



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

H Y D E R A B A D

Computational Linguistics 1

Project Report

Submitted by:-

Aamir Farhan (20161078)

Ojaswi Binnani (20161006)

Title

Verb Argument Patterns in Urdu

Introduction

The way through which we share our contents or feelings have always great importance in understanding and processing of language. Parsing is the most suited approach in identifying and scanning what the available sentences express. Parsing is the process in which syntactic structure of sentence is identified using grammatical tags. Phrase and Dependency are two main structure formalisms for parsing natural language sentences. We explored that dependency parsing is more appropriate for Urdu and other free word order languages.

The Shakti Standard Format (SSF) format is generated with the help of a morph analyser and chunking. SSF is a highly readable representation for storing language analysis. It is designed to be used as a common format or common representation on which all modules of a system operate.

Problem Statement

Extracting verb argument patterns from Urdu treebanks.

Method Followed

The project was divided into three phases which were as follows:-

Phase 1 - We built our own parser in python which takes the folder containing ssf files and extracts the verb frames and their frequency corresponding to each verb.

Phase 2 - We used the ssfAPI and we ran it on the ssf files to extract the TAM information for each verb present in the folder.

Phase 3 - Now, using the frequency data generated in phase 1, we give score to each verb frame and from this score we predict the dependencies of the verb used at some other instance.

Discussion

The sample SSF format from an Urdu treebank is shown below :-

<pre><body> <tb number="1" segment="no" bullet="no"> <foreign language="select" writingsystem="RTL"> </foreign> <text> <Sentence id='1'> 1 ((NP <fs name='NP' drel='k2:VGF'></pre>		<pre> <'>=adj,any,any,,d,,,' posn='10' name,سابق=<fs af <'>=adj,any,any,,o,,,' posn='20' name,عالمی=<fs af <'>=n,m,sg,3,o,0,0,' posn='30' name,نمبر=<fs af <'>=num,any,any,,any,,,' posn='40' name,ایک=<fs af <'>=n,m,sg,3,o,0,0,' posn='50' name,فیڈرر=<fs af <'>=psp,,,,,,,' posn='60' name,کو=<fs af</pre>		JJ	سابق	1.1		
				JJ	عالمی	1.2		
				NN	نمبر	1.3		
				QC	ایک	1.4		
				NNP	فیڈرر	1.5		
				PSP	کو	1.6		
				((((
2 ((NP <fs name='NP2' drel='r6:NP3'>				<pre> <'>=n,m,sg,3,d,0,0,' posn='70' name,ٹینس=<fs af <'>=n,f,sg,3,o,0,0,' posn='80' name,تاریخ=<fs af <'>=psp,m,sg,,d,,,' posn='90' name,کا=<fs af</pre>		NN	ٹینس	2.1
				NN	تاریخ	2.2		
				PSP	کا	2.3		
				((((
3 ((NP <fs name='NP3' drel='k2s:VGF'>				<pre> <'>=adj,any,any,,d,,,' posn='100' name,بہترین=<fs af <'>=n,m,sg,3,d,0,0,' posn='110' name,کھلاڑی=<fs af</pre>		JJ	بہترین	3.1
				NN	کھلاڑی	3.2		
				((((
4 ((JJP <fs name='JJP' drel='pof:VGF'>				<pre> <'>=adj,any,any,,,,,' posn='120' name,قرار=<fs af</pre>		JJ	قرار	4.1
				((((
5 ((VGF <fs name='VGF' style='declarative' voicetype='active'>				<pre> <'>=yA' posn='130' name,ہا,ہا,v,m,sg,any,,=<fs af <'>=wA' posn='140' name,تا,تا,v,any,any,any,,=<fs af <'>=hE' posn='150' name,ہے,ہے,v,any,sg,3,,=<fs af <'>=punc,,,,,,,' posn='160' name,,=<fs af</pre>		VM	ہا	5.1
				VAUX	ہا	5.2		
				VAUX	ہے	5.3		
				SYM	.	5.4		
				((((
</Sentence>								
</text>								
</tb>								
</body>								

The above is the ssf format of the Urdu sentence “*Sabik aalami number ek Federer ko tennis tarikh ka behtareen khiladi karar diya jata hai*”

In the above sentence, the main verb is “diya”. We ran the ssfAPI on this sentence and we got the TAM information about the verb, which is, “jata hai” (VAUX).

Next, we ran a self coded parser on a folder containing the ssf files and we extracted the verb frames and its frequency corresponding to a particular verb. Now, we assigned a score value to each of these frames and we can use these values to predict the dependencies this verb would have at any other instance. Thus, the verb argument patterns are used in form of the verb frame frequencies.

Conclusion

We successfully extracted the Verb frames, its frequencies and the TAM information from the given Urdu tree banks in form of ssf files. We can use this frequency data and scores assigned to each verb frames for predicting the dependencies a verb will have at an instance.

Papers Read and References

SSF: Shakti Standard Format Guide, LTRC, IIIT-H.
Formal specifications of Morphology, HCU
A review on Urdu Parsing, Muhammad Javed Arslan Ali Raza, Gomal University, Pakistan.