LexStat

Evaluation

LexStat: Automatic Detection of Cognates in Multilingual Wordlists

Johann-Mattis List*

* Institute for Romance Languages and Literature Heinrich Heine University Düsseldorf

April 24, 2012

Structure of the Talk



Keys to the Past

Identification of Cognates





Keys to the Past

Identification of Cognates

LexStat

Evaluation



Keys to the Past



AAGTCTAGGTTACCCGTAT



LexStat

Evaluation

Charles Lyell on Languages

LexStat

Charles Lyell on Languages



LexStat

Charles Lyell on Languages



LexStat

Evaluation

Uniformitarianism and Abduction

LexStat

Evaluation

Uniformitarianism and Abduction

Uniformitarianism

LexStat

Evaluation

Uniformitarianism and Abduction

Uniformitarianism

 "Universality of Change" – Change is independent of time and space

Uniformitarianism and Abduction

Uniformitarianism

- "Universality of Change" Change is independent of time and space
- "Graduality of Change" Change is neither abrupt nor chaotic

Uniformitarianism and Abduction

Uniformitarianism

- "Universality of Change" Change is independent of time and space
- "Graduality of Change" Change is neither abrupt nor chaotic
- "Uniformity of Change" Change is not heterogeneous

Uniformitarianism and Abduction

Uniformitarianism

- "Universality of Change" Change is independent of time and space
- "Graduality of Change" Change is neither abrupt nor chaotic
- "Uniformity of Change" Change is not heterogeneous

Abduction

Uniformitarianism and Abduction

Uniformitarianism

- "Universality of Change" Change is independent of time and space
- "Graduality of Change" Change is neither abrupt nor chaotic
- "Uniformity of Change" Change is not heterogeneous

Abduction

- Present Events or Patterns
- + Known Laws
- => Abduction of Historical Facts

Uniformitarianism and Abduction

Uniformitarianism

- "Universality of Change" Change is independent of time and space
- "Graduality of Change" Change is neither abrupt nor chaotic
- "Uniformity of Change" Change is not heterogeneous

Abduction

- Present Events or Patterns
- + Known Laws
- => Abduction of Historical Facts

Similarities Between Languages

- + Language Change
- => Inference of Proto-Languages



LexStat

The Comparative Method

LexStat

The Comparative Method

Basic Procedure

• Compile an initial list of putative cognate sets.

The Comparative Method

- Compile an initial list of putative cognate sets.
- Extract an initial list of putative sets of sound correspondences from the initial cognate list.

The Comparative Method

- Compile an initial list of putative cognate sets.
- Extract an initial list of putative sets of sound correspondences from the initial cognate list.
- Refine the cognate list and the correspondence list by

The Comparative Method

- Compile an initial list of putative cognate sets.
- Extract an initial list of putative sets of sound correspondences from the initial cognate list.
- Refine the cognate list and the correspondence list by
 - adding and deleting cognate sets from the cognate list, depending on whether they are consistent with the correspondence list or not, and

The Comparative Method

- Compile an initial list of putative cognate sets.
- Extract an initial list of putative sets of sound correspondences from the initial cognate list.
- Refine the cognate list and the correspondence list by
 - adding and deleting cognate sets from the cognate list, depending on whether they are consistent with the correspondence list or not, and
 - adding and deleting correspondence sets from the correspondence list, depending on whether they are consistent with the cognate list or not.

The Comparative Method

- Compile an initial list of putative cognate sets.
- Extract an initial list of putative sets of sound correspondences from the initial cognate list.
- Refine the cognate list and the correspondence list by
 - adding and deleting cognate sets from the cognate list, depending on whether they are consistent with the correspondence list or not, and
 - adding and deleting correspondence sets from the correspondence list, depending on whether they are consistent with the cognate list or not.
- Finish when the results are satisfying enough.

LexStat

Evaluation

The Comparative Method

The Comparative Method

Language-Specific Similarity Measure

 Sequence similarity is determined on the basis of systematic sound correspondences as opposed to similarity based on surface resemblances of phonetic segments.

The Comparative Method

- Sequence similarity is determined on the basis of systematic sound correspondences as opposed to similarity based on surface resemblances of phonetic segments.
- Lass (1997) calls this notion of similarity phenotypic as opposed to a genotypic notion of similarity.

The Comparative Method

- Sequence similarity is determined on the basis of systematic sound correspondences as opposed to similarity based on surface resemblances of phonetic segments.
- Lass (1997) calls this notion of similarity phenotypic as opposed to a genotypic notion of similarity.
- The most crucial aspect of correspondence-based similarity is that it is *language-specific*: Genotypic similarity is never defined in general terms but always with respect to the language systems which are being compared.

The Comparative Method

- Sequence similarity is determined on the basis of systematic sound correspondences as opposed to similarity based on surface resemblances of phonetic segments.
- Lass (1997) calls this notion of similarity *phenotypic* as opposed to a *genotypic* notion of similarity.
- The most crucial aspect of correspondence-based similarity is that it is *language-specific*: Genotypic similarity is never defined in general terms but always with respect to the language systems which are being compared.

Meaning	German	Dutch	English
"tooth"	Zahn [<mark>ts</mark> a:n]	tand [<mark>t</mark> ant]	<i>tooth</i> [<mark>t</mark> υ:θ]
"ten"	zehn [<mark>ts</mark> e:n]	<i>tien</i> [<mark>t</mark> i:n]	ten [<mark>t</mark> ɛn]
"tongue"	Zunge [<mark>ts</mark> ʊŋə]	<i>tong</i> [<mark>t</mark> ວŋ]	<i>tongue</i> [<mark>t</mark> ող]

The Comparative Method

- Sequence similarity is determined on the basis of systematic sound correspondences as opposed to similarity based on surface resemblances of phonetic segments.
- Lass (1997) calls this notion of similarity phenotypic as opposed to a genotypic notion of similarity.
- The most crucial aspect of correspondence-based similarity is that it is *language-specific*: Genotypic similarity is never defined in general terms but always with respect to the language systems which are being compared.

Meaning	Shanghai	Beijing	Guangzhou
"nine"	[<mark>tç</mark> ir ³⁵]	Beijing [<mark>tc</mark> iou ²¹⁴]	[<mark>k</mark> eu ³⁵]
"today"	$\begin{bmatrix} tc i n^{55} ts 2^{21} \end{bmatrix}$	Beijing [<mark>tç</mark> ið ⁵⁵]	[<mark>k</mark> em ⁵³ jet ²]
"rooster"	[koŋ ⁵⁵ tç i ²¹]	Beijing[kuŋ ⁵⁵ tç i ⁵⁵]	[<mark>k</mark> ei ⁵⁵ koŋ ⁵⁵]

LexStat

Evaluation

Automatic Approaches

LexStat

Automatic Approaches

Alignment Analyses

Automatic Approaches

Alignment Analyses

Automatic Approaches

Alignment Analyses





Automatic Approaches

Alignment Analyses



Automatic Approaches

Alignment Analyses



d

Automatic Approaches

Alignment Analyses

Collaboration the algument. In alignment analyses, sequences such a way that corresponding column, while empty cells reg elements are filled with ga

trix in same sponding

LexStat

Evaluation

Automatic Approaches

Sound Classes
LexStat

Evaluation

Automatic Approaches

Sound Classes

LexStat

Automatic Approaches

Sound Classes

k	Ø	p	b
(ť)	¢	ſ	v
t	d		3
θ	ð	s	Z

LexStat

Evaluation

Automatic Approaches

Sound Classes



LexStat

Evaluation

Automatic Approaches

Sound Classes



LexStat

Evaluation

Automatic Approaches

Sound Classes



Automatic Approaches

Sound Classes

Paine the first two consolities of two winds wods ale judged to be cognate offennise not her nach coating her sound classes Sounds which often occur in correspondence relations in genetically related langua can be clustered into cl Concertaint (types). It is assume phonetic correspo inside a 'type' than those b 'types'" (Dolge

S

LexStat

Automatic Approaches

Sound-Class-Based Alignment (SCA)

LexStat

Automatic Approaches

Sound-Class-Based Alignment (SCA)

Sound classes and alignment analyses can be easily combined by representing phonetic sequences internally as sound classes and comparing the sound classes with traditional alignment algorithms.

LexStat

Automatic Approaches

Sound-Class-Based Alignment (SCA)

Sound classes and alignment analyses can be easily combined by representing phonetic sequences internally as sound classes and comparing the sound classes with traditional alignment algorithms.



000000

Automatic Approaches

Sound-Class-Based Alignment (SCA

	et a
Sound-Class-Based Alignment (SCA)	Hondi
Sound classes and alignment analy and the sound classes and alignment analy and the sound classes are the sound classes of the sound cl	alignment
INPUT toxtər doxtər → T 0 INPUT toxtər doxtər → T 0 INPUT toxtər doxtər	N G T
$\begin{array}{c} container = 1\\ container = 1\\ tain the similar \\ th$	$ \begin{array}{c} N \\ x \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \end{array} \begin{array}{c} OUTPUT \\ t \\ z \\ x \\ t \\ z \\ x \\ t \\ z \\ \end{array} \begin{array}{c} output \\ z \\ $

LexStat

Evaluation

Traditional vs. Automatic Approaches

LexStat

Evaluation

Traditional vs. Automatic Approaches

Similarity

Traditional vs. Automatic Approaches

Similarity

Almost all current automatic approaches are based on a language-independent similarity measure, while the comparative method applies a language-specific one. All automatic approaches will therefore yield the same scores for phenotypically identical sequences, regardless of the language systems they belong to.



	LexStat	
	000000	

Working Procedure

	LexStat	
	000000	

Working Procedure

Sequence Input

sequences are read from specifically formatted input files

		LexStat ●oooooo	
Working Proced	ure		
Sequence Ir	nput	sequences are read from specifically fo matted input files	r-

1 Sequence Conversion

sequences are converted to sound classes and prosodic profiles

LexStat •oooooo

Working Procedure

Sequence Input

sequences are read from specifically formatted input files

1 Sequence Conversion

2 Scoring-Scheme Creation

sequences are converted to sound classes and prosodic profiles

using a permutation method, languagespecific scoring schemes are determined

LexStat

Working Procedure

Sequence Input

sequences are read from specifically formatted input files

- 1 Sequence Conversion
- 2 Scoring-Scheme Creation
- 3 Distance Calculation

sequences are converted to sound classes and prosodic profiles

using a permutation method, languagespecific scoring schemes are determined

based on the language-specific scoringscheme, pairwise distances between sequences are calculated

LexStat

Working Procedure

Sequence Input

sequences are read from specifically formatted input files

- 1 Sequence Conversion
- 2 Scoring-Scheme Creation
- 3 Distance Calculation
- 4 Sequence Clustering

sequences are converted to sound classes and prosodic profiles

using a permutation method, languagespecific scoring schemes are determined

based on the language-specific scoringscheme, pairwise distances between sequences are calculated

sequences are clustered into cognate sets whose average distance is beyond a certain threshold

LexStat

Working Procedure

Sequence Input	sequences are read from specifically for- matted input files

- 1 Sequence Conversion
- 2 Scoring-Scheme Creation
- 3 Distance Calculation
- 4 Sequence Clustering

sequences are converted to sound classes and prosodic profiles

using a permutation method, languagespecific scoring schemes are determined

based on the language-specific scoringscheme, pairwise distances between sequences are calculated

sequences are clustered into cognate sets whose average distance is beyond a certain threshold

Sequence Output

information regarding sequence clustering is written to file using a specific format

	LexStat	
	000000	

Implementation

	LexStat o●ooooo	
Implementation		

• LexStat ist implemented as part of the LingPy Python library (see http://lingulist.de/lingpy) for automatic tasks in historical linguistics.

	LexStat o●ooooo	
Implementation		

- LexStat ist implemented as part of the LingPy Python library (see http://lingulist.de/lingpy) for automatic tasks in historical linguistics.
- The current release of LingPy (lingpy-1.0) provides methods for pairwise and multiple sequence alignment (SCA), automatic cognate detection (LexStat), and plotting routines (see the online documentation for details).

	LexStat o●ooooo	
Implementation		

- LexStat ist implemented as part of the LingPy Python library (see http://lingulist.de/lingpy) for automatic tasks in historical linguistics.
- The current release of LingPy (lingpy-1.0) provides methods for pairwise and multiple sequence alignment (SCA), automatic cognate detection (LexStat), and plotting routines (see the online documentation for details).
- LexStat can be invoked from the Python shell or inside Python scripts (examples are given in the online documentation).

LexStat

Evaluation

Input and Output

ID	Items	German	English	Swedish
1	hand	hant	hænd	hand
2	woman	fraʊ	wʊmən	kvina
3	know	kɛnən	ทอช	çɛna
3	know	vīsən	-	veːta

LexStat

Evaluation

Input and Output

ID	Items	German	COG	English	COG	Swedish	COG
1	hand	hant	1	hænd	1	hand	1
2	woman	fraʊ	2	wʊmən	3	kvina	4
3	know	kɛnən	5	ทอซ	5	çɛna	5
3	know	vīsən	6	-	0	veːta	6

	LexStat	
	000000	

Input and Output

Basic Concept: <i>belly</i> (ID: 4)					
CogID	Language	Entry	Aligned Entry		
6	Danish	on∧liw'			
7	German	baux	b au x		
7	Dutch	bœyk	b œy k		
7	Swedish	buk	b u k		
7	Norwegian	b u :k	b u: k		
8	English	bɛlɪ			
9	Swedish	ma:ge	ma:ge		
9	Norwegian	ma : gə	ma:gə		
9	Danish	mæːvə	m æ: v ə		
10	Icelandic	k ^h vī:ðyr			

	LexStat	
	000000	

Internal Representation of Sequences

LexStat

Evaluation

Internal Representation of Sequences

Sound Classes and Prosodic Context

Internal Representation of Sequences

Sound Classes and Prosodic Context

• All sequences are internally represented as sound classes, the default model being the one proposed in List (forthcoming).

Internal Representation of Sequences

Sound Classes and Prosodic Context

- All sequences are internally represented as sound classes, the default model being the one proposed in List (forthcoming).
- All sequences are also represented by *prosodic strings* which indicate the prosodic environment (initial, ascending, maximum, descending, final) of each phonetic segment (List 2012).

LexStat

Internal Representation of Sequences

Sound Classes and Prosodic Context

- All sequences are internally represented as sound classes, the default model being the one proposed in List (forthcoming).
- All sequences are also represented by *prosodic strings* which indicate the prosodic environment (initial, ascending, maximum, descending, final) of each phonetic segment (List 2012).
- The information regarding sound classes and prosodic context is combined, and each input sequence is further represented as a sequence of tuples, consisting of the sound class and the prosodic environment of the respective phonetic segment.

LexStat

Evaluation

Scoring-Scheme Creation

LexStat

Evaluation

Scoring-Scheme Creation

Attested Distribution

LexStat

Scoring-Scheme Creation

Attested Distribution

- carry out global and pairwise alignment analyses of all sequence pairs occuring in the same semantic slot
- store all corresponding segments that occur in sequences whose distance is beyond a certain threshold
Scoring-Scheme Creation

Attested Distribution

- carry out global and pairwise alignment analyses of all sequence pairs occuring in the same semantic slot
- store all corresponding segments that occur in sequences whose distance is beyond a certain threshold

Creation of the Expected Distribution

Scoring-Scheme Creation

Attested Distribution

- carry out global and pairwise alignment analyses of all sequence pairs occuring in the same semantic slot
- store all corresponding segments that occur in sequences whose distance is beyond a certain threshold

Creation of the Expected Distribution

- shuffle the wordlists repeatedly and
 - carry out global and pairwise alignment analyses of all sequence pairs in the randomly shuffled wordlists
 - store all corresponding segments
- average the results

Scoring-Scheme Creation

Attested Distribution

- carry out global and pairwise alignment analyses of all sequence pairs occuring in the same semantic slot
- store all corresponding segments that occur in sequences whose distance is beyond a certain threshold

Creation of the Expected Distribution

- shuffle the wordlists repeatedly and
 - carry out global and pairwise alignment analyses of all sequence pairs in the randomly shuffled wordlists
 - store all corresponding segments
- average the results

Calculation of Similarity Scores

Scoring-Scheme Creation

Attested Distribution

- carry out global and pairwise alignment analyses of all sequence pairs occuring in the same semantic slot
- store all corresponding segments that occur in sequences whose distance is beyond a certain threshold

Creation of the Expected Distribution

- shuffle the wordlists repeatedly and
 - carry out global and pairwise alignment analyses of all sequence pairs in the randomly shuffled wordlists
 - store all corresponding segments
- average the results

Calculation of Similarity Scores

• Calculation of *log-odds scores* from the distributions.

LexStat

Scoring-Scheme Creation

English	German	Att.	Exp.	Score
#[t,d]	#[t,d]	3.0	1.24	6.3
#[t,d]	#[ts]	3.0	0.38	6.0
#[t,d]	#[∫,s,z]	1.0	1.99	-1.5
#[θ,ð]	#[t,d]	7.0	0.72	6.3
#[θ,ð]	#[ts]	0.0	0.25	-1.5
#[θ,ð]	#[s,z]	0.0	1.33	0.5
[t,d]\$	[t,d]\$	21.0	8.86	6.3
[t,d]\$	[ts]\$	3.0	1.62	3.9
[t,d]\$	[∫,s]\$	6.0	5.30	1.5
[θ,ð]\$	[t,d]\$	4.0	1.14	4.8
[θ,ð]\$	[ts]\$	0.0	0.20	-1.5
[θ,ð]\$	[∫,s]\$	0.0	0.80	0.5

LexStat

Scoring-Scheme Creation

English	German	Att.	Exp.	Score
#[t,d]	#[t,d]	3.0	1.24	6.3
#[t,d]	#[ts]	3.0	0.38	6.0
#[t,d]	#[∫ ,s,z]	1.0	1.99	-1.5
#[θ,ð]	#[t,d]	7.0	0.72	6.3
#[θ,ð]	#[ts]	0.0	0.25	-1.5
#[θ,ð]	#[s,z]	0.0	1.33	0.5
[t,d]\$	[t,d]\$	21.0	8.86	6.3
[t,d]\$	[ts]\$	3.0	1.62	3.9
[t,d]\$	[∫,s]\$	6.0	5.30	1.5
[θ,ð]\$	[t,d]\$	4.0	1.14	4.8
[θ,ð]\$	[ts]\$	0.0	0.20	-1.5
[θ,ð]\$	[∫,s]\$	0.0	0.80	0.5

LexStat

Evaluation

Scoring-Scheme Creation

	Initial	Final
English	<i>town</i> [taun]	hot [hət]
German	Zaun [tsaun]	heiß [haɪs]
English	thorn [θວːn]	<i>mouth</i> [maυθ]
German	Dorn [dərn]	Mund [mont]
English	dale [deɪl]	head [hɛd]
German	<i>Tal</i> [ta:1]	Hut [hu:t]

LexStat

Evaluation

Sequence Clustering

	Ger.	Eng.	Dan.	Swe.	Dut.	Nor.
Ger. [frau]	0.00	0.95	0.81	0.70	0.34	1.00
Eng. [wumən]	0.95	0.00	0.78	0.90	0.80	0.80
Dan. [kvenə]	0.81	0.78	0.00	0.17	0.96	0.13
Swe. [kvin:a]	0.70	0.90	0.17	0.00	0.86	0.10
Dut. [vrauv]	0.34	0.80	0.96	0.86	0.00	0.89
Nor. [kvinə]	1.00	0.80	0.13	0.10	0.89	0.00

LexStat

Evaluation

Sequence Clustering

	Ger.	Eng.	Dan.	Swe.	Dut.	Nor.
Ger. [frau]	0.00	0.95	0.81	0.70	0.34	1.00
Eng. [wumən]	0.95	0.00	0.78	0.90	0.80	0.80
Dan. [kvenə]	0.81	0.78	0.00	0.17	0.96	0.13
Swe. [kvin:a]	0.70	0.90	0.17	0.00	0.86	0.10
Dut. [vrauv]	0.34	0.80	0.96	0.86	0.00	0.89
Nor. [kvinə]	1.00	0.80	0.13	0.10	0.89	0.00
Clusters	1	2	3	3	1	3





LexStat

Evaluation •oooo

Gold Standard

LexStat

Evaluation •oooo

Gold Standard

File	Family	Lng.	ltm.	Entr.	Source
GER	Germanic	7	110	814	Starostin (2008)
ROM	Romance	5	110	589	Starostin (2008)
SLV	Slavic	4	110	454	Starostin (2008)
PIE	Indo-Eur.	18	110	2057	Starostin (2008)
OUG	Uralic	21	110	2055	Starostin (2008)
BAI	Bai	9	110	1028	Wang (2006)
SIN	Sinitic	9	180	1614	Hóu (2004)
KSL	varia	8	200	1600	Kessler (2001)
JAP	Japonic	10	200	1986	Shirō (1973)

LexStat

Evaluation

Evaluation Measures

LexStat

Evaluation Measures

Set Comparison

LexStat

Evaluation

Evaluation Measures

Set Comparison

Precision, Recall, and F-Score are calculated by comparing the cognate sets proposed by the method with the cognate sets in the gold standard (see Bergsma & Kondrak 2007).

LexStat

Evaluation

Evaluation Measures

Set Comparison

Precision, Recall, and F-Score are calculated by comparing the cognate sets proposed by the method with the cognate sets in the gold standard (see Bergsma & Kondrak 2007).

Pair Comparison

LexStat

Evaluation

Evaluation Measures

Set Comparison

Precision, Recall, and F-Score are calculated by comparing the cognate sets proposed by the method with the cognate sets in the gold standard (see Bergsma & Kondrak 2007).

Pair Comparison

Pair comparison is based on a pairwise comparison of all decisions present in testset and goldstandard.

	Evaluation
	00000

	Evaluation

 Sound Classes – matching sound classes without alignment (based on Turchin et al. 2010)

- Sound Classes matching sound classes without alignment (based on Turchin et al. 2010)
- Simple Alignment normalized edit-distance (Levenshtein 1966)

- Sound Classes matching sound classes without alignment (based on Turchin et al. 2010)
- Simple Alignment normalized edit-distance (Levenshtein 1966)
- SCA language-independent distance scores derived from sound-class-based alignment analyses (List 2012)

- Sound Classes matching sound classes without alignment (based on Turchin et al. 2010)
- Simple Alignment normalized edit-distance (Levenshtein 1966)
- SCA language-independent distance scores derived from sound-class-based alignment analyses (List 2012)
- LexStat language-specific distance scores

LexStat

Evaluation

General Results

LexStat

Evaluation

General Results

Score	LexStat	SCA	Simple Alm.	Sound Cl.
Identical Pairs	0.85	0.82	0.76	0.74
Precision	0.59	0.51	0.39	0.39
Recall	0.68	0.57	0.47	0.55
F-Score	0.63	0.55	0.42	0.46

LexStat

Evaluation

General Results



LexStat

Evaluation

LexStat

Evaluation

Specific Results

• Pairwise decisions were extracted from the KSL dataset and compared with the Gold Standard.

- Pairwise decisions were extracted from the KSL dataset and compared with the Gold Standard.
- 72 borrowings were explicitly marked along with their source by Kessler (2001).

- Pairwise decisions were extracted from the KSL dataset and compared with the Gold Standard.
- 72 borrowings were explicitly marked along with their source by Kessler (2001).
- 83 chance resemblances were determined automatically by taking non-cognate word pairs with an NED score less than 0.6.

- Pairwise decisions were extracted from the KSL dataset and compared with the Gold Standard.
- 72 borrowings were explicitly marked along with their source by Kessler (2001).
- 83 chance resemblances were determined automatically by taking non-cognate word pairs with an NED score less than 0.6.

	LexStat	SCA	Simple Alm.	Sound Cl.
Borrowings	50%	61%	49%	53%
Chance Resemblances	17%	42%	89%	31%

LexStat

Evaluation 00000





