

Distanz- und Alignmentanalysen in der historischen Linguistik

Johann-Mattis List
Heinrich Heine Universität Düsseldorf

22. Januar 2010

Gliederung

Grundlegendes zur Einführung

- Mengen- vs. Sequenzvergleiche
- Unigrammatische und n -grammatische Segmentierung
- Paarweise und multiple Alinierung
- Der dynamische Programmieralgorithmus

Erweiterung des Alinierungsverfahrens

- Erweiterung des Algorithmus
- Erweiterung der Vergleichsfunktion
- Segmentklassen anstelle "reiner" Segmente

Vorstellung neuerer Ansätze von Sequenzanalysen

- The Automated Similarity Judgment Program
- Covington
- ALINE

Arbeitsstand und Ausblick

Grundlegende Anmerkungen zur Einführung

Mengen- vs. Sequenzvergleiche

Mengenvergleiche Vergleich einer Anzahl ungeordneter, distinkter Elemente

Sequenzvergleiche Vergleich einer Anzahl geordneter Elemente, deren Distinktivität erst durch die Anordnung hergestellt wird

Sequenzvergleiche setzen eine *Alinierung* der Sequenzen voraus, da Sequenzdistanzen nur ermittelt werden können, wenn die korrespondierenden Segmente bestimmt wurden.

Unigrammatische und n -grammatische Segmentierung

Monogrammatisch	θ	i	y	a	t	ϵ	r	a	
Bigrammatische	$-\theta$	θi	$i y$	$y a$	$a t$	$t \epsilon$	ϵr	$r a$	$a-$
Trigrammatische	$--\theta$	$-\theta i$	$\theta i y$	$i y a$	$y a t$	$a t \epsilon$	$t \epsilon r$	$\epsilon r-$	$r--$

Tabelle: Mono- bi und trigram-basierte Segmentierung

Paarweise und multiple Alinierung

d	ɔ:	-	-	t ^h	ə	-	-
θ	i	ɣ	a	t	ɛ	r	a
t ^h	ɔ	x	-	t ^h	ɐ	-	-

Tabelle: Multiple Alinierung von Sequenzen

Grundidee des Algorithmus

- ▶ Erstellen einer Matrix mit allen möglichen Segmententsprechungen zweier Sequenzen
- ▶ Festsetzen von Kosten für die Gegenüberstellung der jeweiligen Segmente durch eine Vergleichsfunktion
 - ▶ Gegenüberstellung von Segmenten (Substitution & Match)
 - ▶ Einfügen von (Null)-Segmenten (Insertion & Deletion)
- ▶ Kumulative Aufrechnung von Kosten für alle möglichen Wege durch die Matrix

Die Vergleichsfunktion der Levenshtein-Distanz

Entscheidung	Bedingung	Kosten
Gegenüberstellung	Identität von Segmenten	0
	Verschiedenheit von Segmenten	1
Einfügen und Ersetzen		1

Tabelle: Die Vergleichsfunktion der Levenshtein-Distanz

Erstellen der Matrix

-	-	-	-	-
h	h	e	r	z
e	e	e	r	z
a	a	e	r	z
r	r	e	r	z
t	t	e	r	z

Tabelle: Vergleichsmatrix für den Wagner-Fischer-Algorithmus

Berechnen der Kosten für jeden Pfad

0 -/-	1 -/h	2 -/e	3 -/r	4 -/z
1 h/-	0 h/h	1 -/e	2 -/r	3 -/z
2 e/-	1 e/-	0 e/e	1 -/r	2 -/z
3 a/-	2 a/-	1 a/-	1 a/r	2 -/z
4 r/-	3 r/-	2 a/-	1 r/r	2 -/z
5 t/-	4 t/-	3 a/-	2 t/r	2 t/z

Tabelle: Vergleichsmatrix nach der Auswertung (mit Kosten)

Transpositionen

Levenshtein	f	r	o	m	a	g	e	-	-
	f	o	r	m	a	g	g	i	o
Damerau-Levenshtein	f	ro		m	a	g	e	-	-
	f	or		m	a	g	g	i	o

Tabelle: Levenshtein und Damerau-Levenshtein

Konsequente Indels

Traditionell	d	i	d	o:	m	i
	d	a	-	-	m	-
Gotoh-Erweiterung	d	i	d	o:	m	i
	-	-	d	a	m	-

Tabelle: Konsequente *indels* in der Sequenzalinierung

Kompressionen und Expansionen

Levenshtein	d	c	-	t ^h	e
	t ^h	c	x	t ^h	e
Oommen-Erweiterung	d	c	t ^h		e
	t ^h	c	xt ^h		e

Tabelle: Kompression und Expansion in der Sequenzalinierung

Lokale Alinierung

Levenshtein	-	ā	p	a	k	o	s	ī	s	-
	w	ā	p	i	k	o	n	ō	h	a
Lokale Alinierung		ā	p	a	k	o	sīs			
	w	ā	p	i	k	o	nōha			

Probleme beim Alinieren phonetischer Daten

θ	i	y	a	t	ε	r	a
d	ɔː	t ^h	ʒ	-	-	-	-

Tabelle: Problem der Alinierung phonetischer Sequenzen

Wie gut muss die Vergleichsfunktion sein?

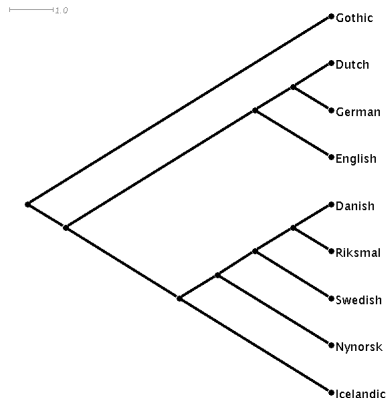


Abbildung: Neighbor-Analyse von Levenshtein-Distanzen

Lautklassen in der Alinierung

- ▶ Zuweisung von Segmenten zu einheitlichen Klassen, innerhalb derer Lautwandel als wahrscheinlicher angenommen wird, als außerhalb der Klassen
- ▶ Umwandlung der Segmente in ihre jeweiligen Klassen
- ▶ Durchführung "traditioneller" Alignmentanalysen

Der Vorteil von lautklassenbasierten Ansätzen liegt in ihrer leicht zu realisierenden Implementierung. Ferner können phonetisch unzureichende Sprachdaten meist relativ einfach in ein Klassenformat überführt werden.

Dolgopolskys Lautklassen

No.	Typ	Beschreibung	Bsp.
1	P	labiale Obstruenten	p,b,f
2	T	dentale Obstruenten	d,t,θ,ð
3	S	alveolare, postalveolare und retroflexe Frikative	s,z,ʃ,ʒ
4	K	velare und postvelare Obstruenten und Affrikaten	k,g,ts,tʃ
5	M	labialer Nasal	m
6	N	übrige Nasale	n,ɲ,ŋ
7	R	Trills, Taps, Flaps und laterale Approximanten	r,l
8	W	stimmhafter labialer Frikativ und initiale gerundete Vokale	v,u
9	J	palataler Approximant	j
10	Ø	Laryngale und initialer velarer Nasal	h,ɦ,ŋ

Tabelle: Dolgopolskys Klassifizierung von Lautwandeltypen

Alinierung mit Hilfe von Dolgopolskys Lautklassen

Interne Alinierung	T	V	K	V	T	V	R	V
	T	V	-	-	T	V	-	-
Ausgabe	θ	i	ʏ	a	t	ε	r	a
	d	ɔː	-	-	t ^h	ʒ	-	-

Tabelle: Interne und externe Darstellung der Dolgopolsky-Alinierung

The Automated Similarity Judgment Program

Sprache	Bedeutung	IPA	ASJP-Code
engl.	'das'	ðis	8is
engl.	'Mund'	mauθ	mau8
engl.	'Zunge'	t ^h əŋ	th~3N
dt.	'Fisch'	fɪʃ	fiS
engl.	'Zahn'	tu:θ	tu8
dt.	'Schwester'	ʃwest ^h ɐ	Swasth~a

Tabelle: Das universale Alphabet des ASJP-Projektes

Covington

Penalty	Conditions
0	Exact match of consonants or glides (w, y)
5	Exact match of vowels
10	Match of two vowels that differ only in length, or i and y, or u and w
30	Match of two dissimilar vowels
60	Match of two dissimilar consonants
100	Match of two segments with no discernible similarity
40	Skip preceded by another skip in the same word
50	Skip not preceded by another skip in the same word

Tabelle: Covingtons 'Evaluationsmetrik' für die Alinierung

Kondraks ALINE

Syllabic	5	Place	40
Voice	10	Nasal	10
Lateral	10	Aspirated	5
High	5	Back	5
Manner	50	Retroflex	10
Long	1	Round	5

Tabelle: Merkmalssalienzen in ALINE

Arbeitsstand und Ausblick

Python-Programm zur umfangreichen Sequenzanalyse

- ▶ Modularer Aufbau, Module für Sequenzanalysen (Alignments und Distanzberechnungen), phonetische Analysen, automatische Kognatenerkennung
- ▶ Eingabe in Form Unicode-kodierter csv-Dateien (comma separated value), deren Struktur sich am Format etymologischer Wörterbücher in der STARLING-Software orientiert
- ▶ Ausgabe von Analysen und Berechnungen in verschiedenen Formaten, die wahlweise durch phylogenetische Softwarepakete weiterverarbeitet werden können (Phylip, Nexus)
- ▶ Skripte für Daten-Ein- und -Ausgabe sind bereits realisiert, ferner sind eine Reihe von Algorithmen/Verfahren zur Sequenzanalyse (erweiterter DPA, Dolgopolsky, ASJP, Covington) bereits verfügbar

Das war's!

