



THESIS SEMINAR MEETING

Philip Georgis
May 26, 2021

OUTLINE OF THESIS

Research Questions:

1) Which types of distance methods are (most) useful for phylogenetic inference?

- Phonetic methods vs. information theory methods (PMI, surprisal) [+ mixed methods]
- Phonemes vs. sound classes

2) At what stage of phylogenetic inference are they useful?

(i.e. are they useful for cognate detection, or only for evaluating already-detected cognates?)

OUTLINE OF THESIS

(Sample) Abstract:

“Jäger (2018) reports that a hybrid method for phylogenetic inference, combining both binary character- and distance-based cognate detection methods, yields the best fit with expert-created gold standard trees. However, Jäger employs only one method for measuring the distance between cognate pairs, namely by means of pointwise mutual information on the basis of global sound class correspondences. **This thesis proposes several additional word pair distance metrics based on phonetic features and information content, in order to determine whether these might also prove useful, either for cognate detection in the first step or for subsequent generation of phylogenetic trees on the basis of cognate distance evaluation metrics.**”

OUTLINE OF THESIS

I) Introduction: Overview of language families, phylogenetic inference, cognates & cognate detection

II) Background / Literature Review

Survey of techniques in (computational) historical linguistics:

- comparative method, glottochronology, Bayesian phylogenetic methods, distance-based approaches

What else?

III) Methodology

Outline of cognate detection procedure

Phonetic and sound class distance metrics

Information-theory distance metrics

Clustering languages / generating phylogenetic trees

IV) Key Studies and Evaluation

Target families: Slavic, Romance, Turkic, Arabic, Sinitic

Cognate detection: evaluate against gold cognate sets

Distance measures: evaluate using tree topography comparison metrics, possibly also mutual intelligibility

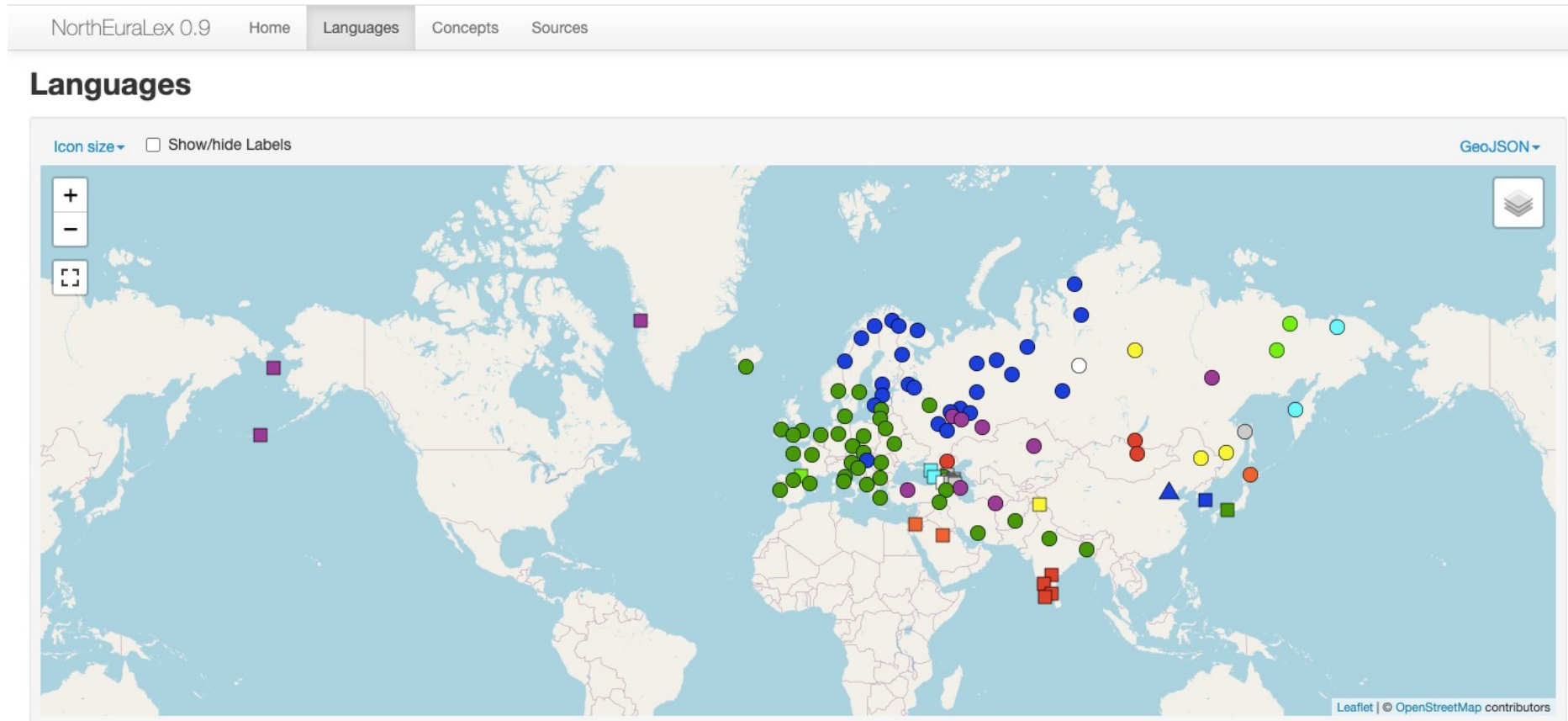
THESIS PROCEDURE

- Decide on input datasets (wordlists) and corresponding gold trees/cognate sets for evaluation
- Develop and implement distance metrics
- Develop/extend/adapt data structures for processing inputs, implementing cognate detection and evaluation steps, and comparing results with gold standard
- Perform cognate detection using each method, evaluate results against gold cognate sets
- Take results from **best-performing** cognate detection method, perform cognate evaluation
- Compare calculated linguistics distances against expert trees and intelligibility data

DATASETS

NorthEuraLex: 107 languages (Indo-European, Uralic, Turkic, Dravidian, Eskimo-Aleut, etc.)

→ use for (Balto-)Slavic, Uralic, *any others?*



DATASETS

NorthEuraLex: 107 languages (Indo-European, Uralic, Turkic, Dravidian, Eskimo-Aleut, etc.)

→ use for (Balto-)Slavic, Uralic, *any others?*

Sinitic (Líu et al., 2007): 19 Sinitic varieties

Turkic (Savelyev & Robeets, 2020): 31 Turkic languages

Arabic (Ratcliffe, 2020): 14 dialects of Arabic

Global Lexicostatistic Database: Romance (58 varieties), and many others...



SINITIC, ROMANCE, NORTHEURALEX DATASETS

- Downloaded and preprocessed all of the data to yield input files in the format needed for my implementation

ID	AUTOID	CHARACTERS	CHARACTERS_IS	CLASSES	COGIDS	CONCEPT	CONCEPTCON_ID	CONCEPTCON_ID	LOOSEID	STRICTID
CONCEPT_CHINESE	DOCULECT	NUMBERS	PAGE	DOCULECT_ID	PROSTRINGS	DUPLICATES	IPA	LANGID	SONARS	
MORPHEMES	STRUCTURE	STRUCTURES	TOKENS	SAMPA	SEGMENTS	VALUE	WEIGHTS	LEXEME_NOTE		
CHARACTER_NOTE	NOTE									
1	999	吐	吐	TY3	1	vomit	1278	嘔吐	Beijing 1	1
	t ^h u ⁵¹	1	487	spit/吐	1.T.C 1.Y.V 1.3.T	77	AXT	t_hu51	t ^h	
u ⁵¹	1 7 8	717	i n t	i n t	t ^h u ⁵¹ t ^h u ⁵¹	1.6 3.0 1.0				
2	999	吐	吐	TY3	1	vomit	1278	嘔吐	Ha_erbin	2
	1	t ^h u ⁵³	7	487	spit/吐	7.T.C 7.Y.V 7.3.T	77	AXT		
t_hu53	t ^h u ⁵³	1 7 8	717	i n t	i n t	t ^h u ⁵³ t ^h u ⁵³	1.6 3.0 1.0			
3	999	吐	吐	TY3	1	vomit	1278	嘔吐	Jinan 3	1
	t ^h u ³¹	8	487	spit/吐	8.T.C 8.Y.V 8.3.T	77	AXT	t_hu31	t ^h	
u ³¹	1 7 8	717	i n t	i n t	t ^h u ³¹ t ^h u ³¹	1.6 3.0 1.0				
4	999	嘔吐	嘔吐	U6_TY6	2 1	vomit	1278	嘔吐	Rongcheng	4
	0	ou ²¹³ +t ^h u ²¹⁴	14	487	nausea/嘔	spit/吐				
14._._	14.T.C 14.Y.V 14.6.T	77	XT_AXT	ou213-35 t_hu214	o u ^{35/213} + t ^h u ²¹⁴					
214	7 8 9 1 7 8	868	n t + i n t	n t + i n t						
ou ²¹³⁻³⁵	t ^h u ²¹⁴	3.0 1.0 0.0 1.6	3.0 1.0							
5	1000	嘔嘔	嘔嘔	TY3_LE6	2 5	vomit	1278	嘔吐	Taiyuan 5	0
	t ^h u ⁵³ +lǎ ⁰	16	487	nausea/嘔	_:PERFECTIVE/了					
16.3.T	16._._ 16.L.C 16.E.V 16.6.T	77	AXT_AXT	t_hu53 lǎ ⁰	t ^h u ⁵³ + lǎ ⁰					
7 8 9 5 7 8	869	i n t + i n t	i n t + i n t	t ^h u ⁵³ + lǎ ⁰	t ^h u ⁵³ lǎ ⁰					
1.6 3.0 1.0 0.0 1.6 3.0 1.0										
6	999	嘔吐	嘔吐	NU3_TY3	2 1	vomit	1278	嘔吐	Xi_an 6	0
	ηou ²¹ +t ^h u ⁵³	18	487	nausea/嘔	spit/吐					

I/我/ŋwai^{3 2}
all/個郎下/ko^{2 1 2} louŋ^{5 3} ŋa^{2 4 2}
and/共/kɔŋ^{2 4 2}
animal/四骸爬/si^{2 1 2} k^ha^{5 5} βa^{5 3}
ash/灰灰/hwoi^{5 5} hwoi^{5 5}
ash/灰/hwoi^{5 5}
back/背/p^hjaŋ^{5 5}
bad/痞/p^hai^{3 2}
bad/呆/ŋai^{5 3}
bark/树皮/ts^hjeu^{2 1 2} p^hwoi^{5 3}
because/因为/iŋ^{5 5} ŋwoi^{2 4 2}
belly/腹老/pu^{2 4} lo^{3 2}
big/大/twai^{2 4 2}
bird/鳥/tsɛu^{3 2}
bite/咬/ka^{2 4 2}

- Downloaded and preprocessed all of the data to yield input files in the format needed for my implementation

ID		AUTOID	CHARACTERS	CHARACTERS_IS	CLASSES	COGIDS	CONCEPT	CONCEPTCON_ID	CON_ID	CON_ID
CONCEPT	CHINESE	DOCULECT	DOCULECT_ID	DOCULECT_ID	COGIDS	CONCEPT	CONCEPTCON_ID	CON_ID	CON_ID	CON_ID
CONCEPT	CHINESE	DOCULECT	DOCULECT_ID	DOCULECT_ID	COGIDS	CONCEPT	CONCEPTCON_ID	CON_ID	CON_ID	CON_ID
puu	nort2646	Ohr::N	استرکه	st3rga	STERKA	STVRKV	validate	Beijing	1	1
puu	nort2646	Auge::N	خوړ	gwaz	GWAS	KWVS	validate	AXT	t_hu51	th
puu	nort2646	Nase::N	پوزه	'poza	PUSA	PVSV	validate			
puu	nort2646	Mund::N	خوړه	xu'la	GYLA	KVRV	validate			
puu	nort2646	Zahn::N	غاش	gas	GAS	KVS	validate			
puu	nort2646	Zunge::N	جبه	'd3aba	j3ba	CEPA	KVPV	Ha_erbin	2	
validate	puu	nort2646	Lippe::N	شونده	'fundu	Sunda	SVNTV	77	AXT	
validate	puu	nort2646	Wange::N	باړخو	ba'xu	barxu	PARGY	1.0		
validate	puu	nort2646	Gesicht::N	مخ	mex	MEG	MVK	Jinan	3	1
validate	puu	nort2646	Stirn::N	تلدي	tan'daj	tandaj	TANTAJ	AXT	t_hu31	th
puu	nort2646	Stirn::N	وچولي	wu'fj'wulai	wu'fj'wulai	wucwulai				
puu	nort2646	WVKWVRV	validate					Rongcheng	4	
puu	nort2646	Haar::N	ويښنه	wes'ta	WESTE	WVSTV	validate	14.U.V	14.6.T	
puu	nort2646	Schnurrbart::N	بريت	bret	PRET	PRVT	validate			
puu	nort2646	Bart::N	يوړه	'zira	SIRA	SVRV	validate			
puu	nort2646	Kinn::N	نښه	'zana	SENA	SVNV	validate			
puu	nort2646	Kiefer[Anatomie]::N	جامه	'd3oma	d3oma	jama	CAMA			
puu	nort2646	Kehle::N	غارو	'goja	GARA	KVRV	validate	Taiyuan	5	0
puu	nort2646	Halss::N	غارو	'goja	GARA	KVRV	validate	16.T.C	16.Y.V	
puu	nort2646	Genick::N	ورښير	wor'maz	w3rmaz	WERMES	WVRMS	th u 53 + l a 0 1		
validate	puu	nort2646	Kopf::N	سر	sar	SVR	validate	th u 53 l a 0 1		
puu	nort2646	Rücken::N	شا	fa	SA	SVS	validate			
puu	nort2646	Bauch::N	نس	nas	NAS	NVS	validate	Xi_an	6	0
puu	nort2646	Bauch::N	گېډه	'gaqa	g3da	KETA	KVTV	validate		
puu	nort2646	Nabel::N	نوم	nom	NUM	NVM	validate	18.U.V	18.3.T	
puu	nort2646	Busen::N	تي	tai	TAI	TV	validate			
puu	nort2646	Brust::N	سینه	si'na	SINA	SVNV	validate			
puu	nort2646	Schulter::N	اوړه	o'za	OZA	VSV	validate			
puu	nort2646	Arm::N	لاس	las	LAS	RVS	validate			
puu	nort2646	Ellenbogen::N	غلغل	fan'gal	fan'gal	cang3l	CANKEL			
puu	nort2646	Hand::N	لاس	las	LAS	RVS	validate			
puu	nort2646	Handfläche::N	ورښو	wor'gawai	w3rx3wai					
WVWVWV	WVWVWV	validate								
WVWVWV	WVWVWV	Finger::N	ښو	'guta	KYTA	KVTV	validate			

I/我/ŋwai^{3 2}
 all/個郎下/ko^{2 1 2}louŋ^{5 3}ŋa^{2 4 2}
 and/共/kɔŋ^{2 4 2}
 animal/四骹爬/si^{2 1 2}kʰa^{5 5}βa^{5 3}
 ash/灰灰/hwoi^{5 5}hwoi^{5 5}
 ash/灰/hwoi^{5 5}
 back/背/pʰjaŋ^{5 5}
 bad/痞/pʰai^{3 2}
 bad/呆/ŋai^{5 3}
 bark/树皮/tsʰjeu^{2 1 2}pʰwoi^{5 3}
 because/因为/iŋ^{5 5}ŋwoi^{2 4 2}
 belly/腹老/puʔ^{2 4}lo^{3 2}
 big/大/twai^{2 4 2}
 bird/鳥/tsɛu^{3 2}
 bite/咬/ka^{2 4 2}

Northern Pashto.txt

disappear/ورکیدل/wrəkēd'əl
dish/خواره/xwəɽə
dishes/لویښی/l'əʃai
distance/واتن/wəɽ'n
disturb/ماتول/məɽəw'əl
dive/غل غوټه وړل/ɣuɽ'əwəh'əl
divide/تقسیمول/taqsīmaw'əl
divide/وېشل/wəʃ'əl
do/کول/kəw'əl
doctor/طبيب/tab'ib
doctor/اکټر/ɔkt'ar
dog/سپی/spai
doll/ګودی/guɽ'əz
door/دروازه/dəɽwəz'ə
down/ښکته/ʃkəɽə
drag/وړل/wɽl'əl
draw/رسمول/rəsmaw'əl
dream/خوب/xəb
drink/غڼل/ɽsɽ'əl
drive/بیول/biɽ'əl
drive/تلل/tl'əl
drop/قطره/qəɽr'ə
drop/لښه لاسه لوېدل/ləl'əsalwəd'əl
drop/غاسګی/ɽs'əɽsakai
drop, descend/کېدل ښکته/ʃk'əɽəkəd'əl
dry/وچېدل/wuɽɽəd'əl
dry/وچ/wəɽɽ
duck/ایلی/il'əz
dust/خاوره/x'əwɽə
dust/ګرد/gəɽd

ARABIC DIALECT DATASET (RATCLIFFE, 2020)

- Can't access paper: https://brill.com/view/journals/ldc/11/1/article-p1_1.xml
- Extracted data table from PDF of paper's appendix
- Working on manually cleaning/reformatting extracted data

1.1 Data

		CA	Mor	Mlt	Cai	Dms	Irq	Skh	AqArb	Cyp	Glf	Ymn	Nig	Bux	Nub
1	ALL	kull	koll	Kollu	kull	kəll	kull	tʃill	kəll, sa:yi:n	kull	kil	kull	tʃat	kullu	kulu
2	ASH	ra ma:d	rʻm adʻ	rmi:d	ram a:dʻ*	rama :d	ruma :d	l	rʻam a:d	rama t r-m-d	rama ad	l	l	ra'matt	ruman
3	BARK	qirf at	qef r'a	ʔoʃr a	ʔiʃr*	ʔəʃər	giʃra	kifr	səvi:y e	l	giʃra	giʃr	li:he, girfe	l	*kokobo lataka, (girifa)
4	BELLY	batʻ n	ker ʃ	zaʔʔ	batʻ n*	batʻə n	batʻi n	batʻn , tʃarʃ	dʒoof	patn b-tʻ-n	batʻn	batʻn , karʃ	kirʃ	batin	batna
5	BIG	kab i:r	kbi r	kibi r	kibii r	kbi:r	tʃbi:r	tʃabi:r	gəbi:r	k-b-r	ʃood	Kabii r	kabi:r	ka'biir	kbir
6	BIRD	tʻa: ʔir	tʻir	ʃasf ur	tʻee r*	tʻeer	tʻeer	tʻeer	ku:tʃ ka:ye	2	tʻeer	tʻajr	tʻe:ra	tayra	ter
7	BITE	ʃad	ʃed	gide	ʃadʻ	ʃadʻd	ʃaðʻð	ʃaðʻð	l	ʃaðð	ʃadʻd	lugus	adʻdʻa	yaʃazz	adi

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	Index	Gloss	CA	Mor	Mlt	Cai	Dms	Irq	Skh	AqArb	Cyp	Glf	Ymn	Nig	Bux	Nub
2	1	ALL	kull	koll	Kollu	kull	kəll	kull	tʃill	kəll	kull	kil	kull	tʃat	kullu	kulu
3										sa:yi:n						
4	2	ASH	rama:d	rʻmadʻ	rmi:d	rama:dʻ	rama:d	ruma:d		rʻama:d	ramat	ramaad			ra'matt	ruman
5	3	BARK	qirfat	qefr'a	ʔoʃra	ʔiʃr*	ʔəʃər	giʃra	kifr			giʃra	giʃr	girfe		girifa
6																*kokobo
7										savi:ye						
8														li:he		
9																lataka
10	4	BELLY	batʻn			batʻn*	batʻan	batʻin	batʻn		patn	batʻn	batʻn	batʻn	batn	batna
11				kerʃ					tʃarʃ					kirʃ		
12					zaʔʔ											
13										dʒoof						
14	5	BIG	kabi:r	kbir	kibir	kibii r	kbir	tʃbi:r	tʃabi:r	gəbi:r	k-b-r			Kabir	kabi:r	ka'biir
15																kbir
16	6	BIRD	tʻa:ʔir	tʻir	ʃasfur	tʻeer*	tʻeer	tʻeer	tʻeer			ʃood		tʻe:ra	tayra	ter
17										ku:tʃ		tʻeer	tʻajr			
18										ka:ye						
19	7	BITE	ʃadʻd'a	ʃedʻd'	gidem	ʃadʻd'	ʃadʻd'	ʃaðʻð'	ʃaðʻð'		ʃaðð	ʃadʻd'		adʻd'a	yaʃazz	adi
20													lugus			
21	8	BLACK	ʔaswad	khel	iswed	iswid	ʔaswad	ʔaswad	ʔaswad	aswad	isfet	ʔaswad	ʔaswad	axadʻar	ʔaswad	*asue
22	9	BLOOD	dam	demn	derm	damm	damm	damm	dam	Dam	d-m-y	damm	Dam	damm	dam	dum
23																
24	10	BONE	ʃaðʻm	ʃedʻem	ʃadma	ʃadma	ʃadʻme	ʃaðʻma	ʃaðʻm	ʃadʻam		ʃaðʻm	ʃaðʻma	adʻum	ʃazəm	ladim
25																
26	11	BREAST	sʻadr	bezzula	Sider	sidir*	sʻadʻar	nahid	ʃidi(f.)sʻidr		pzaz9pl)jb-	deed(f)jsʻa	sʻadr	de:d	sadar;buʃu	*sudur
27																
28	12	BURN	haraqa	hreq	haraʔ	haraʔ	haraʔ	hirag	harak	h-r-q	h-r-q	harag	harag	hirig,tajʃa	ʃalag	haragu
29																
30	13	CLAW	ðʻufr	dʻferʻ	Difer	dʻifir*	dʻafar	ðʻifir	iðʻfir	dʻafar	ðʻ-f-r	mixdab	ðʻufrin**	xumfur		*dufuru
31																
32	14	CLOUD	yama:mat	ymama	shaba	sahaab	yeeme	yeema	yeemi		y-y-r	sihaab,yeem		dʻa(ha)wij	yeem,deen	*soobura
33																
34	15	COLD	ba:rid	bared	bi:red	bard	barəd	barda:n	bi:rid	beerəd	*peret(brd	baarid	baarid	barda:n	bard	bari
35																
36	16	COME	dʒa:ʔa	ʒa	dʒi:	ga	ʔadʒa	dʒa:	atsa	dʒa:	dʒ-y-y	dʒa	dʒaaʔ	dʒa	dʒakki	dʒa
37																
38																
39	17	DIE	ma:ta	mat	mi:t	maat	ma:t	ma:t	ma:t	ma:t	m-w-t	maat	maat	ma:t	maat	mutu
40																
41	18	DOG	kalb	kelb	kelb	kalb	kalb	tʃalib	tʃalb	kalb	k-l-b	tʃalb	kalb	kalb	kalb	keli

FEATURAL REPRESENTATION OF SOUNDS

- Goal: encode any IPA sequence as a vector or sequence of vectors representing the distinctive phonological features of the sound(s) in question
- Challenges/Considerations
 - Should be able to handle all possible IPA characters and any relevant combinations
 - Should be able to deal with different equivalent transcriptions, e.g. $[p^h] = [p^{hj}]$, $[\widehat{tj}] = [tj]$
 - Ideally: no redundancy in features, i.e. use as few features as possible for all distinct sounds to have distinct representations
- Source
 - PHOIBLE feature set: <https://phoible.org/parameters>

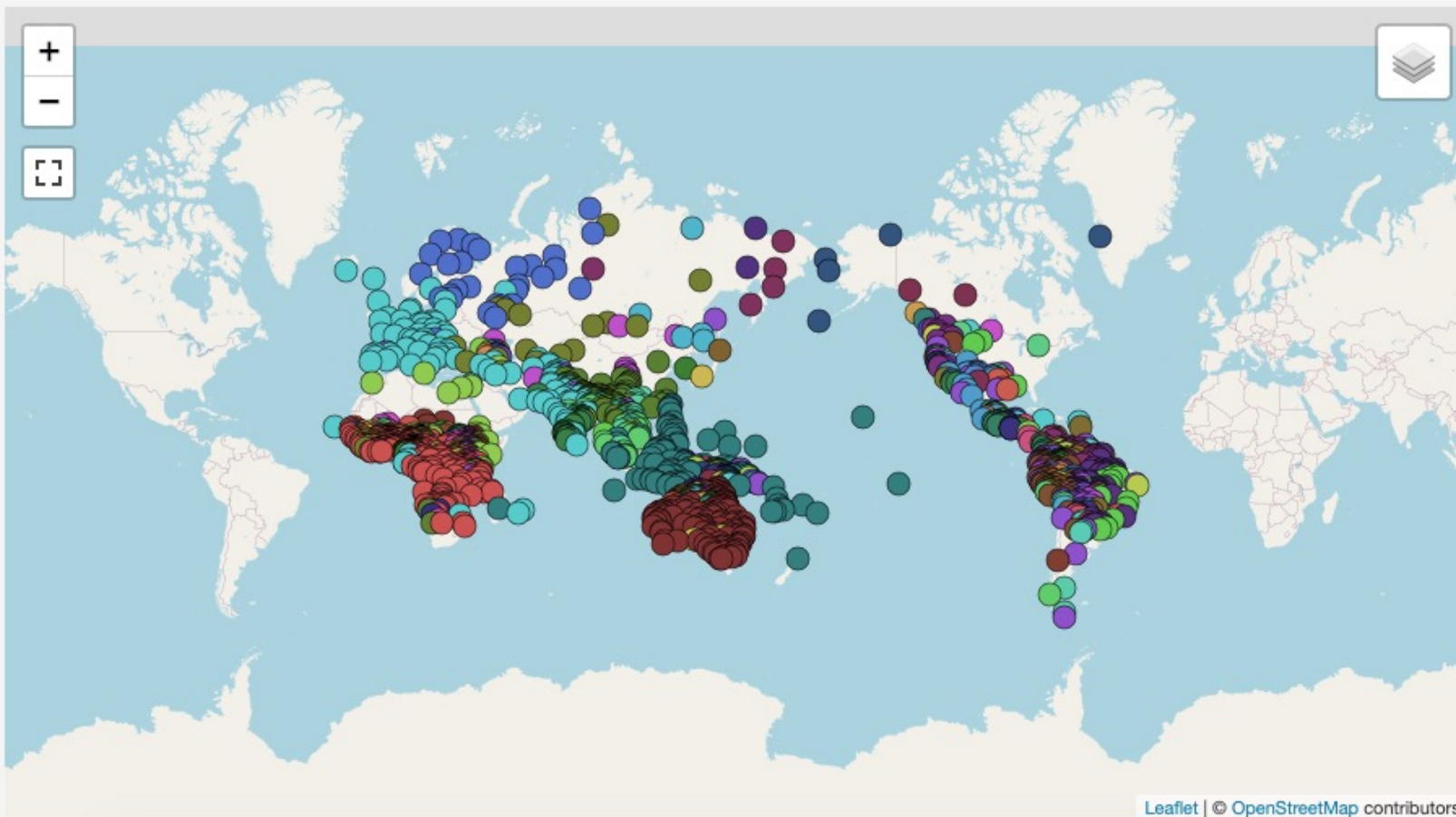
Consonant p W



LATIN SMALL LETTER P

Icon size ▾ ☐ Show/hide Labels

GeoJSON ▾



Features

advanced tongue root	0
anterior	0
approximant	-
back	0
click	-
consonantal	+
constricted glottis	-
continuant	-
coronal	-
delayed release	-
distributed	0
dorsal	-
epilaryngeal source	-
fortis	-
front	0
high	0
labial	+
labiodental	-
lateral	-
long	-
low	0

EXPLORING SEGMENTS/FEATURES

- Examined phonetic symbols used in NorthEuraLex, Sinitic, and Romance datasets, plus set of PHOIBLE features
- All standard IPA characters used in all 4 datasets already covered by my earlier implementation of phonetic distance
 - Only a few diacritics and characters specific to Sinitic transcription conventions missing
 - easy to add or replace with equivalent characters, e.g. /ɲ/ = /z̞/
- 20 of 38 PHOIBLE features have direct equivalents to features in my implementation
 - Some of the remaining PHOIBLE features may be redundant, e.g. “labial” and “labiodental” features are not needed as they can equivalently be expressed by combinations of other features

EXPLORING SEGMENTS/FEATURES

- Strategy for dealing with phonetic diacritics (e.g. \sim , ^h, j, w, :, etc.)
 - PHOIBLE: each character with diacritics has a separate entry in feature matrix
 - e.g. /p/ has a separate entry from /p^h/ and another for /p^h:/, etc.
 - 2162 segment entries
- My method: each base character has a unique featural entry; diacritics associated with specific features and values, modify the base feature vector
 - e.g. /p/ has unique representation; adding /^h/ adds the +SPREAD GLOTTIS feature; /:/ adds +LONG feature
 - Only 116 segment entries, plus diacritics

FEATURAL REPRESENTATION OF SOUNDS

- **Plan: adapt PHOIBLE segment/feature dataset**
 - Use PHOIBLE feature representation for basic IPA characters (i.e. segments without diacritics)
 - Remove any redundant features
 - Extend current method for allowing diacritics to modify base characters' feature vectors
 - Eliminates need for an entry in feature matrix for every combination of IPA characters and diacritics
 - Allows for more flexibility in transcription conventions

MEASURING PHONETIC DISTANCE

- Values of distinctive features can be:

- + the sound has this feature
- the sound does not have this feature
- 0 this feature is not relevant for this sound

		diacritic examples				Coronal Obstruents [+cons, -son, +cor]								Palatal Obstruents [+cor + dors]								Non-coronal Obstruents [+cons, -son, -cor]										Laryngeals [-cons, -son]							
		tʰ	tʰ	t̚	t	d	s	z	ʈ	ʂ	θ	ð	ʃ	ʒ	c	ɟ	ç	j	p	b	f	v	ɸ	β	k	g	x	ɣ	q	ʁ	χ	ʁ	h	ʕ	h	ɦ	ʔ		
Class features	cons	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	
	son	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	syll	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
Place features	labial	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	
	round	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	-	-	-	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	coronal	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	ant	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	dist	-	-	+	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	dorsal	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	-	-	-	-	-	-	+	+	+	+	+	+	+	+	-	-	-	-	-	
	high	0	0	0	0	0	0	0	0	0	0	0	0	0	0	+	+	+	+	0	0	0	0	0	0	+	+	+	+	-	-	-	-	0	0	0	0	0	
	low	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	0	0	0	0	0	0	-	-	-	-	-	-	-	-	0	0	0	0	0	
	back	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	0	0	0	0	0	0	+	+	+	+	+	+	+	+	0	0	0	0	0	0
	tense	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	0	0	0	0	0	0	-	-	-	-	-	-	-	-	0	0	0	0	0	
Place features	pharyngeal	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	-	-	-	-		
	ATR	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	0	0	0	

MEASURING PHONETIC DISTANCE

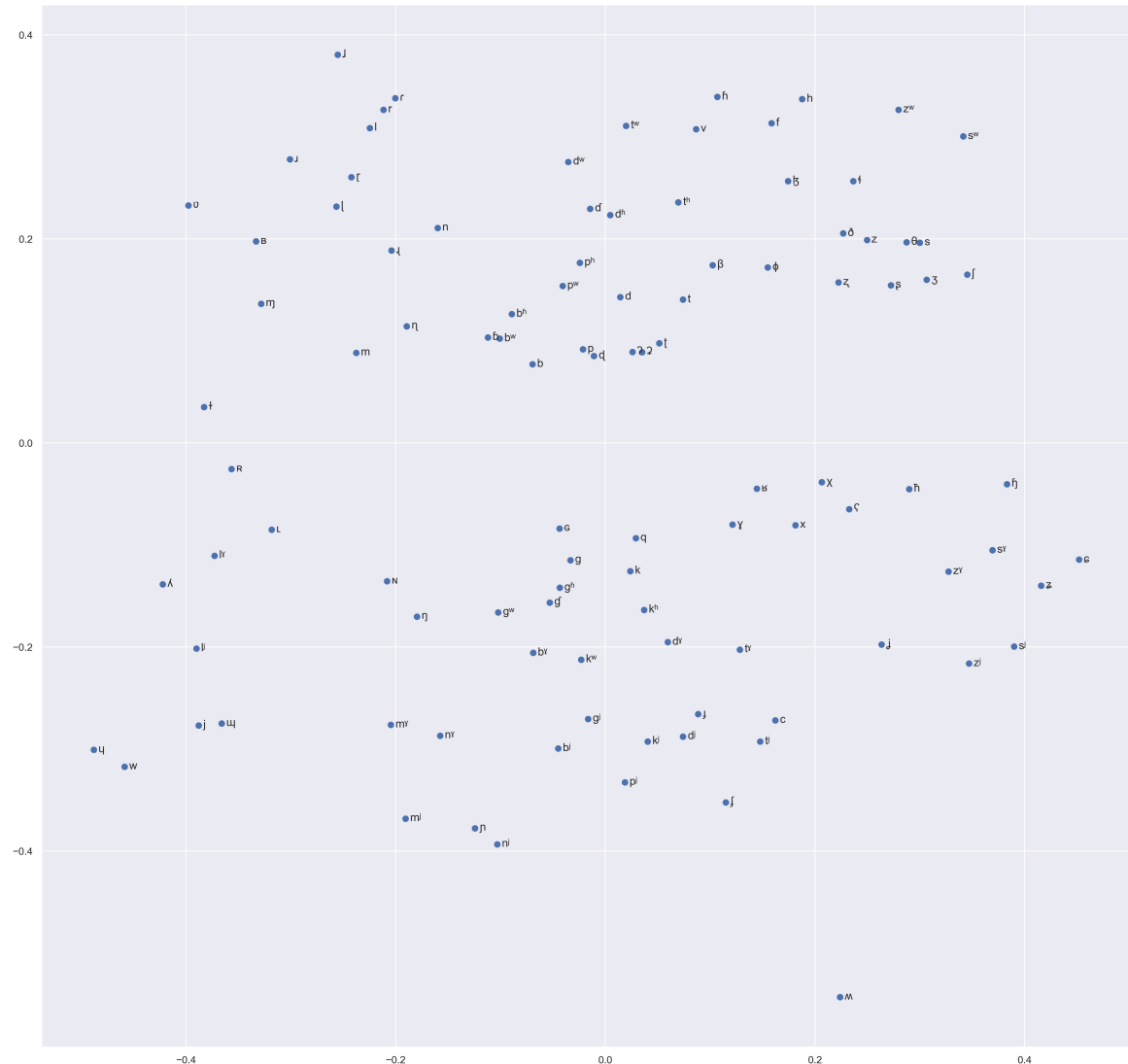
- Values of distinctive features can be:
 - + the sound has this feature
 - - the sound does not have this feature
 - 0 this feature is not relevant for this sound

		Coronal Obstruents												Palatal Obstruents				Non-coronal Obstruents												Laryngeals										
		diacritic examples		[+cons, -son, +cor]												[+cor + dors]				[+cons, -son, -cor]												[-cons, -son]								
		tʰ	tʰ	t̥	t	d	s	z	t̪	ʈ	θ	ð	ʃ	ʒ	c	ɟ	ç	j	p	b	f	v	ɸ	β	k	g	x	ɣ	q	ɢ	χ	ʁ	ħ	ʕ	h	ɦ	ʔ			
Class features	cons	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	
	son	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	syll	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
Place features	labial	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	round	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	-	-	-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	coronal	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	ant	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	dist	-	-	+	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	dorsal	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	-	-	-	-	-	-	+	+	+	+	+	+	+	+	+	-	-	-	-	-	
	high	0	0	0	0	0	0	0	0	0	0	0	0	0	0	+	+	+	+	0	0	0	0	0	0	0	+	+	+	+	-	-	-	-	0	0	0	0	0	
	low	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	0	0	0	0	0	0	0	-	-	-	-	-	-	-	-	0	0	0	0	0	
	back	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	0	0	0	0	0	0	0	+	+	+	+	+	+	+	+	0	0	0	0	0	0
	tense	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	-	-	0	0	0	0	0	0	0	-	-	-	-	-	-	-	-	0	0	0	0	0	0
pharyngeal		-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	-	-	-	-		
	ATR	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	-	0	0	0	0

MEASURING PHONETIC DISTANCE

- Values of distinctive features can be:
 - + the sound has this feature
 - – the sound does not have this feature
 - 0 this feature is not relevant for this sound
- Problem: how to code the '0' (unspecified) value when quantifying differences?
 - **Strategy 1:** Omit 0-valued features entirely when comparing a pair of sounds
Rationale: doesn't make sense to compare using a feature that is unspecified
Problem: different pairs of sounds would be compared using vectors of different lengths
 - **Strategy 2:** Encode 0-valued features the same as – [sound doesn't have this feature]
Rationale: 0-valued features are actually sub-features of other features; if the main feature has value –, then the sub-features would by default have the same – value
Problem: what is actually the difference between value 0 and value –?
 - **Strategy 3:** Create 3-way distinction in feature values and consider 0 different from both + and –
Rationale: avoids problem of either ignoring 0 features or needing to re-encode them
Problem: is it a valid comparison between 0 and +/–?

- How to pick the appropriate method for dealing with 0 values?
- Compare MDS projections of segments based on features (once adapted)



WORKS CITED

Dellert, J., Daneyko, T., Münch, A. et al. Lang Resources & Evaluation (2019). <https://doi.org/10.1007/s10579-019-09480-6> (version 0.9).

Jäger, G. (2018). Global-scale phylogenetic linguistic inference from lexical resources. *Scientific Data*, 5(180189), 1–16. <https://doi.org/10.1038/sdata.2018.189>

Líu, L.; Wáng, H.; Bǎi, Y. (2007): Xiàndài Hànyǔ fāngyán héxīncí, tèzhēng cíjí
现代汉语方言核心词·特征词集 [Collection of basic vocabulary words and characteristic
dialect words in modern Chinese dialects]. Nánjīng: Fèngguāng.

Moran, Steven & McCloy, Daniel (eds.) 2019. PHOIBLE 2.0. Jena: Max Planck Institute for the Science of Human History. <http://phoible.org>.

Ratcliffe, R. R. (2020). The glottometrics of Arabic: Quantifying linguistic diversity and correlating it with diachronic change. *Language Dynamics and Change*, 11(1), 1–29.

Savelyev, A., & Robbeets, M. (2020). Bayesian phylolinguistics infers the internal structure and the time-depth of the Turkic language family. *Journal of Language Evolution*, 5(1), 39–53.