





WORDS BONES GENES TOOLS Tracking Linguistic, Cultural, and Biological Trajectories of the Human Past





A Comparison between Morphological Complexity Measures: Typological Data vs. Language Corpora

Christian Bentz, University of Tübingen Tatyana Ruzsics, University of Zürich Alexander Koplenig, Institute for German Language Mannheim Tanja Samardžić, University of Zürich Contact: chris@christianbentz.de



Introduction

We investigate the degree to which morphological complexity measures are mutually correlated in a sample of more than 500 languages of 101 language families. We use human expert judgements from the World Atlas of Language Structures (WALS), and compare them to four quantitative measures automatically calculated from language corpora.

Measures

Typological

C_WALS: Average complexity value based on up to 28 features on morphology of the World Atlas of Language Structures (Dryer et al. 2013)



f_i= value per feature n=number of features



Depending on the number of WALS features included (1-27), Spearman correlations between C_WALS and the corpus-based measures range from 0.3 to 0.9 (figure above)

The correlations hold between different language families (below)

(Cor : 0.437	Cor : 0.362	Cor : 0.318	Cor : 0.402
2-	Atlantic-C.: 0.604	Atlantic-C.: 0.567	Atlantic-C.: 0.54	Atlantic-C.: 0.632
ALE	Austron.: 0.512	Austron.: 0.273	Austron.: 0.375	Austron.: 0.498
5 1-	Indo-Euro.: 0.349	Indo-Euro.: 0.369	Indo-Euro.: 0.133	Indo-Euro.: 0.357
0-	Other: 0.348	Other: 0.315	Other: 0.256	Other: 0.307
0.8 -		Cor : 0.876	Cor : 0.799	Cor : 0.918
0.6 -		Atlantic-C.: 0.882	Atlantic-C.: 0.897	Atlantic-C.: 0.941
0.4		Austron.: 0.672	Austron.: 0.732	Austron.: 0.887
		Indo-Euro.: 0.865	Indo-Euro.: 0.765	Indo-Euro.: 0.838

Corpus-based

C_H: The word entropy as calculated from parallel texts (Mayer & Cysouw, 2014)

C_D: The difference in character entropy before and after word internal regularities have been masked (Koplenig et al., forthcoming)

C A: normalized difference in word alignments from a fixed source language to a target language



OneToOne

#ManyToOne - #OneToMany

#AllAlignments

OneToMany

because

potomu chto



Implications

1. All corpus-based automated measures display strong correlations between each other, i.e. strong agreement on which languages are morphologically complex. This is the case despite the conceptual differences between automated methods.

C_TTR: type-token ratio for parallel texts



ManyToOne

will make

V= number of word types fr= token frequency of ith word type

2. Given enough feature values, the expert judgements of the WALS also converge with the automated corpus-based methods. If our sole objective is to rank languages on a morphological complexity scale, then automated methods can support human expert rating.

References

Matthew S. Dryer & Martin Haspelmath. 2013. World Atlas of Language Structures online. Max PlanckDigital Library, Munich. Thomas Mayer and Michael Cysouw. 2014. Creating a massively parallel bible corpus. In Nicoletta Calzolari et al., Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014), Reykjavik, Iceland, May 26-31. Alexander Koplenig, Peter Meyer, Sascha Wolfer, and Carolin Mueller-Spitzer. 2016. The statistical tradeoff between word order and word structure: large-scale evidence for the principle of least effort. ArXiv *preprint,* arXiv:1608.03587.

