

Dialektklassifikation auf der Grundlage aggregierter Ausspracheunterschiede

John Nerbonne und Christine Siedle*
Alfa-informatica, Universität Groningen
Postbus 716, 9700 AS Groningen, Niederlande
{nerbonne,siedle}@let.rug.nl

28. Januar 2005

Zusammenfassung

Die Dialektgeographie betrachtet Sprachdaten normalerweise auf einem kategorialen (nominalen) Niveau, wobei die Einteilung in Dialekte auf der Grundlage einzelner sprachlicher Phänomene geschieht. Wenn aber die geographischen Verteilungen verschiedener sprachlicher Phänomene nun nicht übereinstimmen, dann erhält der Begriff 'Dialektgebiet' keine befriedigende analytische Basis— dies ist das Problem der nicht zusammenfallenden Isoglossen. Der vorliegende Beitrag schlägt eine Lösung dieses Problems vor durch die Verwendung eines Messinstruments für Ausspracheähnlichkeit, mit dem sich große Datenmengen systematisch verarbeiten lassen. Das Messinstrument ist numerisch (anstatt kategorial) und deswegen additiv, so dass man Unterschiede nicht nur auf dem Niveau des Wortes oder Phonems charakterisieren kann, sondern auch auf dem aggregierten Niveau der 'Varietät' – z.B. des Dialektes eines Dorfes. Auf dieser Basis lassen sich Dialektgebiete bestimmen, bei denen die Grenzen fließender sind als bei den Gebieten, die in den Dialekthandbüchern besprochen werden. Das Verfahren wird anhand von Daten 186 deutscher Varietäten illustriert.

1 Einleitung

Im Fokus des vorliegenden Berichts ist die Anwendung eines quantitativen Maßes für Ausspracheähnlichkeit auf eine relativ große Stichprobe phonetischer Transkriptionen von Dialekten aus dem Gebiet der heutigen Bundesrepublik. Die zentrale Frage unserer Arbeit lautet, ob zum Erkennen einer kohärenten Sprachgeographie tatsächlich die subjektiv gelenkte Auswahl von Spracheigenschaften der traditionellen Dialektologie (siehe Bloomfield (1933, 340-341)) notwendig ist oder ob die von uns verwendete quantitative Methode eine echte Alternative sein kann. Mit Coseriu (1975, 50) suchen wir eine Alternative zum drohenden Atomismus in der Dialektforschung.

In den folgenden Abschnitten behandeln wir die anzuwendende Methode, die Datensammlung, die Resultate und schließlich die künftigen Möglichkeiten einer auf Ausspracheabständen basierenden Dialektologie.

*Dieser Bericht entstand in einem Projekt in Zusammenarbeit mit Prof.Dr. Hermann Niebaum. Peter Kleiweg hat die Programme implementiert, und Frits Steenhuizen und Elwin Koster haben die Konvertierung der geographischen Daten vorgenommen. Prof.Dr. Jan Goossens und Prof.Dr. Jürgen-Erich Schmidt haben wertvolle Kommentare gegeben. Wir danken diesen Forschern und den Teilnehmern des 1. Kongresses der Internationalen Gesellschaft für Dialektologie des Deutschen in Marburg, 5.–8.3.2003 für ihre Verbesserungsvorschläge.

2 Methode

In diesem Abschnitt beschreiben wir den Hintergrund des vorliegenden Experiments, die Technik, die ihm zu Grunde liegt, und den experimentellen Entwurf mitsamt den zu erwartenden Resultaten.

2.1 Hintergrund

Der erste Befürworter des Einsatzes quantitativer Techniken in der Dialektologie war Séguy (1971), wogegen GOEBL (1982, 1984) viele neue Techniken untersucht und verfeinert hat. In der germanistischen Dialektologie waren Herrgen & Schmidt (1989) und Hummel (1993) frühe Anwender der dialektometrischen Techniken, die eine zunehmende Popularität gerade in den letzten Jahren genießen (siehe u.a. Lenz (2002), Lameli (2004), und Möller (2003)).

GOEBL betont, dass bezüglich bestimmter Fragen der Dialektgeographie, insbesondere der Definition von Dialektgebieten, das Problem der nicht-übereinstimmenden Isoglossen nie eine befriedigende Antwort erhalten hat. Dies bedeutet wiederum, dass man Methoden entwickeln muss, um mit den durch Isoglossen indizierten widersprüchlichen Tendenzen auszukommen. Sowohl SÉGUY als auch GOEBL behandeln Sprachdaten allerdings auf einem kategorialen Niveau, so dass Wahrnehmungen (Dialektdaten) letztendlich als gleich oder ungleich bewertet werden, nicht aber als *zu einem bestimmten Grad ähnlich bzw. unähnlich*. Nerbonne, Heeringa & Kleiweg (1999) zeigen, dass sich für das Messen von Ähnlichkeit bzw. Abstand von Dialekten anhand (transkribierter) Aussprachebeispiele besonders gut das *Levenshteinsche Sequenzabstandsmaß* eignet. Nerbonne & Heeringa (1998) wenden die Technik auf niederländische Dialektatlasdaten an, und Bolognesi & Heeringa (2002) erweitern diese Technik durch den Gebrauch von statistisch zuverlässigeren Stichproben, die aus sardinischen Aussprachedaten gewonnen werden. Heeringa & Gooskens (2003) verarbeiten mit derselben Methode auch akustische Daten aus dem norwegischen Sprachraum; die Ergebnisse vergleichen sie anschließend mit perceptionsbasierten Einschätzungen der Dialektabstände durch die Dialektsprecher selber. Für einen umfassenden Überblick über Technik und Möglichkeiten des Levenshteinschen Sequenzabstandsmaßes sei auf Heeringa (2004) verwiesen. Die vorliegende Arbeit ist die erste, die den Vergleich von deutschen Dialekten mittels dieser Technik zum Ziel hat.

2.2 Technik

Grundlage für das Berechnen von (Aussprache-)Distanzen zwischen Dialekten ist das Berechnen der Distanzen zwischen Wörtern, welches wiederum auf dem Vergleich der Phone basiert. Bei der einfachsten und bei hinreichender Datenmenge dennoch zuverlässigen Methode findet hier ein simpler Identitätsabgleich statt; d.h. sind die Phone identisch, beträgt der Abstand 0, sind sie nicht-identisch, ist er gleich 1. Um jedoch der Tatsache Rechnung zu tragen, dass z.B. ein [y] einem [i] ähnlicher ist als ein [p], experimentieren wir auch mit einigen weitestgehend artikulatorisch motivierten Merkmalsystemen, wie z.B. VIEREGGE & CUCCHIARINI oder ALMEIDA & BRAUN, die beide entwickelt wurden, um die Qualität von phonetischen Transkriptionen zu beurteilen (Vieregge, Rietveld & Jansen 1984, Almeida & Braun 1986, Heeringa & Braun 2003), und die beide zu ähnlichen Ergebnissen kommen. Wir illustrieren das Verfahren zum Bestimmen des Abstandes zwischen Segmenten anhand des Systems für Vokale von VIEREGGE in Abbildung 1. Wir haben bereits Experimente auf der Grundlage einiger Merkmalsysteme durchgeführt, bei der Validierung unserer Ergebnisse sind wir jedoch zu dem Schluss gelangt, dass die Gewichtung der Merkmale für unsere Belange noch sorgfältig modifiziert werden muss, so dass die hier präsentierten Ergebnisse auf dem einfacheren, fürs erste aber zuverlässigeren Identitätsabgleich der Phone basieren. Neben der Anpassung der Merkmalsysteme wäre es sicherlich interessant, zum einen noch das Merkmalsystem von HERRGEN & SCHMIDT zu implementieren, das eigens dazu aufgestellt wurde,

	[i]	[e]	[u]	d([i],[e])	d([i],[u])
Lage	2 (vorn)	2 (vorn)	6 (hinten)	0	4
Höhe	4 (hoch)	3 (mittelhoch)	4 (hoch)	1	0
Länge	3 (kurz)	3 (kurz)	3 (kurz)	0	0
Rundung	0 (nicht-rund)	0 (nicht-rund)	1 (rund)	0	1
				1	5

Abbildung 1: Der Segmentabstand wird auf der Basis von Merkmalsbeschreibungen der Segmente errechnet. Die Distanz d zwischen zwei Lauten entspricht der Summe der Abstände zwischen einzelnen Merkmalsausprägungen (z.B. vorn – hinten). Innerhalb des Systems von VIEREGGE ist $d([i],[e]) = 1$ (d.h der Abstand zwischen [i] und [e] beträgt insgesamt eins) und $d([i],[u]) = 5$.

Dialektalität zu messen (Herrgen & Schmidt 1989), und zum anderen direkt mit akustisch ermittelten Lautabständen zu arbeiten, womit Heeringa (2004) für das Niederländische noch bessere Ergebnisse erzielt hat als mit den artikulatorischen Merkmalssystemen. Nicht unbedingt geeignet sind phonologisch inspirierte Systeme, bei denen man sich vor einer Übergewichtung von Merkmalen hüten sollte, deren Motivation vor allem auf die Knappheit der Regeldarstellung zurückgeht. So ist z.B. das Merkmal $[\pm$ gespannt], das CHOMSKY und HALLE gebrauchen, um extrem vordere und extrem hintere Vokale von mittleren Vokalen zu unterscheiden, *prima facie* keine gute Basis für ein Abstandsmaß (Chomsky & Halle 1968).

Wie auch immer man nun die Abstände zwischen den Lautsegmenten ermittelt hat, sie bilden die Basis für den Levenshteinschen Sequenzvergleich. Dieser definiert den Abstand zwischen zwei Sequenzen (hier: Wörtern) als die Summe der *Kosten* aller Operationen, die nötig sind, um die eine in die andere Sequenz zu transformieren. Muss hierzu beispielsweise ein [b] durch ein [p] ersetzt werden, entsprechen die Kosten dieser Operation dem vorher definierten Segmentabstand zwischen den beiden Lauten. Ist nur diese eine Operation nötig, um die Wörter zu transformieren, entspricht der Segmentabstand auch dem Sequenzabstand. Sind mehrere Operationen nötig, werden die jeweiligen Kosten der Operationen addiert. Um eine willkürliche Sequenz in andere zu transformieren, braucht man neben dem Ersetzen von Lauten auch die Operationen 'Tilgen' und 'Hinzufügen'. Die Kosten dieser Operationen leitet man ebenfalls aus dem Segmentabstand ab, wozu man eine Definition des Segments 'Stille' benötigt. Bei der Verwendung von Merkmalssystemen wird hierfür die etwas modifizierte Merkmalsdefinition eines artikulatorisch 'neutralen' Vokals bzw. Konsonanten ([ə] bzw. [ʔ]) verwendet. Die Kosten entsprechen dann dem Segmentabstand zwischen der so definierten 'Stille' und dem zu tilgenden bzw. hinzuzufügenden Segment. Das Beispiel in Abbildung 2 soll die Idee veranschaulichen. Hierin wird die Realisierung des Wortes *Durst* im Dialekt von Aachen [tʊəf] in die korrespondierende Realisierung des Dialekts von Vielbrunn [tʊft] überführt (der Anschaulichkeit halber wird hier auf die Verwendung von Diakritika verzichtet). Die Kosten der Einzeloperationen, die hier auf dem Merkmalssystem von VIEREGGE basieren, werden zum

tʊəf	ersetze ʊ mit ɔ	1.5
tʊəf	tilge ə	1
tʊf	füge t hinzu	5
tʊft		
		7.5

Abbildung 2: Levenshteinabstand entspricht den Kosten für die effizienteste Reihe von Operationen, die eine Sequenz in die andere transformiert.

Sequenzabstand addiert.

Da die drei Arten von Operationen zur Transformation von Sequenzen beliebig oft und in beliebiger Reihenfolge angewendet werden können, gibt es meist mehrere Transformationsreihen, die sich durchaus im 'Preis' unterscheiden können. Als Levenshteinabstand werden jedoch nur die Gesamtkosten der 'kostengünstigsten' Reihe von Operationen bezeichnet. Der so definierte Abstand entspricht allen mathematischen Definitionselementen des Abstands, d.h. Symmetrie, Nichtnegativität, Dreiecksungleichheit sowie 0 genau dann wenn zwei Sequenzen gleich sind. Wichtig für unsere Zwecke ist, dass man die Resultate der Levenshteinmessung als echte Abstände analysieren kann. Kruskal (1999) bietet eine zugängliche Einführung in die Technik des Sequenzvergleichs.

Es sei noch kurz bemerkt, dass es inzwischen effiziente Algorithmen gibt, um den Levenshteinabstand sehr schnell zu berechnen (siehe u.a. <http://www.let.rug.nl/~kleiweg/lev> für ein public-domain-Paket für Levenshteinmessungen).

Obwohl ein ausführlicher Vergleich zwischen dieser Methode und Alternativen den Rahmen dieses kurzen Beitrags sprengen würde, könnten einige Bemerkungen wohl nützlich sein. Die meisten dialektometrischen Ansätze behandeln die Daten im Grunde genommen auf einem nominalen Niveau, d.h. als identisch mit anderen Daten oder nicht-identisch (SEGUY 1971; GOEBL 1982, 1984; MÖLLER 2003). Das innovative Element unserer Methode liegt darin, dass Ausspracheunterschiede als metrische Daten betrachtet werden, die erstens um einiges sensibler sind, um die sich zweitens deswegen den numerischen Analysetechniken erschließen lassen. Denselben Vorteil dürften HERRGEN & SCHMIDT (1989) und andere, die sich ebenfalls von ihrer Methode bedienen, für sich behaupten (HUMMEL 1993; LENZ 2002; LAMELI 2004), allerdings mit der Einschränkung, dass ihr Ansatz eines manuellen Eingriffs bedarf, um vergleichbare Laute zu betimmen, wogegen die hier beschriebene Methode vollautomatisch auf Aussprachelisten von vergleichbaren Wörtern angewendet wird. Wir erwarten nicht, dass die Ergebnisse der zwei Methoden radikal anders sein werden. Allerdings möchten wir betonen, dass unsere Methode auf Material aus Dialektatlanten angewendet werden kann, die aus anderen Gründen entworfen wurden, ohne einen zusätzlichen Aufwand für die Selektion der Daten. Es scheint uns daher gerechtfertigt, von einer objektiven – oder, vorsichtiger ausgedrückt, objektivierbaren – Methode zu sprechen.

2.3 Das aktuelle Experiment

Ziel des aktuellen Experiments ist es zu untersuchen, ob ein systematischer Vergleich vieler Wortrealisierungen verschiedener Dialekte im Aggregat zu einer zufriedenstellenden alternativen Basis für die Einteilung der deutschen Dialekte führen kann, die geeignet ist, die Probleme der subjektiv ausgewählten und letztendlich nie perfekt übereinstimmenden Isoglossen zu überwinden. Um eine Antwort liefern zu können, werden wir geeignetes Material von unterschiedlichen Erhebungsorten mit dem oben beschriebenen Sequenzvergleichsalgorithmus analysieren. Unsere Hypothese ist, dass die Prozedur zeigen wird, dass Orte mit ähnlicher Aussprache einander auch geographisch nah sind. Dies muss sich keineswegs zwangsläufig ergeben. Wenn die Technik nicht ausreichend valide oder sensibel ist, oder wenn die Dialektgeographie des Deutschen letztlich ohne eine subjektiv gelenkte Auswahl von Unterscheidungskriterien nicht auskommen kann, werden sich die als aussprachenah ermittelten Erhebungsorte willkürlich über den Erhebungsraum verteilen. Wir werden also untersuchen, ob die Ergebnisse erstens geographische Kohärenz zeigen, und ob sie zweitens dem wissenschaftlichen Konsensus im Großen und Ganzen entsprechen.

Abend	acht	Affe	alle	als	alte	alter
ändern	anfängt	Apfelbäume	Äpfelchen	Äpfeln	auch	auf
Augenblick	austrinken	bald	bauen	Bauern	bei	beißen
Berge	besser	bestellt	Blätter	bleib	böse	bösen
braune	Brot	Bruder	das	deiner	dich	die
Dienstag	Donnerstag	Dorf	drei	dreschen	du	durchgelaufen
dürft	Durst	Ei	Eier	eingebrochen	eingeschlafen	eins
Eis	Eise	elf	erzählt	fängt	Feld	Felde
fest	Feuer	Fleisch	fliegen	Frau	Freitag	fünf
für	Gänse	Garten	geblieben	gebracht	gebrannt	gefahren
gefallen	gefunden	gekannt	gekommen	genug	geschlafen	geschmolzen
gestern	gestohlen	gestorben	glaube	gleich	groß	größer
gut	hat	Haus	Hause	Häuser	heiß	heiße
Herz	Herzen	hoch	höher	hört	ich	isst
ist	ja	kalte	kein	Kind	Kochlöffel	Kohlen
Kühe	lauter	Leute	Leuten	liebes	liegen	Luft
machen	Mann	mein	meinem	Milch	mit	mitgehen
Mittwoch	Montag	müde	muss	müsst	Mutter	nähen
neue	neun	oben	Ochsen	Ofen	ohne	Ohren
Pfeffer	Pfund	recht	rot	roten	Salz	Samstag
schlechte	schlechten	Schnee	schneien	schöne	schwarz	Schwester
sechs	Seife	sein	selbst	sich	sie	sieben
so	sollen	Sonntag	Stückchen	Tisch	Tochter	tot
treiben	trinken	tu	tut	über	um	und
uns	unsere	unserem	unten	verkaufen	versteht	viel
vier	von	vor	wachsen	wäre	was	Wasser
weh	Wein	weiße	wem	werden	Wetter	wieviel
will	Winter	wo	Wochen	wollte	Wort	Wurst
zehn	Zeiten	zum	zwei	zwölf		

Tabelle 1: Die 201 vorstehenden Wörter aus den Wenkersätzen wurden für den Kleinen Deutschen Lautatlas – Phonetik (siehe Text) erhoben und bilden die Grundlage für die Messungen in dieser Arbeit.

3 Daten

Um den Ausspracheabstand innerhalb der deutschen Dialekte zu bestimmen, wurden die Wortabstände zwischen phonetischen Transkriptionen von 201 Wörtern aus 186 Erhebungsorten auf dem Gebiet der heutigen Bundesrepublik gemessen. Die Summe dieser 201 Wortabstände¹ gilt als das Maß des Dialektabstands zwischen den jeweiligen Erhebungsorten.

Die Daten wurden im Rahmen des am Forschungsinstitut für deutsche Sprache “Deutscher Sprachatlas“ in Marburg initiierten Projekts “Kleiner Deutscher Lautatlas – Phonetik“ (auch bekannt als “Phonetischer Atlas der Bundesrepublik Deutschland“) erhoben (Göschel 1992). Weil diese qualitativ hochwertigen Daten bis jetzt allerdings relativ unbekannt geblieben sind, ist es sinnvoll, sie hier näher zu beschreiben.

Die Transkriptionen wurden auf der Basis von Aufnahmen gemacht, die unter der Leitung von JOACHIM GÖSCHEL zuerst in den 1960er und 1970er Jahren in der alten BRD erhoben worden waren (Göschel 1992, S. 64-70). Nach der Wende konnte man vergleichbare Erhebungen auch auf dem Gebiet der ehemaligen DDR durchführen, so dass das Netz von Erhebungsorten innerhalb der heutigen Staatsgrenzen relativ vollständig ist. Die elizitierten Wörter (siehe Tabelle 1) stammen aus den für die Aufnahmen verwendeten ‘Wenker-Sätzen’, d.h. den Sätzen, die GEORG WENKER 1879–87 durch die Lehrer von ca. 40.000 Schulorten des damaligen Deutschen Reichs in den jeweiligen Ortsdialekt übertragen ließ. Der Rückgriff auf die Wenker-Sätze war durch den Wunsch motiviert, die historische Entwicklung

¹Tatsächlich sind es jeweils etwas weniger, da i.d.R. nicht für jeden Dialekt alle Wortrealisierungen vorliegen. Zur Berechnung des Abstandes zwischen zwei Dialekten werden ausschließlich die Wörter verwendet, für die in beiden Dialekten Realisierungen vorhanden sind.

über ein Jahrhundert hinweg näher zu studieren. Die Erhebungsorte sind in Abbildung 3 eingetragen; die Motivation für ihre Auswahl scheint jedoch undokumentiert geblieben zu sein.

Die Daten wurden von einer Gruppe in Marburg in IPA transkribiert (IPA 1949), wobei alle Transkriptionen dem Konsens der Mitarbeiter entsprechen mussten. Die Transkribenten waren ANTONIO ALMEIDA, ANGELIKA BRAUN, RAPHAELA LAUF und KLAUS MONTERMANN. Im Jahr 2002 haben ROGIER NIEUWEBOER, SAAKJE VAN DELLEN und BERTHIEN MARKVOORT im Rahmen einer Kooperation zwischen den Universitäten Groningen und Marburg² die Transkriptionen digitalisiert, wobei eine etwas modifizierte Version von X-SAMPA als maschinenlesbare Kodierung verwendet wurde (für Informationen über X-SAMPA, siehe <http://coral.lili.uni-bielefeld.de/LangDoc/EGA/Formats/Sampa/sampa.html>).

Das Material wird vom Deutschen Sprachatlas in Marburg (<http://www.uni-marburg.de/dsa/>) verwaltet und ist bislang nicht öffentlich zugänglich.

4 Resultate

Das unmittelbare Resultat der Abstandsbestimmung ist eine Abstandstabelle, in der die ermittelten Ausspracheabstände in den Zellen stehen (anstatt von z.B. Kilometerangaben in Tabellen für Städteentfernungen). Weil diese Tabelle aber viel zu groß ($\binom{186}{2} = 17.205$ Abstände) ist, um hier abgedruckt zu werden, bedienen wir uns einer Visualisierung, um die Resultate global zu verstehen. Abbildung 4 zeigt im Prinzip *alle* phonetischen Abstände dadurch, dass kleine phonetische Abstände zwischen Orten dunklen und große Abstände hellen Linien entsprechen. Hierbei fallen vier größere Gebiete ins Auge. Dies wären zum einen ein jeweils relativ homogenes südliches, östliches und nördliches Dialektgebiet und zum anderen ein offensichtlich extrem heterogenes westliches Gebiet.

Die 'Linienkarte' liefert also offenbar eine Antwort auf die Frage nach der geographischen Kohärenz. Wenn die geographisch entfernteren Orte nach dem Levenshteinschen Maß nicht auch phonetisch unähnlicher wären, würden wir u.a. viele dunklere Linien zwischen relativ weit auseinanderliegenden Orten sehen. Wie wir aber Abbildung 4 entnehmen können, kommt signifikantere phonetische Ähnlichkeit praktisch nur für Paare von Orten vor, die auch geographische Nachbarn sind.

Der zweite Aspekt der Untersuchungsfrage gilt dem Vergleich der aktuellen Ergebnisse mit der traditionellen Dialekteinteilung durch die heutige Dialektwissenschaft. Weil diese vor allem mit Begriffen wie 'Dialektgebiet' arbeitet, ist es sinnvoll, die Gebiete, die implizit in der Tabelle der phonetischen Abstände vorhanden sind und die sich bereits in der Linienkarte andeuten, klarer zum Vorschein zu bringen. Eine hierfür intuitiv einleuchtende Methode ist das sogenannte *Clustering* (genauer: hierarchisches agglomeratives Clustering) (Jain & Dubes 1988). Beim Clustering identifiziert man innerhalb einer Abstandstabelle stets die beiden Elemente mit dem geringsten Abstand, die man anschließend zusammenfügt und damit eine kleinere Tabelle erstellt. Wir illustrieren das Verfahren in Tabelle 2.

Dem Clustering wird zwar vorgeworfen, dass es mathematisch nicht 'stabil' sei, d.h., dass kleine Unterschiede in den Ausgangsdaten mitunter zu großen Unterschieden in den Ergebnissen führen können. Dafür ist diese Methode anschaulich und führt oft zu sehr respektablen Ergebnissen. Es gibt, in Abhängigkeit von der Bestimmung der restlichen Abstände zum neu zusammengeführten Element (Aldenderfer & Blashfield 1984, Heeringa 2004), sehr viele Varianten des Clusterings, die ebenfalls zu deutlich unterschiedlichen Ergebnissen kommen können. Wir haben aus diesem Grund mehrere Clustering-Verfahren getestet und sind dabei zu dem Schluss gelangt, dass mit dem Clustering durch gewogenes arithmetisches Mittel (*Weighted arithmetic average clustering*) die am einfachsten zu interpretierenden

²Die Kooperation wurde von HERMANN NIEBAUM initiiert; ANGELIKA BRAUN, JÜRGEN-ERICH SCHMIDT, JOHN NERBONNE und WILBERT HEERINGA sind bislang die weiteren Teilnehmer.



Abbildung 3: Verteilung der 186 Aufnahmeorte

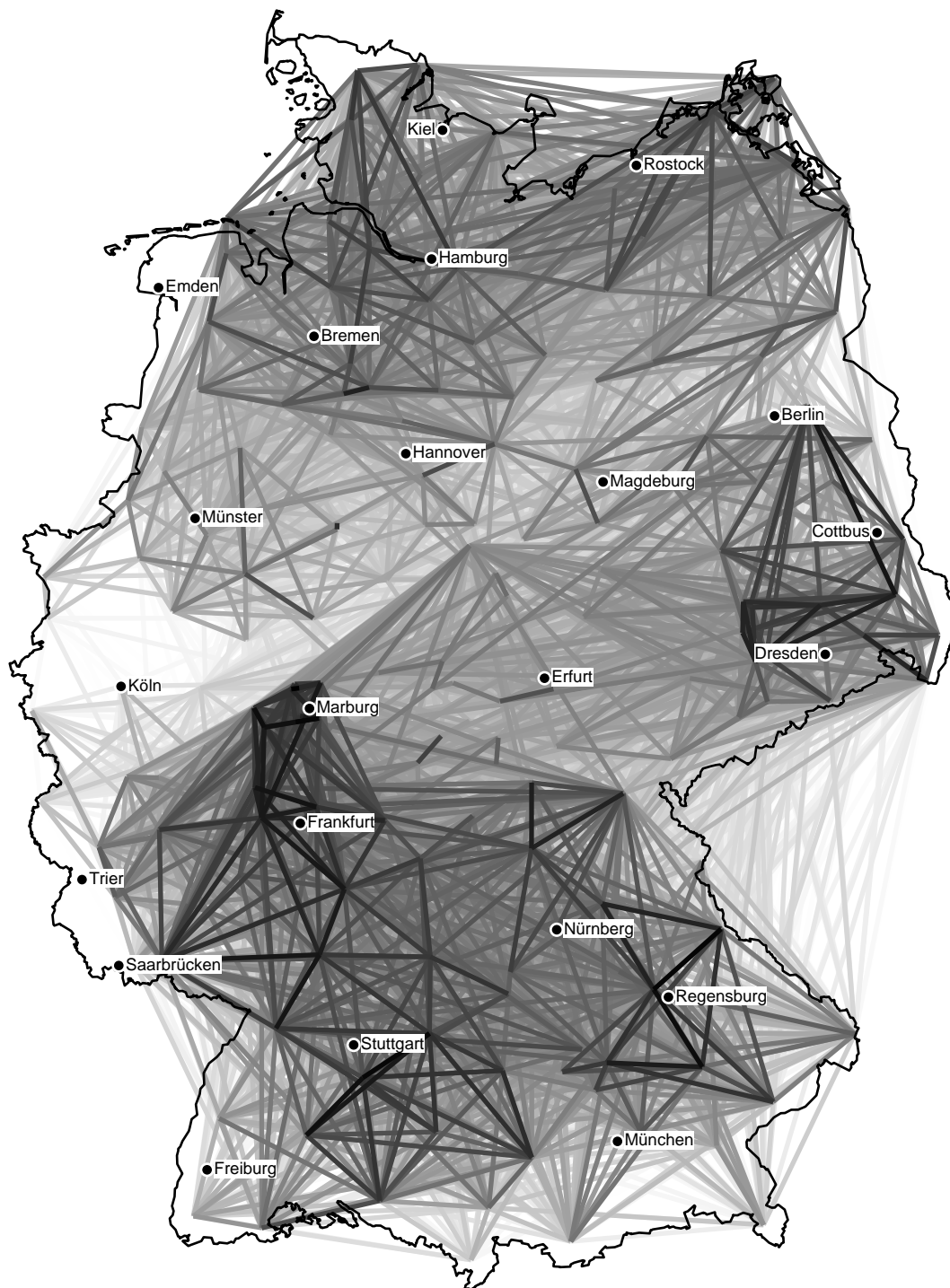


Abbildung 4: Ausspracheabstände zwischen *allen* Orten, wobei gilt: Je ähnlicher die beiden Dialekte, desto dunkler ihre Verbindungslinie. Zur Orientierung sind einige größere Städte eingezeichnet, von denen jedoch (mit Ausnahme von Köln und Nürnberg) keine Aufnahmen existieren.

	Jever	Oberau	Köln	Aachen	Hagen
Jever	0	73	64	67	79
Oberau		0	81	74	68
Köln			0	43	91
Aachen				0	86
Hagen					0

	Jever	Oberau	Köln-Aachen	Hagen
Jever	0	73	65.5	79
Oberau		0	77.5	68
Köln-Aachen			0	88.5
Hagen				0

Tabelle 2: Clustering identifiziert Gruppen in einer Abstandstabelle durch wiederholtes Zusammenfügen der jeweils nächsten Elemente. In der ersten Abstandstabelle weisen Köln und Aachen die geringsten Ausspracheunterschiede auf, weswegen sie in der unteren Tabelle zusammengefügt werden. Die Abstände der anderen Dialekte zu dem neuen Element Köln-Aachen entsprechen in diesem Beispiel einfach dem Durchschnitt der Abstände zu den Komponenten Köln und Aachen.

Resultate zu erzielen sind; diese Methode wird daher im Folgenden verwendet. Beim Clustering durch gewogenes arithmetisches Mittel wird die Anzahl der Elemente in den Clustern nicht berücksichtigt, weswegen die einzelnen Elemente in einem Cluster mit vielen Elementen weniger stark gewichtet werden als in einem Cluster mit wenigen Elementen. Dies mag auf den ersten Blick als wenig sinnvoll erscheinen, ist für unsere Zwecke jedoch von Bedeutung, da die Verteilung der Aufnahmeorte nicht völlig gleichmäßig ist, so dass sich deren Dichte in den späteren Clustern erheblich unterscheiden kann. Würde die Anzahl der Elemente in die Cluster-Analyse einfließen, würden Gebiete mit einer hohen Aufnahmeortdichte (z.B. der Raum Marburg) anders verarbeitet als Gebiete mit einer niedrigen Aufnahmeortdichte (z.B. der Raum Freiburg oder insbesondere das Gebiet nördlich von Köln), was zu einer unrealistischen Klassifikation führen könnte.

Wenn wir die Resultate des Clusterings auf einer Karte abbilden (Abb. 5), so erhalten wir nähere Informationen zu unserer Untersuchungsfrage, ob eine systematische Messung von aggregierten Ausspracheunterschieden im Stande ist, die Dialektunterschiede im Deutschen zu erfassen. Wie sich bereits in der Linienkarte andeutete, finden wir eine geographisch kohärente Dreiteilung, die beim Clustering deutlich zu Tage tritt, sowie ein Gebiet im Westen Deutschlands, das offenbar durch eine große Heterogenität gekennzeichnet ist. Auf die Gefahr hin, uns zu wiederholen, sei hier bemerkt, dass das Clustering in dieser Anwendung keineswegs geographisch inspiriert ist. Wenn die Levenshteinsche Methode nicht geeignet wäre oder der Dialektsystematik des Deutschen keine geographische Struktur zu Grunde läge, wären die Schattierungen in Abbildung 5 willkürlich im Raum verteilt.

Ein zweiter Aspekt der Untersuchungsfrage war, das Verhältnis zwischen den Ergebnissen der Levenshteinschen Methode und der Sicht der aktuellen Dialektwissenschaft zu vergleichen. Abbildung 5 zeigt eine Verteilung der großen Dialektgebiete, die nur wenig von der der Standardwerke abweicht (vgl. z.B. König (1994) oder Niebaum & Macha (1999)). Auch die 'Benrather Linie' ist anhand typischer Merkmale (z.B. der 'Berliner Trichter' oder die Harzer Dialektenklave) eindeutig zu erkennen, und auch die große Heterogenität der Dialekte im Westen Deutschlands findet sich in der traditionellen Dialektologie als 'Rheinischer Fächer' wieder. Die Hierarchie der Dialektgebiete, die im Dendrogramm (Abb. 6)



Abbildung 5: Ergebnis des Clusterings als Karte dargestellt. Erkennbar sind drei Hauptgebiete, die sich im Wesentlichen mit den Verteilungen des Nieder-, Ostmittel- und Oberdeutschen (Cluster 1, 4 und 5) nach traditioneller Einteilung decken, sowie ein heterogenes Gebiet im Westen, das in etwa Ripuarisch (Cluster 3) und Niederrheinisch-Westmünsterländisch (Cluster 2) entspricht.

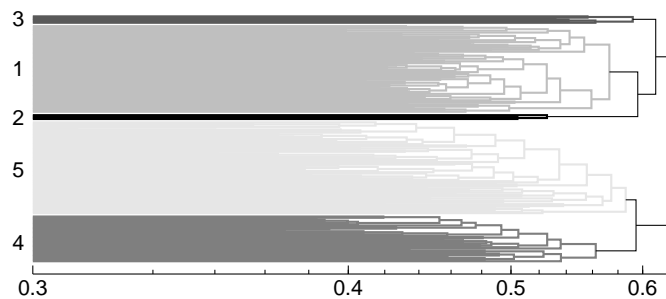


Abbildung 6: Resultat des Clusterings als Dendrogramm dargestellt. Die Grautöne und Clusternummern stimmen mit denen in Abbildung 5 überein. Die Grenze zwischen den nieder- und hochdeutschen Dialekten ist dominant. Das ripuarische Cluster wird mit dem hier verwendeten Clustering durch gewichtetes arithmetisches Mittel dem Niederdeutschen zugeordnet (Diskussion hierzu im Text).

dargestellt ist, deckt sich ebenfalls mit der herkömmlichen Klassifikation. Die ostmittel- und oberdeutschen Dialekte (Cluster 4 und 5) werden auch durch das Clustering zu 'Hochdeutsch' zusammengefasst; die Grenze zu den niederdeutschen Dialekten (Cluster 1-3) ist ausgeprägt. Der wesentliche Unterschied zur traditionellen Einteilung besteht darin, dass auf der Basis unserer bisherigen Resultate ein kohärentes westmitteldeutsches Dialektgebiet nicht auszumachen ist. Einerseits werden die traditionell als 'Westmitteldeutsch' klassifizierten Dialekte wie Hessisch, Mosel- und Rheinfränkisch zwar dem Hochdeutschen zugeordnet, sie werden jedoch nicht, oder zumindest nicht im gleichen Ausmaß vom Oberdeutschen abgegrenzt wie die ostmitteldeutschen Dialekte. Andererseits wird Ripuarisch (Cluster 3) als Niederdeutsch klassifiziert (s. Dendrogramm), wobei allerdings anzumerken ist, dass dieser Zusammenschluss auf einer relativ hohen Ebene in der Hierarchie stattfindet, was darauf hindeutet, dass dieses Gebiet eine Sonderstellung zwischen den hoch- und niederdeutschen Dialekten einnimmt. Verwendet man ein anderes geeignetes Cluster-Verfahren, wie z.B. Clustering durch ungewichtetes arithmetisches Mittel, wird das ripuarische Cluster dem Hochdeutschen zugeordnet. Es ist jedoch nicht unwahrscheinlich, dass die Verwendung einer anderen Methode der Lautabstandsbestimmung, die die Konsonanten stärker gewichtet als dies der einfache Phon-Identitätsabgleich tut, was sich natürlich auf die Dialektabstände auswirkt, bei einem Gebiet mit einer so kritischen Stellung in der Hierarchie durchaus zu einer anderen Klassifikation führen kann. Generell sollte die Verwendung einer anderen Methode zur Bestimmung der Lautabstände allerdings eher zu quantitativen als qualitativen Unterschieden führen; d.h. es kann zu kleinen Veränderungen im Grenzverlauf kommen, was in einigen Gebieten auch noch notwendig erscheint. Zu den Grenzen allgemein sei noch angemerkt, dass ihr Verlauf natürlich auch durch die Anzahl an Messpunkten mitbestimmt wird. Um die Punkte herum werden Polygone definiert, also kleine Areale, für die derselbe Dialekt angenommen wird, der in dem jeweiligen Referenzort gesprochen wird. Bei abnehmender Messpunkt-Dichte kommt es demnach zu einer zunehmend ungenauen Abbildung der realen Verteilung der Dialekte. Unser Messpunkt-Raster ist relativ grob, d.h. Feinheiten können im Allgemeinen nicht wiedergegeben werden. Für unsere Zwecke entscheidend ist jedoch, dass es wie demonstriert fein genug ist, um die wesentlichen Charakteristika der Verteilung wiederzugeben.

Die klaren Grenzen in den Cluster-Karten dienen vor allem der Verdeutlichung, entsprechen jedoch nicht unbedingt der Realität. Wie bereits erwähnt, verlaufen die Grenzen zwischen den Dialekten tatsächlich oft deutlich fließender, so dass man insgesamt eher von einem Dialektkontinuum sprechen sollte. Mit Hilfe des Levenshteinschen Sequenzabstandsmaßes kann man nun auch dieses Kontinuum darstellen. Wir bedienen uns hierfür der sogenannten *multidimensionalen Skalierung* (Tourgerson 1952, Heeringa 2004). Zur Illustration dieser Methode stelle man sich z.B. einen Quader vor, in dem jede Ecke durch ihren Abstand zu allen anderen Ecken definiert ist. Stünden einem nun nur die Abstände zur Verfügung, könnte man den Quader allein aus diesen Angaben rekonstruieren. Der Unterschied bei der Anwendung auf Dialekte liegt hauptsächlich darin, dass diese nicht von vorneherein dreidimensional sind. Nichtsdestotrotz werden sie bei einer multidimensionalen Skalierung so in drei Dimensionen³ angeordnet, dass ihre jeweiligen Abstände zu allen anderen Dialekten möglichst genau wiedergegeben werden. Es findet also eine, natürlich nicht verlustfreie, Datenkompression statt. Ordnet man nun jeder Dimension eine Farbe zu (hier rot, grün und blau), erhält jeder Dialekt auf Grund seiner Lage in den drei Dimensionen eine für ihn typische Farbe. Ähnliche Dialekte, die im dreidimensionalen Raum nah beieinander liegen, erhalten demnach auch ähnliche Farben. Diese lassen sich wiederum in einer (zweidimensionalen) Karte darstellen. Abbildung 4 zeigt solch eine Regenbogenkarte als Ergebnis einer multidimensionalen Skalierung. Nun liegt es in der Natur der Methode, dass Grenzen kaum auszumachen sind, dennoch lassen sich die wichtigsten Unterscheidungen (Oberdeutsch grün, Ostmitteldeutsch

³Die Beschränkung auf drei Dimensionen ist nicht theoretisch begründet. Möglich wären auch zwei oder mehr Dimensionen, wobei gegen zwei Dimensionen jedoch die schlechteren Ergebnisse und gegen mehr Dimensionen die Probleme bei Darstellung und Verständlichkeit sprechen.



Abbildung 7: Visualisierung des Dialektkontinuums mittels multidimensionaler Skalierung. Die Farbzeichnungen sind u.U. diskussionsbedürftig, um jedoch Missverständnissen im Text vorzubeugen, legen wir uns an dieser Stelle auf folgende Termini fest: Oberdeutsch entspricht grün, Ostmitteldeutsch hellblau, Ostniederdeutsch dunkel-lila und Westniederdeutsch weitestgehend rötlich-violett, wobei auch Westniederdeutsch in Richtung des rheinischen Fächers zunehmend heterogener wird, für den wir die Bezeichnung 'rötlich) bunt' verwenden.

hellblau, Niederdeutsch dunkel-lila bis rötlich-violett, in Richtung des rheinischen Fächers zunehmend rötlich-bunt) relativ deutlich erkennen. Tatsächlich erweist sich nur die Abgrenzung der Dialektgebiete im Westen als wirklich unmöglich. Innerhalb der großen Gebiete sind weitere Farbabstufungen auszumachen, die Grenzen verlaufen hier allerdings noch um einiges fließender. So ist in dem von uns als 'Oberdeutsch' definierten grünen Gebiet ein etwas heller gefärbter Bereich im Westen zu erkennen, der mit Hessisch, Mosel- und Rheinfränkisch korrespondiert sowie ein tendenziell dunklerer Bereich im Osten, der in etwa Ostfränkisch und Bairisch entspricht. Zudem sind einige leicht anders gefärbte 'Ausreißer' zu erkennen, die meist in der Nachbarschaft größerer Städte oder im unmittelbaren Grenzbereich zu finden sind. Dies müsste allerdings noch näher untersucht werden. Im niederdeutschen Dialektraum lässt sich ein eher dunkel-lila gefärbter ostniederdeutscher von dem ins Rötlich-Violette tendierenden westniederdeutschen Bereich abgrenzen. Letzterer weist allerdings vor allem in Richtung Süden zum rheinischen Fächer hin eine so große Heterogenität auf, dass in der gesamten Region keine klaren Grenzen auszumachen sind.

Wie wir zeigen konnten, sind die mit dem Levenshteinschen Sequenzabstandsmaß gewonnenen Ergebnisse nicht nur geographisch kohärent, sie decken sich alles in allem auch mit der Einteilung der traditionellen Dialektologie. Natürlich ist Letzteres kein Qualitätskriterium an sich, es fördert jedoch das Vertrauen in eine neue Methode, dass man eine Kongruenz mit bekannten Techniken und Resultaten sehen kann. Insgesamt gelangen wir daher zu dem Schluss, dass ein systematischer oder gar mechanischer Vergleich aggregierter Ausspracheunterschiede eine zufrieden stellende alternative Basis für die Dialektgeographie darstellt, wodurch nicht nur Resultate erzielt werden können, die den mit traditionellen Ansätzen gewonnenen nahezu ebenbürtig sind, sondern zusätzlich, wie mit der Darstellung des Dialektkontinuums mittels multidimensionaler Skalierung demonstriert, bislang unzugängliche Möglichkeiten eröffnet werden.

5 Ausblick

Diese Arbeit ist ein Bericht aus einem laufenden Projekt. Eine detaillierte Betrachtung vieler Alternativen innerhalb der Messtechnik ist sinnvoll, wie auch nähere Untersuchungen zum Einfluss der phonetischen Merkmale, der Auswahl der Wörter sowie der Auswahl der Erhebungsorte. Weiterhin wäre es zweckmäßig, die Konsistenz und Validität der Messungen genauer zu betrachten (Heeringa, Nerbonne & Kleiweg 2002), und auch die traditionelle Einteilung der Dialektgebiete aus diesem Blickwinkel nochmals zu inspizieren.

Es wäre jedoch falsch zu vermuten, dass die bloße Bestätigung der zeitgenössischen Sicht die Ziele dieser Forschungsrichtung erschöpft. Wie Nerbonne & Heeringa (1998) und Heeringa & Nerbonne (2002) in Ansätzen zeigen, eröffnet die Verwendung des Levenshteinschen Sequenzabstandsmaßes eine Reihe neuer Möglichkeiten. So ist es offensichtlich, dass diese Methode bei geeignetem Material auch auf Fragen der Varietätslinguistik außerhalb der Dialektgeographie angewendet werden könnte. Vielleicht noch interessanter dürfte allerdings der Aspekt sein, dass sich die Dialektologie mit diesem quantitativen Maß nun die Techniken der numerischen Analyse erschließt, welche entscheidend dazu beitragen könnten, Antworten auf Fragen nach z.B. dem Einfluss außersprachlicher Variablen auf die sprachliche Varietät zu finden.

Literatur

- Aldenderfer, Mark S. & Roger K. Blashfield (1984), *Cluster Analysis*, Quantitative Applications in the Social Sciences, Sage, Beverly Hills.
- Almeida, Almerindo & Angelika Braun (1986), ‘Richtig’ und ‘Falsch’ in phonetischer Transkription: Vorschläge zum Vergleich von Transkriptionen mit Beispielen aus deutschen Dialekten’, *Zeitschrift für Dialektologie und Linguistik* **LIII**(2), 158–172.
- Bloomfield, Leonard (1933), *Language*, Holt, Rhinehart and Winston, New York.
- Bolognesi, Roberto & Wilbert Heeringa (2002), ‘De invloed van dominante talen op het lexicon en de fonologie van Sardische dialecten’, *Gramma/TTT: Tijdschrift voor Taalwetenschap* **9**(1), 45–84.
- Chomsky, Noam A. & Morris Halle (1968), *The Sound Pattern of English*, Harper and Row, New York.
- Coseriu, Eugenio (1975), *Die Sprachgeographie*, Gunter Narr, Tübingen. ¹1956.
- Goebel, Hans (1982), *Dialektometrie: Prinzipien und Methoden des Einsatzes der Numerischen Taxonomie im Bereich der Dialektgeographie*, Österreichische Akademie der Wissenschaften, Wien.
- Goebel, Hans (1984), *Dialektometrische Studien: Anhand italo-romanischer, rätoromanischer und galloromanischer Sprachmaterialien aus AIS und ALF. 3 Vol.*, Max Niemeyer, Tübingen.
- Göschel, Joachim (1992), Das Forschungsinstitut für Deutsche Sprache “Deutscher Sprachatlas”, Wissenschaftlicher Bericht, Das Forschungsinstitut für Deutsche Sprache, Marburg.
- Heeringa, Wilbert (2004), *Measuring Dialect Pronunciation Differences using Levenshtein Distance*, PhD thesis, Rijksuniversiteit Groningen.
- Heeringa, Wilbert & Angelika Braun (2003), ‘The use of the Almeida-Braun system in the measurement of Dutch dialect distances’, *Computers and the Humanities* **37**(3), 257–271. Special Issue on Computational Techniques in Dialectometry, ed. by John Nerbonne and William Kretzschmar.
- Heeringa, Wilbert & Charlotte Gooskens (2003), ‘Norwegian dialect examined perceptually and acoustically’, *Computers and the Humanities* **37**(3), 293–315.
- Heeringa, Wilbert & John Nerbonne (2002), ‘Dialect areas and dialect continua’, *Language Variation and Change* **13**, 375–400.
- Heeringa, Wilbert, John Nerbonne & Peter Kleiweg (2002), Validating dialect comparison methods, in W.Gaul & G.Ritter, eds, ‘Proceedings of the 24th Annual Meeting of the Gesellschaft für Klassifikation’, Springer, Heidelberg, pp. 445–452.
- Herrgen, Joachim & Jürgen Erich Schmidt (1989), Kontrastive Dialektgeographie, in W.Putsche, W.Veith & P.Wiesinger, eds, ‘Dialektgeographie und Dialektologie’, Vol. 90 of *Deutsche Dialektgeographie*, N.G.Elwert Verlag, Marburg, pp. 304–346.
- Hummel, Lutz (1993), *Dialektometrische Analysen zum kleinen deutschen Sprachatlas (KDSA)*, Max Niemeyer Verlag, Tübingen.
- IPA (1949), *The Principles of the International Phonetic Association*.

- Jain, K. & R. C. Dubes (1988), *Algorithms for Clustering Data*, Prentice Hall, Englewood Cliffs, New Jersey.
- König, Werner (1994), *DTV Atlas zur deutschen Sprache*, Deutscher Taschenbuch Verlag, München.
¹1978.
- Kruskal, Joseph (1999), An overview of sequence comparison, in D.Sankoff & J.Kruskal, eds, 'Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison', CSLI, Stanford, pp. 1–44. ¹1983.
- Lameli, Alfred (2004), *Standard und Substandard: Regionlismen im diachronen Längsschnitt*, Steiner, Stuttgart.
- Lenz, Alexandra (2002), Struktur und Dynamik des Substandards. Eine Studie zum Westmitteldeutschen (Wittlich/Eifel), PhD thesis, Marburg.
- Möller, Robert (2003), 'Zur diatopischen Gliederung des alltagsprachlichen Wortgebrauchs. eine dialektometrische Auswertung von Jürgen Eichhoff: Wortatlas der Deutschen Umgangssprache, Bd. 1/4 (Bern/München)', *Zeitschrift für Dialektologie und Linguistik* **LXX**(3), 259–297.
- Nerbonne, John & Wilbert Heeringa (1998), 'Computationele vergelijking en classificatie van dialecten', *Taal en Tongval* **50**(2), 164–193.
- Nerbonne, John, Wilbert Heeringa & Peter Kleiweg (1999), Edit distance and dialect proximity, in D.Sankoff & J.Kruskal, eds, 'Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison, 2nd ed.', CSLI, Stanford, CA, pp. v–xv.
- Niebaum, Hermann & Jürgen Macha (1999), *Einführung in die Dialektologie des Deutschen*, Niemeyer, Tübingen. ¹1983.
- Séguy, Jean (1971), 'La relation entre la distance spatiale et la distance lexicale', *Revue de Linguistique Romane* **35**, 335–357.
- Tourgerson, Warren S. (1952), 'Multidimensional scaling I', *Psychometrika* **17**, 401–419.
- Vieregge, Wilhelm H., A.C.M. Rietveld & Carel Jansen (1984), A distinctive feature based system for the evaluation of segmental transcription in Dutch, in M. P.van den Broecke & A.Cohen, eds, 'Proc. of the 10th International Congress of Phonetic Sciences', Dordrecht, pp. 654–659.