

Clustering of Analogies for Inter-Language Similarities

Software project - 5th presentation

Justine Diliberto, Cindy Pereira, Anna Nikiforovskaja

Université de Lorraine, IDMC

10.12.2021



UNIVERSITÉ
DE LORRAINE



Institut des
sciences du Digital
Management & Cognition

Summary of the project

Subject: Analogies between morphological rules

Main goal: find out about the closeness of languages and if they have common rules

Final product: Predict if two languages will transfer well, based on the rules they share

What was done before

- Listed possibly close and far pairs of languages
 - Started to study some rules in these pairs
 - Trained and ran the model on 6 languages
-
- Looked for new ways to simplify rule extraction
 - Experimented with 2 tools
 - Trained the multilingual model on old data and transferred to the new data

Rule extraction - issues

- Very time consuming
- Complex to do manually

n ⁵ ndʔɛʔ ¹²	bá ⁵ tʔɛʔ ¹² há ³¹	V;IPFV;PL;1+EXCL;PRS
n ⁵ ndí ⁵ tzhwenʔ ³	n ⁵ ndí ⁵ tzhwenʔ ³	V;IRR;SG;3;FUT
n ⁵ ndí ⁵ khé ⁵	kwí ⁵ khé ⁵ há ³	V;SBJV;SG;1
n ⁵ ndí ⁵ chiʔ ³⁵	chi ⁵³	V;IPFV;SG;1;PRS
n ⁵ ndí ⁵ ʔ ⁵ han ³	kwí ⁵ ʔ ⁵ han ³ ʔ ^{u3}	V;SBJV;SG;2
n ⁵ ndí ³ nkiʔ ³ xken ³	n ⁵ ndí ³ nkiʔ ³ nkên ^{1o3}	V;IRR;PL;3;FUT
n ⁵ ndí ⁵ hnduáʔ ⁵	tí ⁵ hnduáʔ ⁵	V;PFV;SG;3;PST
n ⁵ ndí ⁵ tzhaʔ ⁵³ xken ³	to ³ ndí ⁵ tzhaʔ ⁵³ nkên ³¹	V;PROG;PL;1+EXCL;PST
n ⁵ ndyio ³ xú ⁵	tyio ³ ʔ ^{o3} xú ⁵	V;PFV;PL;3;PST
n ⁵ ndí ⁵ be ³⁴ n ⁵ no ³	n ⁵ ndí ⁵ be ³⁴ n ⁵ no ⁵³	V;IRR;SG;1;FUT

Figure: Some morphological inflections of Amuzgo

Rule extraction - ideas

- Heard about ALEA and Lepage during meeting
- Read papers about morphological rules extraction
- Found 2 tools to try:
 - AutoLEX webpage
 - NLG module (Lepage)

AutoLEX Project:

- Create descriptive grammars automatically
- Focus on "agreement" rules between head and dependent token
- All languages from Universal Dependencies
- Still in progress

AutoLEX Language Descriptions Explorer ¹

Output: agreement + case marking + word order

Early results:

- Only 8 languages are completely extracted
- 14 languages from our dataset are partially done
- Rest is to be filled

¹<https://neulab.github.io/autolex/index.html>

Rule extraction - AutoLEX

Examples of rules on the Finnish language:

```
Tense relation,head-pos,child-pos
conj,VERB,VERB conj,AUX,VERB conj,VERB,AUX conj,AUX,AUX mod,AUX,VERB mod,VERB,VERB
mod,VERB,AUX mod,AUX,AUX
```

Figure: Some agreement rules about the Tense

```
Person relation,head-pos,child-pos
conj,VERB,VERB subj,VERB,PRON conj,AUX,AUX subj,AUX,PRON conj,AUX,VERB parataxis,AUX,VERB
parataxis,AUX,AUX mod@relcl,VERB,VERB comp:aux,AUX,VERB mod@relcl,VERB,AUX
```

Figure: Some agreement rules about the Person

- Few data available
- About word relations
- Maybe combine with NLG?

Rule extraction - NLG

Creates analogical grids:

walk : walk^s : walk^{ing} : walk^{ed}
show : show^s : show^{ing} : show^{ed}
open : open^s : open^{ing} :
study : : study^{ing} :
play : : play^{ing} : play^{ed}

Figure: Example of an analogical grid obtained by NLG - example taken from their article

Rule extraction - NLG

- Creates vectors with morphosyntactic description as features

$n^5ndí^5chi^?^3^5$ chi^5^3 $V;IPFV;SG;1;PRS$

Figure: Example of morphosyntactic description in Amuzgo

- Extracts analogical rules using these vectors

$kwí^5ntyén^5?u^3 : kwí^5sì^1?u^3 :: to^3ndí^5ntyén^5hâ^3^1 : to^3ndí^5sì^1hâ^3^1$
 $kwí^5ndà^?^1hâ^3^1 : n^5ngo^3ndí^5ndà^?^1o^3 :: kwí^5kí^5chì^1hâ^3^1 : n^5ngo^3ndí^5kí^5chì^1?o^3$
 $to^3nchhe^3o^?^3 : to^3hndò^?^1^2b?i^1^2 :: nchhe^3o^?^3 : hndò^?^1^2b?i^1^2$

Figure: Example of analogical rules in Amuzgo

Rule extraction - NLG

Problem: not always readable and easy to analyse

- Print analogical rules with other format

```
n5ngo3nchhe3      : ki3nchhe3  
n5ngo3ndui3 4?u3   : ki3ndui3 4?u3  
n5ngo3ndui?o?o?   : ki3ndui?o?o?  
n5ngo3tyio3 4     : ki3tyio3 4
```

- Better but still some problems
 - Many gaps in these rules
 - Need manual analysis and comparison

Running multilingual model

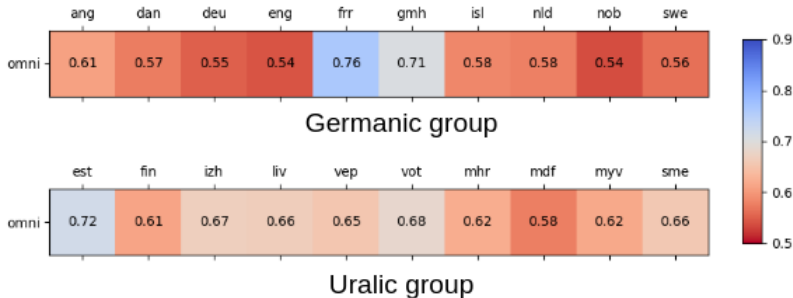
- Multilingual model trained on old data:
 - Arabic, Finnish, Georgian, German, Hungarian, Japanese, Maltese, Navajo, Russian, Spanish, Turkish
- Transfer to new languages.
 - Germanic group, Uralic group

Multilingual model: transfer results

Full transfer, 50000 analogies.

$$F_1 = 2 \cdot \frac{p \cdot n}{p + n}, \text{ where:}$$

- p — accuracy on positive analogies
- n — accuracy on negative analogies



Observations of the transfer

Three languages with the best results:

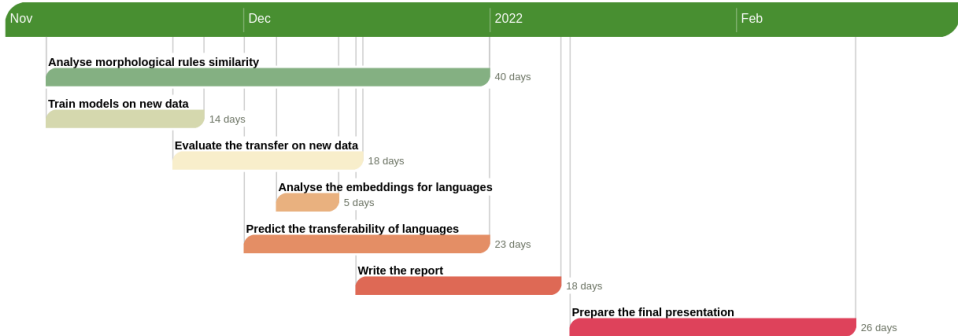
- North Frisian – small dataset
- Middle High German – small dataset
- Estonian
 - Close to Finnish from the old dataset
 - Why did not transfer well to the new Finnish?
 - New Finnish dataset: many adjectives
 - Estonian: only nouns and verbs
 - Old Finnish dataset: mainly noun and verbs

Observations of the transfer

In general – Uralic group performed better, than Germanic group of languages.

Probable reason: Finnish and Hungarian in the old dataset, while only German from Germanic group. Data imbalance.

Timeline



Thank you for your attention!