

Growing Trees from Big Data

Bayesian Phylogeny for Historical Linguistics

Gereon Kaiping

2017-07-18

1 Solutions to All Your Problems!

- Crunching Numbers
- Bayesian ...
- ... Phylogenetics

2 [Examples]

- Austronesian: Branches and times
- Bantu: Phylogeography
- Indo-European: Ancient written sources

3 And why actually not.

4 Conclusions

- Further Reading

Problem

Using the comparative method is hard and limited, because

- it is a lot of painstaking work,
- we don't know how to weigh the evidence,
- is a mix of hypothesis generation and validation,
- loan words and chance resemblances make our lives more difficult,
- cognates may have changed meanings (but how far?).

And then it doesn't even give us dates, just "not before" or "not after" if we are lucky.

Problem

Using the comparative method is hard and limited, because

- it is a lot of painstaking work,
- we don't know how to weigh the evidence,
- is a mix of hypothesis generation and validation,
- loan words and chance resemblances make our lives more difficult,
- cognates may have changed meanings (but how far?).

And then it doesn't even give us dates, just "not before" or "not after" if we are lucky.

Solution

Tree reconstruction methods from **Bioinformatics**

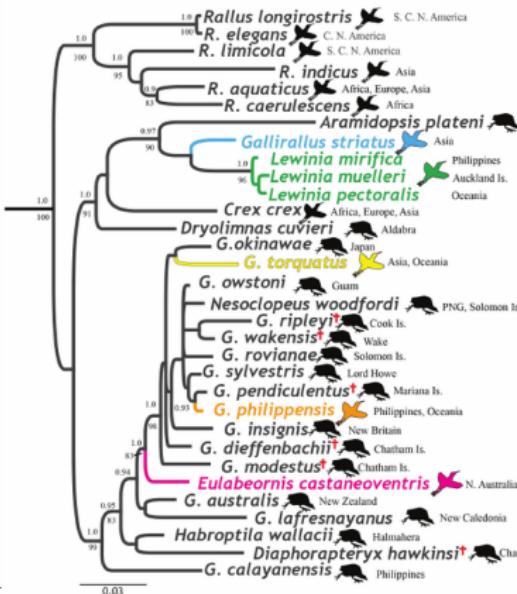
cctccacgcc aacggagcct cattcttctt
cctccacgcc aacaaagcct cattctt---
cctccacgcc aacaaaggcct cattcttctt →
cctccacgcc aacggagcct caggtgtctt
cctccactcc aacggagcct caggtgtctt

¹Garcia-Ramirez et al. 2015

Solution

Tree reconstruction methods from Bioinformatics

```
cctccacgcc aacggagcct catttttctt  
cctccacgcc aacaaaggct catttt---  
cctccacgcc aacaaaggct catttttctt →  
cctccacgcc aacggagcct caggtgtctt  
cctccactcc aacggagcct caggtgtctt
```

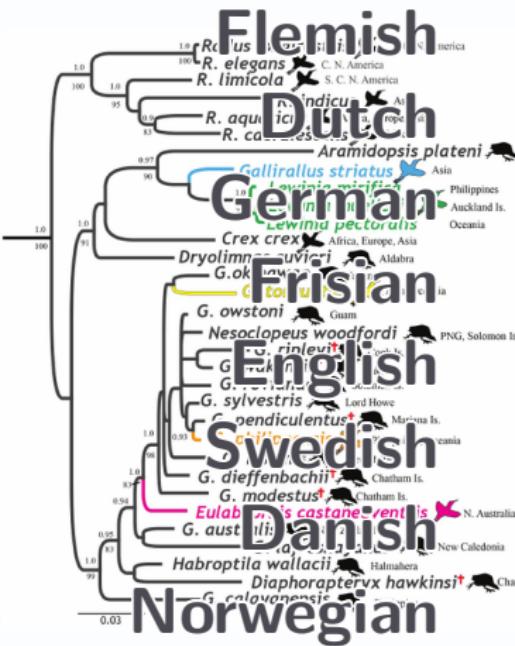


¹Garcia-Ramirez et al. 2015

Solution

Tree reconstruction methods from Bioinformatics

h o n - e r - t _ v c u - t
 h o n d e - - t _ v c : - - t
 h o n d e - - t _ β æ - - t →
 h e n d - u e d _ w ɔ: - - d
 h o n d - r a - - o: - - d



Evolution as a random process

Idea Evolution = a random process on a tree².

Generate a tree and some language history using dice rolls.



Example

- Start with TATATA
- Roll a dice: On a 5 or 6, split the tree here.
- Roll a dice: Replace that letter in the word with E.

“Likelihood”: $P(\text{Data} \mid \text{Model})$

²or a network or a population, as long as we can formalize it.

Evolution as a random process

Idea Evolution = a random process on a tree².

Generate a tree and some language history using dice rolls.



Example

- Start with TATATA
- Roll a dice: On a 5 or 6, split the tree here.
- Roll a dice: Replace that letter in the word with E.

“Likelihood”: $P(\text{Data} \mid \text{Model})$

²or a network or a population, as long as we can formalize it.

Evolution as a random process

Idea Evolution = a random process on a tree².

Generate a tree and some language history using dice rolls.



Example

- Start with TATATA
- Roll a dice: On a 5 or 6, split the tree here.
- Roll a dice: Replace that letter in the word with E.

“Likelihood”: $P(\text{Data} \mid \text{Model})$

²or a network or a population, as long as we can formalize it.

Going back from data

Is given data compatible with this model?

- TATETA
- TATATA
- TATUTA
- TATETA, TAEATA
- TAEETA, TETEEA, EATEEA

How compatible is this data with this or that model?

Probabilities = confidence of belief. Not: repeatable random experiment.

“Posterior probability”: $P(\text{Model} \mid \text{Data})$

Going back from data

Is given data compatible with this model?

- TATETA
- TATATA
- TATUTA
- TATETA, TAEATA
- TAEETA, TETEEA, EATEEA

How compatible is this data with this or that model?

Probabilities = confidence of belief. Not: repeatable random experiment.

“Posterior probability”: $P(\text{Model} \mid \text{Data})$

Going back from data

Is given data compatible with this model?

- TATETA
- TATATA
- TATUTA
- TATETA, TAEATA
- TAEETA, TETEEA, EATEEA

How compatible is this data with this or that model?

Probabilities = confidence of belief. Not: repeatable random experiment.

“Posterior probability”: $P(\text{Model} \mid \text{Data})$

Going back from data

Is given data compatible with this model?

- TATETA
- TATATA
- TATUTA
- TATETA, TAEATA
- TAEETA, TETEEA, EATEEA

How compatible is this data with this or that model?

Probabilities = confidence of belief. Not: repeatable random experiment.

“Posterior probability”: $P(\text{Model} \mid \text{Data})$

Going back from data

Is given data compatible with this model?

- TATETA
- TATATA
- TATUTA
- TATETA, TAEATA
- TAEETA, TETEEA, EATEEA

How compatible is this data with this or that model?

Probabilities = confidence of belief. Not: repeatable random experiment.

“Posterior probability”: $P(\text{Model} \mid \text{Data})$

Going back from data

Is given data compatible with this model?

- TATETA
- TATATA
- TATUTA
- TATETA, TAEATA
- TAEETA, TETEEA, EATEEA

How compatible is this data with this or that model?

Probabilities = confidence of belief. Not: repeatable random experiment.

“Posterior probability”: $P(\text{Model} \mid \text{Data})$

Going back from data

Is given data compatible with this model?

- TATETA
- TATATA
- TATUTA
- TATETA, TAEATA
- TAEETA, TETEEA, EATEEA

How compatible is this data with this or that model?

Probabilities = confidence of belief. Not: repeatable random experiment.

“Posterior probability”: $P(\text{Model} \mid \text{Data})$

Bayes' Theorem

$$P(\text{Model} \mid \text{Data}) \propto P(\text{Data} \mid \text{Model}) \times P(\text{Model}) \quad (3)$$

“What did the language history look like?”

=

“What trees are compatible with the data and my idea of language change?”

=

“Weighted by how ‘strange’ they are, how well does each tree explain my data?”

Bayesian inference may look complicated, but it is

- model-based
- can incorporate prior knowledge
- outputs result uncertainty
- gives implicit weights from first principles



Bayes' Theorem

$$P(\text{Model} \mid \text{Data}) \propto P(\text{Data} \mid \text{Model}) \times P(\text{Model}) \quad (3)$$

“What did the language history look like?”

=

“What trees are compatible with the data and my idea of language change?”

=

“Weighted by how ‘strange’ they are, how well does each tree explain my data?”

Bayesian inference may look complicated, but it is

- model-based
- can incorporate prior knowledge
- outputs result uncertainty
- gives implicit weights from first principles



Computational Phylogenetics

“Roll dice to generate trees, but only keep the good ones”

Need:

- simple stochastic model(s) of language evolution, with parameters
- intuition (“prior”) of how parameters look like
- large dataset of model-compatible data

Data:

- Swadesh lists: models based on semantic change (like Glottochronology – but without its problems)
- Geography: various models
- Phonetic alignments: still in infancy
- Typological data: Some approaches, problems with universals/pathways
- Morphosyntax: ??????

Computational Phylogenetics

“Roll dice to generate trees, but only keep the good ones”

Need:

- simple stochastic model(s) of language evolution, with parameters
- intuition (“prior”) of how parameters look like
- large dataset of model-compatible data

Data:

- Swadesh lists: models based on semantic change (like Glottochronology – but without its problems)
- Geography: various models
- Phonetic alignments: still in infancy
- Typological data: Some approaches, problems with universals/pathways
- Morphosyntax: ??????

Computational Phylogenetics

“Roll dice to generate trees, but only keep the good ones”

Need:

- simple stochastic model(s) of language evolution, with parameters
- intuition (“prior”) of how parameters look like
- large dataset of model-compatible data

Data:

- Swadesh lists: models based on semantic change (like Glottochronology – but without its problems)
- Geography: various models
- Phonetic alignments: still in infancy
- Typological data: Some approaches, problems with universals/pathways
- Morphosyntax: ??????

So much the theory.

- Austronesian: Branches and times
- Bantu: Phylogeography
- Indo-European: Ancient written sources

Example 1: Austronesian

abvo.org ~ Austronesian ~ Crown ~ (Beta)

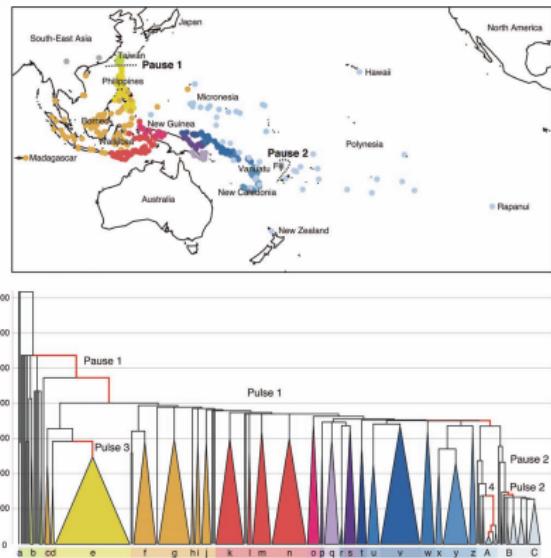
Austronesian Basic Vocabulary Database

Word: hand

Entries for "hand":

ID	Language	Item	Annotate	Category	Classification	Lean
Hand						
217959	Nevali (Bengt)	re-varas				
191500	Ngaju (Djambi)	re-varas				
33064	Mentawai Telingan	kamay				
259818	Proto-Oceanic	*kamay				
227485	Osage (Kaweenah)	inseŋat				
217958	Nevali (Bengt)	re-varas				
20936	Proto-Mon-Khmer	*t*ka?				
209327	Proto-Mon-Khmer	*t*ka?				
186016	Chewung	caet				
186017	Chewung	caet				
177236	Hut	tut				
289824	Monokorene, Cebuano	tae				
246690	Hang (Amping)	teris				
205324	Malagasy	tsia				
204427	Batak	tu				
247574	Buginese (Nalati)	taess				
247575	Buginese (Makassar)	taess				
123700	Malagasy	tae				
308716	Sarawak	ta				
205	Proto-Austronesian	*tae?ima	1	Austronesian		
184446	Proto-Austronesian	*kanday	1, 2	Austronesian		
384949	Proto-Austronesian (Crown)	*g*la-irang	1, 2	Austronesian		
216	Malay - Cleft USA	pasar	2	K-A-K-CVf		
109621	Malay - Cleft USA	qala?	2	K-A-K-CVf		
205323	Malay - Cleft USA	qla?	2	K-A-K-CVf		
71434	Malay - Soqali Fath	qba?	2	K-A-K-Sqaf		
71425	Malay - Soqali Fath	rapa?	2	K-A-K-Sqaf		
71426	Malay - Soqali Fath	kaas	2	K-A-K-Hqk		
71427	Malay - Soqali Fath	ketem	2	K-A-K-Hqk		
71428	Malay - Soqali Fath	avet	2	K-A-K-Sqaf		
205324	Malay - Cleft USA	baap?	2	K-A-K-Sqaf		
235	Sundanese (Sukarela)	bawo?	51	K-A-S		
207320	Sundanese (Pasean)	bawo?	51	K-A-S		
207320	Sundanese (Tatah)	bawo	51	K-A-S		
207320	Sundanese (Husua)	bawo	51	K-A-S		
207320	Sundanese (Husua)	bawo	51	K-A-S		
71419	Burun Pit, Southern	inser	1	A-Burun		
71420	Burun Pit, Southern	lapald		A-Burun		
38630	Burun (Tabukulan LBS)	reñat?	1	A-Burun		
203720	Burun (Tabukulan LBS)	reñat?	1	A-Burun		
203946	Burun (Tabukulan)	reña	1	A-Burun		
204355	Burun (Tabukulan)	reñe	1	A-Burun		
60822	Arme (Central)	kañay	13	K-S-FitS		
203946	Burun (Tabukulan)	reñay	1	A-Burun		
203946	Burun (Tabukulan)	reñay	1	A-Burun		
203941	Arme (Central)	zulma?	1	A-Z-MitS		
203941	Arme (Central)	zulma?	1	A-Z-MitS		
210373	Reuey (Yehmo)	chrea	1	A-Z-Densey		
210373	Reuey (Yehmo)	chrea	1	A-Z-Densey		
210595	Reuey (Yehmo)	chrea	1	A-Z-Densey		
63899	Kawaleo (Fid)	lins?	1	K-S-N-Kawaleo		
63899	Kawaleo (Fid)	gralp	3	K-A-N-Kawaleo		

Example 1: Austronesian



"The invention of the outrigger canoe and its sail may have enabled the Austronesians to move across this channel before spreading rapidly over the 7000 km from the Philippines to Polynesia (4). This is supported by linguistic reconstructions showing that the terminology associated with the outrigger canoe complex can only be traced back to Proto-Malayo-Polynesian and not Proto-Austronesian (41)."

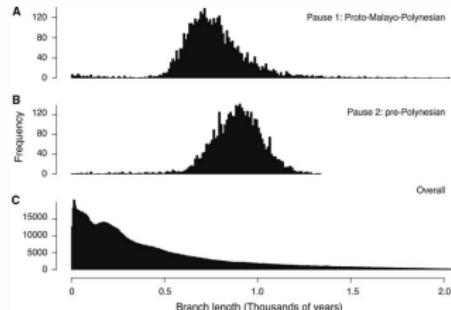


Fig. 3. Histograms of the branch length distributions. (A) The distribution of the Proto-Malayo-Polynesian pause, (B) the distribution of the pre-Polynesian pause, and (C) the overall branch-length distribution.

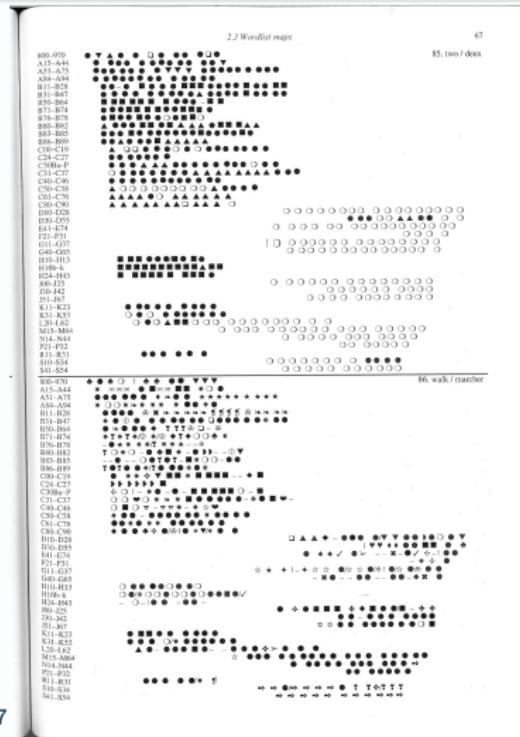
⁶Gray, Drummond & Greenhill 2009

Example 1: Austronesian – Critique

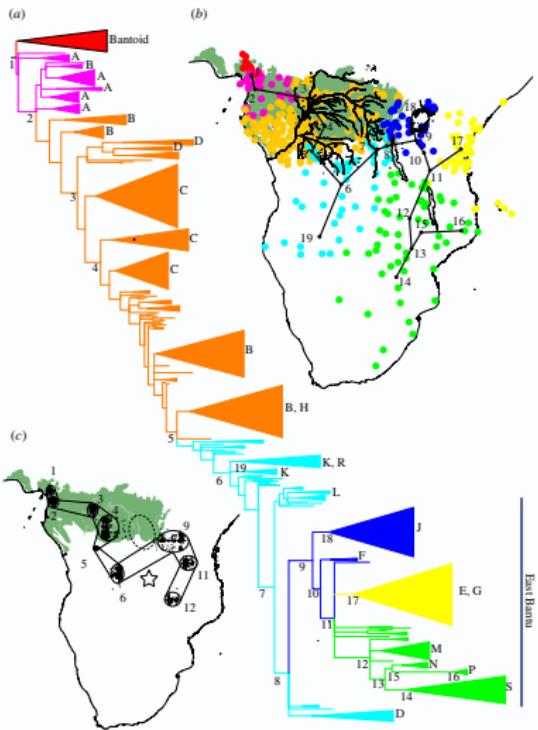
- Pauses and pulses appear with high posterior probability
- Prior? Do the results follow from data or original guess?
- Some subgroupings not linguistically supported – Data contains sociogeography
- How realistic is binary covarion?

Example 2: Bantu

- 2908 cognate classes for 90 meanings in 542 varieties of Bantu/Bantoid, with geographical point-coordinates
 - Binary covariation with 6 (empirical) rate categories
 - Brownian motion ancestral state reconstruction of latitudes and longitudes on 500 best trees
 - Branch-dependent speed of movement and lexical change
 - Other statistical analyses



Example 2: Bantu



⁸Currie et al. 2013

"These debates have implications regarding the origin and spread of important cultural innovations, such as metallurgy and cattle-keeping."

"[...] explicit mapping of ancestral locations to make inferences about the specific route taken during the dispersal of Bantu languages. The results clearly support the 'pathway through the rainforest' scenario for the expansion of Bantu through much of sub-Saharan Africa. There is no support in these analyses for an early, deep split between East and West Bantu languages and a movement by one branch north of the rainforest."

Example 2: Bantu – Critique

- Several robustness checks of parameters
- Prior?
- How good is Brownian motion as model for language spread?
Language shift and post-split contact might affect geographic inference.
(Though the fundamental results look robust.)

Example 3: Indo-European – Data and Prologue

a 066 HAND

b

066 73 Ossetic
066 59 Gujarati

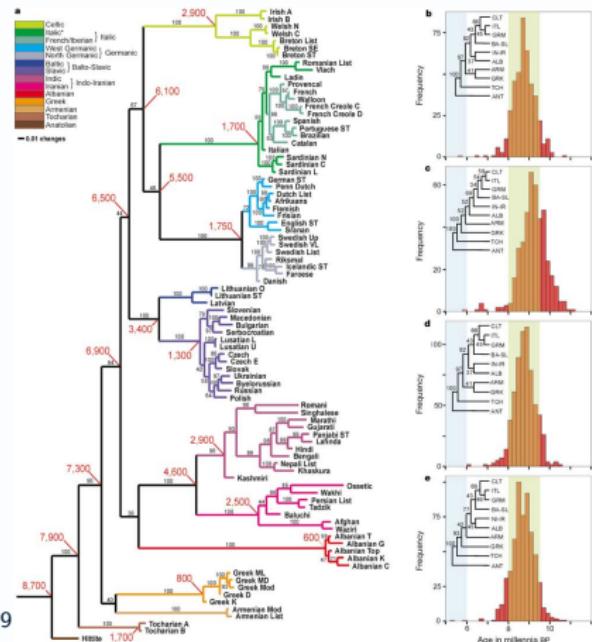
001

K"YX
NATH

b

066 17 Sardinian N	MANU
066 18 Sardinian L	MANU
066 09 Vlach	MYNE
066 22 Brazilian	MAO
066 21 Portuguese ST	MAO
066 15 French Creole C	LAME
066 13 French	MAIN
066 16 French Creole D	LAME
066 14 Walloon	MIN
066 12 Provencal	MAIN
066 20 Spanish	MANO
066 23 Catalan	MA
066 10 Italian	MANO
066 19 Sardinian C	MANU
066 11 Ladin	MAUN
066 08 Rumanian List	MINA

002



⁸Dyen 1997

⁹Gray & Atkinson 2003

Example 3: Indo-European

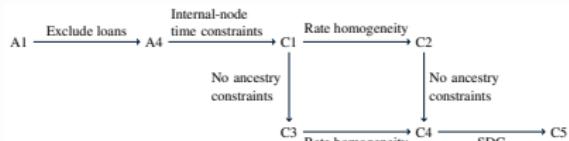


IELex

hand

login

ID	Language	Source Form	Phonological Form	Notes	Cognate Class
114	Proto-Indo-European	*mon-u-			E
114	Proto-Indo-European	*θ̥e̥s-r(o)-,			C
114	Proto-Indo-European	*mar-			E
80	Hittite	keššar			C
133	Luvian	ıssaris			C
134	Lycian	izre			C
81	Tocharian A	tsar			C
82	Tocharian B	ṣar			C
88	Albanian	dorë		A singularised neut. plural PAIb ...	C
143	Standard Albanian	dorë			C
2	Albanian Sicily	dorë		A singularised neut. plural PAIb ...	C
4	Albanian Corinth	dorë		A singularised neut. plural PAIb ...	C
3	Albanian Gheg	dorë		A singularised neut. plural PAIb ...	C
6	Albanian Tosk	dorë		A singularised neut. plural PAIb ...	C
173	Mycenaean Greek	ke ⁰	kʰer-	Attested as an element in ...	C
110	Ancient Greek	χείρ	kʰé:r	G.sg. χειρός	C
152	Tsakonian	χερά			C
32	Greek	χεπλ	'çeri		C
31	Greek Lesbos	CHERI			C
129	Classical Armenian	ձեռն	jeřn		C
8	Armenian Eastern	ձեղ	dzerkʰ		C
7	Armenian Western	ձեղ	ðz'erkʰ		C
10	Avestan	zastō			C



Starting from a replication of previous work¹¹, improve

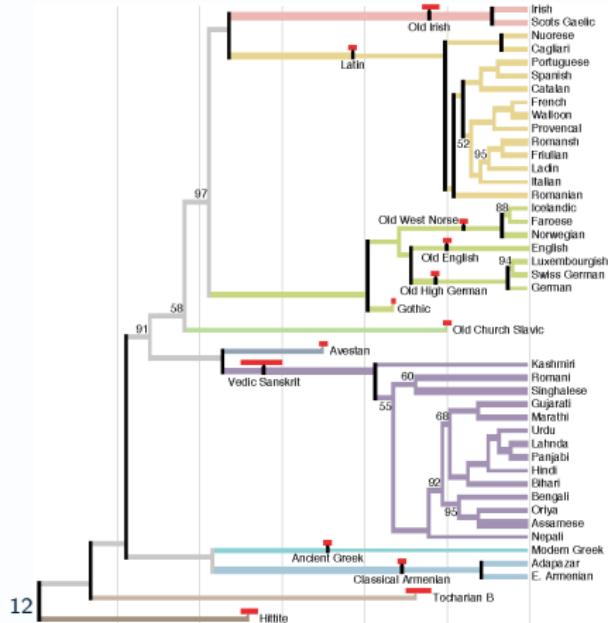
- data
- methodology
- tree prior
- post-processing

comparing each step.

¹⁰Dunn 2015

¹¹Bouckaert et al. 2012

Example 3: Indo-European



"Here we present a phylogenetic analysis in which ancestry constraints permit more accurate inference of rates of change, based on observed changes between ancient or medieval languages and their modern descendants, and we show that the result strongly supports the steppe hypothesis."

"Because previous statistical phylogenetic research supported the Anatolian hypothesis, linguists who find that hypothesis implausible for other reasons may dismiss statistical analyses that purport to determine ancestral chronology. [...] statistical phylogenetic analysis can yield reliable information about pre-historic chronology, at least where all of the available data is taken into consideration."

¹²Chang et al. 2015

Example 3: Indo-European – Critique

- Very explicit about methodology (small steps, driver files available)
- Careful description of data coding
- Would someone have been this careful if the original results *had* matched the linguists' expectations?
- Ancestral constraints are very strong, and somewhat artificial in the model.

It will not solve all problems

Some papers

- disregard prior knowledge
- use models that don't fit their data
- don't show their priors

Not even in a better world.

State of the art models

- can only build trees, no language contact
- only support cognate data (baby steps towards phonetic data and typology), no distinction between innovation and retention
- are not realistic / not calibrated / have biases
- can't actually decide high-level relationships

I think we will always

- need domain experts
- have problems modelling morphosyntax
- have qualitative data that are hard to integrate

Not even in a better world.

State of the art models

- can only build trees, no language contact
- only support cognate data (baby steps towards phonetic data and typology), no distinction between innovation and retention
- are not realistic / not calibrated / have biases
- can't actually decide high-level relationships

I think we will always

- need domain experts
- have problems modelling morphosyntax
- have qualitative data that are hard to integrate

Not even in a better world.

State of the art models

- can only build trees, no language contact
- only support cognate data (baby steps towards phonetic data and typology), no distinction between innovation and retention
- are not realistic / not calibrated / have biases
- can't actually decide high-level relationships

I think we will always

- need domain experts
- have problems modelling morphosyntax
- have qualitative data that are hard to integrate

Not even in a better world.

State of the art models

- can only build trees, no language contact
- only support cognate data (baby steps towards phonetic data and typology), no distinction between innovation and retention
- are not realistic / not calibrated / have biases
- can't actually decide high-level relationships

I think we will always

- need domain experts
- have problems modelling morphosyntax
- have qualitative data that are hard to integrate

Not even in a better world.

State of the art models

- can only build trees, no language contact
- only support cognate data (baby steps towards phonetic data and typology), no distinction between innovation and retention
- are not realistic / not calibrated / have biases
- can't actually decide high-level relationships

I think we will always

- need domain experts
- have problems modelling morphosyntax
- have qualitative data that are hard to integrate

Not even in a better world.

State of the art models

- can only build trees, no language contact
- only support cognate data (baby steps towards phonetic data and typology), no distinction between innovation and retention
- are not realistic / not calibrated / have biases
- can't actually decide high-level relationships

I think we will always

- need domain experts
- have problems modelling morphosyntax
- have qualitative data that are hard to integrate

Conclusions

- Computer models can help make sense of large data sets
- The computer only tests consistency or helps build intuition, it does not replace expertise
- Very few language-appropriate models so far
- Building a *good* inference is hard!
- Mathematical models can handle and combine new types of data for new *types* of results

If you disagree with results, *what parameters or choices do you disagree with?*

Sources and Further Reading I

-  Bastin, Yvonne, André Coupez & Michael Mann. 1999. *Continuity and divergence in the bantu languages: perspectives from a lexicostatistic study.* Musée royal de l'Afrique centrale.
-  Bouckaert, Remco, Philippe Lemey, Michael Dunn, Simon J. Greenhill, Alexander V. Alekseyenko, Alexei J. Drummond, Russell D. Gray, Marc A. Suchard & Quentin D. Atkinson. 2012. Mapping the Origins and Expansion of the Indo-European Language Family. *Science* 337(6097). 957–960. <https://doi.org/10.1126/science.1219669>. <http://www.sciencemag.org/content/337/6097/957> (3 December, 2014).
-  Chang, Will, Chundra Cathcart, David Hall & Andrew Garrett. 2015. Ancestry-constrained phylogenetic analysis supports the Indo-European steppe hypothesis. *Language* 91(1). 194–244. <https://doi.org/10.1353/lan.2015.0005>. <http://www.linguisticsociety.org/files/news/ChangEtAlPreprint.pdf> (27 February, 2015).

Sources and Further Reading II

-  Currie, Thomas E., Andrew Meade, Myrtille Guillon & Ruth Mace. 2013. Cultural phylogeography of the bantu languages of sub-saharan africa. *Proceedings of the Royal Society of London B: Biological Sciences* 280(1762). <https://doi.org/10.1098/rspb.2013.0695>. <http://rspb.royalsocietypublishing.org/content/280/1762/20130695>.
-  Dunn, Michael. 2015. *IELex – Indo-European Lexical Cognacy Database*. <http://ielex.mpi.nl/> (12 March, 2015).
-  Dyen, Isidore. 1997. *COMPARATIVE INDOEUROPEAN DATABASE COLLECTED BY ISIDORE DYEN*. <http://www.wordgumbo.com/ie/cmp/iedata.txt>.
-  Garcia-Ramirez, Juan C, Graeme Elliott, Kath Walker, Isabel Castro & Steven A Trewick. 2015. Trans-equatorial range of a land bird lineage (aves: rallidae) from tropical forests to subantarctic grasslands. *Journal of Avian Biology*.

Sources and Further Reading III

-  Gray, R. D., A. J. Drummond & S. J. Greenhill. 2009. Language phylogenies reveal expansion pulses and pauses in pacific settlement. *Science* 323(5913). 479–483. <https://doi.org/10.1126/science.1166858>. <http://science.sciencemag.org/content/323/5913/479>.
-  Gray, Russell D. & Quentin D. Atkinson. 2003. Language-tree divergence times support the Anatolian theory of Indo-European origin. *Nature* 426(6965). 435–439. <https://doi.org/10.1038/nature02029>. <http://www.nature.com/nature/journal/v426/n6965/abs/nature02029.html> (27 November, 2014).
-  Greenhill, S.J., R Blust & R.D. Gray. 2008. The austronesian basic vocabulary database: from bioinformatics to lexomics. *Evolutionary Bioinformatics*. 271–283.
-  McMahon, April & Robert McMahon. 2005. *Language classification by numbers*. Oxford University Press. 285 pp.

Sources and Further Reading IV

-  Michael, Lev, Natalia Chousou-Polydouri, Keith Bartolomei, Erin Donnelly, Sérgio Meira, Vivian Wauters & Zachary O'Hagan. 2015. A Bayesian Phylogenetic Classification of Tupí-Guaraní. *LIAMES: Línguas Indígenas Americanas* 15(2). 193–221.
<https://doi.org/10.20396/liames.v15i2.8642301>.
<https://periodicos.sbu.unicamp.br/ojs/index.php/liames/article/view/8642301> (29 August, 2017).
-  Verkerk, Annemarie. 2017. Phylogenies: Future, not fallacy. *Language Dynamics and Change* 7(1). 127–140.
<https://doi.org/10.1163/22105832-00601013>.
<http://booksandjournals.brillonline.com/content/journals/10.1163/22105832-00601013> (9 October, 2017).