



Under The Guidance Of :
Dr. Pushpak Bhattacharyya
Dr. Manish Srivastava



Center For Indian Language Technology

Indian Institute of Technology, Bombay



Presentation Outline

- Research Problem Statement
- Motivation
- Approach
- Results
- References
- Future Scope
- Internship Work Timeline



Center For Indian Language Technology

Indian Institute of Technology, Bombay



Problem Statement-

“IdentifyingThe Gender Category Of Nouns (Common/Proper) In Hindi Corpus”

Eg. - राम, जो मिठाई बेचता है, आज घर चला गया ।

↓
M

▼
F

↓
M



Center For Indian Language Technology

Indian Institute of Technology, Bombay

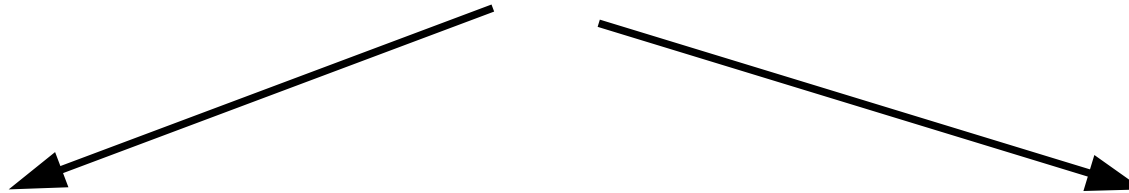


MOTIVATION:

- > Gender of proper and common nouns cannot be identified unless they are specifically assigned a gender in the predefined lexicon list
- > It's tedious to make a list of all the existing proper nouns and manually assign their gender category to them
- > Failure of precise morphology analysis of nouns leads to overall decrease in efficiency of the morphology analyzer techniques and loss of information related to attributes of existing nouns



Can Gender Category be recognised at Morphology / POS Tagging Level ?



Ending Sound / Suffix (If Any)(If Root Form Can Be Identified)

- >>Dependent on the root word list
And suffix replacement rules
- >>Ending Sound (Not Always Right)
अ - पुस्तक, छत(f) , घर, हाथ (m)
ई - लड़की (f) , माली (f)
- >>Not All nouns have suffixes

Number (For Common Nouns)

- >> Word Form Along With
Plurals.
लड़का - लड़कें
क्रीड़ा - क्रीडायें
But ऋतु, औं - feminine
शत्रु, औं - masculine
- >> Also plural to singular is
possible but not vice versa



Center For Indian Language Technology

Indian Institute of Technology, Bombay



Gender of Proper Nouns / Common Nouns cannot be identified at these levels with notable efficiency.



Center For Indian Language Technology

Indian Institute of Technology, Bombay



Taking The Approach To Next Level Of Chunking And Parsing ----

Sentence = Noun Phrase + Verb Phrase

Nouns- Proper/Common



Pronouns



Adjective



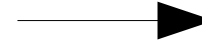
Ordinals



Main Verb



Auxillary Verb



Adverb



Taking The Approach To Next Level Of Chunking And Parsing ----

Sentence = Noun Phrase + Verb Phrase

Nouns- Proper/Common

Pronouns

Adjective

Ordinals

Main Verb

Auxillary Verb

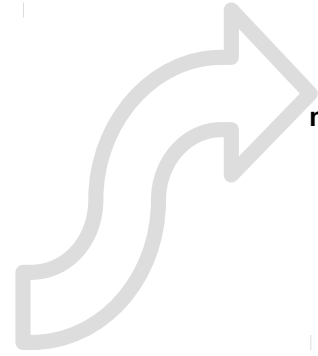
Adverb



--- Categories which change according to gender

Machine was Trained For
Gender Identification Of
Nouns via
Stanford Classifier 3.4
using The Following
Features:-

- 1.useNGrams=true
- 1.usePrefixSuffixNGrams=true
- 1.maxNGramLeng=10
- 1.minNGramLeng=1
- 1.binnedLengths=10,20,30
- 2.useString=true
- 3.usePrefixSuffixNGrams=true
- 4.useString=true
- 5.realValued=true
- 6.useString=true
- 7.useString=true
- 8.useNGrams=true
- 8.usePrefixSuffixNGrams=true
- 8.maxNGramLeng=6
- 8.minNGramLeng=1
- 9.useString=true



m	उत्पादन	n	NP	head	3	-	उत्पादन	0_में	sg
-	में	psp	NP	child	1	-	में	-	-
-	आई	v	VGNF	head	4	f	आ	या	sg
m	गिरावट	n	NP2	head	17	-	गिरावट	0_का_वजह_से	sg
-	की	psp	NP2	child	4	f	का	-	-
f	वजह	n	NP2	child	4	-	वजह	-	-
-	से	psp	NP2	child	4	-	से	-	-
m	गेहूँ	n	NP3	head	10	-	गेहूँ	0_का	sg
-	की	psp	NP3	child	8	f	का	-	sg
f	कीमत	n	NP4	head	17	-	कीमत	0_में	sg
-	में	psp	NP4	child	10	-	में	-	-
m	उछाल	n	NP5	head	14	-	उछाल	0_का	sg
-	की	psp	NP5	child	12	f	का	-	pl
f	आशंकाओं	n	NP6	head	17	-	आशंका	0_को	pl
-	को	psp	NP6	child	14	-	को	-	-
-	खारिज	adj	JJP	head	17	any	खारिज	-	any
-	करते	v	VGNF2	head	21	m	कर	ता_हो+या	sg
-	हुए	v	VGNF2	child	17	m	हो	या	sg
f	सरकार	n	NP7	head	21	-	सरकार	0_ने	sg
-	ने	psp	NP7	child	19	-	ने	-	-
-	कहा	v	VGf	head	0	m	कह	या_है	sg
-	है	v	VGf	child	21	any	है	है	sg
-	कि	avy	CCP	head	21	-	कि	-	-
-	वर्तमान	adj	NP8	child	25	any	वर्तमान	-	any
f	जरूरतों	n	NP8	head	28	-	जरूरत	0_को	pl
-	को	psp	NP8	child	25	-	को	-	-
-	पूरा	adj	JJP2	head	28	m	पूरा	-	sg
-	करने	v	VGNN	head	37	any	कर	ना_के_लिए	any
-	के	psp	VGNN	child	28	-	के	-	-
-	लिए	psp	VGNN	child	28	-	लिए	-	-
m	देश	n	NP9	head	37	-	देश	0_में	sg
-	में	psp	NP9	child	31	-	में	-	-
m	खाद्यान्नों	n	NP10	head	36	-	खाद्यान्न	0_का	pl
-	का	psp	NP10	child	33	m	का	-	sg

SAMPLE TRAINING DATA



Center For Indian Language Technology

Indian Institute of Technology, Bombay



RESULTS OBTAINED

Corpus Identity	No. Of Lines		Accuracy For Nouns
Hindi Tree Bank – LTRC, IIIT Hyderabad	Training Set 12,041	Testing Set 1233	74.24%
Random Hindi Corpus (Siva Reddy's (IIIT hyd) Hindi Dependency Parser and POS Tagger)	1466	227	60.3%

Baseline - 50%

91.3%

84.32%-Unlabelled
78.92%-Labelled



Center For Indian Language Technology

Indian Institute of Technology, Bombay



References :

- [1] Anshuman Tripathi & Manaal Faruqui : Gender prediction of Indian names. In Proceedings of the IEEE Students' Technology Symposium (TechSym) 2011, Kharagpur, India.
- [2] Nivre, Joakim, Johan Hall, and Jens Nilsson. "Maltparser: A data-driven parser-generator for dependency parsing." *Proceedings of LREC*. Vol. 6. 2006.
- [3] Husain, Samar, Prashanth Mannem, Bharat Ram Ambati, and Phani Gadde. "The ICON-2010 tools contest on Indian language dependency parsing." *Proceedings of ICON-2010 Tools Contest on Indian Language Dependency Parsing, ICON 10 (2010): 1-8*.
- [4] Bharati, Akshar, Rajeev Sangal, Dipti Misra Sharma, and Lakshmi Bai. "Anncorra: Annotating corpora guidelines for pos and chunk annotation for Indian languages." *LTRC-TR31 (2006)*
- [5] Begum, Rafiya, Samar Husain, Arun Dhawaj, Dipti Misra Sharma, Lakshmi Bai, and Rajeev Sangal. "Dependency Annotation Scheme for Indian Languages." In *IJCNLP*, pp. 721-726. 2008.



Center For Indian Language Technology

Indian Institute of Technology, Bombay



Future Scope :

- 1 Sentiment Analysis Can Be Done For Determining The Gender Category Of Unrelated Nouns Existing In Any Sentence**
- 2 Assignment Of Gender Category via Sonorants as mentioned in the last referenced paper**
- 3 Efficient Hindi Dependency Parser**

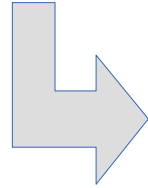


Center For Indian Language Technology

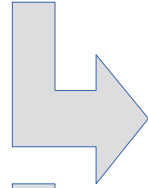
Indian Institute of Technology, Bombay



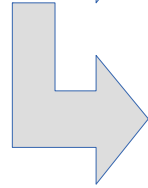
Progress Achieved -



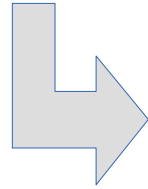
Studying Basic Theory



Debugging and Adding Features (Gender, Number) in Hindi Morphology Analyzer



Solution to Gender of Nouns In Hindi



Final Presentation And Report Submission

