms of time and space complexity as opposed to stemmers. But its lucidity in results and its performance obviously outweighs accuracy of stemmers. As compared to other Indo-Iranian family of languages, Sanskrit is comparatively structured and morphologically rich language. Day by day the importance of Sanskrit language in domain of linguistic research and traditional usage is increasing. Sanskrit is one of the 22 Scheduled languages listed in Constitution of India and it is official language of state of Uttarakhand, India. Hence realizing the importance and widespread roots imbibed in works carried out in computational linguistics domain, developing Sanskrit lemmatizer would be helpful for related tasks. Here to increase efficiency and accuracy of rule based Sanskrit morphological analyzer, the results of lemmatizer is injected before the mainstream logic of morphological analyzer. Here morphological analyzer is responsible for POS tagging of Sanskrit words in core grammatical category like indeclinable, pronoun, noun and verbs. In Sanskrit, nouns and adjective almost shares similar inflectional morphological rules, hence adjectives are not segregated. Sanskrit being free word order language, the word itself reflects complete grammatical function as syntactic-semantic relation within sentence. Morphological analysis provides ample information of word with its syntactic and semantic role played in a sentence. Here initially after retrieving lemma, lemmatizer also provides surface level grammatical POS information to morphological analyzer as it is concurrently attached with lexical resource. There after grammatical information like gender, number, person, tense, etc is marked through the inflectional affix rules [6].

## 2. Related Work

A self learning context aware lemmatizer for German language is built which is lexicon and noun based. Lexicon based approaches are quite costly in terms of maintenance but accurate. An important feature of the system is, it automatically updates lexicon from processed documents which leads to better accuracy and widened scope[1]. A memory based Lemmatizer for Greek language is built which uses lexicon and Frog framework for POS tagging. Frog is open source NLP framework suitable to be trained for all the languages. The word is first POS tagged with Frog's NLP framework, then looked up in lexicon for lemma[2]. The Bengali lemmatizer "BenLem" is developed and evaluated especially

for WSD. It handles inflectional and derivational Bengali morphology. BenLem's architecture has two main pillars - A valid suffix list of the language and a dictionary [3]. The rule base Lemmatizer for Kannada is built by Pratibha R J, et al. The designed lemmatizer uses affix rule set and root-dictionary approach with approximately 94% efficiency [4]. Similarly noun lemmatization technique for Bangla is developed by Alok P, et al. The TDIL corpus of Bangla language is used to tag noun words using POS tagger. Using longest suffix stripping methodology in decreasing order through a suffix list it extracted a Bangla lemma [5].

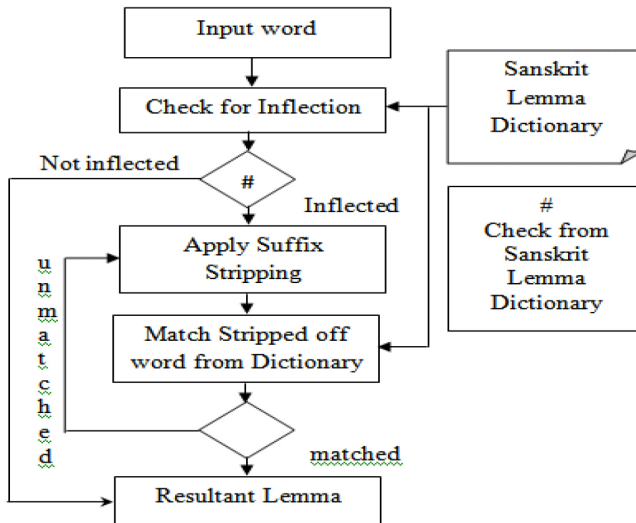
Syntax dependency parser is developed by Pawan G, et al to analyzing Sanskrit sentences[7]. Sanskrit parser algorithm based on rich case endings is developed by Shashank S and Raghav A. Using paradigm table and DFA, the root word along with its attributes are retrieved by checking against avyaya, pronoun, verb and noun tree sequentially [8]. A tool to to analyze inflection forms of Sanskrit words by FST is built by Amba K and Devanand S [9]. A rule based POS tagger is built by Namrata T and Suresh J [10] where rules are stored in the database and the word was compared to database after suffix stripping. Akshar B, et al built morphological analyzer using modular approach of programming paradigm and included modules for Sandhi - Samasa analyzer and formation, Subanta, Tinanta and Kridanta Analyzer [11].

### 3. Approach and Implementation

The availability of digital linguistic resources for the language is scarce. Machine Learning techniques sucks enormous amount of clean and structured corpus for accurate results. Hence a rule based approach is implemented to build the same. The generated results definitely forms basis for understanding surface structure of sentence and can be utilized for further improvement of related systems like Information Retrieval, Part of Speech taggers, Machine Translation etc.

#### 3.1 The Lemmatizer

The Lemmatizer accepts Sanskrit word in Devanagari UTF-8 Unicode encodings. The word is matched repetitively with Sanskrit lexicon containing approximately 28500 words, by stripping off suffixes until the correct lemma and first level POS information is achieved. The algorithm is depicted in form of flowchart through figure 1. The collection of word lemma are gathered from several printed dictionaries, but the vital resource was *Bruhat kosh* [16].



**Figure 1**  
**The Lemmatizer Flowchart**

### 3.2 The Morphological Analyzer

The surface level POS information, lemma and inflection information achieved from lemmatizer is feed into morphological analyzer. The lemma without noun and verb is matched with indeclinable and pronoun datastore. The pronoun datastore contains all the grammatical information like Gender (Masculine, Feminine, Neuter), Number (Singular, Dual, Plural) and Case/Vibhakti (Nominative, Accusative, Instrumental, Dative, Ablative, Genitive, Locative). The lemma with noun and verb POS information is matched with verb and noun affix rules defined in Table 1 and Table 2 respectively. The algorithmic flowchart of morphological analyzer is mentioned in figure 2.

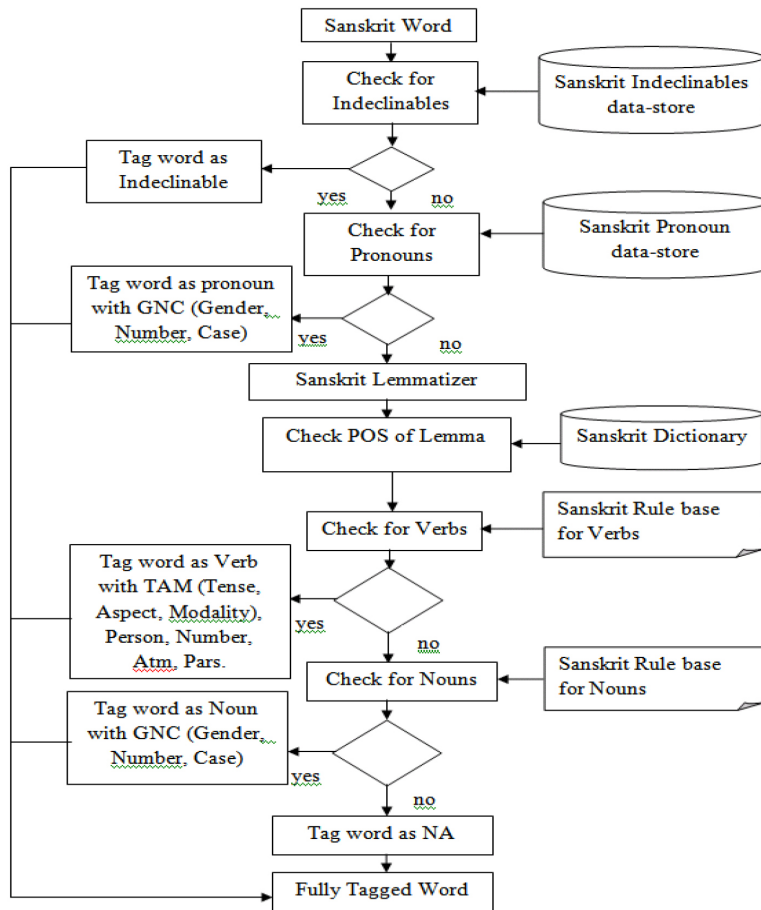


Figure 2  
The Morphological Analyzer Flowchart

Table 1

## Verbal Inflections for Sanskrit [13][14][15]

[FP-First Person, SP-Second Person, TP-Third Person, S-Singular, D-Dual, P-Plural]

Present Tense [Parasmaipada]			
	S	D	P
FP	मि	वः	मः
SP	सि/षि	थः	थ
TP	ति	तः	न्ति

PresentTense [Atmanepada]			
	S	D	P
FP	े	वहे	महे
SP	से	ेथे	ध्वे
TP	ते	ेते, ते	न्ते, ते

Imperfect (Past) [Parasmaipada]			
	S	D	P
FP	अ..म्	अ..व	अ..म
SP	अ..ः	अ..तम्	अ..त
TP	अ..त्	अ..ताम्	अ..न्

Imperfect (Past) [Atmanepada]			
	S	D	P
FP	अ..े, अ..ि	अ..वहि	अ..महि
SP	अ..थाः	अ..थाम्	अ..ध्वम्
TP	अ..त	अ..ताम्	अ..न्त, अ..त

Imperative Mood [Parasmaipada]			
	S	D	P
FP	ानि	ाव	ाम
SP	-	तम्	त
TP	तु	ताम्	न्तु

Imperative Mood [Atmanepada]			
	S	D	P
FP	ै	ावहै	ामहै
SP	स्व	थाम्	ध्वम्
TP	ताम्	ेताम्, ाताम्	न्ताम्, ताम्

Potential Mood [Parasmaipada]			
	S	D	P
FP	ेयम्, याम्	ेव, याव	ेम, याम
SP	ेः, थाः	ेतम्, यातम्	ेत, यात
TP	ेत, यात्	ेताम्, याताम्	ेयुः

Potential Mood [Atmanepada]			
	S	D	P
FP	ेय	ेवहि	ेमहि
SP	ेथाः	ेथाम्	ेध्वम्
TP	ेत	ेयाताम्	ेरन्

First / Periphrastic Future [Parasmaipada]			
	S	D	P
FP	तास्मि	तास्वः	तास्मः
SP	तासि	तास्थः	तास्थ
TP	ता	तारौ	तारः

First / Periphrastic Future [Atmanepada]			
	S	D	P
FP	ताहे	तास्वहे	तास्महे
SP	तासे	तासाथे	ताध्वे
TP	ता	तारौ	तारः

Second Future [Parasmaipada]			
	S	D	P
FP	ष्यामि	ष्यावः	ष्यामः
SP	ष्यसि	ष्यथः	ष्यथ
TP	ष्यति	ष्यतः	ष्यन्ति
ष्य can also have forms like क्ष्य / स्य			

Second Future [Atmanepada]			
	S	D	P
FP	ष्ये	ष्यावहे	ष्यामहे
SP	ष्यसे	ष्येथे	ष्यध्वे
TP	ष्यते	ष्येते	ष्यन्ते
ष्य can also have forms like क्ष्य / स्य			

Conditional Mood [Parasmaipada]			
	S	D	P
FP	अ..ष्यम्	अ..ष्याव	अ..ष्याम
SP	अ..ष्यः	अ..ष्यतम्	अ..ष्यत
TP	अ..ष्यत्	अ..ष्यताम्	अ..ष्यन्
ष्य can also have forms like क्ष्य / स्य			

Conditional Mood [Atmanepada]			
	S	D	P
FP	अ..ष्ये	अ..ष्यावहि	अ..ष्यामहि
SP	अ..ष्यथाः	अ..ष्येथाम्	अ..ष्यध्वम्
TP	अ..ष्यत	अ..ष्येताम्	अ..ष्यन्त
ष्य can also have forms like क्ष्य / स्य			

Perfect (Past) [Parasmaipada]			
	S	D	P
FP	-	व	म
SP	थ	थुः	-
TP	-	तुः	तुः

Perfect (Past) [Atmanepada]			
	S	D	P
FP	े	वहे	महे
SP	षे	थे	ध्वे
TP	े	ते	रे

Aorist (Past)[Parasmaipada]			
	S	D	P
FP	अ..म्	अ..व	अ..म
SP	अ..ः	अ..तम्	अ..त
TP	अ..त्	अ..ताम्	अ..न्

Aorist (Past)[Atmanepada]			
	S	D	P
FP	अ..म्	अ..वहि	अ..महि
SP	अ..ः	अ..थाम्	अ..ध्वम्
TP	अ..त्	अ..ताम्	अ..न्

Benedictive [Parasmaipada]			
	S	D	P
FP	यासम्	यास्व	यास्म
SP	याः	यास्तम्	यासुत
TP	यात्	यास्ताम्	यासुः

Benedictive [Atmanepada]			
	S	D	P
FP	षीय	षीवहि	षीमहि
SP	षीष्ठाः	षीयास्थाम्	षीध्वम्
TP	षीष्ट	षीयास्ताम्	षीरन्
षी can also have forms like सी			



Table 2  
Nominal Inflections for Sanskrit [13][14][15]

Karaka [Case]	Masculine			Feminine			Neuter		
	Singular	Dual	Plural	Singular	Dual	Plural	Singular	Dual	Plural
Nominative [प्रथमा] [Subject]	ः, ा, ूः, ौः, ान्	ौ, ी, ू ः	ः, नः, ः, यः, वः	ः, ूः, ीः, ौः, ः	े, यौ, ौ, ू ः	ः, यः, ठः, वः	म्, ि, ठः, ु	े, ी, िणी, णी, नी	ानि, ि, शीणि, नि, न्ति, ि
Accusative [द्वितीया] [Object]	म्	ौ, ी, ू ः	ान्, न्, ीन्, ून्, ः	म्	े, यौ, ू, ौ ः	ः, ूः	म्, ि, ु	े, िणी, नी	ानि, ि, ीणि, नि, न्ति, ि
Instrumental [तृतीया] [by, with]	ेन, ेण, ा, ना,	ान्याम्, भ्याम्, याम्	ैः, भ्यः, भिः	या, याम्, ता, ा	भ्याम्, याम्	भिः, िः	ैन, णा, ना	भ्याम्, भ्याम्	ैः, भिः
Dative [चतुर्थी] [for]	ाय, े, ये	ान्याम्, भ्याम्, याम्	ेभ्यः, भ्यः, यः	ये, ये, वै, ते, े	भ्याम्, याम्	भ्यः	ाय, णे, ने	भ्याम्, भ्याम्	भ्यः
Ablative [पञ्चमी] [from]	ात्, ः, ैः, ौः	ान्याम्, भ्याम्, याम्	ेभ्यः, भ्यः, यः	या, ैः, ौः, वा, नोः	भ्याम्, भ्याम्	भ्यः	ात्, णि, नः	भ्याम्	ेभ्यः, भ्यः
Genitive [षष्ठी] [shows possession]	स्य, ः, ैः, ौः, नः	यो, ोः, वो, नोः	ानाम्, ाणाम्, णाम्, नाम्, ीनाम्, ाम्	या, वाः	यो, वोः	नाम्, ाम्	स्य, णः, नः	यो, णोः, नोः	ानाम्, नाम्, णाम्
Locative [सप्तमी] [in, on]	े, ि, ौ, नि, व	यो, ोः, वो, नोः	ेषु, षु, क्षु, सु	याम्, ाम्, ि	यो, वोः, ि	सु, षु, क्षु	े, नि, णि	यो, णोः, नोः	षु, क्षु

**Table 3**  
**Lemmatizer Evaluation Results**

Total no. of Words inputted in Lemmatizer	No. of words with incorrect lemma or missed out lemma	No. of words correct lemma / root
8240	737	7503
Accuracy = 91.06 %		

#### 4. Results

Sanskrit text from various internet resources were collected, also text was digitized from Sanskrit textbook[15]. The lemmatizer and morphological analyzer is tested on 8240 words. The results obtained from the system were verified through human linguist expert and renowned portal “spokensanskrit.org”. The results of Lemmatizer is described below in Table 3.

Similarly morphological analyzer was tested with same set of words and the results are described in Table 4,5,6 and 7. The accuracy in indeclinables and pronoun category is absolute due to direct lexicon matching, the incorrectness observed in this category is mostly due to incorrect spelling, and typo error. The untagged word count is 389 is mainly due to typing and incorrect spelling errors.

#### Formula 1.

Accuracy of Individual POS tags =

$$\frac{\text{No. of correctly tagged} \ll \text{POS category} \gg}{\text{Total no. of} \ll \text{POS category} \gg \text{ in corpus}} \times 100$$

<<POS category>> values are Indeclinables, Pronouns, Verbs and Nouns.

**Table 4**  
**Evaluating Morphological Analyzer**

Total Words in Corpora	8240
No. of Sentences	1451
Total Identified (correctly tagged) Words	7261
Incorrectly tagged words + Unidentified words	979 (590+389)
Total Incorrectly tagged Words	590
Completely Unidentified Words (untagged words)	389

**Table 5**  
**POS category wise Evaluation Results**

POS Category	Status	No. of Words	No. of Words	Individual POS category Accuracy %
Indeclinables	Correctly tagged	325	333	97.60
	Incorrectly tagged	8		
Pronouns	Correctly tagged	657	674	97.48
	Incorrectly tagged	17		
Verbs	Correctly tagged	1842	1984	92.84
	Incorrectly tagged	142		
Nouns (Subanta)	Correctly tagged	4437	4860	91.30
	Incorrectly tagged	423		
Completely Unidentified Words (untagged words)		389		
Total Words		8240		

**Table 6**  
**Accuracy of tagged words**

Total no. of Words inputted in Morphological Analyzer	Total no. of tagged words	No. of words with invalid (590) + missed POS tag (389)	No. of words with valid POS tag
8240	7851	979	7261
$\text{Accuracy} = \frac{\text{No. of words with valid POS tag (7261)}}{\text{Total no. of tagged words (7851)}} \times 100$			
Accuracy of Morphological Analyzer = 92.49 %			

**Table 7**  
**Overall accuracy of Morphological Analyzer**

Total no. of Words inputted in Morphological Analyzer	No. of words with invalid and missed POS tag	No. of words with valid POS tag
8240	979	7261
$\text{Accuracy} = \frac{\text{No. of words with valid POS tag (7261)}}{\text{Total no. of Words in Corpus (8240)}} \times 100$		
Overall Accuracy of Morphological Analyzer = 88.11 %		

The sample output of the implemented algorithm for morphological analyzer and lemmatizer is depicted in table-8 and table-9 respectively.

**Table 8**  
**Sample output of Morphological Analyzer**

Sr. No	Input	Output
1	शशकः	#शशकः, Noun,Nom,Mas,Sin
2	वृक्षस्य	#वृक्षस्य, Noun,Gen,Neu,Sin
3	अपतत	#अपतत, Verb,1,Im_Past,Pars,SP,Plu
4	वानरः	#वानरः, Noun,Nom,Mas,Sin
5	प्रणम्य	#प्रणम्य, Verb,1,Imperative,Pars,SP,Sin
6	प्राप्तः	#प्राप्तः, Verb,1,Present,Pars,TP,Dua
7	इति	#इति, Indec
8	निवसति	#निवसति, Verb,1,Present,Pars,TP,Sin
9	उत्पठिति	#उत्पठिति, Verb,1,Present,Pars,TP,Sin
10	अहं	#अहं, Pro,Nom,-,Sin

[Nom-Nominative, Gen-Genitive, Mas-Masculine, Neu-Neuter, Sin-Singular, Im-Imperfect, Pars-Parasmaipada, Atm-Aatmanepada, FP-First Person, SP-Second Person, TP-Third Person]

**Table 9**  
**Sample output of Sanskrit Lemmatizer**

Sr No	Inflected Sanskrit Form with ITRANS	Sanskrit Lemma
1	शशकः (shashakaH)	शशकः (shashaka)
2	वृक्षस्य (vRRikShasya)	वृक्षस्य (vRRikSha)
3	अपतत (apatata)	अपतत (apatata)
4	वानरः (vAnaraH)	वानरः (vAnara)
5	प्रणम्य (praNamya)	प्रणम्य (praNam)
6	प्राप्त (prAptaH)	प्राप्त (prApta)
7	निवसति (nivasati)	निवसति (nivasa)
8	उपतिष्ठति (upatiShThati)	उपतिष्ठति (upatiShTha)
9	शृणोतु (shRRiNotu)	शृणोतु (shRRiNa)
10	तिष्ठति (tiShThati)	तिष्ठति (tiShTha)

## 5. Conclusion

Working with the most ancient and structured language like Sanskrit is always challenging and invigorating. The performance of lemmatizer is the result of rich lexicon whose results again transformed morphological analyzer with acceptable accuracy. The outcome will definitely form the basis for other NLP related activities for the language. The morphological analyzer performance drop backs a bit when compound words are encountered hence Sandhi splitter would indeed balance the tradeoffs.

## References

- [1] Perera Praharshana and Witte Rene, "A Self-Learning Context-Aware Lemmatizer for German", Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP), pg 636-643, Oct 2005.
- [2] Corien Bary, et al, "A Memory-Based Lemmatizer for Ancient Greek", In Proceedings of DATeCH2017, Gottingen, Germany, June 1-2, 2017.
- [3] Chakrabarty Abhisek, Garain Utpal, "BenLem (A Bengali Lemmatizer) and Its Role in WSD", ACM Transaction, Asian Low-Resour. Lang. Inf. Process Vol. 15, No. 3, Article 12, Feb 2016.

- [4] Prathibha R J and Padma M C, "Design of Rule based Lemmatizer for Kannada Inflectional Words", International Conference on Emerging Research in Electronics, Computer Science and Technology, 2015.
- [5] Alok R P, Niladri S D, et al, "An Innovative Lemmatization Technique for Bangla Nouns by using Longest Suffix Stripping Methodology in Decreasing Order", International Conference on Computing and Network Communications, Dec 2015.
- [6] Raulji Jaideepsinh, Jatinder kumar saini, "Morphological Analyzer for Sanskrit Language", International Conference on "The Journey of Indian Languages : Perspective on Culture and Society" , 2017.
- [7] Goyal Pawan, Gerard Huet, Kulkarni Amba, Peter Scharf and Ralph Bunker, " A Distributed Platform for Sanskrit Processing", Proceedings of COLING 2012 : pp 1011-1028,
- [8] Saxena Shashank and Agrawal Raghav, "Sanskrit as a Programming Language and Natural Language Processing", *Global Journal of Management and Business Studies*, Vol 3, No. 10 , pp 1135-1142, 2013.
- [9] Kulkarni Amba, and Shukl Devanand, "Sanskrit morphological analyser: Some issues." *Indian Linguistics* 70.1-4 169-177, 2009.
- [10] Tapaswi Namrata and Jain Suresh, "Treebank Based Deep Grammar Acquisition and Part of Speech Tagging for Sanskrit Sentences", 6th International Conference on Software Engineering, Sept 2012.
- [11] Akshar Bharati, Kulkarni Amba, Sheeba V, "Building a Wide Coverage Sanskrit Morphological Analyzer : A Practical Approach", The First National Symposium on Modelling and Shallow Parsing of Indian Languages, 2006.
- [12] Sampada S. Wazalwar & Urmila Shrawankar (2017), "Interpretation of sign language into English using NLP techniques, *Journal of Information and Optimization Sciences*, 38:6, 895-910.
- [13] Kale M R, "A Higher Sanskrit Grammar", Motilal Banarasidas Publishers pvt ltd, 11th Reprint Delhi, 2016.
- [14] Patel Purushottam, "Sanskrit Vyakran Prabodh", Saraswati Pustak Bhandar, Ahmedabad, Edition 2012-13.
- [15] Bhatt Vasantkumar, "Sanskrit Vakyasanrachana", Saraswati Pustak Bhandar, Ahmedabad, Edition 2014.
- [16] Naik Ratilal, "Bruhat Kosh", Akshara Prakashan, Ahmedabad, Edition 2012.