

AN AGGREGATE ANALYSIS OF PRONUNCIATION IN THE GOEMAN-TAEDEMAN-VAN REENEN- PROJECT DATA

Abstract

Contemporary Dutch dialects are compared using the Levenshtein distance, a measure of pronunciation difference. The material consists of data from the most recent Dutch dialect source available: the Goeman-Taeldeman-Van Reenen-Project (GTRP). This data consists of transcriptions of 1876 items for 613 localities in the Netherlands and Belgium gathered during the period 1980 – 1995. In addition to presenting the analysis of the GTRP, we compare the dialectal situation it represents to the *Reeks Nederlands(ch)e Dialectatlassen* (RND), in particular to the 350-locality sample studied by Heeringa (2004), noting areas of convergence and divergence. Although it was not the purpose of the present study to criticize the GTRP, we nonetheless note that transcriptions from Belgian localities differ substantially from the transcriptions of localities in the Netherlands, impeding the comparison between the varieties of the two different countries. We therefore analyze the developments in the two countries separately.

1. Introduction

The Goeman-Taeldeman-Van Reenen-Project (GTRP; Goeman & Taeldeman 1996) is an enormous collection of data collected from the Dutch dialects, including transcriptions of over 1800 items from over 600 localities, all collected over a relatively brief, and therefore, unproblematic time interval (15 years, 1980 – 1995). The GTRP is the first large-scale collection of Dutch dialect data since Blancquaert & Peé's *Reeks Nederlands(ch)e Dialectatlassen* (RND; 1925 – 1982), and complements it as a more recent and temporally more limited set. The GTRP provides a rich and attractive database, designed by the leading experts in Dutch dialectology, who likewise collaborated in obtaining, transcribing, and organizing its information. The GTRP rivals the RND in being fully available digitally and being designed with an eye toward contemporary questions in phonology, morphology and variationist linguistics (Van Oostendorp, to appear).

We present the GTRP and the RND in more detail in Section 2.

The present paper provides an aggregate analysis of the pronunciation variation in this collection, using the same techniques for analysis which Nerbonne et al. (1996) first applied, and which Heeringa (2004) lays out in full detail. The aggregate analysis proceeds from a word-by-word measurement of pronunciation differences, which has been shown to provide consistent probes into dialectal relations, and which correlates strongly ($r > 0.7$) with lay dialect speakers' intuitions about the degree to which non-local dialects sound "remote" or "different" (see Heeringa 2004: Chapter 7; and Heeringa et al. 2006 for rigorous discussions of the consistency and validity of the measures). The aggregate analysis differs from analyses based on a small number of linguistic variables in providing a global view of the relations among varieties, allowing more abstract questions to be posed about these relations. We sketch the necessary technical background for the measurement of pronunciation differences in Section 3 below.

For various technical reasons, we restrict our analysis to 562 items in the GTRP, which is nonetheless notably large compared to other analyses. We present the results of this analysis in Sections 4.1 and 4.2 below.

A second, related goal of this paper is to examine what has changed between the time of the RND and that of the GTRP. For this purpose we focus our attention on 224 localities which are common to the GTRP and the RND varieties analyzed by Heeringa (2004). To allow interpretation to be as exact as possible, we also focused on the 59 words which were common to the GTRP and the RND. Since the two projects differed in methodologies, especially transcription practice, we approach the comparison indirectly, via regression analyses. We are able to identify several areas in which dialects are converging (relatively), and likewise several in which they are diverging. The results of the comparison are the subject of Section 4.3 below.

It was not originally a goal of the work reported here to examine the GTRP with respect to its selection and transcription practices, but several preliminary results indicated that the Belgian and the Dutch collaborators had not been optimally successful in unifying these practices. We follow these indications up, and conclude in Section 4.1 that caution is needed in interpreting aggregate results unless one separates Dutch and Belgian material. We further suggest that these problems are likely to infect other, non-aggregating approaches as well. At the end

of Section 4.2 we discuss some clues that fieldworker and transcription practices in the Netherlands may be confounding analyses to some degree. Also Hinskens & Van Oostendorp (2006) reported transcriber effects in the GTRP data.

2. Material

In this study two Dutch dialect data sources are used: data from the Goeman-Taeldeman-Van Reenen-Project (GTRP; Goeman & Taeldeman 1996) and data from the *Reeks Nederlands(ch)e Dialectatlassen* (RND; Blancquaert & Peé 1925 – 1982) as used by Heeringa (2004).

2.1. GTRP

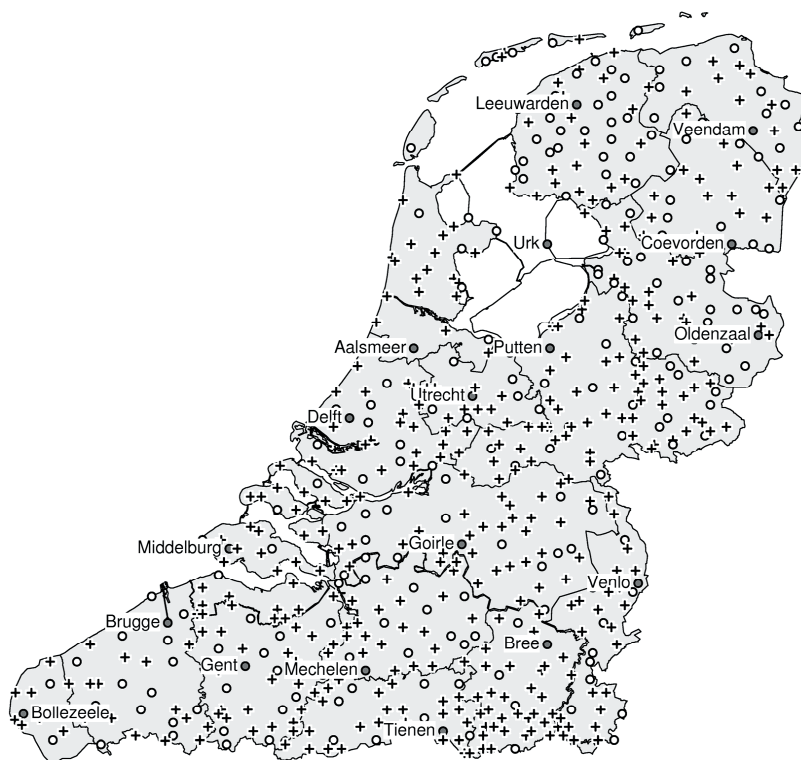
The GTRP consists of digital transcriptions for 613 dialect varieties in the Netherlands (424 varieties) and Belgium (189 varieties; see Figure 1 for the geographical distribution). All data was gathered during the period 1980 – 1995, making it the most recent broad-coverage Dutch dialect data source available. The GTRP is moreover available digitally, making it especially useful for research. For every variety, a maximum of 1876 items was narrowly transcribed according to the International Phonetic Alphabet. The items consisted of separate words and word groups, including nominals, adjectives and nouns. A more specific overview of the items is given in Taeldeman and Verleyen (1999).

The recordings and transcriptions of the GTRP were made by 25 collaborators, but more than 40% of all data was transcribed by only two individuals who created reliable transcriptions (Goeman, 1999). In most cases there were multiple transcribers operating in a single region, ranging from 1 (Drenthe) to 13 (Zuid-Holland). In general the dialectal data of one variety was based on one dialect speaker.

Our analyses are conducted on a subset of the GTRP items. Because the Levenshtein distance is used to obtain dialect distances, we only take single words into account (like Heeringa 2004). Unfortunately, word boundaries are not always clearly identified in the transcriptions (primarily for Belgian dialect varieties), making segmentation very hard. For this reason, we restrict our subset to items consisting of a single word. Because the singular nouns are (sometimes, but not always) preceded by an article (*'n*) these will not be included. The first-person plural is the only verb form not preceded by a pronoun and therefore the only verb form which is included. Finally, no items are included where multiple

lexemes are possible.

Figure 1: The geographic distribution of the 613 GTRP localities. The 224 localities



marked with a circle appear both in the GTRP and in the 360-element sample of the RND studied by Heeringa (2004). Localities marked by a '+' occur only in the GTRP. See the text for further remarks.

The GTRP was compiled with a view to documenting both phonological and morphological variation (De Schutter et al. 2005). Because our purpose here is the analysis of variation in pronunciation, we ignore many items in the GTRP whose primary purpose was presumably the documentation of morphological variation. If we had included this material directly, the measurements would have confounded pronunciation and morphological variation. Differently inflected forms of one word (e.g., base and comparative forms of an adjective) are very similar and therefore are not both selected in the subset to keep the distance measurement focused on pronunciation.

The following forms are included in the subset:

- The plural nouns, but not the diminutive nouns (the singular nouns are preceded by an article and therefore not included)
- The base forms of the adjectives instead of the comparative forms
- The first-person plural verbs (the transcriptions of other verb forms include pronouns and therefore not included)

The complete list of the 562 remaining items used in our analysis is displayed in Table 1.

2.2. RND

We will compare the results obtained on the basis of the GTRP with results obtained on the basis of an earlier data source, the *Reeks Nederlands(ch)e Dialectatlassen* (RND). The RND is a series of atlases covering the Dutch language area. The Dutch area comprises the Netherlands, the northern part of Belgium (Flanders), a smaller northwestern part of France and the German county Bentheim. The RND contains 1956 varieties, which can be found in 16 volumes. The first volume appeared in 1925, the last in 1982. The first recordings were made in 1922, the last ones in 1975. E. Blancquaert initiated the project. When Blancquaert passed away before all the volumes were finished, the project was finished under the direction of W. Peé. In the RND, the same 141 sentences are translated and transcribed in phonetic script for each dialect.

aarde	daken	gebruiken	juist	leren	over	schuw	treffen	wegen
aardig	damp	geel	kaas	leugens	paarden	simpel	treinen	wegen
acht	dansen	gehad	kaf	leunen	padden	slaan	trouwen	weinig
achter	darmen	geld	kalm	leven	paden	slapen	tussen	weken
adem	deeg	geloven	kalveren	lezen	Pasen	slecht	twaaft	wensen
af	denken	genoeg	kamers	licht	pekel	slijm	 twee	werken
anders	derde	geraken	kammen	liederen	pelln	slijpen	tweetde	weten
appels	deuren	gerst	kammen	liggen	peper	slim	twijfel	wieden
arm	dienen	geven	kanten	lijken	peren	sluiten	twintig	wijd
armen	diep	geweest	karren	likken	piepen	smal	uilen	wijn
auto's	dieven	gewoon	kasten	lomp	pijpen	smeden	vader	wijven
baarden	dik	gisteren	katten	lopen	planken	smelten	vallen	wild
bakken	dingen	glazen	kennen	lucht	pleinen	smeren	vals	willen
barsten	dinsdag	god	kermis	lui	ploegen (wrktg)	sneeuw	vangen	winnen
bedden	dochters	goed	kersen	luiden	potten	sneeuwen	varen	wippen
beenderen	doeken	goud	kervel	luisteren	proeven	soep	vast	wit
beginnen	doen	gouden	keuren	maandag	proper	spannen	vaten	woensdag
benen	dol	gras	kiezen	maanden	raar	sparen	vechten	wol
beren (wild)	donder	graven	kijken	maart	raden	spartelen	veel	wonen
best (bijw)	donderdag	grijs	kinderen	magen	recht	spelden	veertig	woorden
beurzen	donker	groen	klaver	mager	reddn	spelen	ver	worden
beven	doof	grof	kleden	maken	regen	sport (spel)	verf	wrijven
bezems	doeien	groot	klederen	marmer	rekken	spreken	vers	zacht
bezig	door	haast	klein	maten	ribben	springen	vesten	zakken
bidden	dopen	haastig	kloppen	mazelen	riet	spuiten	vet	zand
bier	dorsen	haken	kloppen	meer	rijden	staan	veulens	zaterdag
bij (vz)	dorst	halen	knechten	mei	rijk	stallen	vier	zee
bijen	draaien	half	kneden	meid	rijp	stampen	vieren	zeep
bijten	draden	handen	knien	melk	rijst	steken	vijf	zeggen
binden	dragen	hanen	koeien	menen	ringen	stelen	vijftig	zeilen
bitter	dreigen	hangen	koel	merg	roepen	stenen	vijgen	zeker
bladen	drie	hard	koken	metselen	roeren	sterven	vinden	zelf
bladeren	drinken	haver	komen	meubels	rogge	stijf	vingers	zes
blauw	dromen	hebben	kommen	missen	rokken	stil	vissen	zetten
blazen	droog	heel	konijnen	modder	rond	stoelen	vlaggen	zeven
bleek	dubbel	heet	koorts	moe	rondes	stof (huisvuil)	vlas	zeventig
blijven	duiven	heffen	kopen	moes	rood	stokken	vlees	ziek
blind	duizend	heilig	koper	moeten	rook	stom	vliegen	ziektes
bloeden	dun	helpen	kort	mogelijk	ruiken	stout	vloeken	zien
bloeien	durven	hemden	koud	mogen	runderen	straten	vlooiën	zijn

blond	duur	hemel	kousen	morgen (demain)	ruzies	strepn	voegen	zilveren
blozen	duwen	hengsten	kraken	mossels	sap	strooien	voelen	zitten
bokken	dweilen	heren	kramp	muizen	saus	sturen (zenden)	voeten	zoeken
bomen	echt	heten	kreupel	muren	schade	suiker	vogels	zoet
bonen	eeuwen	hier	krijgen	naalden	schapen	taai	vol	zondag
boren	eieren	hoeden	krimpen	nat	schaven	taarten	volgen	zonder
boter	eigen	hoesten	krom	negen	scheef	tafels	volk	zonen
bouwen	einde	hol	kruipen	negers	scheel	takken	voor	zorgen
boven	elf	holen	kwaad	nieuw	scheiden	tam	vragen	zout
braaf	engelen	honden	laag	noemen	schepen	tanden	vreemd	zouten
braden	enkel	honger	laat	nog	scheppen	tangen	vriezen	zuchten
branden	eten	hoog	lachen	noorden	scheren	tantes	vrij	zuigen
breed	ezels	hooi	lam	noten	scherp	tarwe	vrijdag	zuur
breien	fel	hoop (espoir)	lammeren	nu	schieten	tegen	vrijen	zwaar
breken	fijn	hopen	lampen	ogen	schimmel	tellen	vroeg	zwart
brengen	flauw	horen	lang	om	schoenen	temmen	vuil	zwellen
broden	flessen	horens	lastig	ons	scholen	tenen	vuur	zwem- men
broeken	fruit	houden	laten	oogst	schoon	tien	wachten	zwijgen
broers	gaan	huizen	latten	ook	schrijven	timmeren	wafels	
bruin	gaarne	jagen	leden	oosten	schudden	torens	warm	breder
buigen	gal	jeuken	ledig	op	schuiven	traag	wassen	
buiten	ganzen	jong	leem	open	schuld	tralies	weer	
dagen	gapen	jongen	leggen	oud	schuren	trams	weg	

Table 1: List of all 562 words in the GTRP subset. The 59 words in boldface are used for RND – GTRP comparison (see Section 4.3). The word *breder* is included in the set used for comparison with the RND, but not in the base subset of 562 words (due to the presence of *breed*).

The recordings and transcriptions of the RND were made by 16 collaborators, who mostly restricted their activities to a single region (Heeringa, 2004). For every variety, material was gathered from multiple dialect speakers.

In 2001 the RND material was digitized in part. Since digitizing the phonetic texts is time-consuming, a selection of 360 dialects was made and for each dialect the same 125 words were selected from the text. The words represent (nearly) all the vowels (monophthongs and diphthongs) and consonants. Heeringa (2001) and Heeringa (2004) describe the selection of dialects and words in more detail and discuss how differences introduced by different transcribers are processed.

Our set of 360 RND varieties and the set of 613 GTRP varieties have 224 varieties in common. Their distribution is shown in Figure 1. The 125 RND words and the set of 562 GTRP words share 58 words. We added one extra word, *breder* ‘wider’, which was excluded from the set of 562 GTRP words since we used no more than one morphologic variant per item and the word *breed* ‘wide’ was already included. So in total we have 59 words, which are listed in boldface in Table 1. The comparisons between the RND and GTRP in this paper are based only on the 224 common varieties and the 59 common words.

3. Measuring linguistic distances

In 1995 Kessler introduced the Levenshtein distance as a tool for measuring linguistic distances between language varieties. The Levenshtein distance is a string edit distance measure, and Kessler applied this algorithm to the comparison of Irish dialects. Later the same technique was successfully applied to Dutch (Nerbonne et al. 1996; Heeringa 2004: 213–278), Sardinian (Bolognesi & Heeringa 2002), Norwegian (Gooskens & Heeringa 2004) and German (Nerbonne & Siedle 2005).

In this paper we use the Levenshtein distance for the measurement of pronunciation distances. Pronunciation variation includes phonetic and morphologic variation, and excludes lexical variation. Below, we give a brief explanation of the methodology. For a more extensive explanation see Heeringa (2004: 121–135).

The Levenshtein algorithm provides a rough, but completely consistent measure of pronunciation distance. Its strength lies in the fact that it can be implemented on the computer, so that large amounts of dialect material can be compared and analyzed. The usage of this computational technique enables dialectology to be based on the aggregated comparisons of millions of pairs of phonetic segments.

3.1. Levenshtein algorithm

Using the Levenshtein distance, two varieties are compared by comparing the pronunciation of words in the first variety with the pronunciation of the same words in the second. We determine how one pronunciation might be transformed into the other by inserting, deleting or substituting sounds. Weights are assigned to these three operations. In the simplest form of the algorithm, all operations have the same cost, e.g., 1. Assume *melk* ‘milk’ is pronounced as [mœlkə] in

the dialect of Veenwouden (Friesland), and as [mɛlək] in the dialect of Delft (Zuid-Holland). Changing one pronunciation into the other can be done as follows (ignoring suprasegmentals and diacritics):

mɔəlɔə	delete ə	1
mɔlkə	subst. ɔ/ɛ	1
mɛlkə	delete ə	1
mɛlk	insert ə	1
mɛlək		
		4

In fact many sequence operations map [mɔəlɔə] to [mɛlək]. The power of the Levenshtein algorithm is that it always finds the cost of the cheapest mapping.

To deal with syllabicity, the Levenshtein algorithm is adapted so that only vowels may match with vowels, and consonants with consonants, with several special exceptions: [j] and [w] may match with vowels, [i] and [u] with consonants, and central vowels (in our research only the schwa) with sonorants. So the [i], [u], [j] and [w] align with anything, the [ə] with syllabic (sonorant) consonants, but otherwise vowels align with vowels and consonants with consonants. In this way unlikely matches (e.g., a [p] with an [a]) are prevented. In our example we thus have the following alignment:

m	ɔ	ə	l		k	ə
m	ɛ		l	ə	k	
	1	1		1		1

In earlier work we divided the sum of the operations by the length of the alignment. This normalizes scores so that longer words do not count more heavily than shorter ones, reflecting the status of words as linguistic units. However, Heeringa et al. (2006) showed that results based on raw Levenshtein distances approximate dialect differences as perceived by the dialect speakers better than results based on normalized Levenshtein distances. Therefore we do not normalize the Levenshtein distances in this paper but use the raw distances, i.e. distances which give us the sum of the operations needed to transform one pronunciation into another, with no transformation for length.

3.2. Operation weights

The example above is based on a notion of phonetic distance in which phonetic overlap is binary: non-identical phones contribute to phonetic distance, identical ones do not. Thus the pair [i,ɒ] counts as different to the same degree as [i,i]. In earlier work we experimented with more sensitive versions in which phones are compared on the basis of their feature values or acoustic representations. In that way the pair [i,ɒ] counts as more different than [i,i].

In a validation study Heeringa (2004) compared results of binary, feature-based and acoustic-based versions to the results of a perception experiment carried out by Charlotte Gooskens. In this experiment dialect differences as perceived by Norwegian dialect speakers were measured. It was found that generally speaking the binary versions approximate perceptual distances better than the feature-based and acoustic-based versions. The fact that segments differ appears to be more important in the perception of speakers than the degree to which segments differ. Therefore we will use the binary version of Levenshtein distance in this article, as illustrated in the example in Section 3.1. All substitutions, insertions and deletions have the same weight, in our example the value 1.

3.3. Diacritics

We do not process suprasegmentals and diacritics. Differences between the way in which transcribers transcribe pronunciations are found especially frequently in the use of suprasegmentals and diacritics (Goeman 1999). The RND transcribers, instructed by (or in the line of) Blancquaert, may have used them differently from the GTRP transcribers. To make the comparison between RND and GTRP results as fair as possible, we restrict our analyses to the basic phonetic segments and ignore suprasegmentals and diacritics.

3.4. Dialect distances

When comparing two varieties on the basis of n_w words, we analyze n_w word pairs and get n_w Levenshtein distances. The dialect distance is equal to the sum of n_w Levenshtein distances divided by n_w . When comparing n_d varieties, the average Levenshtein distances are calculated between each pair of varieties and arranged in a matrix which has n_d rows and n_d columns.

To measure the consistency (or reliability) of our data, we use Cronbach's α (Cronbach 1951). On the basis of variation of one single word (or item)

we create a $n_d \times n_d$ distance matrix. With n_w words, we obtain n_w distance matrices, for each word one matrix. Cronbach's α is a function of the number of linguistic variables and the average inter-item correlation among the variables. In our case it is a function of the number of words n_w and the average inter-word correlation among the n_w matrices. Its values range between zero and one, higher values indicating greater reliability. As a rule of thumb, values higher than 0.7 are considered sufficient to obtain consistent results in social sciences (Nunnally 1978).

4. Results

4.1. GTRP data of all varieties

To find the distance between two pronunciations of the same word, we use the Levenshtein distance. The dialect distance between two varieties is obtained by averaging the distances for all the word pairs. To measure data consistency, we calculated Cronbach's α for the obtained distance measurements. For our results, Cronbach's α is 0.99, which is much higher than the accepted threshold in social science (where $\alpha > 0.70$ is regarded as acceptable). We conclude that our distance measurements are highly consistent.

Figure 2 shows the dialect distances geographically. Varieties which are strongly related are connected by darker lines, while more distant varieties are connected by lighter lines. Even where no lines can be seen, very faint (often invisible) lines implicitly connect varieties which are very distant.

When inspecting the image, we note that the lines in Belgium are quite dark compared to the lighter lines in the Netherlands. This suggests that the varieties in Belgium are more strongly connected than those in (the north of) the Netherlands. Considering that the northern varieties in the Netherlands were found to have stronger connections than the southern varieties (Heeringa 2004: 235), this result is opposite to what was expected.



Figure 2: Average Levenshtein distance among 613 GTRP varieties. Darker lines connect close varieties, lighter lines more distant ones. We suggest that this view is confounded by differences in transcription practice. See the text for discussion, and see Figure 5 (below) for the view we defend.

We already indicated that the data of varieties in Belgium hardly contained any word boundaries (see Section 2.1), while this was not true for varieties in the Netherlands. Although unimportant for our subset containing only single word items, this could be a clue to the existence of a structural difference in transcription method between Belgium and the Netherlands.

We conducted a thorough analysis of the GTRP dialect data, which showed large national differences in the number of phonetic symbols used to transcribe the items. Table 2 indicates the number of unused phonetic symbols in both countries, four neighboring provinces and two neighboring cities. For completeness, the

number of unused tokens for all 1876 items for both countries is also included. Figure 3 gives an overview of the phonetic tokens which are not used in Belgium (for the complete set of 1876 items).

562 items			
The Netherlands	1 (1.077.169)	Belgium	33 (469.155)
Noord-Brabant	12 (130.324)	Antwerp	40 (86.257)
Limburg	15 (80.535)	Belgian Limburg	38 (110.294)
Goirle (NB)	39 (2.553)	Poppel (Ant)	49 (2.687)
1876 items			
The Netherlands	0 (4.790.266)	Belgium	27 (2.128.066)

Table 2: Total number of distinct phonetic symbols in boldface (out of 83) which do not occur in the transcriptions. The total size (number of phonetic symbol tokens) of the dialect data for each region is given between parentheses.

voice	lab	inter/lab dent	alv	palalv	alvpal	pal	vel	uv	phar	gl
-	p		t	tʰ		ɸ	k	q		ʔ plosive
+	b		d			g²ʰ	gʷ	g⁸G		plosive
-		θ								fricative
+		ð								fric
-	f		s	s²ʰ	s³ɸ	x²ɸ	x	x⁷ɸ		fric
+	v		z	z²ʰ	z³ɸ	j²ʰ	g³ʷ	g⁷ʰ		fric
-		ɸ					w³M		hʲh	'no' fric
+		β	w²U				w		h⁸ʰ	'no' fric
+			r	r⁵ʰ		r²ʰ		r⁷R		fric
+				r³ʰ				r⁹B		'no' fric
+			r⁴ʰ			j				semi-vowel
+			l			l³ʰ	l²ʰ			low fric
+	y⁴ʰ		l⁴ʰ							semi-vowel
+	m		n			n²ʰ	n, ŋ	n⁷N		nasal
+	m, ŋ									nasal

	spread front	rounded mid	spread back	rounded back
closed	i i	y y	ɨ ɨ	u u
half-closed	ɪ ɪ	ʏ ʏ	ɘ ɘ	ɤ ɤ
				ɔ ɔ
half-closed	e e	ø ø	ɛ ɛ	ɜ ɜ
half-open	ɛ ɛ	ɔ ɔ	ɶ ɶ	ɷ ɷ
open	æ æ	ɶ ɶ	ɑ ɑ	ɶ ɶ

Figure 3: All 83 Keyboard-IPA symbols used in the GTRP data (without diacritics). Symbols on a black background are not used in Belgian transcriptions. Original image: Goeman, Van Reenen & Van den Berg (Meertens Instituut).

Table 3 illustrates some transcription differences between two neighboring places near the border of Belgium and the Netherlands (see Figure 4). For this example, note that the phonetic symbols unused in Belgium include \mathfrak{p} , \mathfrak{l} , \mathfrak{z} , \mathfrak{u} and \mathfrak{m} .

Dutch	English	Goirle (NL)	Poppel (BEL)
baarden	beards	$\mathfrak{b}\mathfrak{o}\mathfrak{r}\mathfrak{d}\mathfrak{\partial}$	$\mathfrak{b}\mathfrak{o}\mathfrak{r}\mathfrak{d}\mathfrak{\partial}$
bij (vz.)	at	$\mathfrak{b}\mathfrak{e}\mathfrak{i}$	$\mathfrak{b}\mathfrak{e}\mathfrak{i}$
blond	blonde	$\mathfrak{b}\mathfrak{o}\mathfrak{n}\mathfrak{t}$	$\mathfrak{b}\mathfrak{l}\mathfrak{o}\mathfrak{n}\mathfrak{t}$
broeken	pants	$\mathfrak{b}\mathfrak{r}\mathfrak{u}\mathfrak{k}\mathfrak{\partial}$	$\mathfrak{b}\mathfrak{r}\mathfrak{u}\mathfrak{k}\mathfrak{\partial}\mathfrak{n}$
donker	dark	$\mathfrak{d}\mathfrak{o}\mathfrak{ŋ}\mathfrak{k}\mathfrak{\partial}$	$\mathfrak{d}\mathfrak{o}\mathfrak{ŋ}\mathfrak{k}\mathfrak{\partial}$
hard	hard	$\mathfrak{h}\mathfrak{a}\mathfrak{r}\mathfrak{t}$	$\mathfrak{h}\mathfrak{a}\mathfrak{r}\mathfrak{t}$
haver	oats	$\mathfrak{h}\mathfrak{o}\mathfrak{v}\mathfrak{\partial}\mathfrak{R}$	$\mathfrak{h}\mathfrak{o}\mathfrak{v}\mathfrak{\partial}\mathfrak{R}$
kamers	rooms	$\mathfrak{k}\mathfrak{o}\mathfrak{m}\mathfrak{\partial}\mathfrak{\partial}\mathfrak{s}$	$\mathfrak{k}\mathfrak{o}\mathfrak{m}\mathfrak{\partial}\mathfrak{\partial}\mathfrak{s}$
kinderen	children	$\mathfrak{k}\mathfrak{e}\mathfrak{n}\mathfrak{d}\mathfrak{\partial}\mathfrak{R}$	$\mathfrak{k}\mathfrak{e}\mathfrak{n}\mathfrak{d}\mathfrak{\partial}\mathfrak{R}$
kloppen	knock	$\mathfrak{k}\mathfrak{t}\mathfrak{o}\mathfrak{p}\mathfrak{\partial}$	$\mathfrak{k}\mathfrak{l}\mathfrak{o}\mathfrak{p}\mathfrak{\partial}$
luisteren	listen	$\mathfrak{l}\mathfrak{a}\mathfrak{s}\mathfrak{t}\mathfrak{\partial}\mathfrak{\partial}\mathfrak{\partial}$	$\mathfrak{l}\mathfrak{\partial}\mathfrak{s}\mathfrak{t}\mathfrak{\partial}\mathfrak{\partial}\mathfrak{\partial}$
missen	miss	$\mathfrak{m}\mathfrak{i}\mathfrak{s}\mathfrak{\partial}$	$\mathfrak{m}\mathfrak{i}\mathfrak{s}\mathfrak{e}$
simpel	simple	$\mathfrak{s}\mathfrak{i}\mathfrak{m}\mathfrak{p}\mathfrak{\partial}\mathfrak{t}$	$\mathfrak{s}\mathfrak{e}\mathfrak{m}\mathfrak{p}\mathfrak{\partial}\mathfrak{l}$
sneeuw	snow	$\mathfrak{s}\mathfrak{n}\mathfrak{o}\mathfrak{u}\mathfrak{m}$	$\mathfrak{s}\mathfrak{n}\mathfrak{e}\mathfrak{a}\mathfrak{w}$
tralies	bars	$\mathfrak{t}\mathfrak{r}\mathfrak{o}\mathfrak{l}\mathfrak{i}\mathfrak{s}$	$\mathfrak{t}\mathfrak{r}\mathfrak{o}\mathfrak{l}\mathfrak{i}\mathfrak{s}$
twalf	twelve	$\mathfrak{t}\mathfrak{w}\mathfrak{\partial}\mathfrak{l}\mathfrak{\partial}\mathfrak{f}$	$\mathfrak{t}\mathfrak{w}\mathfrak{o}\mathfrak{l}\mathfrak{\partial}\mathfrak{f}$
vogels	birds	$\mathfrak{v}\mathfrak{o}\mathfrak{v}\mathfrak{\partial}\mathfrak{s}$	$\mathfrak{v}\mathfrak{o}\mathfrak{u}\mathfrak{v}\mathfrak{\partial}\mathfrak{s}$
vriezen	freeze	$\mathfrak{v}\mathfrak{r}\mathfrak{i}\mathfrak{z}\mathfrak{\partial}$	$\mathfrak{v}\mathfrak{r}\mathfrak{i}\mathfrak{z}\mathfrak{\partial}\mathfrak{n}$
woensdag	Wednesday	$\mathfrak{w}\mathfrak{u}\mathfrak{n}\mathfrak{s}\mathfrak{d}\mathfrak{\partial}\mathfrak{x}$	$\mathfrak{w}\mathfrak{u}\mathfrak{n}\mathfrak{z}\mathfrak{d}\mathfrak{\partial}\mathfrak{x}$
zeggen	say	$\mathfrak{z}\mathfrak{e}\mathfrak{z}\mathfrak{\partial}$	$\mathfrak{z}\mathfrak{e}\mathfrak{z}\mathfrak{\partial}\mathfrak{n}$

Table 3: Phonetic transcriptions of Goirle (NL) and Poppel (BEL) including Dutch and English translations. Even though phonetic transcriptions are of comparable length and complexity, the Dutch sites vary consistently use a much wider range of phonetic symbols, confounding measurement of pronunciation distance.

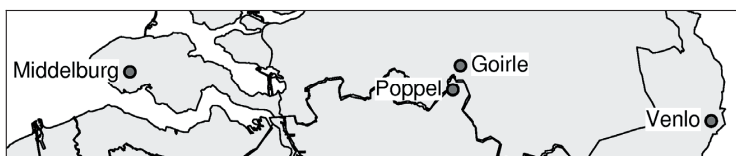


Figure 4: Relative locations of Poppel (Belgium) and Goirle (the Netherlands).

Transcriptions using fewer phonetic symbols are likely to be measured as more similar due to a lower degree of possible variation. Figure 2 shows exactly this result. Because of these substantial transcriptional differences between the two countries (see also Van Oostendorp, to appear; Hinskens & Van Oostendorp 2006) it is inappropriate to compare the pronunciations of the two countries directly. Therefore, in what follows, we analyze the transcriptions of the two countries separately, and also discuss their pronunciation differences separately.

4.2. GTRP data, the Netherlands and Belgium separately

The data was highly consistent even when regarding the countries individually. Cronbach's α was 0.990 for dialect distances in the Netherlands and 0.994 for dialect distances in Belgium.

In Figure 5, the strong connections among the Frisian varieties and among the Groningen and Drenthe (Low Saxon) varieties are clearly shown. The dialect of Gelderland and western Overijssel can also be identified below the dialect of Drenthe. South of this group a clear boundary can be identified, known as the boundary between Low Saxon (northeastern dialects) and Low Franconian (western, southwestern and southern dialects). The rest of the map shows other less closely unified groups, for example, in Zuid-Holland and Noord-Brabant as well as less cohesive groups in Limburg and Zeeland.

Just as was evident in Figure 2, Belgian varieties are tightly connected in both the varieties of Antwerp as well as in those of West Flanders (see Figure 5). A lot of white lines are present in Belgian Limburg, however, indicating more dissimilar varieties in that region. Note the weak lines connecting to the Ghent variety (indicating it to be very different from the neighboring varieties); they appear to be masked by lines of closer varieties in the surrounding area.

By using multidimensional scaling (MDS; see Heeringa 2004: 156–163) varieties can be positioned in a three-dimensional space. The more similar two varieties

are, the closer they will be placed together. The location in the three-dimensional space (in x-, y- and z-coordinates) can be converted to a distinct color using red, green and blue color components. By assigning each variety its own color in the geographical map, an overview is obtained of the distances between the varieties. Similar varieties have the same color, while the color differs for more distant varieties. This method is superior to a cluster map (e.g., Heeringa, 2004: 231) because MDS coordinates are assigned to individual collection sites, which means that deviant sites become obvious, while clustering reduces each site to one of a fixed number of groups. Hence, clustering risks covering up problems.⁽¹⁾



Figure 5: Average Levenshtein distance between 613 GTRP varieties. Darker lines connect close varieties, lighter lines more distant ones. The maps of the Netherlands (top) and Belgium (bottom) must be considered independently.

⁽¹⁾ We discuss apparently exceptional sites at the end of this section, and we note here that these exceptions are indeed obvious in clustering as well.

Because we are reducing the number of dimensions in the data (i.e. the dialect differences) to three by using the MDS technique, it is likely that some detail will be lost. To get an indication of the loss of detail, we calculate how much variance of the original data is explained by the three-dimensional MDS output. For the Netherlands, the MDS output explains 87.5% of the variance of the original dialect differences. For Belgium a comparable value is obtained: 88.1%. We therefore conclude that our MDS output gives a representative overview of the original dialect differences in both countries.

In Figure 6 (see page 113) and 7 (see page 114) the MDS color maps of the Netherlands and Belgium are shown. The color of intermediate points is determined by interpolation using Inverse Distance Weighting (see Heeringa 2004: 156–163). Because the dialect data for Belgium and the Netherlands was separated, the maps should be considered independently. Varieties with a certain color in Belgium are not in any way related to varieties in the Netherlands having the same color. Different colors only identify distant varieties within a country.

To help interpret the color maps, we calculated all dialect distances on the basis of the pronunciations of every single word in our GTRP subset. By correlating these distances with the distances of every MDS dimension, we were able to identify the words which correlated most strongly with the distances of the separate MDS dimensions.

For the Netherlands we found that the dialect distances on the basis of the first MDS dimension (separating Low Saxon from the rest of the Netherlands) correlated most strongly ($r = 0.66$) with distances obtained on the basis of the pronunciation of the word *moeten* ‘must’. For the second MDS dimension (separating the north of the Netherlands, most notably Friesland, from the rest of the Netherlands) the word *donderdag* ‘Thursday’ showed the highest correlation ($r = 0.59$). The word *schepen* ‘ships’ correlated most strongly ($r = 0.49$) with the third MDS dimension (primarily separating Limburg from the rest of the Netherlands). For Belgium we found that the dialect distances obtained on the basis of the pronunciation of the word *wol* ‘wool’ correlated most strongly ($r = 0.82$) with the first MDS dimension (separating eastern and western Belgium). The word *schrijven* ‘write’ correlated most strongly ($r = 0.63$) with the second MDS dimension (separating the middle part of Belgium from the peripheral eastern and western parts), while the word *vrijdag* ‘Friday’ showed the highest correlation with the third MDS dimension (primarily separating Ghent and

eastern Belgium from the rest). Figure 6 and 7 also display these words and corresponding pronunciations in every region.

On the map of the Netherlands, varieties of the Frisian language can clearly be distinguished by the blue color. The town Frisian varieties are purpler than the rest of the Frisian varieties. This can be seen clearly in the circle representing the Leeuwarden variety. The Low Saxon area can be identified by a greenish color. Note that the dialect of Twente (near Oldenzaal) is distinguished from the rest of Overijssel by a less bluish green color. The Low Franconian dialects of the Netherlands can be identified by their reddish tints. Due to its bright red color, the dialect of Limburg can be identified within the Low Franconian dialects of the Netherlands.

For the Belgian varieties, the dialects of West Flanders (green) and Brabant (blue) can be clearly distinguished. In between, the dialects of East Flanders (light blue) and Limburg (red) can also be identified. Finally, the distinction between Ghent (pink) and its surrounding varieties (greenish) can be seen clearly.

Apparent dialect islands

A careful examination of Figure 6 reveals a few sites whose MDS dimensions (and therefore colors) deviate a great deal from their surroundings. For example, there are two bright points around Twente (above the Oldenzaal label) which might appear to be dialect islands. Upon inspection it turns out that these points both used transcriptions by the same fieldworker, who, moreover, contributed almost only those (four) sets of transcriptions to the entire database. We therefore strongly suspect that the apparent islands in Twente are “transcriber isoglosses”. Also reported Hinskens & Van Oostendorp (2006) the existence of transcriber effects in the GTRP data.

But these are not the only apparent dialect islands. What can we do about this? Unfortunately, there are no general and automated means of correcting deviant transcriptions or correcting analyses based on them. At very abstract levels we can correct mathematically for differences in a very small number of transcribers (or fieldworkers), but we know of no techniques that would apply in general to the GTRP data. It is possible to compare analyses which exclude suspect data to analyses which include it, but we should prefer not to identify suspect data only via its deviance with respect to its neighbors.

4.3. GTRP compared to RND

Our purpose in the present section is to examine the GTRP against the background of the RND in order to detect whether there have been changes in the Dutch dialect landscape. We employ a regression analysis (below) to detect areas of relative convergence and divergence. The regression analysis identifies an overall tendency between the RND and GTRP distances, against which convergence and divergence may be identified: divergent sites are those for which the actual difference between the RND and GTRP distances exceeds the general tendency, and convergent sites are those with distances less than the tendency.

We are not analyzing the rate of the changes we detect. Given the large time span over which the RND was collected, it would be illegitimate to interpret the results of this section as indicative of the rate of pronunciation change. This should be clear when one reflects first, that we are comparing both the RND and the GTRP data at the times at which they were recorded, and second, that the RND data was recorded over a period of fifty years. One could analyze the rate of change if one included the time of recording in the analysis, but we have not done that.

We verify first that the regression analysis may be applied, starting with the issue of whether there is ample material for comparison.

In Section 2.2 we mentioned that the comparisons between the RND and GTRP in this paper are based only on the 224 common varieties and the 59 common words. Although one might find this number of words quite small, we still obtained consistent results. When we use the RND data, Cronbach's α is 0.95 for the data from the Netherlands and 0.91 for the data from Belgium. For the GTRP data we found Cronbach's α values of 0.91 and 0.95 respectively.

We correlated the RND distances with the GTRP distances and found a correlation of $r = 0.83$ for the Netherlandic distances, and a correlation of $r = 0.82$ for the Belgian distances. These correlations are significant ($p < 0.001$) according to the Mantel test, a test which takes into account the fact that the distances within a distance matrix are not fully independent (see Heeringa 2004: 74–75 for a brief explanation of this test). The correlations indicate a strong, but not perfect relationship between the old RND dialect distances and the newer GTRP dialect distances. In the sections below we will examine these differences.

4.3.1. Comparison of transcriptions and distances

In Section 3 we described how we have measured pronunciation distances. The RND and the GTRP distances are measured in the same way, but the measurements are based on different kinds of transcriptions. As shown in Section 4.1, these differences may be reflected in the number of different phonetic symbols used in the transcriptions. Therefore we counted the number of different speech segments in the set of common varieties and common words for both the RND and the GTRP. Ignoring suprasegmentals and diacritics we found the following results:

	RND original	RND modified	GTRP
Netherlands	43	40	73
Belgium	42	40	44

In the column ‘RND original’ counts are given on the basis of the original, unchanged transcriptions. When calculating Levenshtein distances, we used a modified version of the transcriptions in which some of the different notations used by different transcribers, are normalized (see Heeringa 2004: 217–226). Counts on the basis of these modified transcriptions are given in the column ‘RND modified’.

If we wished to compare pronunciation directly between the RND and the GTRP, it would be important to verify that measurements were calibrated, i.e. that they were using the same scale. The table above shows that the number of different segments is about the same in all cases, except for the Netherlandic GTRP data which has a much higher number of different tokens (73). We now examine whether the number of different tokens influences our Levenshtein distance measurements. For both countries within each data source we calculated the mean and the standard deviation of the average Levenshtein distances of all pairs of varieties. Remember that each dialect distance represents the average number of operations needed to transform one pronunciation into another. We found the following results:

	RND mean	GTRP mean	RND st. dev.	GTRP st. dev.
Netherlands	1.58	2.03	0.51	0.44
Belgium	1.47	1.64	0.36	0.52

When comparing the means with the corresponding number of different tokens in the table above, we find the expected tendency that a lower number of distinctive tokens corresponds to lower distances. We do not find a clear relationship between the standard deviations and the number of different tokens.

We compared the RND dialect distances to corresponding GTRP dialect distances by means of a matched-pairs *t*-test. It turns out that GTRP distances are significantly higher than the RND distances ($p < 0.001$ for both the Netherlands and Belgium). We emphasize that we do not interpret this as evidence that the Dutch dialects are diverging from one another in general for reasons we turn to immediately.

The differences in the number of different tokens on the one hand, and the differences in distances on the other hand, show that the results of the GTRP cannot be directly compared to the results of the RND. We will therefore use regression analysis to compare the results of the two different data sources.

4.3.2. Linear regression analysis

The idea behind regression analysis is that a dependent variable can be predicted by an independent variable. A linear regression procedure finds a formula which defines a linear relationship between the independent variable and the dependent variable. Because the relationship will usually not be perfectly linear, the values predicted by the formula on the basis of the independent variable will differ from the values of the dependent variable. The differences between the predicted values and the real observations of the dependent variable are called residues.

Since the past may influence the present but not *vice versa*, we regard the RND distances as the independent variable, and the GTRP distances as dependent. With regression analysis we obtain differences between the predicted GTRP distances and the real GTRP distances, i.e. the residues. These residues can be either positive or negative. A positive residue means that the real GTRP distance is larger than the GTRP distance predicted on the basis of the corresponding

RND distance. A negative residue means the real GTRP distance is lower than expected on the basis of the corresponding RND distance.

As mentioned above, we cannot directly compare GTRP distances with RND distances. This means that we cannot determine whether varieties have converged or diverged in absolute terms. But residues tell us whether and to what degree some varieties have become relatively closer, and others relatively more distant. ‘Relatively’ means: in relation to distances of the other dialect pairs. We will refer to this as ‘relative convergence’ and ‘relative divergence’.

For instance, assume that variety A converged to variety B, variety C converged to variety D, and variety E converged to variety F. The varieties A and B converged more strongly than varieties C and D, and varieties E and F converged less strongly than varieties C and D. We are not able to detect the overall pattern of convergence, but we are able to detect that the relationships among the dialect pairs have changed with respect to their relative distances. Ignoring the overall pattern, we would find that varieties A and B have relatively converged, and varieties E and F have relatively diverged.

Figure 8 shows the residues (see page 115). Varieties which have relatively converged are connected by blue lines, and varieties which have relatively diverged are connected by red lines. When we consider the Netherlands, we find that the Frisian dialects in the northwest, and the dialects in the eastern part of the province of Noord-Brabant (north of Goirle) and those in the province of Limburg (north and south of Venlo) have converged relative to one another.

The Frisian dialects are known to be very homogeneous. Therefore it is striking that the dialects became—relatively—even more similar to each other. The Frisian dialects have not converged toward the surrounding dialects, for example toward the Noord-Holland dialects, which are relatively close to standard Dutch. The internal convergence could be the result of the influence of standard Frisian in which case these dialects have become more standardized, i.e. closer to standard Frisian.

In contrast, the Limburg dialects are known to be very heterogeneous and relatively distant from standard Dutch. The strong relative convergence of Limburg and eastern Noord-Brabant dialects may be explained by convergence toward

standard Dutch, which makes them more similar to each other and to some surrounding dialects which are relatively similar to standard Dutch. This idea is supported by a slight relative convergence toward dialects north of Brabant, in the south of the province of Gelderland.

Strong relative divergence is found among the Twente varieties, the area including and west of Oldenzaal. We have no good dialectological explanation for this. However, there were a large number of transcribers (6) in this small region and it could be that the divergence is caused by transcriber problems (e.g., see Section 4.2).

When examining Flanders in Figure 8, we find relative convergence in most provinces, probably again as the result of convergence towards standard Dutch. One clear exception is the variety of Ghent. Phonologically the variety of Ghent differs strongly from the surrounding varieties. For instance, all vowels in the variety of Ghent are longer than in the surrounding varieties. Since the data of Ghent was gathered and transcribed by the same field worker who collected and transcribed the data of other varieties in East and West Flanders, we would conjecture that the variety of Ghent has been influenced much less by standard (Flemish) Dutch, making the contrast to the surrounding dialects larger.

4.3.3. Principal component analysis

Principal component analysis (PCA) is a technique used for reducing multiple dimensions of a dataset to a lower number of dimensions. Dimensions which show similar patterns across the items, thus having high correlations, are fused to a single component. The PCA procedure transforms the data to a lower number of components so that the greatest variance is placed on the first principal component, the second greatest variance on the second component, and so on. The number of dimensions is reduced so that characteristics of the dataset that contribute most to its variance are retained (Tabachnik and Fidell 2001: Chapter 13).

With linear regression analysis we obtained a residue for each pair of varieties. When we have n_d varieties, each variety has a residue with respect to each of the other $n_d - 1$ varieties and to itself (which is always 0). In this way a variety is defined as a range of n_d values, i.e. there are n_d dimensions. Each dimension shows a pattern of relative convergence and divergence among the varieties.

Because we have 164 varieties in the Netherlands, they are represented by 164 dimensions. The SPSS PCA procedure reduces the number of dimensions to 14 components. The first component represents 34.9% of the variance in the original data, the second component represents 13.6%, the third component 11.5%, etc. The 60 Belgian varieties represent 60 dimensions. With PCA the number of dimensions is reduced to 7 components. The first component represents 41.8% of the variance in the residues, the second one represents 22.5%, the third one 8%, etc. As we mentioned above, the greatest variance is placed on the first component. In Figure 9 the values of the first principal component are geographically visualized (see page 116). Higher values are represented by lighter greytone.

Considering the Netherlands, we find a sharp distinction between the Frisian area which is nearly white, and the rest of the Netherlands which is colored more darkly. White colors signify dialects which behave similar to Frisian, and in this case, this is only Frisian. White thus means that varieties have a strong relative convergence towards Frisian. Black represents varieties without any pattern which converge or diverge to all other varieties to the same degree. So the main finding is that Frisian dialects converged with respect to each other, but not with respect to other dialects. Especially striking is the dark area found between Oldenzaal and Putten. This area is geographically close to the border between the northeastern Low Saxon dialects and the southern Low Franconian dialects. In our analysis these varieties do not converge or diverge more strongly with respect to some dialects as compared to others, but we hasten to add the dark area seen there may also be influenced by the transcriber problems noted in Section 4.2.

When looking at Flanders, we see a clear east-west division. The east is colored nearly white, especially the province of Antwerp (north of Mechelen). The western part is colored more darkly. White means that varieties have a strong relative convergence to dialects in the east (Brabant, Antwerp, and Limburg). Dark represents varieties that strongly converged toward dialects in the west (French Flanders and West Flanders). So the main pattern is that western varieties and eastern varieties both converge internally, even while they do not converge toward each other.

5. General discussion

In this paper we have provided an aggregate analysis of the pronunciation in contemporary Dutch dialects as these are sampled in the GTRP. The sheer scale of the GTRP guarantees the basis for a reliable analysis, which in turn demonstrates that the Dutch-speaking landscape is still richly contoured with Friesland, Limburg and Low Saxony as the most distinct areas.

In order to protect our analysis from potential, perhaps subtle differences in measurement scale due to transcription differences between the RND and the GTRP, we used the residues of a regression analysis in order to identify the most dynamic areas of convergence and divergence. The comparison between the situation in roughly the mid-twentieth century (as documented in the RND) and the current situation (as documented by the GTRP) revealed that Friesland, Flemish Brabant, West Flanders, and Limburg are areas of dynamic convergence, while Ghent and the southeastern part of Low Saxony are areas of divergence. We also qualified this general characterization, noting that the RND was collected over a fifty year period, which prevents us from drawing conclusions with respect to the rate of pronunciation change.

We extracted the first principal component from the residues of the regression analysis, which revealed that Friesland and eastern Flanders are behaving coherently. We would like to emphasize that the use of regression analysis, including the application of PCA to its residues, is an innovation in dialectometric technique.

In addition, we examined an apparent discrepancy in the degree of phonetic discrimination provided by GTRP transcriptions for the Netherlands as opposed to that provided for transcriptions for Belgium. After further examination, we concluded that the discrepancy is genuine, and that care is required in analyses involving subsamples of the GTRP involving sites in both countries. An aggregate analysis such as ours is certainly prone to confounding due to discrepancies in data sampling, recording and transcription, but let us hasten to add that single variable analyses are by no means immune to these problems.

This line of work suggests several further investigations. First, it would be interesting to attempt to interpret the second and third principal components of the relative changes, an undertaking which would require more space than we have at our disposal here. Second, we are interested in potential means of

correcting for the sorts of transcription differences noted. Are there automatic means of “reducing” the more detailed transcriptions to less detailed ones? Or must we settle for purely numeric corrections, which would mean that we have little to no opportunity to interpret the “corrections” linguistically? A project which would interest us, but which could only be undertaken in collaboration with the “stewards” of the GTRP, would be to map the more complex Dutch transcription system onto the simpler Flemish one. This could, of course, turn out to involve too many decisions about individual sounds to be feasible, but it could also turn out to be straightforward.

Third, in discussing the Netherlandic part of the GTRP we noted clues that fieldworker and transcription practices may be confounding analyses to some degree (see also Hinskens & Van Oostendorp 2006). This potential confound is bothersome, and it would be rewarding to eliminate it. The most rewarding, but the most difficult strategy would be to try to analyze pronunciation difference purely acoustically, eliminating the need for transcriptions. Perhaps more realistic would be to develop strategies to identify clues that transcriptions are being produced differently and also to quantify the degree to which different transcription might distort measurements. But even in the absence of general techniques, it would be useful to know where transcriber differences may exist in the GTRP.

A more exciting, and more promising opportunity suggests itself in the rich sample of morphological variation represented in the GTRP, which, after all, is the basis of the *Morfologische Atlas van de Nederlandse Dialecten* (MAND; De Schutter et al. 2005). Although Seguy (1973) and Goebel (1984) include morphological variables in their dialectometric work, the morphological material is analyzed at a categorical level, i.e. in which only “same” and “different” are distinguished. The development of a measure of morphological distance reflecting not only the potentially differing exponence of common morphological categories (which after all are already reflected in pronunciation difference), but also reflecting the effect of non-coincidental categories (such as the second Frisian infinitive), would be a rewarding challenge.

Acknowledgements

We are grateful to Peter Kleiweg, whose L04 package was used extensively to carry out and visualize the analyses we present, and also to an anonymous

Taal en Tongval reviewer, whose remarks on transcriber problems in Twente led to several qualifications in our interpretation of data vis-à-vis transcribers above. We thank the Meertens Instituut for making the GTRP data available for research and especially Boudewijn van den Berg for answering our questions regarding this data.

References

- BLANQUAERT, E., & PEÉ, W. (EDS.)
1925 - 1987, *Reeks Nederlans(ch)e Dialectatlassen*. De Sikkel, Antwerpen.
- BOLOGNESI, R., & HEERINGA, W.
2002, De invloed van dominante talen op het lexicon en de fonologie van Sardische dialecten. In: *Gramma/TTT: tijdschrift voor taalwetenschap*, 9(1), 45–84.
- CRONBACH, L.
1951, Coefficient alpha and the Internal Structure of Tests. *Psychometrika*, 16, 297–334.
- GOEBL, H.
1984, *Dialectometrische Studien. Anhand italo-romanischer, rätoromanischer und galloromanischer Sprachmaterialien aus AIS und ALF*, volume 191–193 of *Beihefte zur Zeitschrift für romanische Philologie*. Max Niemeyer Verlag, Tübingen.
- GOEMAN, A.
1999, *T-deletie in Nederlandse dialecten. Kwantitatieve analyse van structurele, ruimtelijke en temporele variatie*. The Hague: Holland Academic Graphics/Thesus.
- GOEMAN, A., & TAEDEMAN, J.
1996, Fonologie en morfologie van de Nederlandse dialecten. Een nieuwe materiaalverzameling en twee nieuwe atlasprojecten. In: *Taal en Tongval*, 48, 38–59.
- GOOSKENS, C., & HEERINGA, W.
2004, Perceptive Evaluation of Levenshtein Dialect Distance Measurements using Norwegian Dialect Data. In: *Language Variation and Change*, 16(3), 189–207.
- HEERINGA, W.
2001, De selectie en digitalisatie van dialecten en woorden uit de Reeks Nederlandse Dialectatlassen. In: *TABU, Bulletin voor Taalwetenschap*, 31(1/2), 61–103.
- HEERINGA, W.
2004, *Measuring Dialect Pronunciation Differences using Levenshtein Distance*. Doctoral dissertation. University of Groningen.

- HEERINGA, W., ET AL.
 2006, Evaluation of String Distance Algorithms for Dialectology. In: J. Nerbonne & E. Hinrichs (eds.), *Linguistic Distances Workshop at the joint conference of International Committee on Computational Linguistics and the Association for Computational Linguistics*, Sydney, July 2006, 51–62.
- HINSKENS, F., & VAN OOSTENDORP, M.
 2006, De palatalisering en velarisering van coronale nasaal-plosief clusters in GTR. Talige, dialectgeografische en onderzoekerseffecten. In: *Taal en Tongval*, 58, 103–122.
- KESSLER, B.
 1995, Computational Dialectology in Irish Gaelic. In: *Proceedings of the 7th Conference of the European Chapter of the Association for Computational Linguistics*, Dublin. EACL, 60–67.
- NERBONNE, J., ET AL.
 1996, Phonetic Distance between Dutch Dialects. In: G. Durieux, W. Daelemans, & S. Gillis (eds.) *CLIN VI: Proceedings of the Sixth CLIN Meeting*. Antwerp, Centre for Dutch Language and Speech (UIA), 185–202.
- NERBONNE, J., & SIEDLE, C.
 2005, Dialektklassifikation auf der Grundlage aggregierter Ausspracheunterschiede. In: *Zeitschrift für Dialektologie und Linguistik*, 72(2), 129–147.
- NUNNALLY, J.
 1978, *Psychometric Theory*. McGraw-Hill, New York.
- OOSTENDORP, M. VAN
 Kenmerkeconomie in de GTRP-database. To appear in *Taal en Tongval*.
- SCHUTTER, G. DE, ET AL.
 2005, *Morfologische atlas van de Nederlandse dialecten – deel 1*. Amsterdam University Press, Meertens Instituut – KNAW.
- SEGUY, J.
 1973, La Dialectométrie dans l'Atlas linguistique de la Gascogne. In: *Revue de linguistique romane*, 37, 1–24.
- TABACHNIK, B., & FIDELL, L.
 2001, *Using Multivariate Statistics*. Boston: Allyn & Bacon: 4th edition.
- TAELEDEMAN, J., & VERLEYEN, G.
 1999, De FAND: Een kind van zijn tijd. In: *Taal en Tongval*, 51, 217–240.

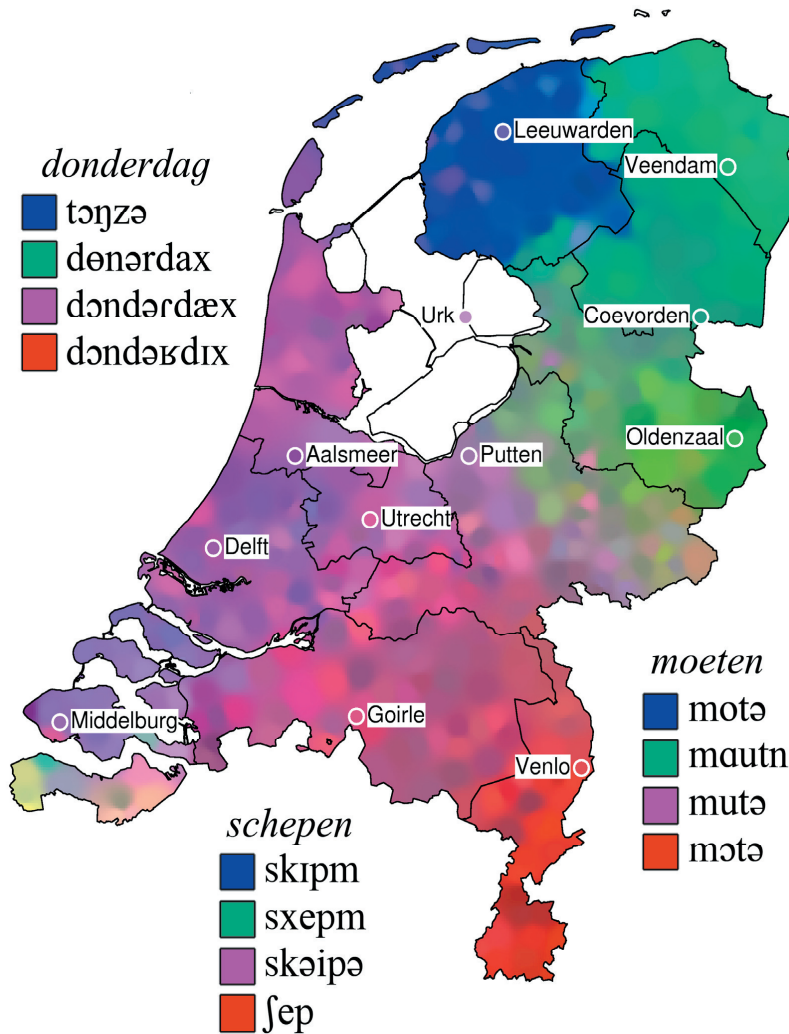


Figure 6: The GTRP data of the Netherlands reduced to its three most important dimensions via MDS (accounting for roughly 88% of dialect variation). Pronunciations of the word *moeten* ‘must’, *donderdag* ‘Thursday’, and *schepen* ‘ships’ correlate most strongly with the first, second and third MDS dimensions respectively.

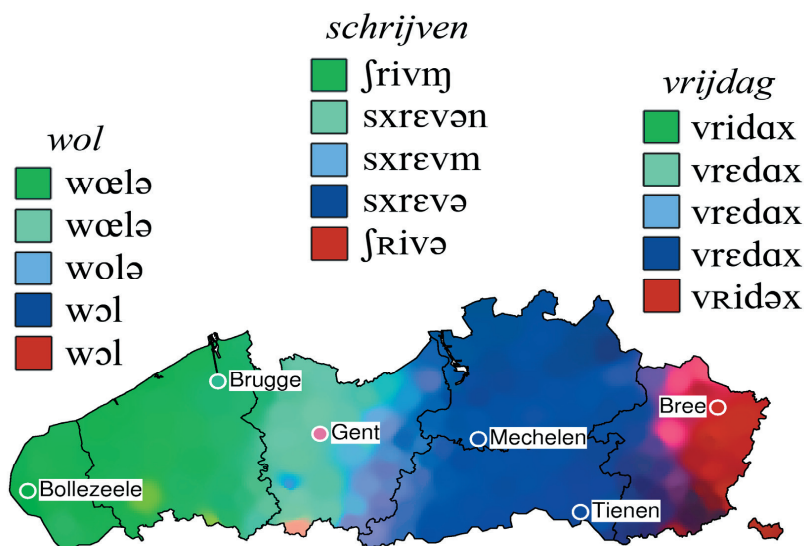


Figure 7: The GTRP data of Belgium reduced to its three most important dimensions via MDS (accounting for roughly 88% of dialect variation). Pronunciations of the word *wol* ‘wool’, *schrijven* ‘write’ and *vrijdag* ‘Friday’ correlate most strongly with the first, second and third MDS dimension respectively.

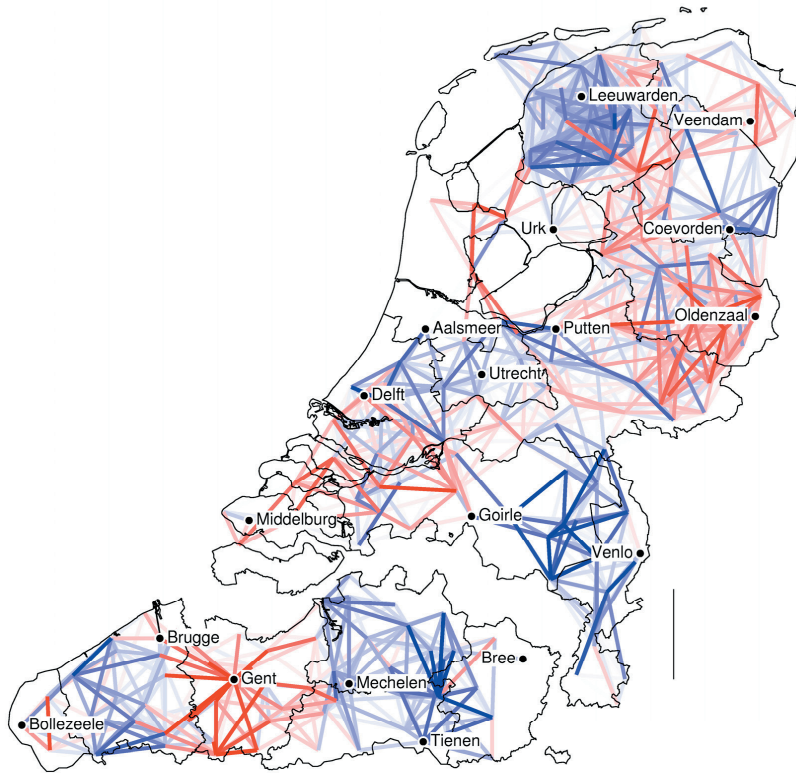


Figure 8: Relative convergence and divergence among dialects. Relative convergence means that dialects have become closer in relation to distances of the other dialect pairs and relative divergence means that dialects have become more distant in relation to distances of the other dialect pairs. The intensity of blue (red) represents the degree of relative convergence (divergence), and lines longer than the black vertical line in the lower right corner are not shown.



Figure 9: Grey tones represent values of the first component obtained with principal component analysis applied to the residues shown in Figure 8. Varieties which have the same pattern of relative convergence and divergence with respect to other varieties show similar grey tones. Thus Friesland and East Flanders house groups of dialects which have developed similarly within the two countries, and in fact, convergently. The maps of the Netherlands and Belgium should be interpreted independently from each other.