

Sequence Alignment

Gerhard Jäger

ESLLI 2016

Sequence alignment: Motivation

Motivation

- suppose we have no information except word lists
- goals:
 - estimate distances between languages
 - estimate cognate classes
 - track individual sound changes

Example

Meaning	Italian	English	cognate
<i>few</i>	'pɔko	fju:	1
<i>rub</i>	fre'gare	rʌb	0
<i>dull</i>	ot'tuzo	dʌl	0
<i>hunt</i>	kat'tʃare	hʌnt	0
<i>year</i>	'anno	jɪə	0
<i>this</i>	'kwesto	ðɪs	0
<i>fish</i>	'peʃfe	fɪʃ	1
<i>rotten</i>	'martʃo	'rɒtən	0
<i>right</i>	'dʒusto	raɪt	0
<i>when</i>	'kwando	wen	1
<i>drink</i>	'bere	dɪŋk	0
<i>heavy</i>	pe'sante	'heɪvɪ	0
<i>heavy</i>	'grave	'heɪvɪ	0
<i>egg</i>	'wɔvo	ɛg	1
<i>earth</i>	'terra	ɜ:θ	0
<i>dust</i>	'polvere	dʌst	0
<i>laugh</i>	'ridere	lɑ:f	0
<i>grass</i>	'erba	grɑ:s	0
<i>sharp</i>	taʎ'ʎente	ʃɑ:p	0
<i>wash</i>	la'vare	wɒʃ	0

Motivation

- suppose we have no information except word lists
- goals:
 - estimate distances between languages
 - estimate cognate classes
 - track individual sound changes

Example

Meaning	Italian	English	cognate
<i>few</i>	'pɔko	fju:	
<i>rub</i>	fre'gare	rʌb	
<i>dull</i>	ot'tuzo	dʌl	
<i>hunt</i>	kat'tʃare	hʌnt	
<i>year</i>	'anno	jɪə	
<i>this</i>	'kwesto	ðɪs	
<i>fish</i>	'peʃʃe	fɪʃ	
<i>rotten</i>	'martʃo	'rɒtən	
<i>right</i>	'dʒusto	raɪt	
<i>when</i>	'kwando	wen	
<i>drink</i>	'bere	dɪŋk	
<i>heavy</i>	pe'sante	'heɪvɪ	
<i>heavy</i>	'grave	'heɪvɪ	
<i>egg</i>	'wɔvo	ɛɡ	
<i>earth</i>	'terra	ɜ:θ	
<i>dust</i>	'polvere	dʌst	
<i>laugh</i>	'ridere	lɑ:f	
<i>grass</i>	'erba	grɑ:s	
<i>sharp</i>	taʎ'ʎente	ʃɑ:p	
<i>wash</i>	la'vare	wɒʃ	

Preprocessing

- IPA is open-ended — 107 letters, 52 diacritics, 4 prosodic marks → 200,000 combinations
- good practice: map IPA strings to a uniform representation with fewer symbols
- common choices:
 - 10 Dolgopolsky sound classes (Dolgopolsky 1986; used i.a. in List 2014)
 - 41 ASJP sound classes
- this course: ASJP

No.	Class	Description	Example
1	P	labial obstruents	p,b,f
2	T	dental obstruents	d,t,θ,ð
3	S	sibilants	s,z,ʃ,ʒ
4	K	velar obstruents, dental and alveolar affricates	k,g,ʈ,ʣ
5	M	labial nasal	m
6	N	remaining nasals	n,ɲ,ŋ
7	R	liquids	r,l
8	W	voiced labial fricative and initial rounded vowels	v,u
9	J	palatal approximant	j
10	ø	laryngeals and initial velar nasal	h,ɦ,ŋ

Preprocessing

- IPA is open-ended — 107 letters, 52 diacritics, 4 prosodic marks → 200,000 combinations
- good practice: map IPA strings to a uniform representation with fewer symbols
- common choices:
 - 10 Dolgopolsky sound classes (Dolgopolsky 1986; used i.a. in List 2014)
 - 41 ASJP sound classes
- this course: ASJP

Meaning	Italian	English
<i>few</i>	poko	fyu
<i>rub</i>	fregare	rob
<i>dull</i>	ottuzo	dol
<i>hunt</i>	kattSare	hunt
<i>year</i>	anno	yi3
<i>this</i>	kwesto	8is
<i>fish</i>	peSSe	fiS
<i>rotten</i>	martSo	rot3n
<i>right</i>	dZusto	rait
<i>when</i>	kwando	wEn
<i>drink</i>	bere	driNk
<i>heavy</i>	pesante	hEvi
<i>heavy</i>	grEve	hEvi
<i>egg</i>	wovo	Eg
<i>earth</i>	tErra	38
<i>dust</i>	polvere	dost
<i>laugh</i>	ridere	lof
<i>grass</i>	Erba	gros
<i>sharp</i>	tallEnte	Sop
<i>wash</i>	lavare	woS

Pairwise alignment

Levenshtein alignment

- related to as *edit distance*
- defines the distance between two strings as the minimal number of *edit operations* to transform one string into the other
- edit operations:
 - deletion
 - insertion
 - replacement
- example: grm. mEnS vs. Cimbrian menEs
 - ① mEnS → menS (replace)
 - ② menS → menES (insert)
 - ③ menES → menEs (insert)
- $d_L(\text{mEnS}, \text{menEs}) = 3$

Levenshtein alignment

- alternative presentation: alignment

m	E	n	—	S
m	e	n	E	s

- distance for a particular alignment is the number of non-identities
- Levenshtein distance is the number of mismatches for the optimal alignment

Computing the Levenshtein Distance

- recursive definition:

- 1 $d_L(\epsilon, \alpha) = d_L(\alpha, \epsilon) = l(\alpha)$

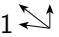
- 2

$$d_L(\alpha x, \beta y) = \min \begin{cases} d_L(\alpha, \beta) + \delta(x, y) \\ d_L(\alpha x, \beta) + 1 \\ d_L(\alpha, \beta y) + 1 \end{cases}$$

- apparently requires exponentially growing number of comparisons \Rightarrow computationally not feasible
- but:
 - if $l(\alpha) = n$ and $l(\beta) = m$, there are $n + 1$ substrings of α and $m + 1$ substrings of β
 - hence there are only $(n + 1)(m + 1)$ many different comparisons to be performed
 - computational complexity is polynomial (quadratic in $l(\alpha) + l(\beta)$)


Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1				
e	2				
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1				
e	2				
n	3				
E	4				
s	5				

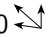
Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0			
e	2				
n	3				
E	4				
s	5				


Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0			
e	2				
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0			
e	2				
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1		
e	2				
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	
e	2				
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2				
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1			
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1		
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3				
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2			
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2		
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4				
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3			
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2		
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5				

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4			

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3		

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

Arrows in the original image point from the cell (n, n) to (n, E) and from (E, E) to (E, n).

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Computing the Levenshtein distance

- Dynamic Programming

	–	m	E	n	S
–	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Computing the Levenshtein distance

- Dynamic Programming

	-	m	E	n	S
-	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

m E n - S
m e n E s

Computing the Levenshtein distance

	-	m	E	n	S
-	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

m E n - S
m e n E s

	-	m	E	n	S
-	0	1	2	3	4
m	1	0	1	2	3
e	2	1	1	2	3
n	3	2	2	1	2
E	4	3	2	2	2
s	5	4	3	3	3

m E n S -
m e n E s

Normalization for length

- grm. mEnS (*Mensch*, 'person') and Hindi manuSya are (partially) cognate
- grm. ze3n (*sehen*, 'see') and Hindi deg are not cognate
- still

$$d_L(\text{mEnS}, \text{manuSya}) = 4$$

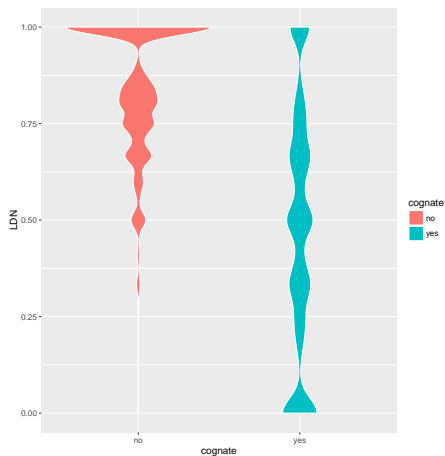
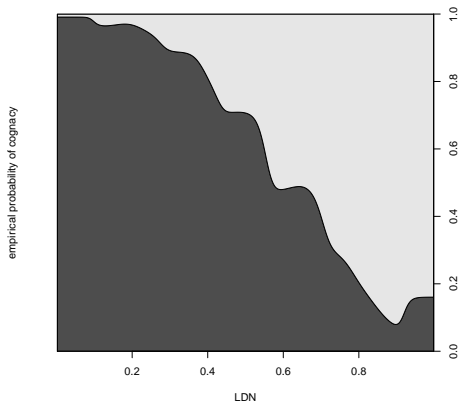
$$d_L(\text{ze3n}, \text{deg}) = 3$$

- normalization: dividing Levenshtein distance by length of longer string:

$$d_{LD}(\text{mEnS}, \text{manuSya}) = 4/7 \approx 0.57$$

$$d_{LD}(\text{ze3n}, \text{deg}) = 3/4 = 0.75$$

How well does normalized Levenshtein distance predict cognacy?



Problems

- binary distinction: match vs. non-match
- frequently genuine sound correspondences in cognates are missed:

c	v	a	i	n	a	z	3	-	-	-	f	i	S
-	-	t	u	n	-	o	s	p	i	s	k	i	s

- corresponding sounds count as mismatches even if they are aligned correctly

h	a	n	t	h	a	n	t
h	E	n	d	m	a	n	o

- substantial amount of chance similarities

Background: probability theory

- Given two sequences: How likely is it that they are aligned?
- More general question: Given some data, and two competing hypotheses, how likely is it that the first hypothesis is correct?

Bayesian Inference!!!

- given:
 - data: d
 - hypotheses: h_1, h_0
 - model: $P(d|h_1), P(d|h_0)$

- wanted:

$$P(h_1|d) : P(h_0|d)$$

Bayesian inference

- Bayes Theorem:

$$P(h|d) = \frac{P(d|h)P(h)}{\sum_{h'} P(d|h')P(h')}$$

- ergo:

$$P(h_1|d) : P(h_0|d) = P(d|h_1)P(h_1) : P(d|h_0)P(h_0)$$

$$P(h_1|d) : P(h_0|d) = \frac{P(d|h_1) P(h_1)}{P(d|h_0) P(h_0)}$$

$$\log(P(h_1|d) : P(h_0|d)) = \log \frac{P(d|h_1)}{P(d|h_0)} + \log \frac{P(h_1)}{P(h_0)}$$

Bayesian inference

- suppose we have many independent data: $\vec{d} = d_1, \dots, d_n$

$$P(\vec{d}|h) = \prod_{i=1}^n P(d_i|h)$$

$$\log P(\vec{d}|h) = \sum_{i=1}^n \log P(d_i|h)$$

$$\log \frac{P(\vec{d}|h_1)}{P(\vec{d}|h_0)} = \sum_{i=1}^n \log \frac{P(d_i|h_1)}{P(d_i|h_0)}$$

$$\log(P(h_1|\vec{d}) : P(h_0|\vec{d})) = \sum_{i=1}^n \log \frac{P(d_i|h_1)}{P(d_i|h_0)} + \log \frac{P(h_1)}{P(h_0)}$$

Bayesian inference

- mein argument against using Bayes' rule: the **prior probabilities** $P(h_1), P(h_0)$ are not known
- there are various heuristics, but no generally accepted way to obtain them
- if n is large though, $\log P(h_1)/P(h_0)$ doesn't matter very much:¹

$$\log(P(h_1|\vec{d}) : P(h_0|\vec{d})) \approx \sum_{i=1}^n \log \frac{P(d_i|h_1)}{P(d_i|h_0)} = \log(P(\vec{d}|h_1) : P(\vec{d}|h_0))$$

- the quantity $\log(P(\vec{d}|h_1) : P(\vec{d}|h_0))$ is called **log-odds**

¹Also, if we choose an *uninformative prior* with $P(h_1) = P(h_0)$, we have $\log P(h_1)/P(h_0) = 0$ anyway.

Log-odds

- log-odds can take any real value
- a positive value indicates evidence for h_1 and a negative value evidence for h_0
- the higher the absolute value, the stronger is the evidence

Weighted alignment

- suppose our data are two aligned sequences \vec{x}, \vec{y}
- for the time being, we assume there are no gaps in the alignment
 - h_1 : they developed from a common ancestor via substitutions
 - h_0 : they are unrelated
- additional assumptions (rough approximation in biology, pretty much off the mark in linguistics): substitutions in different positions occur independently

The null model

- if \vec{x} and \vec{y} are unrelated, their joint probability equals the product of their individual probabilities
- as a start (quite wrong both in biology and in linguistics): let us assume the strings have no “grammar”; each position is independent from all other positions
- then

$$\begin{aligned}
 P(\vec{x}, \vec{y} | h_0) &= P(\vec{x} | h_0) P(\vec{y} | h_0) \\
 &= \prod_i P(x_i | h_0) P(y_i | h_0) \\
 \log P(\vec{x}, \vec{y} | h_0) &= \sum_i \log(P(x_i | h_0) + \log P(y_i | h_0))
 \end{aligned}$$

The null model

- suppose \vec{x} and \vec{y} are generated by the same process (reasonable for DNA and protein comparison, false for cross-linguistic word comparison)
- then $P(x_i|h), P(y_i|h)$ are simply the probabilities of occurrence
- q_a : probability that symbol a occurs in a sequence

$$\log P(\vec{x}, \vec{y}|h_0) = \sum_i \log q_{x_i} + \sum_j \log q_{y_j}$$

- q can be estimated from relative frequencies

The alignment model

- suppose \vec{x} and \vec{y} evolved from a common ancestor via independent substitution mutations
- independence between positions:

$$P(\vec{x}, \vec{y} | h_1) = \prod_i P(x_i, y_i | h_2)$$

- $p_{a,b}$: probability that a position in the latest common ancestor of x and y evolved into an a in sequence \vec{x} and into a b in sequence \vec{y}

$$P(\vec{x}, \vec{y} | h_1) = \prod_i p_{x_i, y_i}$$

$$\log P(\vec{x}, \vec{y} | h_1) = \sum_i \log p_{x_i, y_i}$$

The log-odds score

- taking things together, we have

$$\log(P(\vec{x}, \vec{y}|h_1) : P(\vec{x}, \vec{y}|h_0)) = \sum_i \log \frac{p_{x_i, y_i}}{q_{x_i} q_{y_i}}$$

- $\log \frac{p_{ab}}{q_a q_b}$: **score** of the alignment of a with b
- also goes by the name of **Pointwise Mutual Information** (PMI)
- assembled in a **PMI substitution matrix**

Substitution matrices

- in bioinformatics, several commonly used substitution matrices for nucleotids and proteins
- based on explicit models of evolution and careful empirical testing
- for nucleotids:

	<i>A</i>	<i>G</i>	<i>T</i>	<i>C</i>
<i>A</i>	2	-5	-7	-7
<i>G</i>	-5	2	-7	-7
<i>T</i>	-7	-7	2	-5
<i>C</i>	-7	-7	-5	2

Substitution matrices

- for proteins: different matrices for different evolutionary distances
- for instance: BLOSUM50

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	5	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-3	-1	1	0	-3	-2	0
R	-2	7	-1	-2	-4	1	0	-3	0	-4	-3	3	-2	-3	-3	-1	-1	-3	-1	-3
N	-1	-1	7	2	-2	0	0	0	1	-3	-4	0	-2	-4	-2	1	0	-4	-2	-3
D	-2	-2	2	8	-4	0	2	-1	-1	-4	-4	-1	-4	-5	-1	0	-1	-5	-3	-4
C	-1	-4	-2	-4	13	-3	-3	-3	-3	-2	-2	-3	-2	-2	-4	-1	-1	-5	-3	-1
Q	-1	1	0	0	-3	7	2	-2	1	-3	-2	2	0	-4	-1	0	-1	-1	-1	-3
E	-1	0	0	2	-3	2	6	-3	0	-4	-3	1	-2	-3	-1	-1	-1	-3	-2	-3
G	0	-3	0	-1	-3	-2	-3	8	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4
H	-2	0	1	-1	-3	1	0	-2	10	-4	-3	0	-1	-1	-2	-1	-2	-3	2	-4
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5	2	-3	2	0	-3	-3	-1	-3	-1	4
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5	-3	3	1	-4	-3	-1	-2	-1	1
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6	-2	-4	-1	0	-1	-3	-2	-3
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7	0	-3	-2	-1	-1	0	1
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8	-4	-3	-2	1	4	-1
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10	-1	-1	-4	-3	-3
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5	2	-4	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5	-3	-2	0
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-3	15	2	-3
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	-1
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5

Substitution matrix for the ASJP data

1. identify large sample of pairs of closely related languages (using expert information or heuristics based on aggregated Levenshtein distance)

An.NORTHERN_PHILIPPINES.CENTRAL_BONTOC
An.MESO-PHILIPPINE.NORTHERN_SORSOGON

WF.WESTERN_FLY.IAMEGA
WF.WESTERN_FLY.GAMAWE

Pan.PANOAN.KASHIBO_BAJO_AGUAYTIA
Pan.PANOAN.KASHIBO_SAN_ALEJANDRO

AA.EASTERN_CUSHITIC.KAMBAATA_2
AA.EASTERN_CUSHITIC.HADIYYA_2

ST.BAI.QILIQIAO_BAI_2
ST.BAI.YUNLONG_BAI

An.SULAWESI.MANDAR
An.OCEANIC.RAGA

An.SULAWESI.TANETE
An.SAMA-BAJAW.BOEPINANG_BAJAU

UA.AZTECAN.NAHUATL_HUEYAPAN_TETELA_DEL_VOLCAN
UA.AZTECAN.NAHUATL_CUENTEPEC_TEMIXCO

An.SOUTHERN_PHILIPPINES.KAGAYANEN
An.NORTHERN_PHILIPPINES.LIMOS_KALINGA

An.MESO-PHILIPPINE.CANIPAAN_PALAWAN
An.NORTHWEST_MALAYO-POLYNESIAN.LAHANAN

NC.BANTOID.LIFONGA
NC.BANTOID.BOMBOMA_2

IE.INDIC.WAD_PAGGA
IE.INDIC.TALAGANG_HINDKO

NC.BANTOID.LINGALA
NC.BANTOID.LIFONGA

An.CENTRAL_MALAYO-POLYNESIAN.BALILED0
An.CENTRAL_MALAYO-POLYNESIAN.PALUE

AuA.MUNDA.HO
AuA.MUNDA.KORKU

MGe.GE-KAINGANG.KAYAPO
MGe.GE-KAINGANG.APINAYE

Substitution matrix for the ASJP data

2. pick a concept and a pair of related languages at random
 - languages: Pen.MAIDUAN.MAIDU_KONKAU, Pen.MAIDUAN.NE_MAIDU
 - concept: *one*
3. find corresponding words from the two languages:
 - *nisam, niSem*
4. do Levenshtein alignment

n	i	s	a	m
n	i	S	e	m

5. for each sound pair, count number of correspondences
 - nn: 1; ii: 1; sS; 1; ae: 1; mm: 1

Substitution matrix for the ASJP data

- steps 2-5 are repeated 100,000 times

klem	S3--v	ligini	kulox	Naltir---i	...
klom	S37on	ji---p	Gulox	Naltirtiri	...
		a	a	56,047	: : :
		i	i	33,955	4 8 2
		u	u	23,731	4 a 2
		n	n	21,363	G t 2
		o	o	19,619	i ! 2
		m	m	18,263	G y 2
		t	t	16,975	d ! 2
		k	k	16,773	s G 2
		e	e	12,745	Z 5 2
		r	r	11,601	G s 2
		l	l	11,377	X z 2
		b	b	8,965	! k 2
		s	s	8,245	q 8 2
		d	d	6,829	a ! 2
		p	p	6,681	a ! 2
		w	w	6,613	! y 2
		N	N	6,275	! E 2
		h	h	5,331	j G 2
		y	y	5,321	G i 2
		3	3	5,255	E ! 2

Substitution matrix for the ASJP data

6. determine relative frequency of occurrence of each sound within the entire database

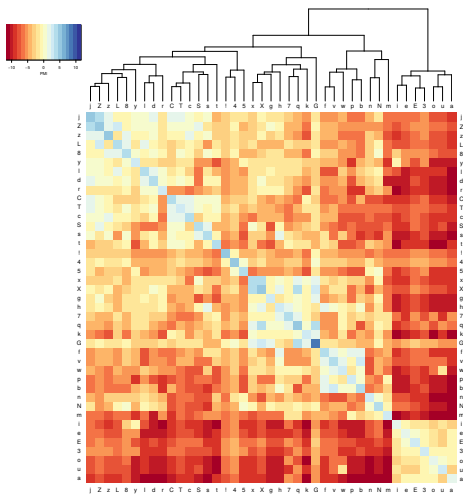
a	0.1479	E	0.0134
i	0.0969	7	0.0124
u	0.0696	C	0.0073
o	0.0626	S	0.0064
n	0.0614	x	0.0062
e	0.0478	c	0.0056
k	0.0478	f	0.0052
m	0.0465	5	0.0049
t	0.0449	v	0.0045
r	0.0346	q	0.0041
l	0.0331	z	0.0035
b	0.0248	j	0.0035
s	0.0243	T	0.0029
w	0.0232	L	0.0027
3	0.0228	X	0.0022
y	0.0222	8	0.0014
d	0.0214	Z	0.0011
h	0.0213	!	0.0009
p	0.0202	4	0.0002
N	0.0201	G	0.0001
g	0.0178		

Substitution matrix for the ASJP data

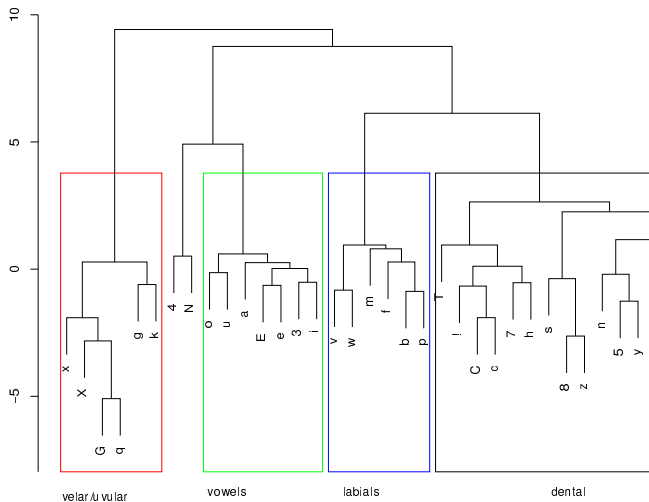
7. estimate p_{ab} as relative frequency of co-occurrence of a with b , q_a, q_b as individual relative frequencies, and determine PMI scores $\log_2 \frac{p_{ab}}{q_a q_b}$

G	G	11.2348	Z	j	4.9386	o	q	-3.2842	
!	!	10.0202	d	d	4.9263	C	a	-3.2893	
4	4	9.1480	g	g	4.8958	j	o	-3.2914	
8	8	8.0650	b	b	4.8906	a	m	-3.2915	
Z	Z	7.9575	s	s	4.8277	E	v	-3.3035	
X	X	7.9375	4	5	4.7508	!	w	-3.3079	
L	L	7.6276	E	E	4.7143	!	u	-3.3087	
z	z	7.2624	w	w	4.6512	5	q	-3.3116	
q	q	7.2542	h	h	4.5819	T	o	-3.3158	
f	f	6.9117	G	x	4.5573	!	k	-3.3526	
v	v	6.8418	Z	z	4.4943	e	z	-3.3763	
5	5	6.7731	y	y	4.4637	!	s	-3.3788	
j	j	6.7587	l	l	4.4037	...	f	q	-3.3942
T	T	6.6580	!	G	4.3760	N	S	-3.3954	
S	S	6.6054	3	3	4.3692	!	b	-3.4077	
c	c	6.5989	r	r	4.3061	L	b	-3.4558	
C	C	6.2439	X	q	4.1200	T	u	-3.4690	
4	G	6.1943	m	m	4.1087	4	i	-3.5529	
x	x	6.1210	t	t	4.1021	5	a	-3.8294	
G	X	5.3342	G	Z	4.0429	C	N	-3.8451	
G	q	5.3017	k	k	3.9046	!	t	-4.2625	
7	7	5.2111	X	x	3.8116	!	e	-4.3534	
p	p	5.0693	T	Z	3.7380	!	i	-4.3712	
N	N	4.9821	8	G	3.6993	!	a	-4.9817	

Evaluation



Evaluation



Gap penalties

- gaps in an alignment correspond either to an insertion or a deletion
- simplified assumption: insertions and deletions are equally likely at all positions; symbols are inserted according to their general frequency of occurrence
- Suppose an item x_i is aligned to a gap. Let α be the probability that an insertion occurred since the latest common ancestor, and β the probability of a deletion

$$\begin{aligned}
 P(x_i, -|h_1) &= \alpha q_{x_i} + \beta q_{x_i} \\
 P(x_i, -|h_0) &= q_{x_i} \\
 \log(P(x_i, -|h_1) : P(x_i, -|h_0)) &= \log(\alpha + \beta) \\
 &= -d
 \end{aligned}$$

- i.e., there is a constant term for each gap
- as $\alpha + \beta < 1$, this term is negative, i.e. there a constant **penalty** for each gap

Affine gap penalties

- deletions/insertions frequently apply to entire blocks of symbols (both in biology and linguistics)
- probability of a gap of length n are higher than the product of probabilities of n individual gaps
- penalty e for **extending** a gap is lower than penalty d for **opening** a gap
- g : length of a gap

$$\gamma(g) = -d - (g - 1)e$$

- no principled way to derive the values of d and e ; have to be fixed via trial and error
- $d = 2.5$ and $e = 1.6$ work quite well for the ASJP data

Weighted alignment

- so far, we assumed that the alignment between \vec{x} and \vec{y} is known
- to assess strength of evidence for h_1 given \vec{x}, \vec{y} , we need to consider all alignments between \vec{x} and \vec{y}
- enumeration is infeasible, because the number of alignments between two sequences of length n is

$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2} \approx \frac{2^{2n}}{\sqrt{\pi n}}$$

- computation is nonetheless possible using *Pair Hidden Markov Models*
- simpler task: find the most likely alignment and determine its log-odds!

The Needleman-Wunsch algorithm

- almost identical to Levenshtein algorithm, except:
 - matches/mismatches are counted not as 1 and 0, but as log-odds scores of the corresponding symbol pair
 - insertions/deletions are counted as gap penalties
 - by convention, the similarity rather than the distance is counted, i.e. we try to find the alignment that maximizes the score
- let \vec{x} have length n , \vec{y} length m , s_{ab} be the log-odds score of a and b , and d/e the gap penalties

The Needleman-Wunsch algorithm

$$F(0,0) = 0$$

$$G(0,0) = 0$$

$$\forall i : \quad 0 < i \leq n$$

$$F(i,0) = F(i-1,0) + G(i-1,0)e + (1 - G(i-1,0))d$$

$$G(i,0) = 1$$

$$\forall j : \quad 0 < j \leq m :$$

$$F(0,j) = F(0,j-1) + G(0,j-1)e + (1 - G(0,j-1))d$$

$$G(0,j) = 1$$

$$\forall i, j : \quad 0 < i \leq n, 0 < j \leq m$$

$$F(i,j) = \max \left\{ \begin{array}{l} F(i-1,j) + G(i-1,j)e + (1 - G(i-1,j))d \\ F(i,j-1) + G(i,j-1)e + (1 - G(i,j-1))d \\ F(i-1,j-1) + s_{x_i y_j} \end{array} \right.$$

$$G(i,j) = 0 \text{ if } \arg \max \left\{ \begin{array}{l} F(i-1,j) + G(i-1,j)e + (1 - G(i-1,j))d \\ F(i,j-1) + G(i,j-1)e + (1 - G(i,j-1))d \\ F(i-1,j-1) + s_{x_i y_j} \end{array} \right\} = 3$$

1 else

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5				
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5				
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Diagram illustrating the dynamic programming table for finding the best alignment. The table shows scores for alignments between characters '–', 'm', 'e', 'n', 'E', and 's'. The scores are: – (0), m (-2.5), E (-4.1), n (-5.7), S (-7.3). Arrows indicate the path from the top-left cell (0) to the cell (m, m) (-2.5), showing the optimal alignment path.

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5				
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13			
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13			
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Arrows point from the cell (m, E) to its left, top, and top-right neighbors.

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13			
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53		
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1				
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53			
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65		
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7				
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03			
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05		
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3				
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47			
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75		
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9				

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97			

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15		

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15	5.1	

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15	5.1	8.84

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15	5.1	8.84

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15	5.1	8.84

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15	5.1	8.84

Arrows in the original image point from the cell (n, n) to (n, E) and from (n, E) to (E, E).

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15	5.1	8.84

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Finding the best alignment

- Dynamic Programming

	–	m	E	n	S
–	0	-2.5	-4.1	-5.7	-7.3
m	-2.5	4.13	1.53	0.03	-1.47
e	-4.1	1.53	5.65	3.05	1.55
n	-5.7	0.03	3.05	9.2	6.6
E	-7.3	-1.47	4.75	6.6	7.62
s	-8.9	-2.97	2.15	5.1	8.84

- memorizing in each step which of the three cells to the left and above gave rise to the current entry lets us recover the corresponding optimal alignment

Evaluation

left: Levenshtein alignment; right: Needleman-Wunsch alignment

-iX ego	iX- ego	-blat folyu	b-lat folyu	han-t manus	han-t manus
du tu	du tu	haut-- -kutis	haut-- k-utis	--brust pektus-	b--rust pektus-
vir nos	vir nos	---blut saNgwis	---blut saNgwis	leb3r yekur	leb3r yekur
ains unus	ain-s -unus	knoX3n --o--s	knoX3n --os--	triNk3n -bibere	triNk3n- -bi-bere
cvai -duo	cvai duo-	horn- kornu	horn- kornu	--ze3n widere	--ze3n widere
---mEnS persona	mEnS--- persona	-au-g3 okulus	a-ug3- okulus	-her3n audire	--her3n audire-
---fiS piskis	fiS--- piskis	na-z3 nasus	naz3- nasus	Sterb3n -mor--i	Sterb3n -mor-i-
hun-t kanis	hun-t kanis	chan dens	chan- d-ens	khom3n wenire	khom3n--- w---enire
-----laus pedikulus	-----laus pedikul-us	-chuN3 liNgwE	chuN--3 -liNgwE	zon3 so-l	zon3 sol-

Evaluation

```
vas3r      --vas3r
-akwa      akwa---
```

```
Stain      Sta-in
lapis      -lapis
```

```
-foia      fo-ia
iNnis      iNnis
```

```
pfat       p-fat
viya       viya-
```

```
bErk       bErk
mons       mons
```

```
n-at       na-t
noks       noks
```

```
---fol     fol----
plenus     p-lenus
```

```
no--i     no-i-
nowus     nowus
```

```
nam-3     nam3-
nomen     nomen
```

German — Swabian

'I': iX i	0.3	'louse': laus laus	15.01	'tongue': chuN3 cuN	9.8	'die': Sterb3n StEab	10.16
'you': du du	8.26	'tree': baum bom	6.57	'knee': kni knui	7.77	'come': khom3n khom	11.84
'we': vir mia	-1.09	'leaf': blat blad	11.92	'hand': hant hEnd	8.6	'sun': zon3 sonE	8.79
'one': ains ois	4.63	'skin': haut haut	14.42	'breast': brust bXuSt	14.81	'star': StErn StEan	16.16
'two': cvai cvoi	16.0	'blood': blut blud	12.88	'liver': leb3r leba	10.01	'water': vas3r vaza	7.8
'person': mEnS mEnZE	12.61	'bone': knoX3n knoXE	16.88	'drink': triNk3n dXiNg	4.99	'stone': Stain Stoi	10.36
'fish': fiS fiS	16.35	'horn': horn hoan	8.75	'see': ze3n se	0.63	'fire': foia fuia	12.43

German — English

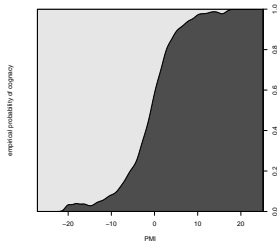
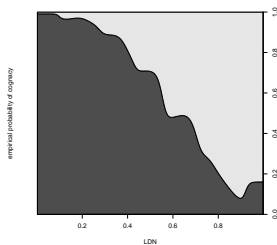
'I': iX Ei	-2.3	'tree': baum tri	-7.83	'tongue': chuN3 t3N	-0.63	'die': Sterb3n dEi	-7.7
'you': du yu	2.34	'leaf': blat lif	-0.47	'knee': kni ni	3.86	'come': khom3n k3m	1.22
'we': vir wi	2.21	'blood': blut bl3d	9.46	'hand': hant hEnd	8.6	'sun': zon3 s3n	1.95
'one': ains w3n	-2.3	'bone': knoX3n bon	-1.36	'breast': brust brest	16.93	'star': StErn star	8.2
'two': cvai tu	-5.25	'horn': horn horn	15.73	'liver': leb3r liv3r	14.65	'water': vas3r wat3r	12.06
'fish': fiS fiS	16.35	'eye': aug3 Ei	-4.1	'drink': triNk3n drink	7.48	'stone': Stain ston	6.75
'dog': hunt dag	-7.46	'nose': naz3 nos	1.63	'see': ze3n si	-3.04	'fire': foia fEir	6.79

German — Latin

'I': -3.87 iX ego	'louse': -0.08 laus pedikulus	'nose': 4.49 naz3 nasus	'see': -4.15 ze3n videre
'you': 3.62 du tu	'tree': -3.85 baum arbor	'tooth': -2.78 chan dens	'hear': -4.24 her3n audire
'we': -5.06 vir nos	'leaf': -3.57 blat folyu	'tongue': -3.4 chuN3 liNgwE	'die': -6.12 Sterb3n mori
'one': 2.39 ains unus	'skin': -0.25 haut kutis	'knee': 0.8 kni genu	'come': -9.25 khom3n wenire
'two': -5.51 cvai duo	'blood': -9.18 blut saNgwis	'hand': 0.73 hant manus	'sun': 0.97 zon3 sol
'person': -4.66 mEnS persona	'bone': -5.72 knoX3n os	'breast': 1.39 brust pektus	'star': 5.72 StErn stela
'fish': 0.29 fiS piskis	'horn': 7.55 horn kornu	'liver': 5.37 leb3r yekur	'water': -5.4 vas3r akwa

How well does PMI similarity predict cognacy?

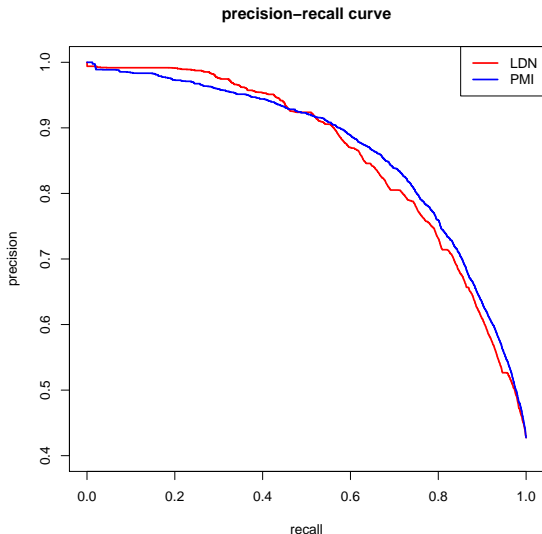
expert cognacy judgments used as gold standard



How well does PMI similarity predict cognacy?

Average Precision

- LDN: 0.847
- PMI: 0.864



Estimating distances from pairwise alignments

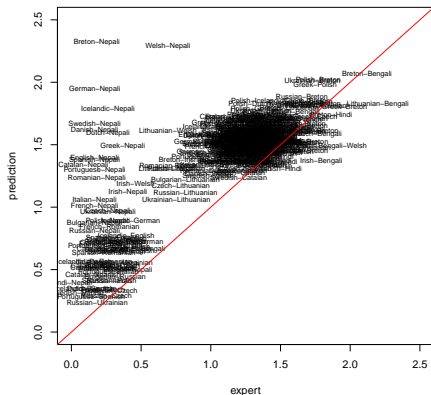
Probability of cognacy

- logistic regression to predict probability of cognacy from PMI similarity

concept	Italian	English	predicted prob.	expert judgment	concept	Italian	English	predicted prob.	expert judgment
<i>sharp</i>	tallEnte	Sop	0.004	0	<i>father</i>	padre	fo83	0.480	1
<i>float</i>	galleddZare	fl3ut	0.004	0	<i>when</i>	kwando	wEn	0.483	1
<i>Kill</i>	ammattsare	kil	0.007	0	<i>night</i>	notte	nait	0.508	1
<i>bark</i>	skordza	bok	0.009	0	<i>and</i>	eed	End	0.518	0
<i>husband</i>	marito	hoz3nd	0.010	0	<i>name</i>	nome	neim	0.519	1
<i>walk</i>	kamminare	wok	0.011	0	<i>worm</i>	vErme	w3m	0.521	1
<i>eat</i>	mandZare	it	0.011	0	<i>round</i>	tondo	raund	0.526	1
<i>bark</i>	kortettSa	bok	0.013	0	<i>many</i>	molti	mEni	0.569	0
<i>know</i>	sapere	n3u	0.015	0	<i>wind</i>	vEnto	wind	0.573	1
<i>come</i>	venire	kom	0.016	1	<i>two</i>	due	tu	0.600	1
<i>swim</i>	nwotare	swim	0.016	0	<i>mother</i>	madre	mo83	0.624	1
<i>back</i>	dosso	bEk	0.018	0	<i>thou</i>	tu	8au	0.629	1
<i>burn</i>	ardere	b3n	0.018	0	<i>child</i>	fantSullo	tSaild	0.638	0
<i>think</i>	pensare	8iNk	0.019	0	<i>long</i>	lungo	loN	0.651	1
<i>dust</i>	polvere	dost	0.019	0	<i>fish</i>	peSSe	fiS	0.659	1
<i>wife</i>	molle	waif	0.020	0	<i>count</i>	kontare	kaunt	0.660	1
<i>swell</i>	gonfyare	swEl	0.021	0	<i>star</i>	stella	sto	0.664	1
<i>sing</i>	kantare	siN	0.022	0	<i>belly</i>	vEntre	bEli	0.679	0
<i>knee</i>	rotElla	ni	0.022	0	<i>sun</i>	sole	son	0.692	1
<i>dry</i>	aSSutto	drai	0.022	0	<i>fly</i>	volare	flai	0.742	0
<i>five</i>	tSinkwe	faiv	0.023	1	<i>three</i>	tre	8ri	0.744	1
<i>skin</i>	pElle	skin	0.024	0	<i>flow</i>	fluire	fl3u	0.759	0
<i>hand</i>	mano	hEnd	0.025	0	<i>heavy</i>	grEve	hEvi	0.769	0
<i>blood</i>	sangwe	blod	0.025	0	<i>person</i>	persona	p3s3n	0.799	1
<i>flow</i>	skorrere	fl3u	0.026	0	<i>animal</i>	animale	Enim3l	0.947	1
<i>wipe</i>	aSSugare	waip	0.026	0	<i>vomit</i>	vomitare	vomit	0.960	1
<i>turn</i>	dZirare	t3n	0.026	0	<i>fruit</i>	frutto	frut	0.966	1

Estimating distances

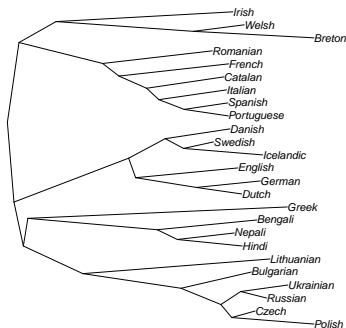
- average or maximal predicted probability of cognacy per concept = expected relative frequency of cognate pairs
 - expected relative frequency of cognate pairs = e^{-t}
- ⇒ distance estimation from raw data
- ⇒ applicable across language families



Estimating distances

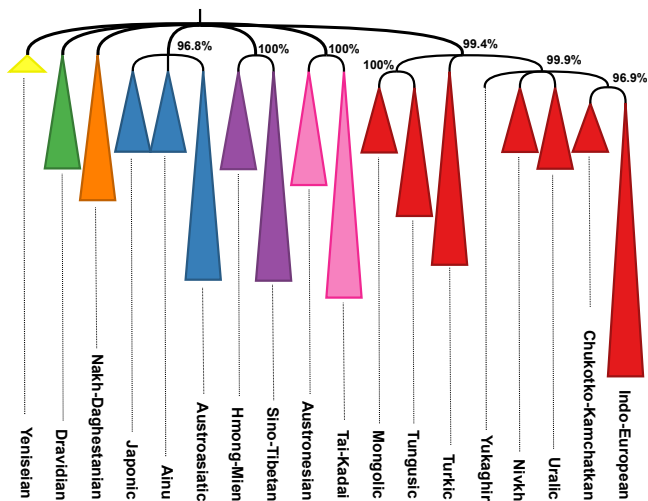
- average or maximal predicted probability of cognacy per concept = expected relative frequency of cognate pairs
- expected relative frequency of cognate pairs = e^{-t}
- ⇒ distance estimation from raw data
- ⇒ applicable across language families

Neighbor Joining tree



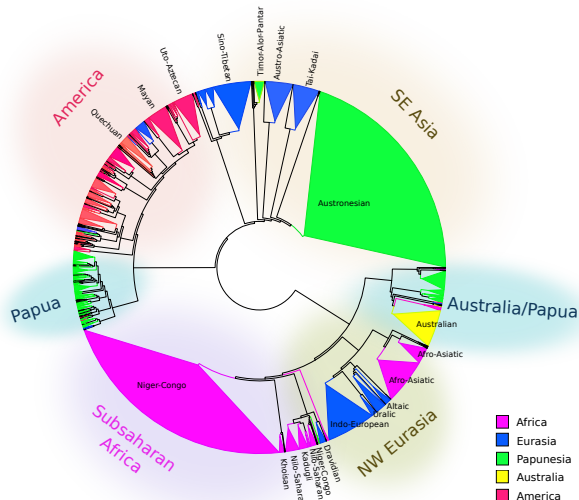
Languages of Eurasia/ ASJP data

cf. Jäger (2015); full tree can be inspected [here](#)



Languages of the World/ ASJP data

cf. Jäger and Wichmann (2016); full tree can be inspected [here](#).



Multiple sequence alignment

Multiple sequence alignment

- Needleman-Wunsch only does pairwise alignment
- desirable: aligning all sequences of a taxon into one matrix
 - necessary for character-based phylogenetic inference
 - improves the quality of the alignment

Multiple sequence alignment

- example: 'one'
 - PIE: oinos
 - Bosian: yedan
 - Kashubian: yEdEn
 - optimal pairwise alignments:

o	i	n	o	s	o	i	n	o	s	y	e	d	a	n
y	e	d	a	n	y	E	d	E	n	y	E	d	E	n

- optimal multiple alignment (maximizing sum of pairwise similarities per column):

y	E	d	E	n	-	-
-	o	-	i	n	o	s
y	e	d	a	n	-	-

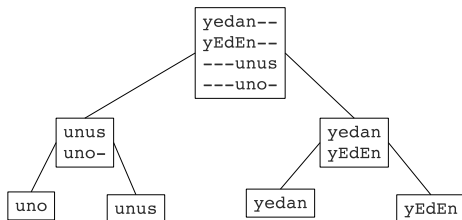
- alignment of all 'n's is etymologically correct

Multiple sequence alignment

- in principle, the Needleman-Wunsch algorithm can be generalized to aligning k sequences
- however, aligning k sequences of length n has complexity $\mathcal{O}(n^{k^2}) \Rightarrow$ computationally intractable
- two strategies
 - heuristic search
 - progressive alignment

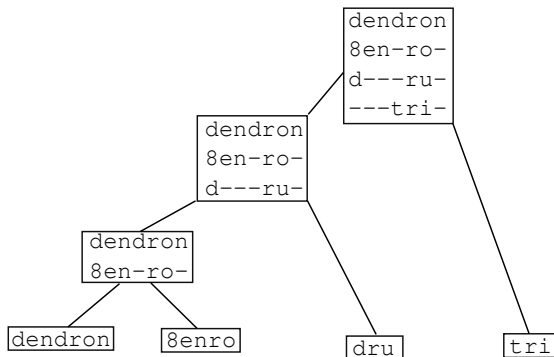
Progressive sequence alignment

- start with a **guide tree** (using some heuristics like pairwise alignment + Neighbor Joining)
- working bottom-up, at each internal node, do pairwise alignment of the block alignments at the daughter node
- complexity is $\mathcal{O}(n^2k^3) \Rightarrow$ computationally feasible



T-Coffee

- progressive alignment only uses (phylogenetically) local information
- erroneous decisions cannot be corrected later



T-Coffee

Tree-based Consistency Objective Function for alignment Evaluation (Notredame et al., 2000)

(slightly adapted for linguistic application)

- 1 pairwise alignment for all word pairs, using PMI scores

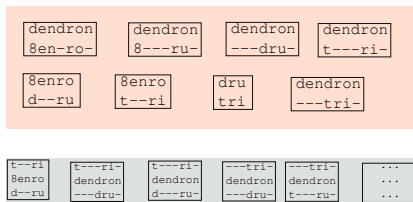
dendron 8en-ro-	dendron 8---ru-	dendron ---dru-	dendron t---ri-
8enro d--ru	8enro t--ri	dru tri	dendron ---tri-

T-Coffee

Tree-based Consistency Objective Function for alignment Evaluation (Notredame et al., 2000)

(slightly adapted for linguistic application)

- 1 pairwise alignment for all word pairs, using PMI scores
- 2 ternary alignments via relation composition

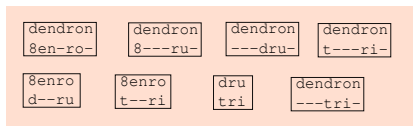


T-Coffee

Tree-based Consistency Objective Function for alignment Evaluation (Notredame et al., 2000)

(slightly adapted for linguistic application)

- 1 pairwise alignment for all word pairs, using PMI scores
- 2 ternary alignments via relation composition
- 3 indirect alignment scores between sound **occurrences**

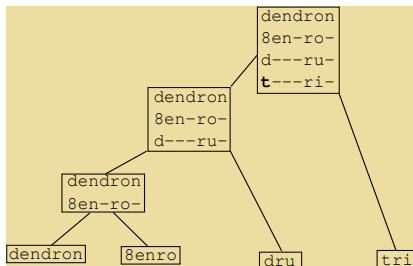
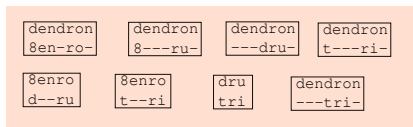


T-Coffee

Tree-based Consistency Objective Function for alignment Evaluation (Notredame et al., 2000)

(slightly adapted for linguistic application)

- 1 pairwise alignment for all word pairs, using PMI scores
- 2 ternary alignments via relation composition
- 3 indirect alignment scores between sound **occurrences**
- 4 progressive alignment using those scores



Examples

cognate class	language	word
<i>one:A</i>	German	-a-i--n-
<i>one:A</i>	Dutch	---e--n-
<i>one:A</i>	English	-w-o--n-
<i>one:A</i>	Danish	---e--n-
<i>one:A</i>	Swedish	---E--n-
<i>one:A</i>	Icelandic	---eidn-
<i>one:A</i>	Irish	---e--n-
<i>one:A</i>	Breton	---i--n-
<i>one:A</i>	French	---E----
<i>one:A</i>	Catalan	---u--n-
<i>one:A</i>	Spanish	---u--no
<i>one:A</i>	Portuguese	---u----
<i>one:A</i>	Italian	---u--no
<i>one:A</i>	Romanian	---u--nu
<i>one:A</i>	Bengali	---E--k-
<i>one:A</i>	Nepali	---e--k-
<i>one:A</i>	Czech	yEdE--n-
<i>one:A</i>	Polish	yEdE--n-
<i>one:A</i>	Ukrainian	-odi--n-
<i>one:A</i>	Russian	-adi--n-
<i>one:A</i>	Bulgarian	-3di--n-

cognate class	language	word
<i>heart:J</i>	German	h-Er-t--s-
<i>heart:J</i>	Dutch	h-or-t----
<i>heart:J</i>	English	h-o--t----
<i>heart:J</i>	Danish	y-Ea-d--3-
<i>heart:J</i>	Swedish	y-E--t--a-
<i>heart:J</i>	Icelandic	S-ar-t--a-
<i>heart:J</i>	French	k-Er-----
<i>heart:J</i>	Catalan	k-or-----
<i>heart:J</i>	Spanish	k-ora8--on
<i>heart:J</i>	Portuguese	k-uras--aw
<i>heart:J</i>	Italian	kwor----e-
<i>heart:J</i>	Hindi	h--r-d--ai
<i>heart:J</i>	Lithuanian	S-ir-dis--
<i>heart:J</i>	Czech	s--r-t-sE-
<i>heart:J</i>	Polish	s-Er-t-sE-
<i>heart:J</i>	Ukrainian	s-Er-t-sE-
<i>heart:J</i>	Russian	s-Erdt-sE-
<i>heart:J</i>	Bulgarian	s-3r-t-sE-
<i>heart:J</i>	Greek	k-ar-8-Sa-

Examples

cognate class	language	word
two:A	German	tsvai-
two:A	Dutch	t-we--
two:A	English	t--u--
two:A	Danish	d--o--
two:A	Swedish	t-vo--
two:A	Icelandic	t-veir
two:A	French	d--e--
two:A	Catalan	d--o-s
two:A	Spanish	d--o-s
two:A	Portuguese	d--oiS
two:A	Italian	d--ue-
two:A	Romanian	d--o-y
two:A	Nepali	d--ui-
two:A	Czech	d-va--
two:A	Polish	d-va--
two:A	Ukrainian	d-wa--
two:A	Russian	d-va--
two:A	Bulgarian	d-va--
two:A	Greek	8-io--

cognate class	language	word
<i>mother:A</i>	German	mu-t--a-
<i>mother:A</i>	Dutch	mu-d--3r
<i>mother:A</i>	English	mo-8--3-
<i>mother:A</i>	Danish	mo----a-
<i>mother:A</i>	Swedish	mu-d--3r
<i>mother:A</i>	Icelandic	mou8--ir
<i>mother:A</i>	French	mE---r--
<i>mother:A</i>	Catalan	ma---r3-
<i>mother:A</i>	Spanish	ma-8-re-
<i>mother:A</i>	Portuguese	ma----i-
<i>mother:A</i>	Italian	ma-d-re-
<i>mother:A</i>	Czech	ma-t-ka-
<i>mother:A</i>	Polish	ma-t-ka-
<i>mother:A</i>	Ukrainian	ma-t--i-
<i>mother:A</i>	Russian	ma-t----
<i>mother:A</i>	Bulgarian	ma-y-k3-
<i>mother:A</i>	Greek	mi-tera-

Examples

cognate class	language	word
<i>tongue:W</i>	German	---tsuN--3
<i>tongue:W</i>	Dutch	---t-oN---
<i>tongue:W</i>	English	---t-oN---
<i>tongue:W</i>	Danish	---d-oN--3
<i>tongue:W</i>	Swedish	---t-3N--a
<i>tongue:W</i>	Icelandic	---t-uNg-a
<i>tongue:W</i>	French	---l-o-g--
<i>tongue:W</i>	Catalan	---l-ENgW3
<i>tongue:W</i>	Spanish	---l-eNgwa
<i>tongue:W</i>	Portuguese	---l-i-gua
<i>tongue:W</i>	Italian	---l-ingwa
<i>tongue:W</i>	Romanian	---l-im-b3
<i>tongue:W</i>	Hindi	---dZi--b-
<i>tongue:W</i>	Czech	ya-z-ik---
<i>tongue:W</i>	Polish	yEwz-3k---
<i>tongue:W</i>	Ukrainian	ya-z-ik---
<i>tongue:W</i>	Russian	yi-z-3k---
<i>tongue:W</i>	Bulgarian	-3-z-ik---

cognate class	language	word
<i>tooth:B</i>	Greek	8-on-di
<i>tooth:B</i>	German	tsan--
<i>tooth:B</i>	Dutch	t-ont-
<i>tooth:B</i>	English	t-u-8-
<i>tooth:B</i>	Danish	d-an--
<i>tooth:B</i>	Swedish	t-and-
<i>tooth:B</i>	Icelandic	t-En--
<i>tooth:B</i>	French	d-o---
<i>tooth:B</i>	Catalan	d-en--
<i>tooth:B</i>	Spanish	dyente
<i>tooth:B</i>	Portuguese	d-e-t3
<i>tooth:B</i>	Italian	d-Ente
<i>tooth:B</i>	Romanian	d-inte
<i>tooth:B</i>	Bengali	d-o-t-
<i>tooth:B</i>	Hindi	d-a-t-

Examples

cognate class	language	word
<i>dog:A</i>	Lithuanian	S-u---o-
<i>dog:A</i>	Ukrainian	s-obaka-
<i>dog:A</i>	Russian	s-abaka-
<i>dog:A</i>	Danish	h-u-n---
<i>dog:A</i>	Swedish	h-3-nd--
<i>dog:A</i>	Icelandic	h-i-ndir
<i>dog:A</i>	German	h-u-nt--
<i>dog:A</i>	Dutch	h-o-nt--
<i>dog:A</i>	Welsh	k-----i-
<i>dog:A</i>	Breton	k-----i-
<i>dog:A</i>	Irish	k-----u-
<i>dog:A</i>	French	S-i---E-
<i>dog:A</i>	Italian	k-a-n-e-
<i>dog:A</i>	Portuguese	k-a---u-
<i>dog:A</i>	Romanian	k-3yn-e-
<i>dog:A</i>	Greek	Tio-n---

cognate class	language	word
<i>tree:C</i>	Danish	d---G-E--
<i>tree:C</i>	Swedish	t---r-Ed-
<i>tree:C</i>	Icelandic	t---ryE--
<i>tree:C</i>	English	t---r-i--
<i>tree:C</i>	Ukrainian	dE--r-Ewo
<i>tree:C</i>	Russian	dE--r-Evo
<i>tree:C</i>	Polish	d---Z-Evo
<i>tree:C</i>	Bulgarian	d3--r--vo
<i>tree:C</i>	Greek	8endr---o

Wrapping up

Important topics not covered

- Bayesian tree estimation \Rightarrow no good introductory texts so far (that I would be aware of); best starting point might be Chen et al. (2014), esp. chapters 1 and 7
- estimation of time depths in years (rather than in “expected number of mutations”) \Rightarrow Chang et al. (2015)
- automatic cognate detection \Rightarrow first part of the Jäger/List-manuscript on the course homepage

Hot research topics

- automatic discovery of regular sound correspondences and sound laws
- automatic reconstruction of proto-forms
- factoring vertical descent from language contact
- integrated probabilistic inference of sequence alignment and phylogenies

- Chang, W., C. Cathcart, D. Hall, and A. Garrett (2015). Ancestry-constrained phylogenetic analysis supports the Indo-European steppe hypothesis. *Language*, **91**(1):194–244.
- Chen, M.-H., L. Kuo, and P. O. Lewis (2014). *Bayesian Phylogenetics. Methods, Algorithms and Applications*. CRC Press, Abingdon.
- Dolgopolsky, A. B. (1986). A probabilistic hypothesis concerning the oldest relationships among the language families of northern eurasia. In V. V. Shevoroshkin, ed., *Typology, Relationship and Time: A collection of papers on language change and relationship by Soviet linguists*, pp. 27–50. Karoma Publisher, Ann Arbor.
- Jäger, G. (2015). Support for linguistic macrofamilies from weighted sequence alignment. *Proceedings of the National Academy of Sciences*, **112**(41):12752–12757. Doi: 10.1073/pnas.1500331112.
- Jäger, G. and S. Wichmann (2016). 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, eds., *The Evolution of Language: Proceedings of the 11th International Conference*

(*EVOLANG11*). Available online:

<http://evolang.org/neworleans/papers/147.html>.

List, J.-M. (2014). *Sequence Comparison in Historical Linguistics*.
Düsseldorf University Press, Düsseldorf.