

Issues and Challenges in Developing Statistical POS Taggers for Sambalpuri

Pitambar Behera, Atul Kr. Ojha and Girish Nath Jha

Jawaharlal Nehru University
Centre for Linguistics, Special Centre for Sanskrit Studies
{pitambarbehera2, shashwatup9k & girishjha}@gmail.com

Abstract

Low-density languages are also known as lesser-known, poorly-described, less-resourced, minority or less-computerized language because they have fewer resources available. Collecting and annotating a voluminous corpus for these languages prove to be quite daunting. For developing any NLP application for a low-density language, one needs to have an annotated corpus and a standard scheme for annotation. Because of their non-standard usage in text and other linguistic nuances, they pose significant challenges that are of linguistic and technical in nature. The present paper highlights some of the underlying issues and challenges in developing statistical POS taggers applying SVM and CRF++ for Sambalpuri, a less-resourced Eastern Indo-Aryan language. A corpus of approximately 121k is collected from the web and converted into Unicode encoding. The whole corpus is annotated under the BIS (Bureau of Indian Standards) annotation scheme devised for Odia under the ILCI (Indian Languages Corpora Initiative) Corpora Project. Both the taggers are trained and tested with approximately 80k and 13k respectively. The SVM tagger provides 83% accuracy while the CRF++ has 71.56% which is less in comparison to the former.

Keywords: low-density language, parts of speech tagger, SVM, CRF++, Sambalpuri, Eastern IA language

1. Introduction

Low-density languages have fewer resources in terms of the availability of voluminous corpus (McEnry et al., 2000) for NLP applications. The unavailability of a corpus for a low-density language proves to bear adverse impacts on its future NLP development. As rightly pointed out by Ostler (1993), languages that lack active participation in the electronic media are doomed to be endangered. Most of these languages are either dialects or languages with no government recognition. As a result, the situations of these languages in South Asia in general and in Indic languages, in particular, are ‘relatively bleak’ (McEnry et al., 2000). Although India is a land of more than 6000 languages with five prominent diverse language families (Abbi, 2001) only 22 are scheduled and the rest are fighting for their survival.

This paper is concerned with demonstrating the issues and challenges in developing statistical parts of speech taggers for a low-resource language, Sambalpuri. The paper has broadly three objectives. Firstly, it highlights the issues in corpus collection with regard to non-uniform orthographic language standards and non-Unicode encodings of the written text. Secondly, it also attempts to bring out the issues in annotation having without any guideline. Finally, it demonstrates the challenges in developing statistical POS taggers for Sambalpuri owing to the typical, unobserved and language-specific linguistic nuances.

2. Sambalpuri: A Low-Density Eastern IA Language

Sambalpuri (ISO 639-3) is an Eastern Indo-Aryan (IA) language is also known as Dom, Kosali, Koshal, Koshali, Western Odia¹. It is spoken in the ten districts of western and south-western Odisha which comprises Bargarh, Bolangir, Kalahandi, Sonapur, Sambalpur, Jharsuguda, Sundargarh, Deogarh, Boudh, Nuapada; and Athmallik

sub-division of Angul district. In comparison to its sister languages such as Maithili, Awadhi, Angika, Bengali, Assamese, Odia and many others, Sambalpuri has not gained much attention; neither from linguists nor from the government. It is really quite obvious to affirm that it shares the genetic affinity with the Indo-Aryan language family by observing some of its linguistic features (Kushal, 2015). Although it has 75-76% lexical similarity with Standard Oriya (Mathai & Kelsall, 2013), it is syntactically a distinct language (Tripathy, 1984, pp. 49).

3. Salient Linguistic Features

Less-described languages have some of the most interesting linguistic features that are typical and language-specific. Some of the features like agglutination, classifiers, serial verbs, multi-words, compounds etc. account for the less accuracy of any of the statistical NLP applications. Some of them are discussed vividly below in four sub-sections: agglutination, classifiers, reduplication and compounds.

3.1. Agglutination

In an agglutinative language words are made of a linear sequence of distinct morphemes each of which corresponds to a definite meaning (SIL International, 2004). In Odia, the categories such as “suffixes, postpositions, and case endings agglutinate with the verbs, nouns, adverbs or pronouns” (Padhy & Mohanty, 2013; Jena et al., 2011). Similar is the case with Sambalpuri language.

For instance,

k^haɛbar-ke ‘to eat’

lɔk-ɔr ‘peoples’

bahar-ke ‘to outside’

mɔr-nɔ ‘from me’

In the examples instantiated above, all the case endings or markers /ke/, /ɔr/, /nɔ/ agglutinate with their head categories verb, noun and pronoun. /ke/ which is

¹ <https://www.ethnologue.com/language/spv>

equivalent to the English infinitive and preposition is alternating here with both verb and spatial noun /k^hæbar/ and /bahar/ respectively.

3.2. Classifiers

Classifier is one of the most prominent phenomena in Indian languages; especially in the Tibeto-Burman languages and Dravidian languages. Besides, in some of the EIA languages like Bengali (Bhattacharya, 1999), Odia (Neukom, 2003 and Behera, 2015), Bhojpuri (Shukla, 1981), Marathi (Baskaran et al., 2008) and so on, it is a dominant linguistic feature in Sambalpuri as well. The classifiers mainly occur either as proper classifiers, attached to numerals or to the quantity word /keṭe/ ‘how many, some’, or as indefinite markers, in combination with the suffix /-e/” (Neukom, 2003) as /te/, /ṭa/, /k^hṣḍe/, /ṣ^hone/, /ṭi/ etc. in Sambalpuri. One of the rarely observed phenomena of Indian languages found in Sambalpuri is that classifiers also occur with post positions.

For instance, /m̄r lek^he-ṭa/ ‘like me’.

3.3. Reduplication

“It is the repetition of a segment, a syllable, or some part or whole of a lexical or phrasal unit leading to a semantic or grammatical modification” (Pandey, 2007). There are two types of reduplication: partial and total. In total reduplication, the whole part of the base is reduplicated and in the partial reduplication, some part is reduplicated (Abbi, 1992). In the following instances, the first two are fully reduplicated while the rest of the following are partial. In the partial reduplication, the final syllables /-na/ of both the words are reduplicated like in the first example whereas the final example contains the reduplications of the initial syllables /hu-/.

For instance,

çık çık ‘shining’
ḡ^hire ḡ^hire ‘slowly’
jona s̄na ‘known’
huḡa huḡi ‘abusing’

When a sequence of verbs occurs in a chronology, they are called serial verbs (Jha et al., 2014) and some of them are reduplicative in nature. In the below-instantiated example, it is quite confusing as to how to annotate the verbal occurrences. Because the initial verb /ḡḡḡ/ is a non-finite verb followed by a verbal reduplication which is behaving like a manner adverb modifying the finite following verb /p̄leiḡla/. The issue here is how to annotate the verbal reduplication.

For example,

se mark̄ri ḡḡḡ ḡḡḡ p̄leiḡla
he V-Nonfinite V-reduplication V-Finite
“He went away beating.”

3.4. Compounds

Compound or Sandhi is one of the most productive linguistic phenomena which is quite typical in most of the worlds’ languages in general and in Indian languages in particular. There are three basic types of compounds: vowel, consonant and visarga. In the following instance, the first word is an adjective and the second is a noun, but when get combined they comprise a nominal category. Since Sambalpuri a head final language, the annotation

label is decided on the basis of the category of the head. Here the head is a nominal element and hence the judgment goes in favor of the category of the label of the head word.

For example,

s̄ḡḡḡ + p̄ḡḡḡ\N_NN = s̄ḡḡḡḡ\N_NN ‘good path’

So, in the above example, the decision whether to annotate the word as an adjective or noun goes for the right-headedness feature of Sambalpuri. This feature is typical to most of the IA languages and the word /s̄ḡḡḡḡ/ is labelled as a noun.

4. Methodology

This section deals with (a) the total corpus collected in four major domains, (b) the BIS annotation guideline adapted for Sambalpuri, (c) size of the corpus for training, testing and development stages and (d) features selection for SVM and CRF++ POS taggers.

4.1. Corpus Size

The tabulated data (see table 1) demonstrates the total corpus size collected for developing the Sambalpuri POS taggers. The whole corpus size comprises of five domains, viz. literature, sports, tourism, entertainment, and miscellaneous. The highest corpus size is registered in the domain of entertainment i.e., approximately 40k while the ‘miscellaneous’ section accounts for the lowest number of data.

Domains	Tokens
Literature	30, 344
Sports	21, 121
Tourism	26, 767
Entertainment	40, 554
Miscellaneous	2, 424
Total	1, 21, 210

Table 1. Total Corpus Size Domain-wise

4.2. Corpus Annotation

The whole Sambalpuri corpus is annotated using the ILCIANN² (Kumar et al., 2012) following the BIS-ILCI tagset (see table 2) devised for Odia language since there is no tagset available for it. The BIS tagset is a hierarchical set designed by the POS Standardization Committee appointed by the Department of Information and Technology, Government of India. It has a total number of 11 categorical labels at the top level and 39 fine-grained labels for the annotation. The tagset is framed keeping in view both the fineness and coarseness or flat and hierarchical structures in view. The table below contains the nomenclatures of all the categories in the second column, annotation labels in the third and categorical IPA examples in the fourth.

Sl. No.	Category	Annotation Labels	Examples of Sambalpuri in IPA
1	Noun	N	
1.1	Common	N_NN	p̄ḡḡḡ, p̄ṭṭṭ, b ^h abna, monus

² <http://sanskrit.jnu.ac.in/ilciann/index.jsp>

1.2	Proper	N_NNP	ram, hīmalāj, gəŋaḍʰər meher bīstəbīḍjalāj, səmbəlpur etc.
1.3	Verbal	N_NNV	pəḍʰa, pəhōra, ḍəga, nacbarṭa
1.4	Spatial & temporal	N_NST	aḡke, pəcʰaḍe, pəre, purbo, etc.
2	Pronoun	PR	
2.1	Personal	PR_PRP	moī, toi, apən, se etc.
2.2	Reflexive	PR_PRF	nṛje etc.
2.3	Relative	PR_PRL	jaḡar, jaḡākər, jeṇmankər
2.4	Reciprocal	PR_PRC	nṛjər bʰiṭre, ḍohe etc.
2.5	Wh-word	PR_PRQ	kie, kaḡar, ken mane, etc.
2.6	Indefinite	PR_PRI	əŋjərə, kəuṇsī, kehī etc.
3	Demonstrati ve	DM	
3.1	Deictic	DM_DMD	ī, se, iḡuṭakə, seḡuṭakə etc.
3.2	Relative	DM_DMR	jeṇḡuṭakə, jaḡar
3.3	Wh-word	DM_DMQ	kaṇa, kenər, keṇḡuṭakər etc.
3.4	Indefinite	DM_DMI	oniā, kənsī etc.
4	Verb	V	
	Main	V_VM	
4.1. 1	Main	V_VM	sū, ḍəuṭ, ḍəkʰ etc.
4.1. 2	Non-finite	V_VNF	kʰar kəri, nacī nacī, etc.
4.1. 3	Infinitive	V_VINF	kʰaɛbarke, kʰarbar laḡi, nacbar etc.
4.1. 5	Gerund	V_VNG	kʰarṭʰibar, kʰarṭʰibar etc.
4.2. 1	Auxiliary	V_VAUX	uṭi, ḍərkar, kəri, ṭʰibar etc.
5	Adjective	JJ	bʰəl, uṭṭəm, sūḍər etc.
6	Adverb	RB	
7	Postposition	PSP	saṇe, lekʰe, laḡi etc.
8	Conjunction	CC	
8.1	Coordinator	CC_CCD	kā helaje ki, karən, au, jəḍi etc.
8.2	Subordinator	CC-CCS	jəḍi ṭebe, jeṭebele seṭebele, je, bəlī etc.
8.3	Quotative	CCS_UT	aare, hæ lb, hæḡp, aḡjā etc.
9	Particles	RP	
9.1	Default	RP_RPD	məḍʰjə, hī, ṭə etc.
9.2	Classifier	RP_CL	ḡuṭe, ḍuṭa, kʰḍḍe etc.
9.3	Interjection	RP_INJ	vaḡ, hæe, ah, vḡv etc.
9.4	Intensifier	RP_INTF	əṭi, kʰub, bəḡuṭ, jəbər etc.
9.5	Negation	RP_NEG	naī, nohe, ni, nīha etc.
10	Quantifiers	QT	
10.1	General	QT_QTF	ṭʰoḍe, besī, ṭike, ḡoḍaḍu etc.
10.2	Cardinal	QT_QTC	ek, ḍə, ṭm, car etc.
10.3	Ordinal	QT_QTO	Pəhela, ḍosra, ṭsra etc.
11	Residuals	RD	
11.1	Foreign words	RD_RDF	languages of the other scripts except Odia
11.2	Symbol	RD_SYM	mathematical and

			other symbols (#, [, {, %, \$, <, >, (,), *, @,)
11.3	Punctuation	RD_PUNC	(, ; : ' ' " " :- etc.)
11.4	Unknown	RD_UNK	Tags that are left undecided
11.5	Echo word	RD_ECH	baḡʰo-pʰaḡo, koṭa- cʰoṭa etc.

Table 2. BIS POS Tagset Adapted for Sambalpuri

4.3. Data Size for the Taggers

The tabulated data (see table 3) represents the different data sets applied to develop the statistical taggers. The total number of training data used for developing the taggers amounts to around 80k. Initially, the tagger is trained with around 50k with manually annotated data and later, the development set consists of 28k which was automatically tagged and manually validated. After the training period, the testing is conducted with a set of approximately 13k corpus size tokens.

Data sets	Tokens
Training	80, 288
Testing	12, 791

Table 3. Training, Development and Testing Data Sets

4.4. Developing POS Taggers

Two statistical taggers are developed for Sambalpuri; the first one is trained with SVM (Joachims, 1999; Giménez & Márquez, 2006) and the second is with CRF++ (Kudo, 2013). So far as the former is concerned, learning phase contains medium verbose (-V 2) and the mode of learning and tagging is set to left-right-left (LRL). The rest of the features like sliding window, feature set, feature filtering, model compression, C parameter tuning, Dictionary repairing and so on are set to the default mode. On the other hand, the latter is trained with the unigram method.

5. Issues and Challenges

This section is divided into three major sub-sections: corpus-related, human annotation-related and tagger-related issues.

5.1. Corpus-related Issues

The issues pertaining to the corpus collection are vividly discussed: corpora collection, unavailability of Unicode encoding, non-standard usage of the language, different writing conventions and Hindi-like constructions.

5.1.1. Corpus Collection

A number of corpora are developed for various languages like English and some European languages. Considering the situations in non-scheduled (lesser-known) Indian languages it is quite unfavorable in comparison to the scheduled languages since some of the Indian institutions have either worked on or are presently developing language resources and technologies for the latter languages only. Because of the indifference of the government towards the lesser-known languages, the former are getting disempowered gradually. The institutions and projects that have worked for the corpus

collection in scheduled Indian languages are IIT-Hyderabad, CIIL-Mysore, ILCI-JNU, and TDIL.

5.1.2. Unavailability of Unicode Encoding

Since low-density languages are less-resourceful or with no resource the software available for them are also less in number. This leads to the non-Unicode encodings which is not favorable for the development of NLP applications. The whole corpus has been converted into UTF-8 encodings using AkruTi Text Converter³. There are different linguistic issues in the corpus itself such as non-standard usage, non-uniform Orthographic forms and Hindi-like constructions.

5.1.3. Non-standard Usage

Sambalpur is not a scheduled Indian language and is written and spoken with varying standards in different regions of the western and south-western Odisha. For example: Sambalpur, Bargadia (spoken in Bargarh), Bolangiri/a (spoken in Bolangir district), Sundargadi/ia (spoken in Sundargarh), Deogarhia (spoken in Deogarh region) etc. There are some dialectical variations among the people of Sambalpur speaking track. The table (see table 4) demonstrates dialectal variations of Sambalpur with reference to negative morpheme ‘no’, adverbs ‘now’ and ‘this way’. Lexical similarity within the varieties of Sambalpur is considerably high which ranges from 90 to 95 percent (Mathai & Kelsall, 2013). This similarity matrix was made by comparing Bargarhi, Bolangiri and Jharsuguda varieties with Sambalpur.

Variety of Sambalpur	Negative Morpheme [nāi] ‘no’	Adverb [ihāḍe] ‘now’	Adverb [iāḍe] ‘this way’
Bargarh	nuhe/nthe	ihāḍe/ε ^h εn	iāḍe/ɪp ^h ale
Bolangir	nī	εk ^h εn	
Deogarh			
Kalahandi	nī	εk ^h εn	ɪbāṭe
Sambalpur	nthe/nuhe		
Sundargarh		ɪgəḍi	

Table 4. Dialectal Variations in Sambalpur (Adapted from Patel)

5.1.4. Different Orthographic Conventions

A large number of words in Sambalpur has different Orthographic conventions; especially the ligatures. In Sambalpur, there are several writing conventions used for a given word form because of the non-uniform usage of language.

For instance, in the following examples two forms are used for one word with two of them having different POS labels with the change of form.

କାଞ୍ଜି N_NN, କାଞ୍ଜି DM_DMQ

³<https://22bc339da9ca3e2462414546a715752e4c2c5e0d.googledrive.com/host/0B5rBGd680WZFemVLA3RxY0pE0/AkrutiUnicode>

କନ୍ଦକର N_NN, କଣ୍ଠକର N_NN

କାନ୍ଦନ N_NN କାଞ୍ଜି N_NNP

This non-standard usage of the words creates issues during both manual and automatic annotation since their POS labels vary with the varying conventions.

5.1.5. Hindi-like Constructions

Sambalpur is more like Hindi than Odia which accounts for the fact that the western region, where it is spoken, is situated just adjacent to Chattisgarh and Jharkhand where the influence of Hindi is largely felt. In the examples instantiated below /bavəjʊḍ/ and /ke/ are postpositions as used in Hindi while the Hindi-like indefinite and reflexive pronouns are also used.

For instance,

/bavəjʊḍ/ PSP

/hər ek/ PR_PRI /ke/ PSP

/əpnə/ PR_PRF /əpnər/ PR_PRF

5.2. Human Annotation-related Issues

One of the prominent challenges is that which pertains to the annotation of the corpus. For annotation of a voluminous corpus and to maintain consistency, one needs a standard tagset. Owing to the fact that a large number of languages like Sambalpur being less-described or less-studied, it is quite daunting to devise a tagset. If one adopts and adheres to the tagset devised for a language of close proximity, then they may either compromise with the saliency of the linguistic data or may end up filling different slots for labels and not researching by delving deep into some interesting structures. For instance, there are large numbers of homophonous words that can neither be included in the reduplicated nor can they be labelled as echo.

5.2.1. Reduplicated Expressions

Generally, in Indian languages the reduplicated expressions follow the meaningful word. Contrastingly, in Sambalpur many of the reduplicated parts precede the meaningful words (see section 3.3). For instance, in the conjunct verb (adjective + finite verb), /c^hɪc^h/ is the meaningless reduplicated part which is preceding the meaningful part /bɪc^h/ ‘scattered’.

For example,

c^hɪc^hrɪRD_ECH bɪc^hrɪJJ heɪc^hɔn ‘have got scattered’

Similarly, in the following verbal reduplication, the meaningless part is preceding the verbal part.

For example,

kʊṭrɪRD_ECH kʊṭerɪV_VM deləV_VM_VF ‘has tickled’

These kinds of constructions pose significant linguistic challenges for the human annotators as to how to label them and so is for the statistical tagger.

5.2.2. Verb-less Constructions

In Sambalpur and many sister languages such as Odia, Bengali, Assamese (Masica, 1993) verb-less constructions or covertly present verbs are commonly used. These constructions are used with adjectives in place of verbs. Therefore, the tagger also labels some of these adjectives as finite verbs because of the annotation of these constructions in the training data.

For example,

saswət̪\N_NNP babər\N_NN canvas\N_NN osar\JJ pɪsar\RD_ECH \.RD_PUNC “Saswat Babu’s canvass is quite large.”

In the above example, /osar\JJ pɪsar\ RD_ECH/ is the reduplicated adjectival phrase which satisfies the need of the verb.

5.2.3. Onomatopoeic Words

Onomatopoeic words are the imitation of a sound associated phonetically with its describing referent. These following expressions are parts of the multi words because individually these words do not have meaning, but when combined they are manner adverbs. As per the ILCI guideline, if we annotate the first sound as noun and the following words as echo-words (RD_ECH), we are missing relevant linguistic information.

For instance,

bʰɛ̃ bʰɛ̃ ‘loudly’
 ʃʰɒ ʃʰɒ ‘heavily’
 ɡʰɔ̃ pɔ̃ ‘gasping’
 bʰɒ bʰɒ ‘bark’

5.2.4. Agglutination of Classifiers with Postpositions

Agglutination (see section 3.2) is one of the common features in Odia (Behera, 2015) and Sambalpuri along with some IA languages like Bengali and Marathi (Baskaran et al., 2008). In Sambalpuri, one of the peculiar constructions with agglutination is that the classifiers and postpositions agglutinate with each other which is also rarely found in the most agglutinating Dravidian languages. Here, to annotate these constructions as classifiers (RP_CL) or postpositions (PSP) is quite difficult.

For instance,

/baɡɪr-ʈa/ ‘as-CL’
 /lekʰeʈa/ ‘like-CL’

Similarly, in the example below, it is quite difficult to decide the annotation labels for both the human and automatic annotation. The reason is the complexity in deciding the head label of the words. The word /ɖɔɪ-ʈa/ comprises of two components, a cardinal and a classifier morpheme. Both of these categories have separate labels in the BIS scheme for Sambalpuri. Therefore, if one annotates the word as cardinal, they are compromising with the other label or linguistic information.

For instance,

ɖɔɪ-ʈa (ɖɔɪ\QT_QTC ʈa\RP_CL)

5.3. Automatic Annotation-related Issues

These issues are mostly pertained to tagger-related ambiguities and some other linguistic errors.

5.3.1. Ambiguity Issues

The data (see table 5) represented below demonstrates that there are different types of ambiguous sets of classes and their accuracy rates. All the ambiguity classes are divided into 244 classes and they are generated automatically by the SVM tool.

Two-label Sets: This section includes the ambiguous words with two conflicting labels. The most commonly ambiguous tags are coordinating-subordinating

conjunction, coordinating conjunction-general quantifier, deictic-interrogative demonstrative and so on.

ɑʊ (CC_CCD or QT_QTF)

For example,

mɔɪ ɑʊ\CC_CCD mɪr bapa “I and my father”
 mɔɪ ɑʊ\QT_QTF kʰana ɖɛrkar “I need some more food”.

In the above examples, the first one suggests that the word /ɑʊ/ is a coordinating conjunction coordinating two noun phrases while the second one states that it is a general quantifier used as a pre-modifier.

Three-label Sets: This section contains the ambiguous words having three labels. The most commonly ambiguous tags are most-expectedly adjective-temporal nouns-finite, negative-main-finite verb and so on.

bahar (V_VM or N_NST or JJ)

For instance,

bahar\V_VM_VF ɡʰəɾɔ\N_NN “Come out of the house”.
 se pələla bahar\JJ ɡʰəɾɔ\N_NN “He went away from the front room”.
 bahar-ke\N_NST as\V_VM_VF “Come to outside”.

In the first instance, the word form /bahar/ has three different POS labels. The first one is annotated as a finite main verb as the sentence is an imperative sentence and the covert subject is the second person pronominal. In the second example, it is labelled as an adjective as it modifies the following noun whereas the third one is a spatial noun as it refers to a location.

More than Three-label Sets: The words having more than three labels are encapsulated in this part. For instance, main-auxiliary-nonfinite-finite verbs, unknown-interjection-default particle and so on.

kəɾɪ (V_VAUX or V_VM or V_VM_VF or V_VM_VNF)

For instance,

kʰaɪ\V_VM kəɾɪ\V_VAUX asle\V_VM_VF
 kam\N_NN kəɾɪ\V_VM asle\V_VM_VF
 kəɾɪ\ɦɪɪ\V_VM_VF
 kʰaɪkəɾɪ\V_VM_VNF asle\V_VM_VF

The verbal word form /kəɾɪ/ has more than three labels in the corpus and which is rightly so. It can be used as main, auxiliary, finite and non-finite verbs as instantiated in the above examples.

Classes of Ambiguity	Label Sets
2 Sets:	CC_CCD_CC_CCS, CC_CCD_QT_QTF, DM_DMD_DM_DMQ
3 Sets:	JJ_N_NST_V_VM_VF, RD_ECH_V_VM_V_VM_VNF, RP_NEG_V_VM_V_VM_VF
More than 3 Sets	V_VAUX_V_VM_V_VM_VF_V_VM_VNF, RP_INJ_RP_NEG_V_VM_V_VM_VNF, RD_UNK_RP_INJ_RP_RPD_V_VM_V_VM_VNF

Table 5. Ambiguity Classes

6. Results and Discussion

In spite of the different issues and challenges, statistical taggers developed for Odia achieves accuracies of 94% and 89% by SVM and CRF++ respectively (Behera, 2015). The results (SVM 83% and CRF++ 71.56%) for

Sambalpuri are comparatively lesser than Odia which accounts for the fact that Sambalpuri has no standardized orthographic convention and hence different regional varieties use the language in their own ways.

The first and foremost point to emphasize for a low-density language is the large-scale writing on the social media using its own script with a soul objective of developing language resources. If there is less availability of the corpus, then one can also take the assistance of mathematical modelling to achieve a higher accuracy rate. So far as the tagset-related issues are concerned, one can take the labels already used by a closely-related language spoken in the region for annotation job by incorporating it on their convenience. With regard to issues in annotation, one can take the exemplary labels from different tagsets developed for Indian languages. For example, we can incorporate WRB label from the IIIT-Hyderabad for the interrogative adverb. The reduplicative expressions need to be considered seriously as they are the most vital parts of the language and they behave quite differently in Sambalpuri. Therefore, it can be averred that labels for reduplication (RD_REDP), possessive pronouns (PR_POS) and demonstratives (DM_POS), interrogative adverbs (WRB) can be introduced. For handling agglutination, a stemmer or a lemmatizer could be used with statistical POS taggers. For punctuations, fine-grained labels should be incorporated based on their functions in a given context as they can be used as coordinators, section headers, list item markers and so on. Finally, the standardization of the language would help solve many of the issues by providing consistency in both the human and statistical annotations.

7. Conclusion

In this paper, we have discussed about different issues and challenges in terms of both corpus collection, annotation and tagger-related issues in detail for a less-resourced language, i.e. Sambalpuri. The results (SVM 83% and CRF++ 71.56%) of the statistical taggers for Sambalpuri in the present study would not only prove to be beneficial for its own future NLP research and development but also would be advantageous for any other morphologically-rich less-resourced language from around the world. At a later stage, we can further use a lemmatizer or stemmer to handle the agglutination issue and incorporate some of the solutions proposed in the research. Furthermore, these POS taggers could be potentially used for developing morph analyzer, chunker, parser and hopefully for enhancing the accuracy of machine translation. For its future development, an online lexical dictionary using Language Explorer, Lexique Pro & Toolbox can also be prepared.

References

- Abbi, A. (1992). *Reduplication in South Asian Languages: An Areal, Typological, and Historical Study*. India: Allied Publishers Pvt. Ltd.
- Abbi, A. (2001). *A Manual of Linguistic Fieldwork and Structures of Indian Languages (Vol. 17)*. Lincom Europa.
- Baskaran, S., Bali, K., Bhattacharya, T., Bhattacharyya, P., & Jha, G. N. (2008). A Common Parts-of-speech Tagset Framework for Indian Languages. In: *LREC 2008*.
- Behera, P. (2015). *Odia Parts of Speech Tagging Corpora: Suitability of Statistical Models*. M. Phil. Thesis. Delhi: Centre for Linguistics, Jawaharlal Nehru University.
- Bhattacharya, T. (1999). *The Structure of the Bangla DP*. Doctoral Dissertation, London: University College.
- Giménez, J. and Màrquez, L. (2006). *Technical Manual v1.3*. Barcelona: Universitat Politècnica de Catalunya.
- Jena, I., Chaudhury, S., Chaudhry, H., & Sharma, D. M. (2011). Developing Oriya Morphological Analyzer Using Lt-toolbox. In: *Information Systems for Indian Languages*. pp. 124-129. Berlin Heidelberg: Springer.
- Jha, G. N., Hellan, L., Beermann, D., Singh, S., Behera, P. and Banerjee, E. (2014). Indian Languages on the TypeCraft Platform– The Case of Hindi and Odia, Iceland. In: *LREC*.
- Joachims, T. (1999). Making Large Scale SVM Learning Practical. Universität Dortmund.
- Kudo, T. (2013). CRF++: Yet another CRF toolkit. Retrieved from: <http://crfpp.sourceforge.net/ptojrcs/crfpp/>. Access date: July 10, 2015.
- Kumar, R., Kaushik, S., Nainwani, P., Banerjee, E., Hadke, S., & Jha, G. N. (2012, March). Using the ILCI Annotation Tool for POS Annotation: A Case of Hindi. In: *13th International Conference on Intelligent Text Processing and Computational Linguistics (CICLing 2012)*, New Delhi, India.
- Kushal, G. (2015). *Case and Agreement in Sambalpuri*. M.Phil. Thesis. Delhi: Centre for Linguistics, Jawaharlal Nehru University.
- Masica, C. P. (1993). *The Indo-Aryan Languages*. Cambridge University Press.
- Mathai, E. K. & Kelsall, J. (2013). Sambalpuri of Orissa, India: A Brief Sociolinguistic Survey. SIL International.
- McEnergy, T., Baker, P., & Burnard, L. (2000). Corpus Resources and Minority Language Engineering. In: *LREC*.
- Mitkov, R. (2005). *The Oxford Handbook of Computational Linguistics*. Oxford University Press.
- Neukom, L., & Patnaik, M. (2003). *A Grammar of Oriya*. Zürich: Seminar für Allgemeine Sprachwissenschaft der University.
- Ostler, N. (1999). Language technology and the Smaller Language. *Elra Newsletter*, 4(2).
- Patel, Kunjabana. (undated). *A Sambalpuri Phonetic Reader*. Sambalpur: Menaka Prakashani.
- Shukla, S. (1981), *Bhojpuri Grammar*. Washington, D. C.: Georgetown University Press.
- Tripathy, B. (1984). *Sambalpuri Semantics*. Graduate Thesis. Sambalpur: Sambalpur University.