

Investigating the Impact of Sample Size on Cognate Detection

Johann-Mattis List

Research Unit Quantitative Language Comparison
Philipps-University Marburg

March 17, 2013

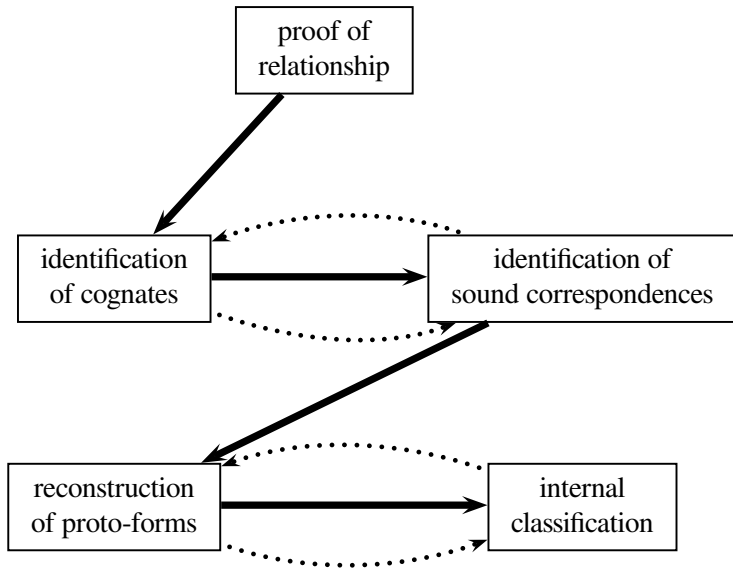
Sanscruta and Italian

Sono scritte le loro scienze tutte in una lingua, che dimandano Sanscruta, che vuol dire bene articolata. [...] et ha la lingua d'oggi molte cose comuni con quella, nella quale sono molti de' nostri nomi, e particolarmente de' numeri il 6, 7, 8 e 9, Dio, serpe, et altri assai.(Sasseti 1855: 415)

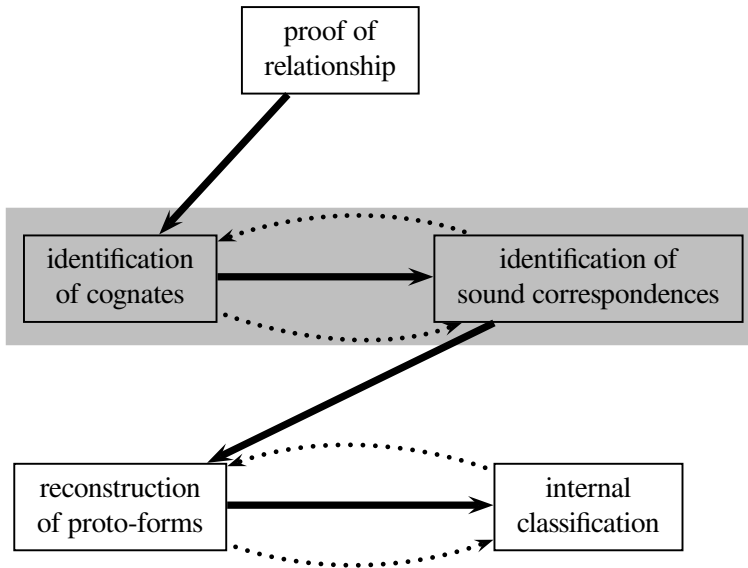
Translation: Everything that is related to science is written in a language which they call “Sanscruta”, meaning as much as “well-articulated”. Our language has much in common with it, among others many of our words, especially the numbers 6, 7, 8, and 9, “God”, “snake”, and many more.

The Comparative Method

Working Procedure



Working Procedure



Cognate Detection

Cognate Detection

Cognate List	
German	<i>dünn</i>
English	<i>thin</i>
German	<i>Ding</i>
English	<i>thing</i>
German	<i>dumm</i>
English	<i>dumb</i>

Cognate Detection

Cognate List		Alignment		
German	<i>dünn</i>	d	ʏ	n
English	<i>thin</i>	θ	ɪ	n
German	<i>Ding</i>	d	ɪ	ŋ
English	<i>thing</i>	θ	ɪ	ŋ
German	<i>dumm</i>	d	ʊ	m
English	<i>dumb</i>	d	ʌ	m

Cognate Detection

Cognate List		Alignment			Correspondence List		
German	<i>dünn</i>	d	y	n	GER	ENG	Frequ.
English	<i>thin</i>	θ	ɪ	n	d	θ	2 x
German	<i>Ding</i>	d	ɪ	ŋ	d	d	1 x
English	<i>thing</i>	θ	ɪ	ŋ	n	n	1 x
German	<i>dumm</i>	d	ʊ	m	m	m	1 x
English	<i>dumb</i>	d	ʌ	m	ŋ	ŋ	1 x

Cognate Detection

Cognate List		Alignment			Correspondence List		
German	<i>dünn</i>	d	ʏ	n	GER	ENG	Frequ.
English	<i>thin</i>	θ	ɪ	n	d	θ	2 x
German	<i>Ding</i>	d	ɪ	ŋ	d	d	1 x
English	<i>thing</i>	θ	ɪ	ŋ	n	n	1 x
German	<i>dumm</i>	d	ʊ	m	m	m	1 x
English	<i>dumb</i>	d	ʌ	m	ŋ	ŋ	1 x
German	<i>Dorn</i>	d	ɔə	n			
English	<i>thorn</i>	θ	ɔ:	n			

Cognate Detection

Cognate List		Alignment			Correspondence List		
German	<i>dünn</i>	d	ʏ	n	GER	ENG	Frequ.
English	<i>thin</i>	θ	ɪ	n	d	θ	3 x
German	<i>Ding</i>	d	ɪ	ŋ	d	d	1 x ?
English	<i>thing</i>	θ	ɪ	ŋ	n	n	2 x
German	<i>dumm</i>	d	ʊ	m	m	m	1 x
English	<i>dumb</i>	d	ʌ	m	ŋ	ŋ	1 x
German	<i>Dorn</i>	d	ɔə	n			
English	<i>thorn</i>	θ	ɔ:	n			

Cognate Detection

Cognate List		Alignment			Correspondence List		
German	<i>dünn</i>	d	y	n	GER	ENG	Frequ.
English	<i>thin</i>	θ	ɪ	n	d	θ	3 x
German	<i>Ding</i>	d	ɪ	ŋ	d	d	1 x
English	<i>thing</i>	θ	ɪ	ŋ	n	n	2 x
German	<i>dumm</i>	d	ʊ	m	m	m	1 x
English	<i>dumb</i>	d	ʌ	m	ŋ	ŋ	1 x
German	<i>Dorn</i>	d	ɔə	n			
English	<i>thorn</i>	θ	ɔ:	n			

Cognate Detection

Cognate List		Alignment			Correspondence List		
German	<i>diinn</i>	d	ɣ	n	GER	ENG	Frequ.
English	<i>thin</i>	θ	ɪ	n	d	θ	3 x
German	<i>Ding</i>	d	ɪ	ŋ	n	n	2 x
English	<i>thing</i>	θ	ɪ	ŋ	ŋ	ŋ	1 x
German	<i>Dorn</i>	d	ɔ̃	n			
English	<i>thorn</i>	θ	ɔ:	n			

Summary

Important Aspects

- language-specific notion of word similarity
- regular sound correspondences
- iterative character

Unspecified Parameters

- number of languages
- semantic similarity of the words
- size of the word lists

Summary

The Problem of the Sample Size

	Albanian	English	French	German
Albanian		0.07	0.10	0.10
English	14		0.23	0.56
French	20	46		0.23
German	20	111	46	

Numbers and proportions of shared cognates in the Swadesh-200 list (Swadesh 1952), taken from Kessler (2001).

Automatic Cognate Detection

Two Types of Similarity

“Phenotypic” Similarity (Lass 1997)

- based on surface resemblances of phonetic segments
- only depends on the words under comparison

“Genotypic” Similarity (ibid.)

- based on sound-correspondences
- depends on the words and the languages under comparison

Two Types of Similarity

German *Mund* [mʊnt]
English *mouth* [mauθ]

Two Types of Similarity

German *Mund* [mʊnt]
 English *mouth* [maʊθ]

German			English		
<i>Milch</i>	[mɪlç]	m	m	[mɪlk]	<i>milk</i>
<i>rund</i>	[rʊnt]	ʊ	au	[raʊnd]	<i>round</i>
<i>anders</i>	[andərs]	n	-	[ʌ(-)θər]	<i>other</i>
<i>südlich</i>	[sytɪç]	t	θ	[sʌθərn]	<i>southern</i>

Language-Independent Approaches

Normalized Edit Distance

- align two words and calculate their hamming distance
- normalize by dividing by the length of the longer word
- assume cognacy for distances beyond a certain threshold

Turchin et al. (2010)

- convert two (or more) words to Dolgopolsky (1966) consonant classes
- assume cognacy if the first two classes match

Language-Independent Approaches

German *Mund* [mʊnt]
English *mouth* [mauθ]

Language-Independent Approaches

German *Mund* [mʊnt]
 English *mouth* [mauθ]

Turchin		NED			
mont	→ M N T	m	ʊ	n	t
mauθ	→ M T	m	au	-	θ
Matches:	x	0	1	1	1
1 match => not cognate		3/4 = 0.75 => not cognate			

Language-Specific Approaches

LexStat (List 2012a)

- represent words as tuples of sound classes and prosodic strings
- use the SCA approach (List 2012b) to guess initial correspondences
- use a Monte-Carlo permutation test to derive language-specific similarity scores
- use the language-specific scores to calculate distance between words
- cluster words into cognate sets using a flat cluster algorithm

LexStat

LexStat

Sound Classes

Sounds which frequently occur in correspondence relations in genetically related languages can be divided in classes (types). It is thereby assumed that “phonetic correspondences inside a ‘type’ are more regular than those between different ‘types’” (Dolgoposky 1986[1966]: 35).

LexStat

Sound Classes

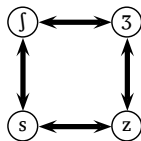
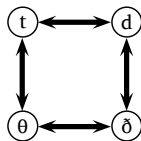
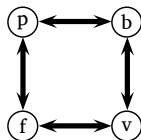
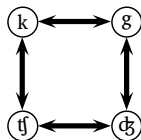
Sounds which frequently occur in correspondence relations in genetically related languages can be divided in classes (types). It is thereby assumed that “phonetic correspondences inside a ‘type’ are more regular than those between different ‘types’” (Dolgoposky 1986[1966]: 35).

Ⓚ	ⓖ	Ⓟ	Ⓟ
Ⓣ	Ⓞ	ⓕ	Ⓥ
Ⓣ	Ⓞ	ⓕ	Ⓥ
Ⓣ	Ⓞ	ⓕ	Ⓥ

LexStat

Sound Classes

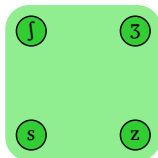
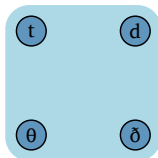
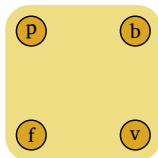
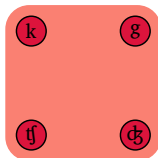
Sounds which frequently occur in correspondence relations in genetically related languages can be divided in classes (types). It is thereby assumed that “phonetic correspondences inside a ‘type’ are more regular than those between different ‘types’” (Dolgoposky 1986[1966]: 35).



LexStat

Sound Classes

Sounds which frequently occur in correspondence relations in genetically related languages can be divided in classes (types). It is thereby assumed that “phonetic correspondences inside a ‘type’ are more regular than those between different ‘types’” (Dolgoposky 1986[1966]: 35).



LexStat

Sound Classes

Sounds which frequently occur in correspondence relations in genetically related languages can be divided in classes (types). It is thereby assumed that “phonetic correspondences inside a ‘type’ are more regular than those between different ‘types’” (Dolgoposky 1986[1966]: 35).

**K****P****T****S**

LexStat

LexStat

Prosodic Strings

Sound change occurs more frequently in weak positions of sound sequences (Geisler 1992). Based on a sonority profile of sound sequences, one can distinguish sound positions according to their prosodic contexts. Prosodic context can be modeled as prosodic string in which different contexts are coded by different symbols.

LexStat

Prosodic Strings

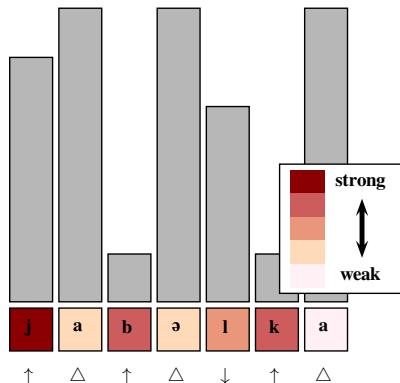
Sound change occurs more frequently in weak positions of sound sequences (Geisler 1992). Based on a sonority profile of sound sequences, one can distinguish sound positions according to their prosodic contexts. Prosodic context can be modeled as prosodic string in which different contexts are coded by different symbols.

j a b ə l k a

LexStat

Prosodic Strings

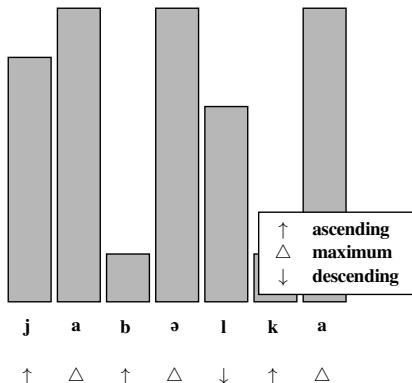
Sound change occurs more frequently in weak positions of sound sequences (Geisler 1992). Based on a sonority profile of sound sequences, one can distinguish sound positions according to their prosodic contexts. Prosodic context can be modeled as prosodic string in which different contexts are coded by different symbols.



LexStat

Prosodic Strings

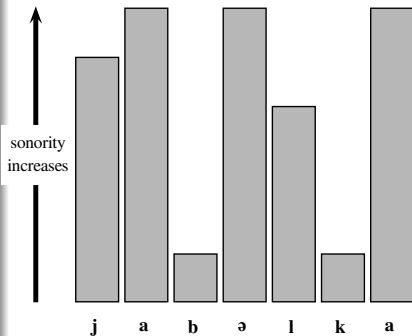
Sound change occurs more frequently in weak positions of sound sequences (Geisler 1992). Based on a sonority profile of sound sequences, one can distinguish sound positions according to their prosodic contexts. Prosodic context can be modeled as prosodic string in which different contexts are coded by different symbols.



LexStat

Prosodic Strings

Sound change occurs more frequently in weak positions of sound sequences (Geisler 1992). Based on a sonority profile of sound sequences, one can distinguish sound positions according to their prosodic contexts. Prosodic context can be modeled as prosodic string in which different contexts are coded by different symbols.



LexStat

Prosodic Strings

Sound change occurs more frequently in weak positions of sound sequences (Geisler 1992). Based on a sonority profile of sound sequences, one can distinguish sound positions according to their prosodic contexts. Prosodic context can be modeled as prosodic string in which different contexts are coded by different symbols.

j a b ə l k a
v C v c C >

LexStat

LexStat

External Representation							
IPA	j	a	b	ə	l	k	a

Internal Representation							
Sound-Class String	J	A	P	E	L	K	A
Prosodic String	#	V	C	V	c	C	>

LexStat

LexStat

Cognate List		Alignment				Correspondence List		
German	<i>Zunge</i>	ts	ʊ	ŋ	ə	GER	ENG	Frequ.
English	<i>tongue</i>	t	ʌ	ŋ	-	ts	t	2 x
German	<i>Zahn</i>	ts	a:	n	-	s	t	2 x
English	<i>tooth</i>	t	u:	-	θ	h	h	1 x
German	<i>heiß</i>	h	ai	s		f	f	1 x
English	<i>hot</i>	h	ɔ	t		n	-	1 x
German	<i>Fuß</i>	f	u:	s	
English	<i>foot</i>	f	ʊ	t				

LexStat

Cognate List		Alignment				Correspondence List		
German	<i>Zunge</i>	ts	ʊ	ŋ	ə	GER	ENG	Frequ.
English	<i>tongue</i>	t	ʌ	ŋ	-	ts	t	2 x
German	<i>Zahn</i>	ts	a:	n	-	s	t	2 x
English	<i>tooth</i>	t	u:	-	θ	h	h	1 x
German	<i>heiß</i>	h	ai	s		f	f	1 x
English	<i>hot</i>	h	ɔ	t		n	-	1 x
German	<i>Fuß</i>	f	u:	s	
English	<i>foot</i>	f	ʊ	t				

LexStat

Cognate List		Alignment				Correspondence List		
German	<i>Zunge</i>	C	U	N	E	GER	ENG	Frequ.
English	<i>tongue</i>	T	A	N	-	C/#	T/#	2 x
German	<i>Zahn</i>	C	A	N	-	S/\$	T/\$	2 x
English	<i>tooth</i>	T	U	-	T	H/\$	H/#	1 x
German	<i>heiß</i>	H	A	S		B/\$	B/#	1 x
English	<i>hot</i>	H	O	T		N/c	-	1 x
German	<i>Fuß</i>	B	U	S	
English	<i>foot</i>	B	U	T				

LexStat

“to dig” (30)			Turchin	NED	LexStat
Albanisch	<i>gërmon</i>	gərmo	1	1	1
Englisch	<i>digs</i>	dɪg	2	2	2
Französisch	<i>creuse</i>	krøze	1	3	3
Deutsch	<i>gräbt</i>	gra:b	1	1	4
Hawaii	<i>‘eli</i>	ʔeli	5	5	5
Navajo	<i>hahashgééd</i>	hahage:d	6	6	6
Türkisch	<i>kazıyor</i>	kaz	7	3	7

Dataset of Kessler (2001)

LexStat

"mouth" (104)			Turchin	NED	LexStat
Albanisch	<i>gojë</i>	goj	1	1	1
Englisch	<i>mouth</i>	mauθ	2	2	2
Französisch	<i>bouche</i>	buʃ	3	3	3
Deutsch	<i>Mund</i>	mund	4	4	2
Hawaii	<i>waha</i>	waha	5	5	5
Navajo	'azéé'	ze:ʔ	6	6	6
Türkisch	<i>ağız</i>	ayz	7	7	7

Dataset of Kessler (2001)

Testing the Impact of Sample Size on Cognate Detection

Gold Standard

IDS-Testset

- 4 languages (German, English, Dutch, French)
- 550 items (glosses)
- translations taken from the IDS (Key & Comrie 2009)
- orthographic entries converted into IPA transcriptions
- cognate judgments follow traditional literature

Subsets of Varying Sample Size

Creating the Subsets

Starting from the basic dataset, subsets of the data were created by

- randomly deleting 5, 10, 15, etc. items from the original dataset, and
- taking 5 different samples for each distinct number of deletions.

This process yielded 550 datasets, covering the whole range of possible sample sizes between 5 and 550 in steps of 5.

Automatic Cognate Detection

Methods for Cognate Detection

- Normalized Edit Distance (NED)
- Turchin et al. (2010, Turchin)
- SCA Distance (List 2012b)
- LexStat (List 2012a)

Implementation

All methods are implemented as part of LingPy-1.0 (see <http://lingpy.org>), a Python library for quantitative tasks in historical linguistics.

Evaluation Measures

B-Cubed Precision and Recall (Amigó et al. 2009)

Given a test (result of an analysis) and a reference (the gold standard),

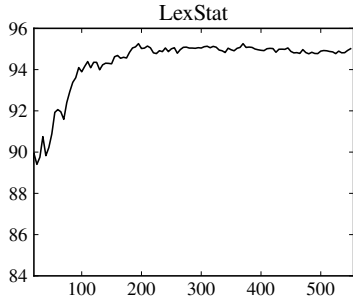
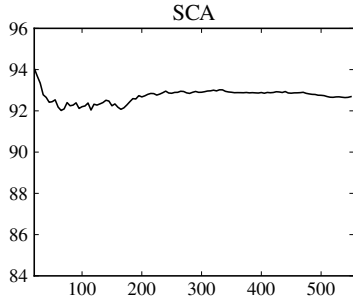
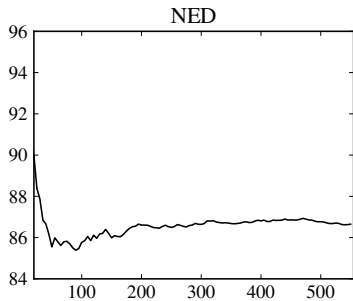
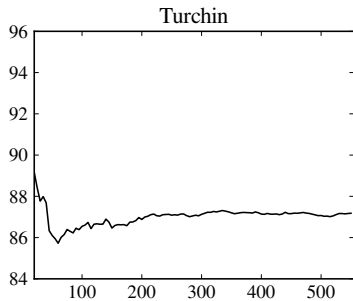
- precision is the proportion of items in the test that also occur in the reference, and
- recall is the proportion of items in the reference that also occur in the test.

Low precision is equivalent to high rates of false positives, low recall is equivalent to high rates of false negatives (missed cognates).

Results

Results

Items	B-Cubed Recall			
	Turchin	NED	SCA	LexStat
50	86.10	85.55	92.44	90.88
100	86.55	85.77	92.20	93.89
200	86.88	86.61	92.68	95.02
300	87.13	86.64	92.90	95.05
400	87.14	86.81	92.89	94.94
500	87.07	86.77	92.75	94.90



Discussion

Discussion

Are 200 words enough?

Although

- the representativity of the data is limited, and
- the number of languages investigated is small,

the test shows that

- sample size has a definite impact on the results of language-specific methods, and
- using 200 words is surely better than using 100 words.

Sanscruta	<i>sarpá-</i>	s	a	r	p	a
Italienisch	<i>serpe</i>	s	ε	r	p	ə
Sanscruta	<i>devá-</i>	d	e	v	a	
Italienisch	<i>Dio</i>	d	i	-	o	
Sanscruta	<i>saptá-</i>	s	a	p	t	a
Italienisch	<i>sette</i>	s	ε	-	t:	ə

Спасибо за Ваше Внимание!