

The Illiterati

Part-of-speech tagging for Magahi and Bhojpuri
without even knowing the alphabet

Thomas Proisl, Peter Uhrig, Philipp Heinrich, Andreas Blombach,
Sefora Mammarella, Natalie Dykes, Besim Kabashi

Chair of Computational Corpus Linguistics
Friedrich-Alexander-Universität Erlangen-Nürnberg

September 11, 2019



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

PHILOSOPHISCHE FAKULTÄT
UND FACHBEREICH THEOLOGIE

NLP Solutions for Under Resourced Languages (NSURL 2019)

Tasks 9 and 10: Part-of-speech tagging for Magahi and Bhojpuri

	Magahi	Bhojpuri
tag set	18	33
training	61.435	94.692
test	8.205	10.582

Our strategies

- use customisable off-the-shelf taggers
- include additional resources (transfer learning)
 - ▶ Brown clusters
 - ▶ word embeddings
 - ▶ tagged corpora of related languages

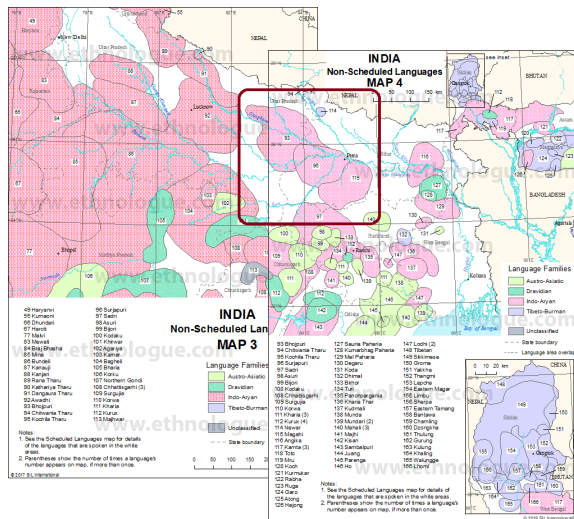
Language Map: Bhojpuri in India



1

¹source:https://commons.wikimedia.org/wiki/File:Bhojpuri_Speaking_Region_in_India.png

Language Map: Magahi and Bhojpuri



²sources: https://www.ethnologue.com/map/IN_03, https://www.ethnologue.com/map/IN_04

About Magahi and Bhojpuri

- principal languages of the Bihari group (West-Eastern Indo-Aryan language)
 - ▶ Magahi
 - ▶ Bhojpuri
 - ▶ Maithili
- Bhojpuri: approx. 51 million native speakers (2011 census), spoken in Western Bihar, Eastern Uttar Pradesh, and in Southwest Nepal
- Magahi: approx. 13 million native speakers (21 million if Khortha, a prominent dialect, is also included), mainly spoken in Southern Bihar

Interesting features w.r.t. POS tagging

- SOV order
- rich verb morphology
- extensive use of postpositions
- Magahi: unusual agreement system (verb has to agree with subject *and* object)

1 Introduction

2 Systems and Strategies

- SoMeWeTa
- BiLSTM-CRF
- Stanford Tagger

3 Results and Error Analysis

- Results
- Error Analysis

4 Conclusion

1 Introduction

2 Systems and Strategies

- SoMeWeTa
- BiLSTM-CRF
- Stanford Tagger

3 Results and Error Analysis

- Results
- Error Analysis

4 Conclusion

Three off-the-shelf POS Taggers

- SoMeWeTa³ (Proisl, 2018)
 - ▶ based on averaged structured perceptron
 - ▶ supports domain adaptation and external resources
- BiLSTM+CRF sequence tagger (Guillaume Genthial, Riedl and Padó (2018)⁴)
 - ▶ based on character and word embeddings
 - ▶ supports transfer learning
- The Stanford Tagger⁵ (Toutanova et al., 2003)
 - ▶ based on a maximum entropy cyclic dependency network
 - ▶ hyperparameter tuning

³<https://github.com/tsproisl/SoMeWeTa>

⁴https://github.com/riedlma/sequence_tagging

⁵<https://nlp.stanford.edu/software/tagger.html>

Additional Resources

freely available resources used in addition to training data

- Hindi UD treebank (HDTB; ca. 352,000 tokens)⁶
- Two Magahi corpora⁷
 - ▶ POS-tagged Magahi corpus (KMI-Mag; ca. 46,000 tokens)
 - ▶ Corpus of untagged Magahi texts (ca. 2.8 million tokens)
- Plain text extracted from Wikimedia dumps⁸
 - ▶ Hindi (ca. 34.7 million tokens)
 - ▶ Bihari (ca. 700,000 tokens)
- Brown clusters (Brown et al., 1992) computed from Wikimedia dumps and untagged Magahi corpus
- Pre-trained fastText embeddings for Hindi and Bihari⁹

⁶https://github.com/UniversalDependencies/UD_Hindi-HDTB

⁷<https://github.com/kmi-linguistics/magahi>

⁸<https://dumps.wikimedia.org>

⁹<https://fasttext.cc/docs/en/crawl-vectors.html>

Experiments Using SoMeWeTa

Focus on Brown clusters and transfer learning

Cross-validation results on training data (selection)

model	accuracy
Bhojpuri (no additional resources)	91.62 (± 0.97)
Bhojpuri (hi)	91.79 (± 1.00)
Bhojpuri (hi+mag)	91.99 (± 0.83)
Bhojpuri (hi+bh+mag)	92.04 (± 0.80)
KMI-Mag \rightarrow Bhojpuri (hi+bh+mag)	92.06 (± 0.94)
Magahi (no additional resources)	88.92 (± 1.24)
Magahi (mag)	89.12 (± 1.23)
Magahi (hi+mag)	89.32 (± 1.15)
Magahi (hi+bh+mag)	89.15 (± 1.17)
KMI-Mag+Bhojpuri \rightarrow Magahi (hi+mag)	89.30 (± 1.14)

Brown clusters beneficial, additional transfer learning not

Experiments Using SoMeWeTa

Focus on Brown clusters and transfer learning

Cross-validation results on training data (selection)

model	accuracy
Bhojpuri (no additional resources)	91.62 (± 0.97)
Bhojpuri (hi)	91.79 (± 1.00)
Bhojpuri (hi+mag)	91.99 (± 0.83)
Bhojpuri (hi+bh+mag)	92.04 (± 0.80)
KMI-Mag \rightarrow Bhojpuri (hi+bh+mag)	92.06 (± 0.94)
Magahi (no additional resources)	88.92 (± 1.24)
Magahi (mag)	89.12 (± 1.23)
Magahi (hi+mag)	89.32 (± 1.15)
Magahi (hi+bh+mag)	89.15 (± 1.17)
KMI-Mag+Bhojpuri \rightarrow Magahi (hi+mag)	89.30 (± 1.14)

Brown clusters beneficial, additional transfer learning not

Experiments Using the BiLSTM-CRF Tagger

Focus on embeddings and transfer learning

Cross-validation results on training data

model	accuracy
Magahi (Hindi embeddings)	88.97 (± 1.14)
Magahi (Bihari embeddings)	89.09 (± 1.00)
HDTB \rightarrow Magahi (Hindi embeddings)	89.85 (± 0.99)
KMI-Mag \rightarrow Magahi (Hindi embeddings)	<i>90.70 (± 0.92)</i>
Bhojpuri (Hindi embeddings)	90.78 (± 0.55)
Bhojpuri (Bihari embeddings)	90.80 (± 0.57)
KMI-Mag \rightarrow Bhojpuri (Hindi embeddings)	<i>91.23 (± 0.68)</i>

Using Hindi embeddings and pretraining on KMI-Mag works best for both languages!

Experiments Using the BiLSTM-CRF Tagger

Focus on embeddings and transfer learning

Cross-validation results on training data

model	accuracy
Magahi (Hindi embeddings)	88.97 (± 1.14)
Magahi (Bihari embeddings)	89.09 (± 1.00)
HDTB \rightarrow Magahi (Hindi embeddings)	89.85 (± 0.99)
KMI-Mag \rightarrow Magahi (Hindi embeddings)	90.70 (± 0.92)
Bhojpuri (Hindi embeddings)	90.78 (± 0.55)
Bhojpuri (Bihari embeddings)	90.80 (± 0.57)
KMI-Mag \rightarrow Bhojpuri (Hindi embeddings)	91.23 (± 0.68)

Using Hindi embeddings and pretraining on KMI-Mag works best for both languages!

Experiments Using the Stanford Tagger: Magahi

parameter	default	value/range
closedClassTags	(none)	ADP AUX CCONJ DET NUM PART PRON SCONJ PUNCT
arch/macro	generic	generic, left3word, bidirectional5words
arch/further unknown-words option	(none)	naacl2003unknowns
arch/unicode shapes for rare words	(none)	(-2,2), (-1,1), (0), (none)
iterations	100	100
learnClosedClassTags	false	false
curWordMinFeatureThresh	2	1..4
minFeatureThresh	5	1..5
rareWordMinFeatureThresh	10	1..10
rareWordThresh	5	1..8
veryCommonWordThresh	250	100, 150, 200, 250

- 76,800 hyperparameter combinations
- 2 runs per parameter combination (first and last 20% as test data)

153,600 training cycles (FAU HPC)

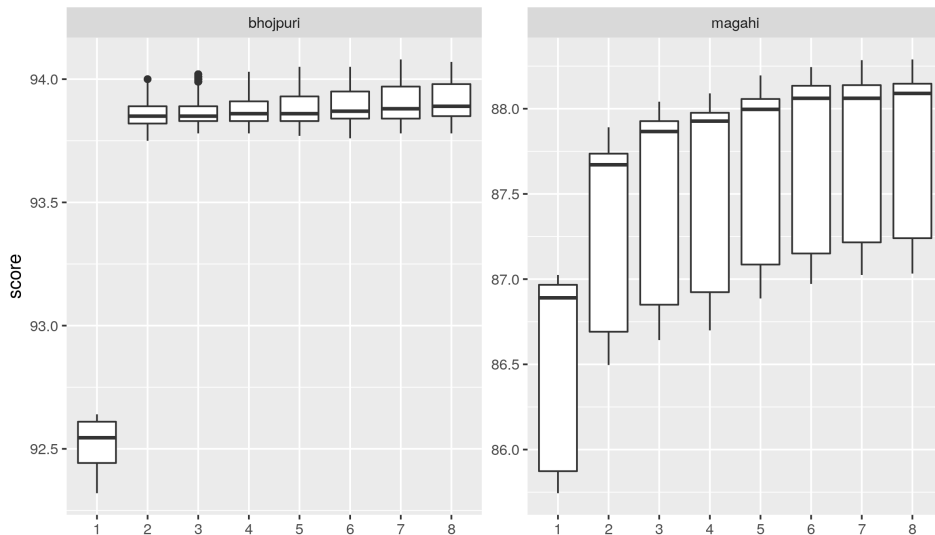
Experiments Using the Stanford Tagger: Bhojpuri

parameter	default	value/range
closedClassTags	(none)	(none)
arch/macro	generic	generic, left3word, bidirectional5words
arch/further unknown-words option	(none)	naacl2003unknowns
arch/unicode shapes for rare words	(none)	(-2,2), (-1,1), (0), (none)
iterations	100	100
learnClosedClassTags	false	true
closedClassTagThreshold	40	40
curWordMinFeatureThresh	2	1..4
minFeatureThresh	5	1..5
rareWordMinFeatureThresh	10	1..10
rareWordThresh	5	1..8
veryCommonWordThresh	250	100, 150, 200, 250

- 76,800 parameter combinations
- full 10-fold crossvalidation

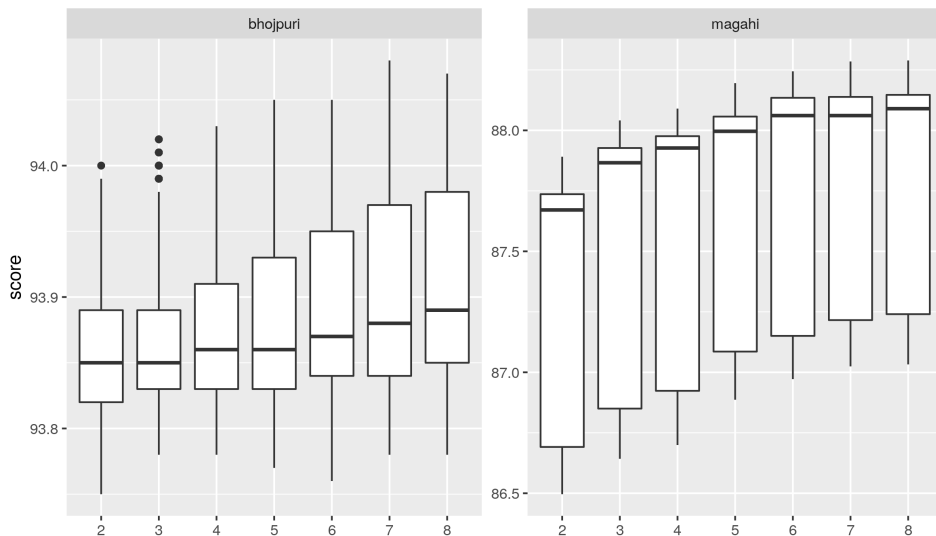
768,000 training cycles (FAU HPC)

Parameter Analysis: rareWordThresh



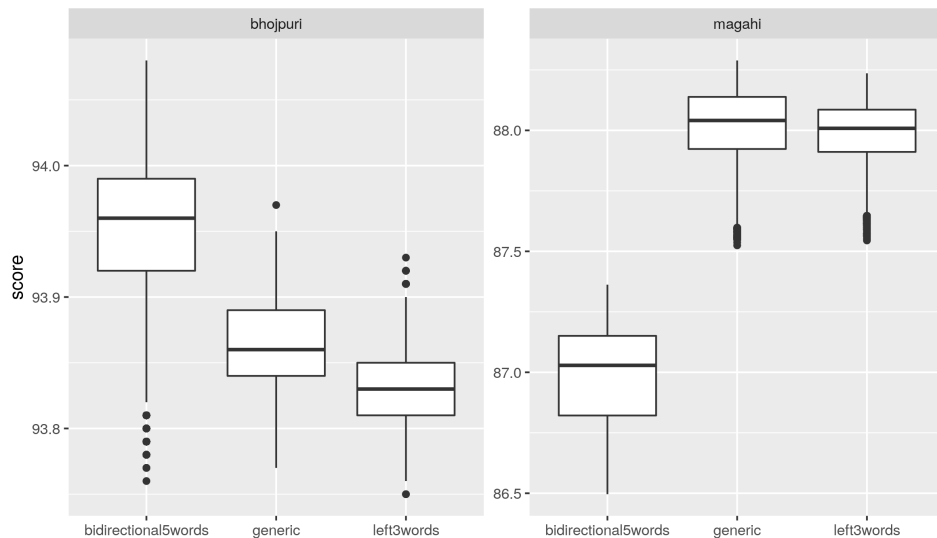
Parameter Analysis: rareWordThresh

exluding rareWordThresh=1



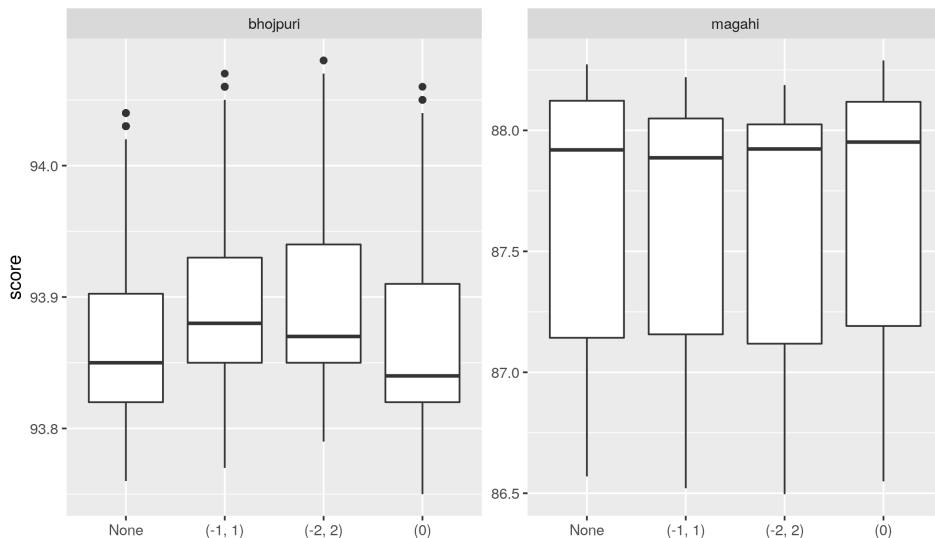
Parameter Analysis: macro

exluding rareWordThresh=1



Parameter Analysis: unicodeshape

exluding rareWordThresh=1



Parameter Analysis: Summary

- performance decreases abruptly when *rareWordThresh* is set to 1 (at least hapax legomena should be treated as rare words)
- performance was insensitive to variation in *veryCommonWordThresh* (this option was ignored by the system)
- macro has most influence
 - ▶ Bhojpuri: `bidirectional5words`
 - ▶ Magahi: `generic` and `left3words` training data annotation?
- *rareWordThresh* explains most of the remaining variation

Official Results

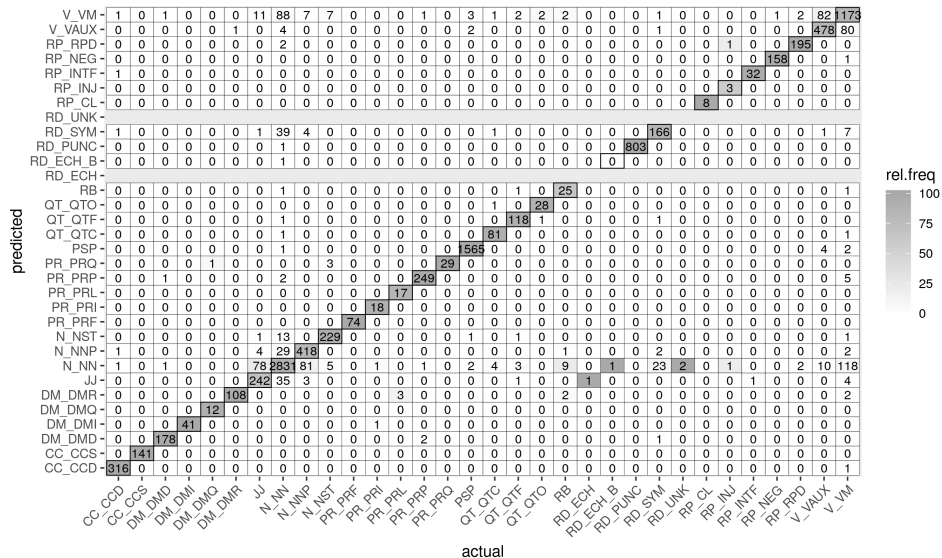
	submission	F_1	acc.
1	Stanford	95	94.78
1	NITK-NLP_SUB1	95	
2	SoMeWeTa	93	92.76
3	BiLSTM-CRF	92	92.01
4	NITK-NLP_SUB2	89	

Table: Results for Bhojpuri

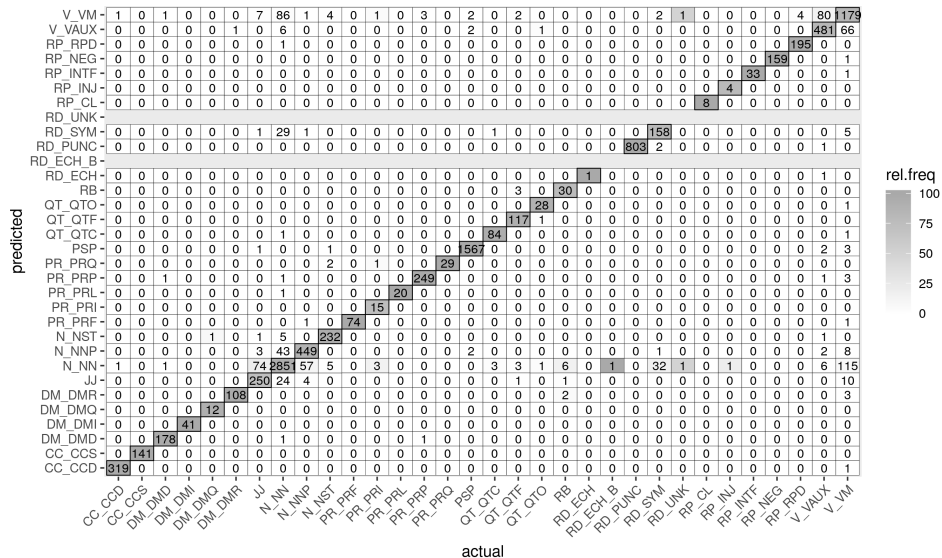
	submission	F_1	acc.
1	NITK-NLP_SUB2	79	
2	BiLSTM-CRF	77	78.86
2	SoMeWeTa	77	78.68
3	Stanford	74	76.57
4	NITK-NLP_SUB1	73	

Table: Results for Magahi

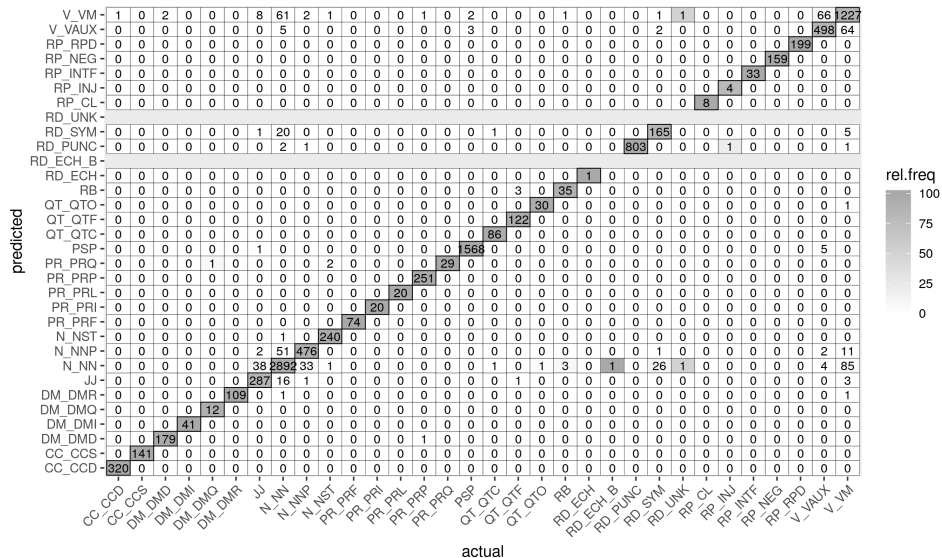
Bhojpuri: Confusion Matrix for BiLSTM



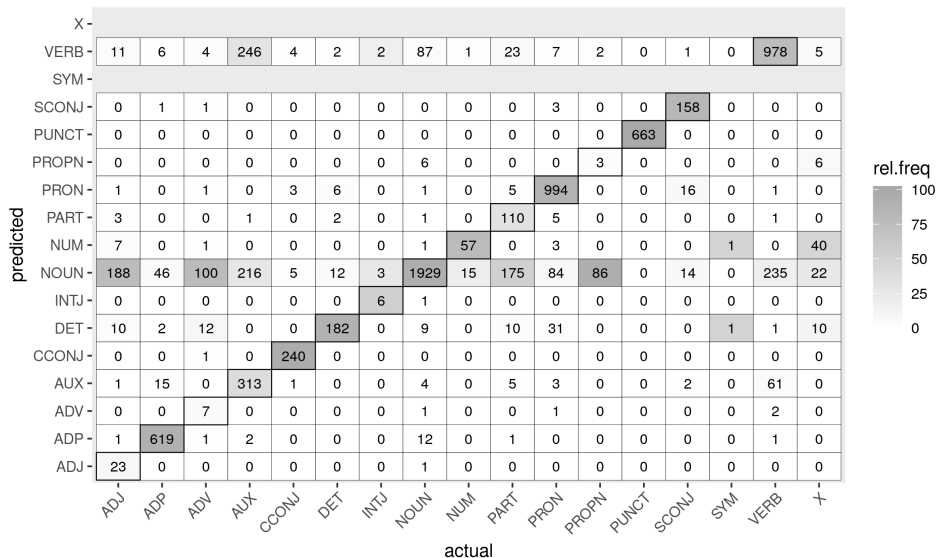
Bhojpuri: Confusion Matrix for SoMeWeTa



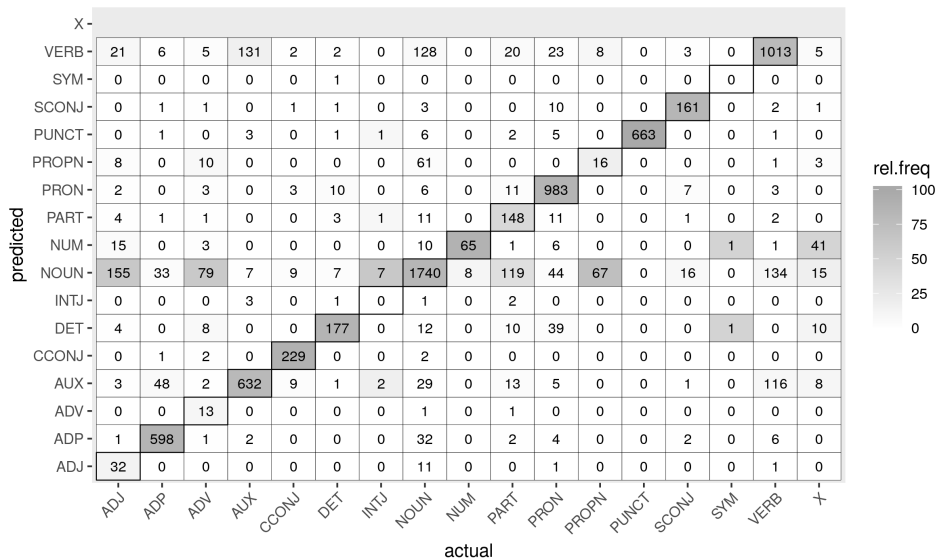
Bhojpuri: Confusion Matrix for Stanford Tagger



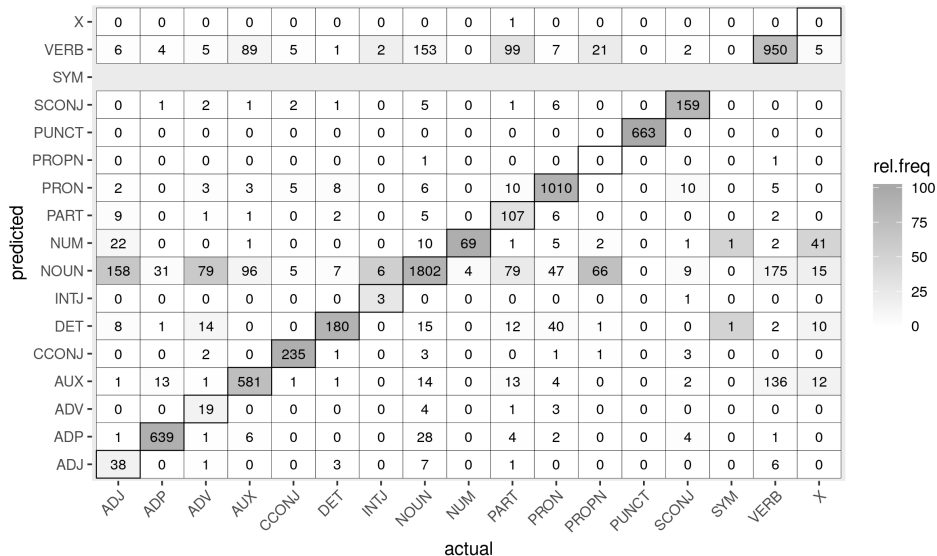
Magahi: Confusion Matrix for Stanford Tagger



Magahi: Confusion Matrix for BiLSTM



Magahi: Confusion Matrix for SoMeWeTa



Error Analysis: Summary

- Bhojpuri ($F_1 \in [89, 95]$)
 - ▶ errors very much what one would expect
 - ▶ rare categories show low recall
 - ▶ frequent tags (N_NN, V_VM) go-to labels for misclassifications
 - ▶ confusion of similar morphosyntactic categories ($V_VM \leftrightarrow V_AUX$, $N_NN \leftrightarrow N_NNP$)
- Magahi ($F_1 \in [73, 79]$)
 - ▶ major problems:
 - ★ recall for PROPN, X
 - ★ ADJ (15.5%), ADV (14.8%), PART (32.5%)
 - ▶ go-to labels NOUN and VERB
 - ▶ obvious confusion: $VERB \leftrightarrow AUX$
 - ▶ different tag distributions in test and training data
 - ★ 602 ADJ and 16777 NOUN in the training set
 - ★ 245 ADJ and 2053 NOUN in the test set

Conclusion

- results for Bhojpuri very satisfying
 - ▶ close to 95% accuracy on a tagset (33 tags, 100,000 tokens training data)
 - ▶ a bit of a downer: mindless parameter-tuning yields best results
 - ▶ differences in system performance probably not significant
- results for Magahi very disappointing
 - ▶ problems with a-priori tag distribution
 - ▶ here: use of additional resources outperforms mere parameter-tuning

Thanks for listening.
Questions?

References

- Peter F. Brown, Vincent J. Della Pietra, Peter V. de Souza, Jennifer C. Lai, and Robert L. Mercer. Class-based n-gram models of natural language. *Computational Linguistics*, 18(4):467–479, 1992.
- Thomas Proisl. SoMeWeTa: A part-of-speech tagger for German social media and web texts. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, pages 665–670, Miyazaki, 2018. European Language Resources Association.
- Martin Riedl and Sebastian Padó. A named entity recognition shootout for German. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL 2018), Volume 2: Short Papers*, pages 120–125, Melbourne, 2018.
- Kristina Toutanova, Dan Klein, Christopher D. Manning, and Yoram Singer. Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of HLT-NAACL 2003*, pages 252–259, Edmonton, 2003.