Data in computational historical linguistics

Gerhard Jäger

ESSLLI 2016

Background

- comparative method strongly focuses on two types of data:
 - morphological paradigms
 - regular sound correspondences
- both are not very suitable for computational approaches, because
 - morphological categories are not easily comparable across languages, especially if we look individual language families
 - also, isolating languages have no morphology
 - identifying regular sound correspondences automatically is a surprisingly hard problem, due to data sparseness
 - currently one of the hot topics, far from resolved (List, 2014; Hruschka et al., 2015; Bouchard-Côté et al., 2013)

Gerhard Jäger Data sources ESSLLI 2016 2 / 25

Background

- what we need (especially if we apply statistical methods):
 - data types which are applicable to all natural languages
 - ideally lots of data
- current practice:
 - word lists + expert annotations about cognacy (currently the dominant paradigm)
 - unannotated word lists in phonetic transcriptions
 - discrete grammatical categorizations (compiled by human experts)

Gerhard Jäger Data sources ESSLLI 2016 3 / 25

Cognate-coded Swadesh lists

Swadesh lists

- collections of 100 200 concepts (there are different versions)
- core vocabulary:
 - not culture dependent
 - diachronically stable, i.e. resistant both against semantic change and aginst borrowing
- proposed by Morris Swadesh (Swadesh, 1955, 1971) to facilitate an early attempt to automatize certain tasks in historical linguistics
- popular among computational historical linguistics because it is a standard
- see (List, 2016) for a thoughtful discussion of the notion of cognacy

Cognates

Cognates are words that have the same origin

Latin filius \Rightarrow French fils, Italian figlio

 traditionally, cognacy excludes loanwords, but terminology among computationalists is sometimes less strict:

Latin $persona \Rightarrow English person$

would also qualify as cognate pair

 on average, the closer two languages are related, the more cognate pairs they share

Cognates

- during language change, the word for a given concept is sometimes replaced by a non-cognate one
- causes: semantic change, borrowing, morphological word formation
 - 'bone': Old High German Bein (cognate to Engl. $bone \Rightarrow New$ High German Knochen
 - Bein is still part of the German lexicon, but it now means leg
- cognate replacement is comparable to a mutation in biological evolution

Cognates

Caveats

- cognacy is not binary, but a matter of degree
 - English woman ← Old English wiff-man
 - first component is cognate to wife, German Weib etc., and second component to man, German Mann etc. Are woman and Weib cognate or not?
- for distantly related languages, experts often disagree about cognacy
 Ancient Greek ὕλη/Latin silva 'woods'

IELex

- Indo-European Lexical Cognacy Database
- freely available online at http://ielex.mpi.nl/
- based on Dyen et al. (1992)
- current version curated by group at MPI Nijmegen
- recently migrated to new MPI Jena; new version not public yet

IELex

- 207-item Swadesh lists for 135 Indo-European languages
- words in orthographic and partially in phonetic transcription (IPA)
- entries are assigned to cognate classes
- sample entries:

language	iso_code	gloss	global_id	local_id	transcription	cognate_class
ELFDALIAN	qov	woman	962	woman	ˈkɛ̀lɪŋg	woman:Ag
DUTCH	nld	woman	962	woman	vrau	woman:B
GERMAN	deu	woman	962	woman	fraŭ	woman:B
DANISH	dan	woman	962	woman	'gʰvenə	woman:D
DANISH_FJOLDE		woman	962	woman	kvin ^j	woman:D
GUTNISH_LAU		woman	962	woman	'kvın:ˌfolk	woman:D
LATIN	lat	woman	962	woman	'mulier	woman:E
LATIN	lat	woman	962	woman	'fe:mina	woman:G
ENGLISH	eng	woman	962	woman	womən	woman:H
GERMAN	deu	woman	962	woman	vaĭp	woman:H
DANISH	dan	woman	962	woman	ˈdɛ:mə	woman:K

Other publicly available cognacy data sources

- Austronesian Basic Vocabulary Database
 http://language.psy.auckland.ac.nz/austronesian/
- ten collections of cognate-coded Swadesh lists from various language families collected by Johann-Mattis List¹
- ten collections of short (40-100 items) cognate-coded Swadesh lists from various language families collected by Søren Wichman and Eric Holman²
- 88 cognate-coded Swadesh lists from Central-Asian languages³

¹List, J.-M. (2014): Data from: Sequence comparison in historical linguistics. GitHub Repository. http://github.com/SequenceComparison/SupplementaryMaterial. Release: 1.0.

²Supplementary material to Wichmann and Holman (2013)

³Supplementary material to Mennecier et al. (2016)

Phonetically transcribed Swadesh lists

The Automatic Similarity Judgment Program

- Project originally hosted at MPI EVA in Leipzig around Søren Wichmann
- since 2009; currently version 17 (2016)
- covers more than 7,000 languages and dialects (4.574 languages with iso code)
- basic vocabulary of 40 words for each language, in uniform phonetic transcription
- freely available at http://asjp.clld.org/

used concepts: *I*, you, we, one, two, person, fish, dog, louse, tree, leaf, skin, blood, bone, horn, ear, eye, nose, tooth, tongue, knee, hand, breast, liver, drink, see, hear, die, come, sun, star, water, stone, fire, path, mountain, night, full, new, name

The Automatic Similarity Judgment Program

Phonetic transcription

- 41 sound classes, all coded as ASCII characters
- various diacritics to capture finer phonetic distinctions, e.g.
 - ph~: aspirated p
 - a*: nasalized a
 - hkw\$: pre-aspirated labalized k

Metadata

- language family, language genus, classification according to Ethnologue and Glottolog
- geographic location
- population size

The Automatic Similarity Judgment Program

ASJP sound classes (from Brown et al. 2013)

ASJP code symbol	Description	IPA symbols
р	voiceless bilabial stop and fricative	р.ф
Ь	voiced bilabial stop and fricative	b. β
f	voiceless labiodental fricative	f
v	voiced labiodental fricative	v
m	bilabial nasal	m
w	voiced bilabial-velar approximant	w
8	voiceless and voiced dental fricative	θ, δ
4	dental nasal	p
t	voiceless alveolar stop	t
d	voiced alveolar stop	d
s	voiceless alveolar fricative	s
z	voiced alveolar fricative	z
c	voiceless and voiced alveolar affricate	ts, dg
n	alveolar nasal	n
r	voiced apico-alveolar flap and all other varieties of "r-sounds"	r, r, R, [
1	voiced alveolar lateral approximant	1
S	voiceless post-alveolar fricative	ſ
Z	voiced post-alveolar fricative	3
Z C	voiceless palato-alveolar affricate	3 tf
j T	voiced palato-alveolar affricate	¢5
	voiceless and voiced palatal stop	c, j
5	palatal nasal	n
У	palatal approximant	j
k	voiceless velar stop	k
g	voiced velar stop	g
×	voiceless and voiced velar fricative	x,
N	velar nasal	ŋ

ASJP code symbol	Description	IPA symbols
q	voiceless uvular stop	q
G	voiced uvular stop	G
X	voiceless and voiced uvular fricative, voiceless and voiced pharyngeal fricative	χ, в, ћ, ۲
h	voiceless and voiced glottal fricative	h, fi
7	voiceless glottal stop	?
L	all other laterals	ι, Į, λ
!	all varieties of "click-sounds"	!. . . ±
i	high front vowel, rounded and unrounded	i, ı, y, y
e	mid front vowel, rounded and unrounded	e, ø
E	low front vowel, rounded and unrounded	æ, ε, œ, Œ
3	high and mid central vowel, rounded and unrounded	i, 9, 9,3, tt, 0, G
a	low central vowel, unrounded	a, e
u	high back vowel, rounded and unrounded	w, u
0	mid and low back vowel, rounded and unrounded	Y, A, G, O, D, D

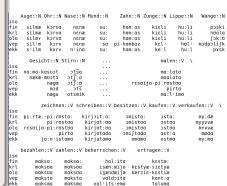
Automated Similarity Judgment Project

concept	Latin	English
1	ego	Ei
you	tu	yu
we	nos	wi
one	unus	w3n
two	duo	tu
person	persona, homo	%pers3n
fish	piskis	fiS
dog	kanis	dag
louse	pedikulus	laus
tree	arbor	tri
leaf	foly \sim u *	lif
skin	kutis	%skin
blood	saNgw \sim is	bl3d
bone	os	bon
horn	kornu	horn
ear	auris	ir
eye	okulus	Ei
nose	nasus	nos
tooth	dens	tu8
tongue	liNgw∼E	t3N

	•	
concept	Latin	English
knee	genu	ni
hand	manus	hEnd
breast	pektus, mama	brest
liver	yekur	liv3r
drink	bibere	drink
see	widere	si
hear	audire	hir
die	mori	dEi
come	wenire	k3m
sun	sol	s3n
star	stela	star
water	akw∼a	wat3r
stone	lapis	ston
fire	iNnis	fEir
path	viya	pE8
mountain	mons	%maunt3n
night	noks	nEit
full	plenus	ful
new	nowus	nu
name	nomen	nem

NorthEuraLex

- Massive data collection effort of the Tübingen EVOLAEMP project
- (currently) translations of 1,017 concepts into 103 (mostly)
 Northern Eurasian languages (cf. Dellert, 2015)
- everything transcriped in IPA
- (so far) no manual cognate coding



Grammatical classifications

Grammatical classification databases

- World Atlas of Language Structure (WALS) http://wals.info/
- Syntactic Structures of the World's Languages (SSWL)
 http://sswl.railsplayground.net/
- collection of syntactic parameters (in the Chomskyan sense) for a few dozen languages collected in the LanGeLin project (Giuseppe Longobardi)

Expert family trees

Expert family trees

- Ethnologue https://www.ethnologue.com/
- Glottolog http://glottolog.org/
 - in many ways improved version of Ethnologue
 - strives to apply uniform standards across all languages
 - rather conservative in accepting family status

Running example

Running example

- 25 living Indo-European languages
- three types of data
 - Swadesh lists in IPA transcription, taken from IELex
 - expert cognate classifications of Swadesh list entries (likewise taken from IELex),⁴ and
 - phonological, grammatical and semantic classifications of languages (taken from WALS)

⁴I only included those entries from IELex where both an IPA transcription and a cognate classification is given.

Running example

sample entries:

language	phonological form (IELex)	cognate class (IELex)	order of subject, object and verb (WALS)
Bengali	-	-	SOV
Breton	-	-	SVO
Bulgarian	mu'rε	sea:B	SVO
Catalan	mar; mar; ma	sea:B	SVO
Czech	'mɔr̞ɛ	sea:B	SVO
Danish	haw/sø?	sea:K/sea:J	SVO
Dutch	ze	sea:J	no dominant order
English	si:	sea:J	SVO
French	mer	sea:B	SVO
German	ze:/'o:tsea:n/me:g	sea:J/sea:E/sea:B	no dominant order
Greek	'θala,sa	sea:F	no dominant order
Hindi	-	-	SOV
Icelandic	ha:v/sjou:r	sea:K/sea:J	SVO
Irish	'fværvji	sea:G	VSO
Italian	'mare	sea:B	SVO
Lithuanian	'ju:re	sea:H	SVO
Nepali	-	-	SOV
Polish	'mɔzɛ	sea:B	SVO
Portuguese	mar	sea:B	SVO
Romanian	'mare	sea:B	SVO
Russian	'mɔrʲɛ	sea:B	SVO
Spanish	mar	sea:B	SVO
Swedish	ha:v/fjø:	sea:K/sea:J	SVO
Ukrainian	'mɔrɛ	sea:B	SVO
Welsh	-	-	VSO

Exercises

- Access the files ielexData.csv and walsData.csv from our running example from http:
 - //www.sfs.uni-tuebingen.de/~gjaeger/esslli2016/data/
 - Are there any WALS feature values exclusively occurring in the Romance languages?
 - Are there any cognate classes exclusively occurring in the Romance languages?
 - Are there any sound shifts (with instances in our data) exclusively occurring in the Romance languages?
 - Answer the same questions for the Slavic languages.

- Bouchard-Côté, A., D. Hall, T. L. Griffiths, and D. Klein (2013). Automated reconstruction of ancient languages using probabilistic models of sound change. *Proceedings of the National Academy of Sciences*, **36**(2):141–150.
- Brown, C. H., E. Holman, and S. Wichmann (2013). Sound correspondences in the world's languages. *Language*, **89**(1):4–29.
- Dellert, J. (2015). Compiling the Uralic dataset for NorthEuraLex, a lexicostatistical database of Northern Eurasia. Proceedings of the First International Workshop on Computational Linguistics for Uralic Languages. January 16, Tromsø, Norway.
- Dyen, I., J. B. Kruskal, and P. Black (1992). An Indoeuropean classification: A lexicostatistical experiment. *Transactions of the American Philosophical Society*, **82**(5):1–132.
- Hruschka, D. J., S. Branford, E. D. Smitch, J. Wilkins, A. Meade,
 M. Pagel, and T. Bhattachary (2015). Detecting regular sound changes in linguistics as events of concerted evolution. *Current Biology*,
 25(1):1–9.

- List, J.-M. (2014). *Sequence Comparison in Historical Linguistics*. Düsseldorf University Press, Düsseldorf.
- List, J.-M. (2016). Beyond cognacy: historical relations between words and their implication for phylogenetic reconstruction. *Journal of Language Evolution*, **1**(1):119–136. Doi: 10.1093/jole/lzw006.
- Mennecier, P., J. Nerbonne, E. Heyer, and F. Manni (2016). A Central Asian language survey: Collecting data, measuring relatedness and detecting loans. *Language Dynamics and Change*, **6**(1). In press.
- Swadesh, M. (1955). Towards greater accuracy in lexicostatistic dating. *International Journal of American Linguistics*, **21**:121–137.
- Swadesh, M. (1971). The Origin and Diversification of Language. Aldine, Chicago.
- Wichmann, S. and E. W. Holman (2013). Languages with longer words have more lexical change. In L. Borin and A. Saxena, eds., *Approaches to Measuring Linguistic Differences*, pp. 249–284. Mouton de Gruyter, Berlin.