

Title: Dissecting the comparative method - comparing traditional reconstruction to computational approaches

Reconstruction is an essential part of historical linguistics. The traditional comparative method (CM) infers states in proto-languages by three core principles: assuming the fewest changes on the tree, plausibility changes and plausibility of combinations of features in proto-languages. However, there is no clear consensus on how to weigh the principles against each other, and many studies are not fully transparent regarding how precisely the principles are applied.

I aim to better understand the comparative method in general and Oceanic proto-language grammar in particular by comparing the reconstruction of grammar by the CM (extracted from Lynch et al 2002, Ross 2004 and more) to computational reconstructions with explicit and transparent mechanisms: **Maximum Parsimony** (MP) and **Marginal Maximum Likelihood** (ML), as implemented in the *R*-packages *castor* and *corHMM* (Louca & Doebeli 2017 and Beaulieu et al 2017). MP is equivalent to CM in that it infers the fewest amounts of changes along the tree. ML on the other hand infers rates of change based on the distribution of values and takes into account branch lengths. Neither method takes into account the plausibility principles, in part because the principles are neither formalised nor sufficiently agreed upon. In addition, I reconstruct proto-language grammar based on which state is the **Most Common** (MC) in the daughter languages, regardless of the tree structure (c.f. Carling & Cathcart 2021, Goldstein 2022). MC, like MP, is an overly simplistic method: while MP doesn't take into account branch lengths and always assumes the slowest rate of change is the most likely, MC ignores tree structure altogether. *A priori*, ML should be the preferred method, but MP should be the most similar to CM.

The results show that CM produces reconstructions that are most similar to those generated by MP or MC. This suggests that, at least for this sample of languages and features, CM does not appear to take branch lengths into account. The methodological drawbacks of using MP and MC imply that we should possibly re-evaluate the way reconstruction is done in historical linguistics and take into account insights from new approaches, such as ML which makes sound assumptions which are likely to give better estimations of the state of proto-languages.

Table 1: Summary table of scores of agreement with CM (numbers range between 0 and 1, with 1 indicating a higher agreement with the results produced by CM).

| | Maximal Parsimony (MP) | | | Maximal Likelihood (ML) | | | Most common in subgroup (MC) |
|-------------|------------------------|------------------|------------------------|-------------------------|------------------|------------------------|------------------------------|
| | Glottolog 4.0 | Gray 2009 (MCCT) | Gray 2009 (posteriors) | Glottolog 4.0 | Gray 2009 (MCCT) | Gray 2009 (posteriors) | |
| Concordance | 0.887 | 0.861 | 0.852 | 0.87 | 0.839 | 0.874 | 0.891 |
| F1-score | 0.882 | 0.842 | 0.853 | 0.859 | 0.817 | 0.86 | 0.85 |

References

- Carling, G., & Cathcart, C. (2021). Reconstructing the evolution of Indo-European grammar: Supplementary material. *Language*, 97(3).
- Goldstein, David (2022). There's no escaping phylogenetics. *Festschrift*.
- Blust, Robert A. 2009. *The Austronesian Languages*. Canberra: Pacific Linguistics.
- Blust, Robert A. & Victoria Chen. 2017. The Pitfalls of Negative Evidence: 'Nuclear Austronesian', 'Ergative Austronesian', and Their Progeny. *Language and Linguistics* 18(4). 577–621.
- Blust, Robert. 2014. Some Recent Proposals Concerning the Classification of the Austronesian Languages. *Oceanic Linguistics* 53(2). 300–391.
- Hammarström, Harald, Robert Forkel & Martin Haspelmath. 2019. Glottolog/Glottolog: Glottolog Database 4.0. URL <https://doi.org/10.5281/zenodo.3260726>.
- Gray, R.D., A.J. Drummond & S.J. Greenhill. 2009. Language Phylogenies Reveal Expansion Pulses and Pauses in Pacific Settlement. *Science* 323. 479–483.
- Beaulieu, Jeremy M., Jeffrey C. Oliver & Brian O'Meara. 2017. Package 'corHMM'. URL <https://CRAN.R-project.org/package=corHMM>.
- Louca, Stilianos & Michael Doebeli. 2017. Efficient Comparative Phylogenetics on Large Trees. *Bioinformatics* 34(6). 1053–1055.
- Lynch, John, Malcolm Ross & Terry Crowley. 2002. Internal Subgrouping. In John Lynch, Malcolm Ross & Terry Crowley (eds.), *The Oceanic Languages Curzon Language Family Series*, 92–120. Richmond: Curzon.
- Ross, Malcolm D. 2004. The Morphosyntactic Typology of Oceanic Languages. *Language and Linguistics* 5(2). 491–541.