PARTS-OF-SPEECH TAGGING FOR KANNADA

Vijayalaxmi F. Patil, Shahid M. Bhat [jeeveshwarip@gmail.com] [shahid.bhat3@gmail.com] LDC-IL, CIIL Mysore

CONTENTS

- **×** Introduction
- × Kannada & Available Language Resources
- Corpus Used In This Work
 - a. Corpus Cleaning
 - b. Corpus normalization
- Tag-set Used In This Work
- × Kannada POS Tagging
- POS Tagging Issues
- Conclusion
- **×** References

INTRODUCTION

- Kannada is the official language of the state Karnataka. Kannada is one of the Dravidian languages with SOV word order.
- It is very important language as it is not only one of the 22 scheduled languages of India but also one of the classical languages.
- It is spoken in Karnataka and its neighboring states like Maharashtra, Tamil Nadu, Andhra, Goa, etc by about 35 million speakers (Wikipedia).

- It is morphologically rich and agglutinative in its nature.
- It shares many morphological features with other Dravidian languages like defective verbs, like alla not and illa (no), some particles like inclusive particle kUDa (also) and some auxiliaries like –koLL (reflexive), paDu (passive) etc which are considered to be one of the type of an auxiliaries.

- Parts-of-Speech tagging refers to the process of assigning a POS tag to the words of a text.
- In other words we can say that it is a process of marking the words in corpus which corresponds to particular parts of speech based on both its definition and the context.
- ▶POS tagging is one of the important level and the ground work for other higher level stages in NLP.

KANNADA & AVAILABLE LANGUAGE RESOURCES

- Like many of the Indian languages, very little work has been done in the area of NLP for Kannada. It is a resource-poor language. Even if resources exist somewhere, they exist without public access (Murthy 2000).
- CIIL has developed a corpus of about 3 million words for Kannada under a project funded by Department of Information Technology (DIT). Further,
- ▶ POS tagger and morphological analyzer have been developed for Kannada under ILMT consortium project. From last few years LDCIL is engaged in creating language resources for Kannada on large scale.

CORPUS USED IN THIS WORK

In the current work we are concerned with POS tagging of text corpus. Text Corpus is a machine readable collection of the text which is generally used as a raw data for various NLP.

➤ We have used 10,000 words of Kannada corpus from a single domain (Aesthetics).

CORPUS USED IN THIS WORK

Category Aesthetics

Literature-Short Stories

Number of words

5654

Aesthetics

Literature-Children's Literature

780

Aesthetics

Literature- Autobiographies

856

Aesthetics

Literature-Essays

6572

Aesthetics

Literature-Biographies

2407

Kannada corpus is not directly used for POS tagging because of various problems that need to be settled before actual tagging.

Whatever we do with corpus to make it fit for tagging is generally called preprocessing. It involves the following two subtasks.

a. Corpus Cleaning

- Corpus usually contains some extra symbols, Sanskrit shlokas and some stanzas of poems, we have removed such elements.
- words and sentences, removed some extra words, sentences and paragraphs according to the text available in the hard copies of the corpus.
- We remain faithful to the text and keep some spelling variations as such which would be considered wrong spellings otherwise.

b. Corpus normalization

- Normalization is sort of tokenization. Since Kannada is highly agglutinative language (with severe fusion of grammatical categories), we need to tokenize corpus so that we can assign POS tags easily.
- In corpus normalization, we tokenize corpus properly by separating punctuations from preceding tokens and by splitting sentences or phrases into their constituent tokens.

For Example:

We segment "hELikoLLuttiruttidda" (had-been-speaking-himself) into 'hELi' (having spoken), 'koLLutta' (himself), 'irutta' (been), and 'idda' (had).

Similarly;

- NOUN: mAtinallIga (now-in-speech) = mAtin-alli (speech-in) + Iga (now)
- PRONOUN: adariMdEnu (with it-what) = adariMda (with it) + Enu (what)
- PRONOUN: nimagArigU (to-you-anyone) = nimage (to-you) + yArigU (anyone)

TAG SET USED IN THIS WORK

In order to assign a tag to a token we must have a tag set according to which we will assign tag to a token. In this work we are using BIS Dravidian tag set. Which has a 11 categories and 35 sub categories. The tag set is summarizes below with examples.

- NOUN: In Kannada we have 3 types of noun namely Common (NN), Proper (NNP) and NIoc (NST).
- Common noun includes words like mara (tree), giDa (plant), manuSya (human) etc. Proper noun includes words like name of places, Institutions, companies like Mysore, CIIL, Bajaj etc. NLOC includes the locations, adverb what we call traditionally adverb of time and place for example mEle (on), keLage(down), Iga (now)
- PRONOUN: In Kannada we have 4 type of pronoun namely Personal (PRP), Reflexive (PRF), Reciprocal (PRC), and Wh-word (PRQ)
- Pronominal includes nAnu (I), nInu (you), avanu (he) etc, Reflexive includes words like tAnu (self), swataha (self)) etc, Reciprocal includes paraspara (each other), obbarigobbaru (each other) etc and Wh-word includes words like yAru (who), Enu (what) etc

- DEMONSTRATIVE: In Kannada we have two types of Demonstratives like Deicic and Wh-word. Deictic includes A (that) and I (this), Wh-word includes yAva (which).
- VERB: In verb we have 2 types of verb like main and auxiliary, main verb has 5 types namely Finite, Non-finite, Infinitive, Gerund and finally participle noun. Auxiliary verb has 3 types namely Finite, non-finite and Infinitive.
- * Finite main verb includes words like baMdanu (has come), hOdanuhas gone) etc, Non-finite forms includes baMdu (having come), hOgi(having gone) etc, Infinitive includes baralu (to come), hOgalu (to go) etc, Gerunds includes baruvudu (coming), hOguvudu(going) etc, Perticiple noun includes baruvavanu (coming person), hOguvavanu (going person) etc. Auxiliary verbs in Dravidian languages occur with the main verb usually. For example bahudu (may), bEku (should) etc.

- * ADJECTIVE: Adjectives has no sub-type and it includes words like suMdaravAda (beautiful), kliSTha (difficult) etc.
- ADVERB: it includes nidhAnavAgi (slowely), jOrAgi (fastly) etc.
- POSTPOSITIONS: it includes locations like mEle (on), keLage (down), hinde (back), muMde (front) etc.

- CONJUNCTIONS: this is divided into 3 three namely coordinator, subordinator and quotative etc.
- Co-ordinator includes words like mattu (and), hAgU (and) etc. Subordinator includes words like AddariMda (therefore), hAgAgi (therefore) etc and quotatives are eMdu (that), anta (that) etc.
- PARTICLES: three sub categories in this section and they are Default, Interjection and Intensifier.
- Default includes kUda (also) etc. Interjections like ayyO, oh etc and Intensifier tuMba (very), bahaLa (many) etc.

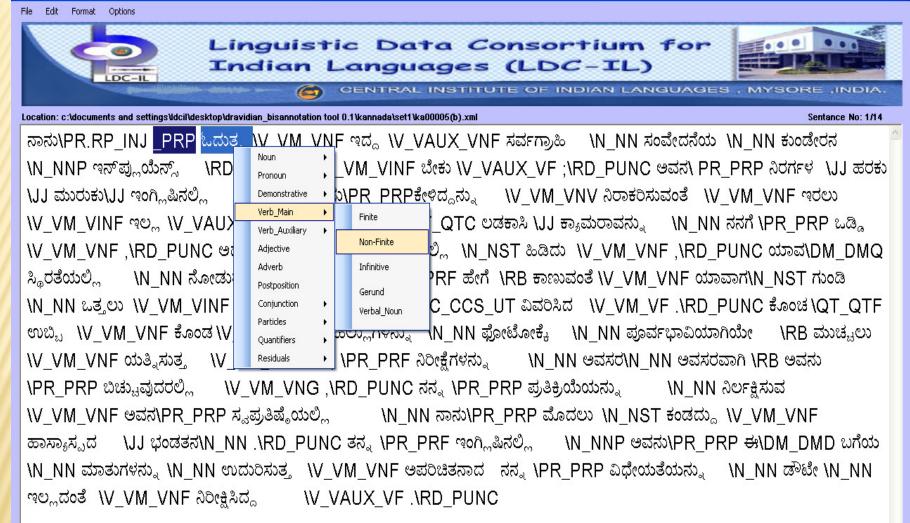
- QUANTIFIERS: we have 3 types in this namely General, Cardinal and Ordinal.
 - General includes ella (all), bahaLa (many) etc, Cardinal includes oMdu(one), eraDu(two) etc and ordinals includes oMdaneya (first), eraDaneya (second) etc.
- RESIDUALS: it includes Foreign, Symbol, Punctuation, Unknown and Echowords.
 - Foreign words usually includes book etc, symbol includes @,& etc. Punctuations like ?, < etc. Unknown includes —— and finally Echowards includes mane gine (houses and other buildings), huli gili (tiger and others) etc.

KANNADA POS TAGGING

LDC-IL has developed annotation tool for POS tagging. It is a customizable manual tool that can be used to implement any tag set.

We have used this customized tool for implementing BIS Dravidian tag set for Kannada. In this work, we have used this tool for tagging the above mentioned preprocessed corpus.





		Sent	/Para No. 1 🔻	< <pre><<pre></pre></pre>	Next>> <	<prev next<="" th="" untag=""><th>t UnTag>> Join Sente</th><th>nce</th><th></th><th></th></prev>	t UnTag>> Join Sente	nce		
🐉 start	Coins of Punj	🇀 AICL 2011 C	Dravidian_BI	Kashmiri Tre	Building a La	Parts-of-spe	. POS SEMINAR	✓ Dravidian-(K	EN 🖁 🖑	4 65

Sentence/Paragraph Selection

POS TAGGING ISSUES

1. Kannada has adverbial suffix which is responsible to make any word into adverb. For example: hasanAgi (cleanly), sukhavAgi (happily), nishcitavAgi (surely), AtmlyavAgi (closely), but there are other cases in Kannada where this Agi can come with the Proper noun like- rudranannu shivanannAgi kANuva kathegaLive (there are some stories where rudra is seen as shivA). If we tag it as adverb the important information like proper noun will be missed out in POS tagging.

2. baMdiruvavanannu (the person who has come +accusative case), baMdu + iruvavanannu here baMdu is main participle noun but iruvavanannu is an auxiliary with png and case marker. We are considering it as auxiliary. When we split this word into baMdu + iruvavanannu, this iru with png and case marker should be an auxiliary participle noun but this auxiliary participle noun is not mentioned in the Dravidian BIS tag set. Similarly the word like baMdiruvudannu we split this word into two tokens to give auxiliary information i.e iruvudannu this token has auxiliary verb with gerundial marker and case marker which is an auxiliary gerund but auxiliary gerund is not there in the Dravidian BIS tag set.

CONCLUSION

In this work we have summarized our experience of POS tagging of 10,000 words of Kannada corpus according to BIS standards. Moreover we have highlighted the problems which Dravidian languages face in general but Kannada in particular at the level of POS tagging because of their agglutinative nature.

REFERENCES

- Schiffman, H. Sptember 1979. Areference grammar of spoken Kannada. Shridhar, S.N. Kannada (Descriptive Grammars).
- T. N. Vikram and Shalini R. Urs. 2007. Development of prototype morphological analyzer for the south Indian language of kannada. In Proceedings of the 10th international conference on Asian digital libraries: looking back 10 years and forging newfrontiers, ICADL'07, pages 109–116, Berlin, Heidelberg. Springer-Verlag.
- B. R Shambhavi, P Ramakanth Kumar, K Srividya, B J Jyothi, Spoorti Kundargi, and G Varsha Shastri. 2011. Kannada morphological analyser and generator using trie. IJCSNS International Journal of Computer Science and Network Security, 11(1), January.
- Kavi Narayana Murthy. 2000. Computer processing of kannada language. Technical report, Computer and Kannada Development, Kannada University, Hampi.
- P. V. S. Avinesh and G. Karthik. 2007. Part-Of-Speech Tagging and Chunking using Conditional Random Fields and Transformation-Based Learning. In Proceedings of the IJCAI and the Workshop On Shallow Parsing for South Asian Languages (SPSAL), pages 21–24.
- P.J. Antony, M Anand Kumar, and K.P. Soman. 2010. Paradigm based morphological analyzer for kannada language using machine learning approach. Advances in Computational Sciences and Technology (ACST), 3(4).
- P.J. Antony and K.P. Soman. 2010. Kernel based part of speech tagger for kannada. In Machine Learning and Cybernetics (ICMLC), 2010 International Conference on, volume 4, pages 2139 –2144, July.

THANK YOU