

Statistical estimation of diachronic stability from synchronic data

Gerhard Jäger

Tübingen University

Cape Town, July 6, 2018



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

ERHARD-SHAW
UNIVERSITÄT
TÜBINGEN



DFG

From the workshop description

“The workshop starts from the null hypothesis that diachronically stable properties are those that appear as the typologically most frequent ones, and that cross-linguistic rarity correlates with diachronic instability.”

Inferring diachronic stability of a feature from its typological frequency is potentially fallacious for three reasons:

1. Processes of different rates may lead to identical equilibrium distributions.
2. Individual languages are not independent random samples, since genetically related languages are likely to have similar typological profiles.
3. The stability of a feature value might depend on the value of other, correlated features.

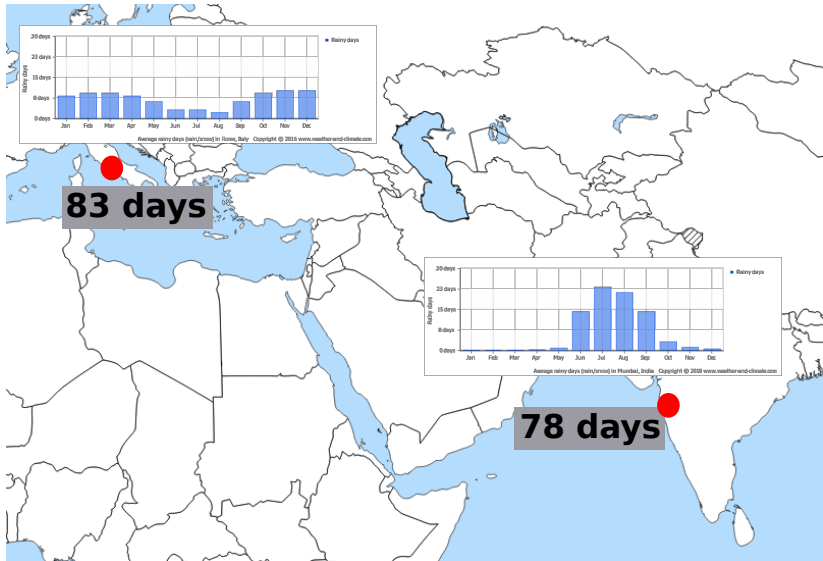
Frequency, stability, and Markov chains

Rainy days per year in Mumbai and Rome



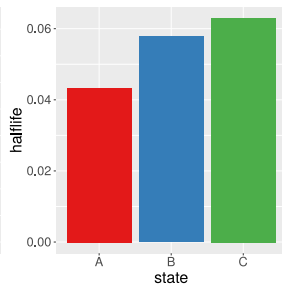
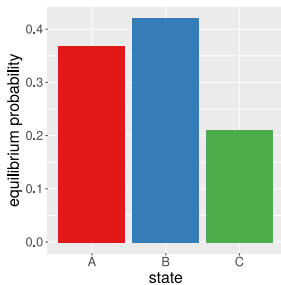
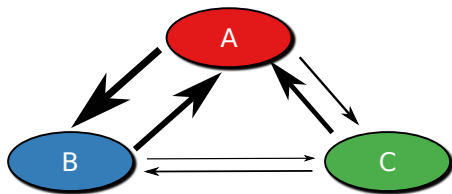
source: <https://weather-and-climate.com>

Rainy days per year in Mumbai and Rome



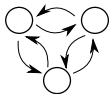
source: <https://weather-and-climate.com>

Markov chains



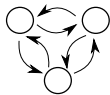
Phylogenetic structure

Markov process



Phylogenetic structure

Markov process

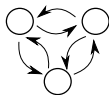


Phylogeny



Phylogenetic structure

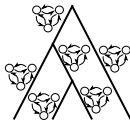
Markov process



Phylogeny

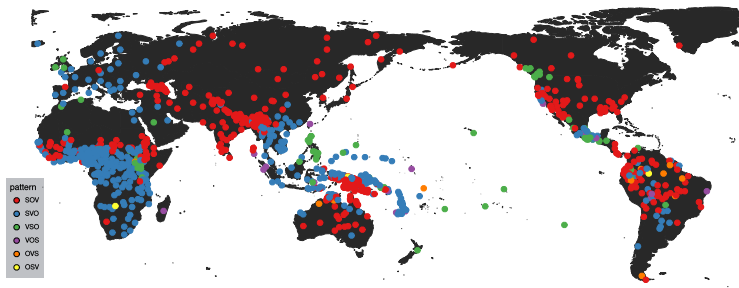


Branching process

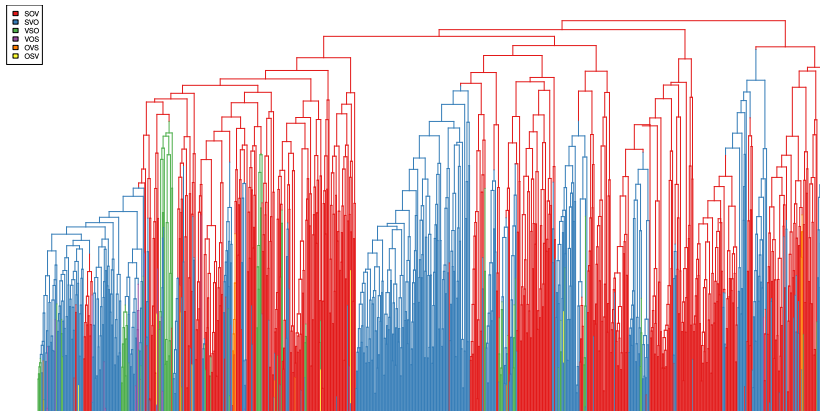


Phylogenetic non-independence

- ▶ languages are phylogenetically structured
 - ▶ if two closely related languages display the same pattern, these are not two independent data points
- ⇒ we need to control for phylogenetic dependencies



Phylogenetic non-independence



Phylogenetic non-independence

Maslova (2000):

“If the A-distribution for a given typology cannot be assumed to be stationary, a distributional universal cannot be discovered on the basis of purely synchronic statistical data.”

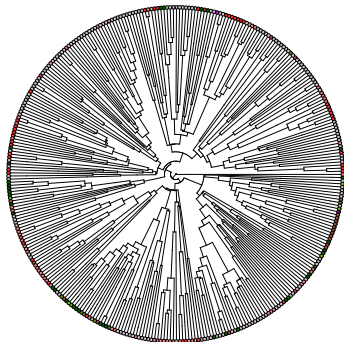
*“In this case, the only way to discover a distributional universal is to **estimate transition probabilities** and as it were to ‘predict’ the stationary distribution on the basis of the equations in (1).”*



The phylogenetic comparative method

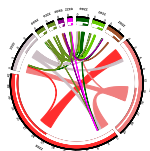
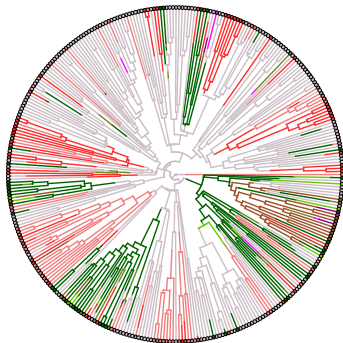
Estimating rates of change

- ▶ if phylogeny and states of extant languages are known...



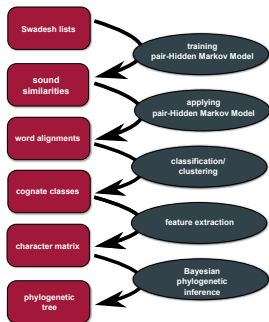
Estimating rates of change

- ▶ if phylogeny and states of extant languages are known...
- ▶ ... transition rates and ancestral states can be estimated based on Markov model

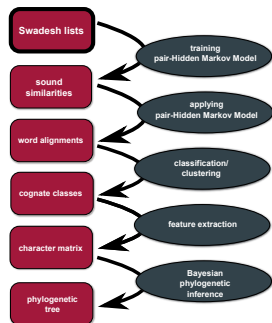


Inferring a world tree of languages

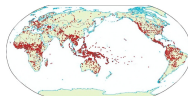
From words to trees



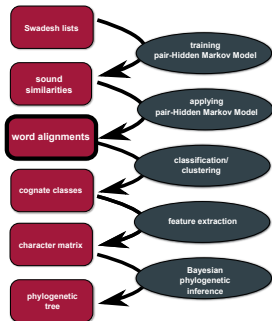
From words to trees



concept	Latin	English
<i>I</i>	ego	Ei
<i>you</i>	tu	yu
<i>we</i>	nos	wi
<i>one</i>	unus	w3n
<i>two</i>	duo	tu
<i>person</i>	persona, homo	pers3n
<i>fish</i>	piskis	fiS
<i>dog</i>	kanis	dag
<i>louse</i>	pedikulus	laus
<i>tree</i>	arbor	tri
<i>leaf</i>	foly~u*	lif
<i>skin</i>	kutis	skin
<i>blood</i>	saNgw~is	bl3d
<i>bone</i>	os	bon
<i>horn</i>	kornu	horn
<i>ear</i>	auris	ir
<i>eye</i>	okulus	Ei

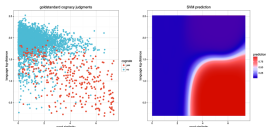
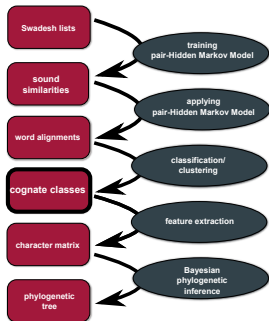


From words to trees



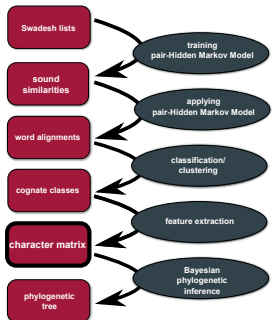
Language	fish:a	tongue:l	snake:l
Abui-Atangmelang	-af-u		
Abui-Fuimelang	-af-u	tal-i-fi-	
Adang	aab--	tal-E-b--	awal--b-a-n-o-7o-
Blagar-Bakalang	-ab--	--j-e-bur-	--ad--b-a-n-aKa-
Blagar-Bama	aab--	teg-e-bur-	-----b-e-n-a-xa-
Blagar-Kuijahi	-ab--	tej-e-bur-	-----b-e-n-aKa-
Blagar-Nule	aab--	tej-e-bur-	--ad--b-e-n-aKa-
Blagar-Tuntuli	aab--	tej-e-bur-	a-adeb-a-n-a-q--
Blagar-Warsalelang	-ab--	tel-e-bur-	a-ad--b-a-n-a-x-
Bunaq			-----b-o-t-o-h-
Deing	haf--		-----buu-n-----
Hamap	7ab--	nar-p-tuU-	-----b-a-n-o-7--
Kabola	hab--	tal-e-b---	awal--b-e-n-e-7o-
Kaera-Padangul	-ab--	tal-e-b---	a-ad--b-e-naa-x-
Kafoa	-afU	tal-i-p---	-----f-o-n-a----
Kamang	-ap-i	nal--p-u-	-----p-u-n-----a-
Kiraman	-Eb--	nal-i-bar-	--ar--b-a-n-o-kan
Klon	-eb-i	gel-E-b--	--ed-ab-o-n-----
Kui	-eb--	tal-i-ber-	--ar--b-o-n-o-k--
Kula	-ap-i	-li-l-p---	-----p--n-eKa-
Nedebang	aaf-i	gel-e-fu-	--ar-ab-o-n-----
Reta	aab--	nal-e-tul-	nal-e-b-o-n-a----
Sar-Adiabang	haf--	--p-e-fal-	--ar--buu-n-----
Sar-Nule	haf--	nal-e-faj-	
Sawila	-ap-i	gal-imuru	-----p-u-n-a-ka-
Teiwa-Madar	xaf--	gel-i-vi-	-----buu-n-----
Wering	-ap-i	nej-e-bur-	--ad--p-u-n-a-k--
Wpantar	hap--	nal-e-tu-	-----b-unn-a----

From words to trees



	English	Spanish	Modern Greek	Standard German
<i>I</i>	Ei:A	yo:B	ego:C	ik:D
<i>you</i>	ya:A	astet:B, ta:C	es:D	du:E
<i>we</i>	wi:A	nosotros:B	esis:C	wir:A
<i>one</i>	w3n:A	uno:B	ena:C, ema:C	eins:D
<i>two</i>	tu:A	dos:B	fy-ou:C, tho:D	zwei:E
<i>person</i>	pez3n:A	persona:A	an3-ropos:B	sthd:C
<i>fish</i>	fi3:A	peakedo:A, pee:A	paari:B	fi3:A
<i>dog</i>	dag:A	perro:B	oTili:C, oTiloo:C	hant:D
<i>come</i>	k3e:A	veni:B	era-3-C	ik3-uch3:A
<i>sun</i>	3n:A	sol:B	ily-ou:C, iloo:C	son3:A
<i>star</i>	star:A	estrella:A	asteri:A, astro:A	StErn:A
<i>water</i>	w33r:A	agu-a:B	sero:C	sero:A
<i>stone</i>	stoo:A	piedra:B	petra:B	Stein:A
<i>fire</i>	f3ir:A	fuego:B	foty-a:C	foia:D
<i>path</i>	p33r:A	senda:B	Bromo3:C	pf-at:A, v3k:D
<i>mountain</i>	mu3nt3n:A	sero:B, mont33n:A	rulo:C, oroo:D	BER3:E
<i>fall</i>	fal:A	yema:B	ymato3-C, plira3:D	fol3:A
<i>tree</i>	tr3:A	nuevo:A	ne33:A, Teary-oo3:B	bol3:A
<i>name</i>	na3:A	nombre:A	onoma:A	nam3:A

From words to trees

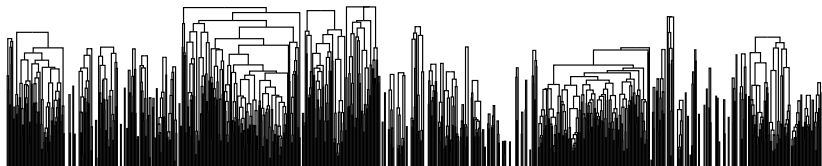


```

TNG. ENGAN. MAIBI
TNG. ENGAN. POLE
TNG. ENGAN. SAU
TNG. ENGAN. YARIBA
TNG. FASU. FASU
TNG. FASU. NAMUMI
TNG. FINISTERRE-HUON. AHARA
TNG. FINISTERRE-HUON. BORONG
TNG. FINISTERRE-HUON. BURUM
TNG. FINISTERRE-HUON. BURUM MIND
TNG. FINISTERRE-HUON. DEDUA
TNG. FINISTERRE-HUON. HUBE
TNG. FINISTERRE-HUON. KATE
TNG. FINISTERRE-HUON. KOMBA
TNG. FINISTERRE-HUON. KOSORONG
TNG. FINISTERRE-HUON. MAPE
TNG. FINISTERRE-HUON. MAPE 2
TNG. FINISTERRE-HUON. MIGABAC
TNG. FINISTERRE-HUON. MINDIK
TNG. FINISTERRE-HUON. MIMOLILI
TNG. FINISTERRE-HUON. NABAK
TNG. FINISTERRE-HUON. NANKINA
TNG. FINISTERRE-HUON. NEK
TNG. FINISTERRE-HUON. NUKNA
TNG. FINISTERRE-HUON. ONO
TNG. FINISTERRE-HUON. SELEPET
TNG. FINISTERRE-HUON. TIMPE
TNG. FINISTERRE-HUON. TOBO
TNG. FINISTERRE-HUON. WANTOAT
TNG. FINISTERRE-HUON. YOPNO
TNG. GOTILALAN. AFOA
TNG. GOTILALAN. KUNIMPAPA
TNG. GOTILALAN. MAFULU
  
```


From tree to forest

- ▶ branch lengths within Glottolog families estimated from lexical data
- ▶ calibration: Proto-Austronesian $\sim 5,000$ years
- ▶ branches above family level effectively set to infinity



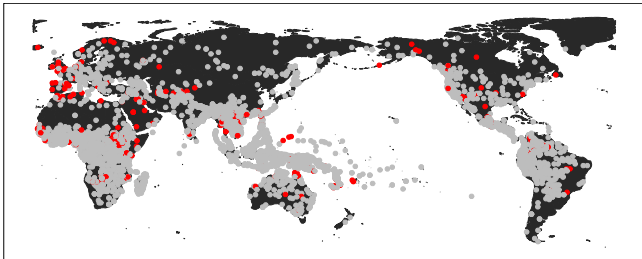
Case study 1: Rare consonants

Synchronic statistics

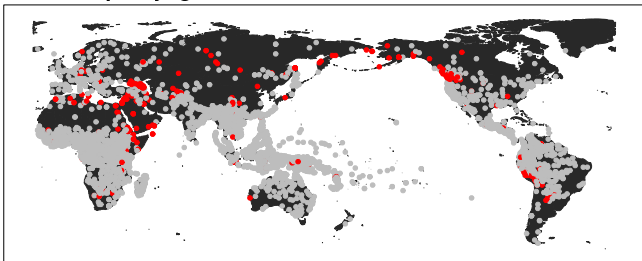
- ▶ data: ASJP word lists (word lists from ca. 6,000 living languages and dialects; Wichmann et al. 2016)
- ▶ variables:
 - ▶ *voiceless and voiced dental fricative* (transcribed as 8)
 - ▶ *voiceless and voiced uvular fricative, voiceless and voiced pharyngeal fricative* (transcribed as X)

	8	X
raw numbers	334	378
average	5.7%	6.6%
weighted by family	14.6	22.2
average	4.6%	7.0%

Dental fricative



Uvular or pharyngeal fricative



Phylogenetic estimates

	8	X
equilibrium probability	5.5%	7.4%
half-life present (kyrs)	1.8	4.6
half-life absent (kyrs)	30.1	58.4

Case study 2: Major word orders

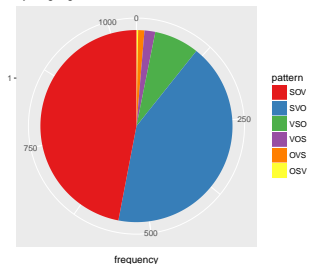
Statistics of major word order distribution

- ▶ data: WALS intersected with ASJP
- ▶ 1,045 languages, 211 lineages

Raw numbers

SOV	SVO	VSO	VOS	OVS	OSV
491	442	79	19	11	3
47.0%	42.3%	7.6%	1.8%	1.1%	0.3%

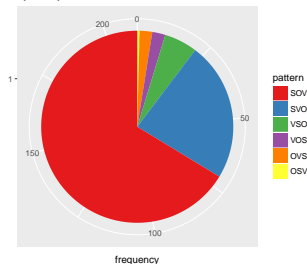
by language



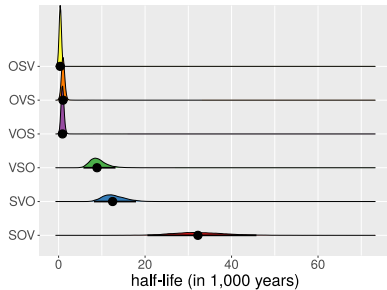
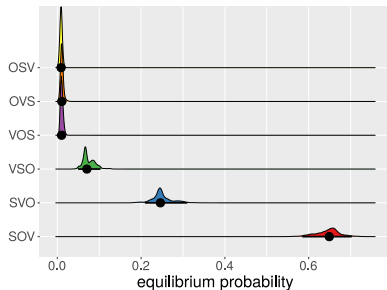
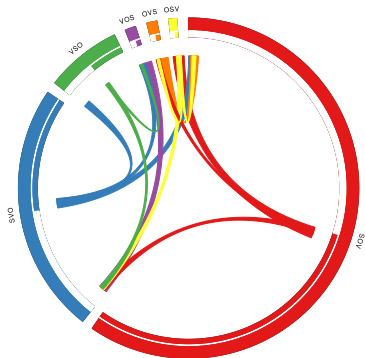
Weighted by lineages

SOV	SVO	VSO	VOS	OVS	OSV
139.1	49.3	11.8	4.7	4.5	0.8
66.3%	23.4%	5.6%	2.2%	2.1%	0.4%

by family



Phylogenetically estimated Markov process



Case study 3: Word order and case

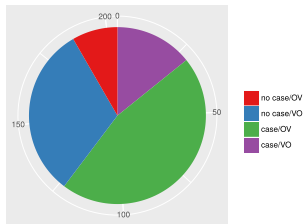
Statistics

- ▶ data: WALS intersected with ASJP
- ▶ 204 languages, 103 lineages

Raw numbers

no case/OV	no case/VO	case/OV	case/VO
17	64	94	29
8.3%	31.4%	46.1%	14.2%

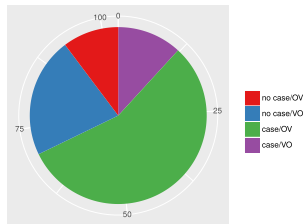
raw frequencies



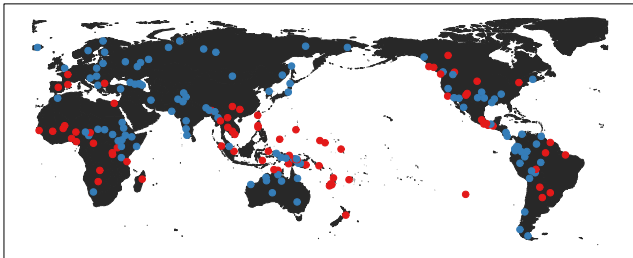
Weighted by lineages

no case/OV	no case/VO	case/OV	case/VO
10.6	22.6	57.7	12.2
10.3%	21.9%	56.0%	11.8%

weighted by family

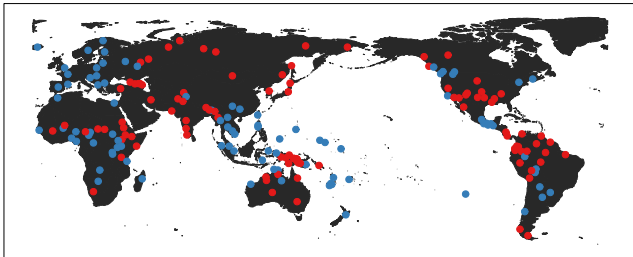


Case



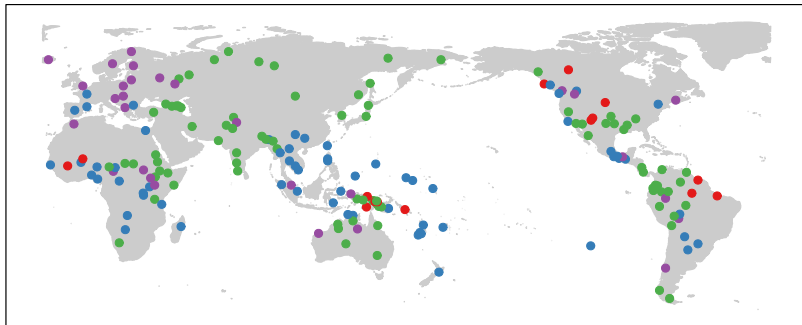
case ● no ● yes

VO vs. OV



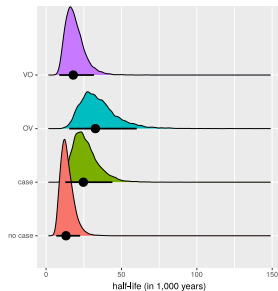
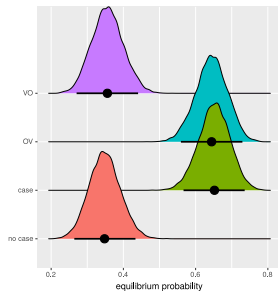
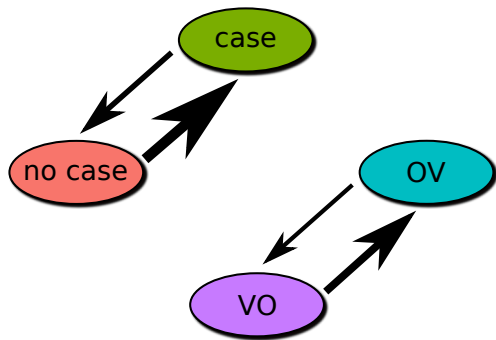
word order ● OV ● VO

case + word order

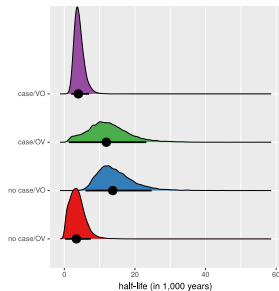
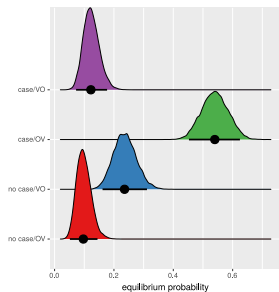
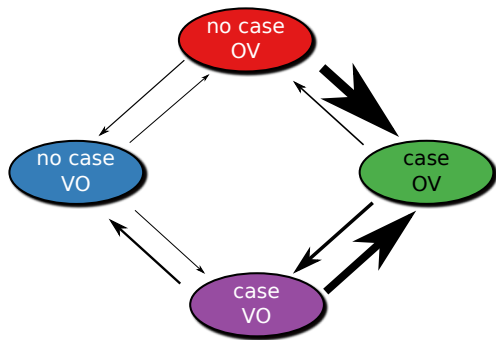


● no case/OV ● no case/VO ● case/OV ● case/VO

Phylogenetically estimated Markov process: features individually



Phylogenetically estimated Markov process: dependent features



Conclusion

Conclusion

- ▶ connection between cross-linguistic frequency and diachronic stability is loose at best
- ▶ to assess diachronic stability, we need information on
 - ▶ phylogenetic structure
 - ▶ branch lengths
- ▶ stability of feature values may depend on other features → potentially complex causal network between typological variables, waiting to be explored
- ▶ todo:
 - ▶ comparison to related but different approaches, such as Bickel's Family Bias Method (Bickel, 2013) or Greenhill et al.'s (2017) approach
 - ▶ factoring in language contact
 - ▶ non-homogeneous Markov chains?

- Balthasar Bickel. Distributional biases in language families. In *Language Typology and Historical Contingency: In honor of Johanna Nichols*, pages 415–444. John Benjamins, Amsterdam, 2013.
- Simon J Greenhill, Chieh-Hsi Wu, Xia Hua, Michael Dunn, Stephen C Levinson, and Russell D Gray. Evolutionary dynamics of language systems. *Proceedings of the National Academy of Sciences*, 114(42):E8822–E8829, 2017.
- Gerhard Jäger. Phylogenetic inference from word lists using weighted alignment with empirically determined weights. *Language Dynamics and Change*, 3(2):245–291, 2013.
- Gerhard Jäger. Support for linguistic macrofamilies from weighted sequence alignment. *Proceedings of the National Academy of Sciences*, 112(41):12752–12757, 2015. doi: 10.1073/pnas.1500331112.
- Gerhard Jäger. Global-scale phylogenetic linguistic inference from lexical resources. arXiv:1802.06079, 2018.
- Gerhard Jäger and Søren Wichmann. Inferring the world tree of languages from word lists. In S. G. Roberts, C. Cuskley, L. McCrohon, L. Barceló-Coblijn, O. Feher, and T. Verhoef, editors, *The Evolution of Language: Proceedings of the 11th International Conference (EVLANG11)*, 2016. Available online: <http://evolang.org/neworleans/papers/147.html>.
- Elena Maslova. A dynamic approach to the verification of distributional universals. *Linguistic Typology*, 4(3): 307–333, 2000.
- Søren Wichmann, Eric W. Holman, and Cecil H. Brown. The ASJP database (version 17). <http://asjp.clld.org/>, 2016.