



# POS (Parts of Speech) Tagging System for Sindhi Language

Ghazala Gul Junejo, Mir Sajjad Hussain Talpur, Taha Nuzhat, and Shakir Hussain Talpur  
Information Technology Centre, Sindh Agriculture University Tandojam  
[mirsajjadhussain@hotmail.com](mailto:mirsajjadhussain@hotmail.com), \*[mirsajjadhussain@sau.edu.pk](mailto:mirsajjadhussain@sau.edu.pk)

## Abstract

Part of Speech (POS) tagging is a fundamental need for any natural language text processing system. However, building such a classifier is quite challenging due to the inherent ambiguity present in the natural languages where the same word may be used as different part of speech in different contexts. Several efforts have been made to build such taggers for many international languages including English, French, German and Arabic. Now, in order to build Sindhi text processing system, a POS tagger for Sindhi language is much needed. Like Arabic, Sindhi POS tagging is more challenging due to its word morphology. In this thesis, we will describe various techniques that are available for POS tagging and discuss why we may or may not opt for them. We will then present a brief survey of the efforts that have been done so far for POS tagging of Sindhi language. In this research we aim to create our own POS tagger for Sindhi by training the famous Stanford POS tagger over a corpus containing more than 5000 Sindhi words. The performance of the trained POS tagger will be tested by using another test corpus containing 2000 Sindhi words. Manual tagging of words (even with the help of semi-automatic tools) for training purpose in such huge corpora is a significant effort in itself and will be retained for later studies.

**Keywords:** Natural language Processing (NLP), Machine learning, core (NLP) library, HMM, Stanford POS taggers.

## I INTRODUCTION

Part-of-Speech Tagger (POS) is a bit of programming that parses message in some language and appoints a tag to each word (and other token)[1] here we use the Sindhi language, for example ضمير, اسم, صفت, etc. and so forth it's a Sindhi grammatical features, albeit for the most part computational applications use all the more fine-grain POS labels like صفت-اسم. This programming is a Java execution of the log-direct grammatical form taggers depicted in these papers [2]. One of such Stanford taggers Corpus which is seen implanted with some NLP devices or applications. (For AI) [3]. These devices use factual taking in calculations to gain from the corpus so as to appoint labels naturally for an archive [4]. This proposal is a survey of measurable methods connected for POS labelling of Sindhi language. Sindhi is a language of more than 25 million individuals tagging module appoints a tag to tokenized word and scan for

equivocal word [5]. The vague words are those words which can go about as a thing and descriptive word in certain unique situation, or go about as a modifier and qualifier in certain specific circumstance. At that point their uncertainty is settled utilizing punctuation rules. Machine interpretation, discourse change. A Part-of-discourse (POS) is the morphological class of a lexical thing (word). Lexical things which share similar POS are accepted to have comparable morphological conduct. Basic POS incorporate Verb, Noun, Adjective, Pronoun, Adverb, Preposition, etc [6] A POS tagger is a framework for consequently deciding the POS tag in a given content, and it ought to most likely recognize morphological classes by appointing labels to words. The given sentence can be tagged as.

Here we use the Sindhi Pos taggers

For Example if we consider two Sindhi sentences

چوڪري آهي راحت سڻي

Here (noun) اسم (noun) راحت



frameworks and sites are highly spreading dynamically. Spell examiners grammatical features labelling frameworks are essential piece of content preparing frameworks. The examination of grammatical features patterns and spell checker advancement is a basic necessity for Sindhi explicit multilingual data preparing frameworks. The exploration contemplate manages the Sindhi content accessible in Sindhi word reference and will comprehend the 90% of Sindhi content accessible in Sindhi writing.

In this part, we mean to make our very own tagger for Sindhi dialects. The postulation centres around executing another Sindhi grammatical forms labelling framework utilizing a Stanford taggers and centre nlp library.

### III METHODOLOGY

The exploration structure of this investigation contains different strides for the improvement of Sindhi tagger dependent on Stanford tagger. The review of the proposed approach has been given in Figure below containing following advances.

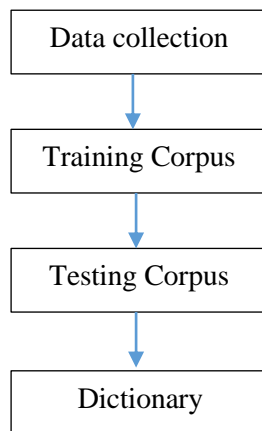


Figure 2: Research Framework

#### 1. Data Collection

The experiments were performed on collected Sindhi text collected from various sources, the text has been collected from Sindhi online books, Sindhi websites, publisher's data and other various sources.

Table 1 represents the list of internet sites and the names of the books where from data is extracted.

Table 1: Sources for gathering data

S.No	Source	URL(s)
1	Daily Kawish	<a href="http://www.thekawish.com">http://www.thekawish.com</a>
2	Sindhi story books	<a href="http://www.majornazaar.blogspot.com">http://www.majornazaar.blogspot.com</a>
3	Sindhi grammar books	<a href="http://www.iqbalkalimatiblog.com">http://www.iqbalkalimatiblog.com</a>
4	Blogs	<a href="http://laari.wordpress.com">http://laari.wordpress.com</a>
5	Literary Writings	<a href="http://voiceofsindh.net">http://voiceofsindh.net</a> <a href="http://sindhssalamat.com">http://sindhssalamat.com</a>

Figure 3 screenshot of Sindhi corpus

سائنسي دور ۾ پهتو آهي. جڏهن انسان اڃا پوريءَ ريت ڳالهائڻ به نه سگهيو هو، ته ان وقت هن کي ڪائنات جي هر شيءِ عجيب ۽ اڃا ئي نظر آئي. جيئن جيئن هن ڳالهائڻ سکي ورتو، تيئن تيئن هن انهن شين جڏهن انسان، ڪ هني. کي پنهنجي سڃاڻپ خاطر نالا ڏيڻ شروع ڪيا. انساني شعور جي ڪنڀل اها پهرين و غارن کي الوداع ڪري پنهنجي هٿن سان جايون اڏي هڪ نئين سماج جي شروعات ڪئي ۽ پوءِ هن اهستي اهستي ڪائنات جي لقائن ۽ مظهرن تي غور ڪرڻ ۽ انهن تي راءِ جوڙڻ جو عمل شروع ڪيو. جو "نفسيات". به انسان کي پنهنجو پاڻ سڃاڻڻ جو هڪ علم مهيا ڪري ٿو "نفسيات" ڪتاب دائرو وسيع آهي، جنهن جو مطلب آهي: هر اها شيءِ جيڪا ساڻهه رکندي هجي، پر هن ڪتاب ۾ انساني جنهن ۾ هو، کي پنهنجو پنهنجو ماحول ڏنو آهي "نفسيات" کي فوڪس ڪيو ويو آهي. خداوند ڪريم هر جيو اهڙي ئي علم تي لکيو ويو آهي. هيءُ هڪ "نفسيات" پنهنجو بچاءُ ۽ واڌ ويجهه ڪري سگهي ٿو. ڪتاب دقيق ۽ ڳوڙهو علم آهي، جنهن تي جيتري شرح لکجي، اها گهٽ آهي. جناب ڊاڪٽر علي احمد قاضي هن ڪتاب جو

#### 2. Training Corpus

The training of the collected corpus is needed to be trained so that the tagger will tag the upcoming data based on the previous training of the corpus. The Stanford tagger provides the various models and built-in APIs to tag the text of various languages ranging from latin to Arabic language text. A particular model is to be used for the particular language.

#### 3. Testing Corpus

The various sets are reserved for the testing. The testing is the checking of the system whether it is working or not. The training corpus is then tested once it was trained on the bases of classifying classes of parts of speech.

#### 4. Explanation of Training and Test Corpus

Corpus of any language can help in taking out and uncovering the examples in content alongside the advancement of measurable information. In the present research, two corpora were considered, a preparation corpus (Rahman, 2010) and a test corpus. Preparing corpus for discovering corpus through, and test corpus for grammatical features labelling utilizing stand passage taggers. The fundamentals for corpus test are like the preparation corpus.

Preparing Corpus incorporates 2.1 million words and it was processed already to recognize remarkable Sindhi words as it were. In the wake of preprocessing the quantity of exceptional Sindhi words was 70576. Preparing corpus is used to build up a word reference (dictionary) for the undertaking of recognizable proof of grammatical forms labelling framework examination of preparing designs. Test corpus incorporates roughly 2453 words and all out number of one of a kind Sindhi words was 2052. The measurements of preparing, corpus and test corpus are appeared Table 2.

Table 2: Data of Training and Test Corpus

	Training Corpus	Test Corpus
<b>Total Number of Characters</b>	168112	3061
<b>Total Number of Words</b>	210000	6453
<b>Total Number of Unique Words</b>	70576	2052

The details of various Sindhi language-specific issues are discussed in the subsequent sections.

Table 3: Valid Sindhi parts of speech table

S.No.	POS(noun)	Word
1	اسم	نالا

2	اسم	شين
3	اسم	شعور
4	اسم	وقت
5	اسم	ڪائنات
6	اسم	انسان
7	اسم	بھادر
8	اسم	مختلف
9	اسم	تاريخ

Table 4: parts of speech

S.No.	POS(pronoun)	Word
1	ضمير	پنهجي
2	ضمير	هن
3	ضمير	ان
4	ضمير	انهن
5	ضمير	اها
6	ضمير	جيڪا
7	ضمير	هو
8	ضمير	هي
9	ضمير	اهڙي

Table 5: parts of speech

S.No.	POS(adjective)	Word
1	صفت	بچاء
2	صفت	دقيق
3	صفت	دقيق
4	صفت	پھريون
5	صفت	اوهان
6	صفت	متن
7	صفت	بن
8	صفت	بھادر

Table 6: parts of speech

S.No.	POS(verb)	Word
1	فعل	ڳالهائڻ
2	فعل	سڪي
3	فعل	ورنو
4	فعل	سوچ
5	فعل	رسيو
6	فعل	آهي
7	فعل	چيائي
8	فعل	عهدي

Table 7: parts of speech

S.No.	POS(adverb)	Word
1	ظرف	اڃ

2	ظرف	جنهن
3	ظرف	پڻ
4	ظرف	نه
5	ظرف	رڳو
6	ظرف	سان
7	ظرف	به
8	ظرف	ان
9	ظرف	جي

Table 8: parts of speech

S.No.	POS(preposition)	Word
1	حرف جر	لاء
2	حرف جر	جي
3	حرف جر	تائين
4	حرف جر	جو
5	حرف جر	تي
6	حرف جر	۾
7	حرف جر	بابت
8	حرف جر	کي
9	حرف جر	سان

Table 9: parts of speech

S.No.	POS(conjunction)	Word
1	حرف جملو	ته
2	حرف جملو	۽
3	حرف جملو	پر
4	حرف جملو	جيڪڏهن
5	حرف جملو	يا

S.No.	POS(interjection)	Word
1	حرف ندا	واه واه
2	حرف ندا	هون
3	حرف ندا	افسوس

T  
a  
b  
l  
e  
1  
0  
:  
p  
a  
r  
t

s of speech

#### 5. Tools used

For the advance of elements of Speech labelling, java language has been chosen. the aim for this determination is that the stage autonomy. The Stanford tagger has been used for the labelling of the Sindhi content. The preparation and testing has been finished by the given genus Apis Machine learning by the Stanford tagger and center IP.

#### IV FINDINGS OF THE STUDY

For English and other western languages various kinds of POS tagging models have been implemented. Resultantly the performance is over 90% marked. Opposite to this, very minute research has been done for Sindhi and some other South Asian languages. The findings of the SPSAL machine learning contest 2006 showed the performance of various taggers that is up to 60 to 78% for Sindhi language. Present study is based on the baseline tagging models which perform well in these intervals for the progress and development of

data provided, without using special device like morphological analyzers and others.

Table 11: Performance of POS Taggers for Sindhi [Test data: 85 sentences, 1000 tokens from the (Protho m- Alo) corpus; Tagset: Level 1 Tagset (14 Tags)]

Token	HMM Accuracy	Unigram Accuracy	Brill Accuracy
0	0	0	0
60	15.4	51.2	50.4
104	18	51.1	44.6
503	34.2	60.7	56.3
1011	42.3	64.2	62.6
2023	45.8	69.1	67.8
3016	49.4	70.1	70.9
4484	45.6	71.2	71.3

#### V DISCUSSION

Many of the languages have their own tagging system. A very little work has been found on regional languages of Pakistan, especially Sindhi language. Automated Sindhi Parts of Speech (POS) tagging system for Sindhi language has been presented. The stand ford taggers and core nlp technique of the Sindhi parts of speech tagging system checker detected a total of 336 discretization and ambiguity words.

#### VI CONCLUSION

It is out of the blue of point by point examination has been completed on the grammatical forms taggers patterns for Sindhi language. Sindhi, being a very homographic language, offers numerous language explicit issues other than customary characters grammatical features like thing, pronoun, descriptive word, verb modifier, and so on. In this proposal, corpora of Sindhi language were breaking down to initially explore the grammatical forms labelling in Sindhi. After the distinguishing proof of these examples, the grammatical features labelling framework checker was created utilizing a stand portage tagger (word reference and factual) approach. With the assistance of lexicon created utilizing preparing corpus and factual investigation on a test corpus, A framework for the labelling of Sindhi has been introduced by taking contribution from the record. The proposed framework takes record

from the framework and afterward labels the sindhi content as indicated by the Parts of Speech and the yield is the labeled content. The framework likewise produces the accuracy and review rate of the content. A significant number of the precedents have been displayed here and the remaining has been left for the reference section. The framework is producing the content as indicated by the fixed information and accessible word references dependent on Stanford tagger and the framework is constrained to given content. The framework might be reached out to different sort of content and the quantity of elements or classes can be expanded. The framework might be reached out to the chatbot and transliteration.

## REFERENCES

- [1] Agrawal, H., & Mani, A. (2006). Part Of Speech Tagging and Chunking Using Conditional Random Fields. In *Proceedings of the NLP/MLcontest workshop, National Workshop on Artificial Intelligence*.
- [2] Antony, P. J., & Soman, K. P. (2011). Machine transliteration for indian languages: A literature survey. *International Journal of Scientific & Engineering Research, IJSER*, 2, 1-8.
- [3] Garg, N., Goyal, V., & Preet, S. (2012, December). Rule based Hindi part of speech tagger. In *Proceedings of COLING 2012: Demonstration Papers* (pp. 163-174).
- [4] Joshi, N., Darbari, H., & Mathur, I. (2013). HMM based POS tagger for Hindi. In *Proceeding of 2013 international conference on artificial intelligence, soft computing (AISC-2013)* (pp. 341-349).
- [5] Kaur, M., Aggerwal, M., & Sharma, S. K. (2014). Improving Punjabi Part of Speech Tagger by Using Reduced Tag Set. *International Journal of Computer Applications & Information Technology*, 7(2), 142.
- [6] Kumar, D., & Josan, G. S. (2010). Part of speech taggers for morphologically rich indian languages: a survey. *International Journal of Computer Applications*, 6(5), 32-41.
- [7] Hasan, F. M., UzZaman, N., & Khan, M. (2007). Comparison of different POS Tagging Techniques (N-Gram, HMM and Brill's tagger) for Bangla. In *Advances and innovations in systems, computing sciences and software engineering* (pp. 121-126). Springer, Dordrecht.
- [8] Mahar, J. A., & Memon, G. Q. (2010, February). Rule based part of speech tagging of sindhi language. In *2010 International Conference on Signal Acquisition and Processing* (pp. 101-106). IEEE.
- [9] Mohnot K, B. N. & Singh S P, K. A. (2014) Hybrid approach for part of speech tagger for hindi language, *International Journal of Computer Technology and Electronics Engineering*, 4, 25-30.
- [10] Shah, S. A. A., Bukhari, S. S. A., Humayun, M., Jhanjhi, N. Z., & Abbas, S. F. (2019, April). Test case generation using unified modeling language. In *2019 International Conference on Computer and Information Sciences (ICCIS)* (pp. 1-6). IEEE.
- [11] Sakiba, S. N., Shuvo, M. M. U., Hossain, N., Das, S. K., Mela, J. D., & Islam, M. A. (2021). A Memory-Efficient Tool for Bengali Parts of Speech Tagging. In *Artificial Intelligence Techniques for Advanced Computing Applications* (pp. 67-78). Springer, Singapore.
- [12] Malema, G., Okgetheng, B., Tebalo, B., Motlhanka, M., & Rammidi, G. (2020, May). Complex Setswana Parts of Speech Tagging. In *Proceedings of the first workshop on Resources for African Indigenous Languages* (pp. 21-24).
- [13] Awwalu, J., Abdullahi, S. E., & Ewwiekpaefe, A. E. (2020). A Corpus Based Transformation-Based Learning for Hausa Text Parts of Speech Tagging. *International Journal of Computing and Digital Systems*, 10, 2-19.
- [14] Awwalu, J., Abdullahi, S. E. Y., & Ewwiekpaefe, A. E. (2020). PARTS OF SPEECH TAGGING: A REVIEW OF TECHNIQUES. *FUDMA JOURNAL OF SCIENCES*, 4(2), 712-721.
- [15] Tucker, B. V., & Wright, R. (2020). Introduction to the special issue on the phonetics of under-documented languages. *The Journal of the Acoustical Society of America*, 147(4), 2741-2744.