

Corpus-based dialectometry – a methodological sketch

Benedikt Szmrecsanyi
Freiburg Institute for Advanced Studies

bszm@frias.uni-freiburg.de

July 22, 2010

This paper introduces methodologies to tap corpora for exploring aggregate linguistic distances between dialects or varieties as a function of properties of geographic space. The paper describes the different steps necessary to obtain an appropriate corpus-based dataset (a so-called ‘distance matrix’), and subsequently discusses several cartographic visualization techniques – network maps, continuum maps, and cluster maps – to project aggregate linguistic relationships to geography. In addition, the paper sketches some statistical methods to quantify these relationships. By way of example, a case study draws on the *Freiburg Corpus of English Dialects*, a major dialect corpus that samples more than thirty traditional English dialects all over Great Britain. With a focus on regional variation in morphosyntax and on the basis of text frequencies of several dozen features, the study probes joint linguistic variability between the dialects sampled in the corpus.

Keywords: corpus linguistics, dialectometry, dialectology, aggregation, morphosyntax, British English dialects

Acknowledgments

I am grateful to Hans Goebel, Wilbert Heeringa, Bernd Kortmann, John Nerbonne, Christoph Wolk, and an anonymous referee for helpful comments and suggestions that have greatly improved this paper. The usual disclaimers apply.

1. Introduction: What is dialectometry?

DIALECTOMETRY is the branch of geolinguistics concerned with measuring, visualizing, and analyzing aggregate dialect similarities or distances as a function of properties of geographic space; for seminal work, see Séguy (1971) (the paper that sparked the dialectometry enterprise); Goebel (1982, 1984, 2006) (The ‘Salzburg School of Dialectometry’); and Nerbonne et al. (1999); Heeringa (2004); Nerbonne (2006) (the ‘Groningen School of Dialectometry’). Whereas practitioners of traditional DIALECTOLOGY are dedicated to the study of ‘interesting’ – typically

phonological or lexical – dialect phenomena, one feature at a time, in a handful of dialects at most, dialectometrical inquiry endeavors to identify “general, seemingly hidden structures from a larger amount of features” (Goebel and Schiltz 1997:13). This means that dialectometricians put a strong emphasis on quantification, cartographic visualization, and exploratory data analysis to infer patterns from feature aggregates. Empirically, the bulk of the dialectometrical literature draws on linguistic atlas material as its primary data source. For example, Goebel (1982) investigates joint variability in 696 linguistic features that are mapped in the *Sprach- und Sachatlas Italiens und der Südschweiz* (AIS), an atlas that covers Italy and southern Switzerland; Nerbonne et al. (1999) analyze aggregate pronunciational dialect distances between 104 Dutch and North Belgian dialects on the basis of 100 word transcriptions provided in the *Reeks Nederlands(ch)e Dialectatlassen* (RND). Some dialectometricians have also relied on dialect dictionaries (for example, Speelman and Geeraerts 2008). Against this backdrop, Leinonen (2008), Heeringa et al. (2009), and Auer et al. (to appear) are rare examples of dialectometrical work which bases claims about aggregate accent differences on the analysis, auditory or acoustic, of actual speech samples. In any case, given that most dialect atlases and dictionaries focus on lexis and pronunciation it should surprise nobody that much of the dialectometrical literature drawing on such material is biased towards lexis and pronunciation at the expense of morphological and, especially, syntactic variation (but see Spruit 2005, 2006; Spruit et al. 2009 for some recent atlas-based yet syntax-centered dialectometrical work).

There is no shortage of corpus-based research on regional variation in morphology and syntax. But then again, while corpus-linguistic methodologies have increasingly found their way into the dialectological toolbox and while more and more dialect corpora are coming on-line (see Anderwald and Szmrecsanyi 2009 for an overview), it is fair to say that the corpus-linguistic community is not exactly drowning in research that marries the qualitative-philological jeweler’s-eye perspective inherent in the analysis of naturalistic corpus data with the quantitative-aggregational bird’s-eye perspective that is the hallmark of dialectometrical research. The present study aims to remedy this shortcoming by discussing a methodology to conduct CORPUS-BASED DIALECTOMETRY (cf. also Szmrecsanyi 2008). As a case study to highlight the empirical potential of the methodology, we shall tap the *Freiburg Corpus of English Dialects*, a naturalistic speech corpus that samples more than three dozen dialects all over Great Britain. On the basis of this corpus, the study calculates a measure of aggregate dialectal distance that is based on joint variability of 57 morphosyntactic dialect features. The investigation subsequently draws on a number of statistical analysis methods and utilizes a range of cartographic projections to geography to aid interpretation of aggregate dialect distances. We wish to emphasize, however, right at the outset that the present paper is methodological talk, prioritizing methodological aspects at the expense of detailed discussions and interpretations of results.

This paper is structured as follows. Section 2 presents two arguments in favor of corpus-based dialectometry. Section 3 introduces the *Freiburg Corpus of English Dialects*. Section 4 discusses the design of an appropriate feature catalogue as the empirical basis for the aggregate analysis. In Section 5, we discuss the feature extraction process and the creation of a so-called ‘frequency matrix’. Section 6 addresses the actual aggregation process, which yields a so-called ‘distance matrix’. In Section 7, we present ways to visually represent, analyze, and interpret aggregate distances and similarities between dialects. Section 8 offers some concluding remarks.

2. Why corpus-based dialectometry?

The marriage of corpus-based variationist research and aggregative-dialectometrical analysis techniques is desirable for two principal reasons.

First, multidimensional objects, such as dialects, call for aggregate analysis techniques. So-called “single-feature-based studies” (Nerbonne 2009:176), with their atomistic focus on typically just one feature, are fine when it is the features themselves that are of analytic interest. They are woefully inadequate, however, when it comes to characterizing multidimensional objects such as dialects or varieties (or relations between them). Outside linguistics, this sort of inadequacy is well-known: Taxonomists, for instance, typically categorize species not on the basis of a single morphological or genetic criterion, but on the basis of many; economists assess the economic climate not on the basis of individual macroeconomic indicators (e.g. unemployment), but also consider inflation, GDP per capita, interest rates, and so on. The problem with single-feature-based studies – in linguistics as well as everywhere else – is that feature selection is ultimately arbitrary (cf. Viereck 1985:94), and that the next feature down the road may or may not contradict the characterization suggested by the previous feature. Thus, there is no guarantee that different dialects will exhibit the same distributional behavior in regard to different features; isoglosses do not necessarily overlap (cf. Bloomfield 1984 [1933]: 329). In addition, individual features may have fairly specific quirks to them that are irrelevant to the big picture. This is why “single-feature studies risk being overwhelmed by noise, i.e., missing data, exceptions, and conflicting tendencies” (Nerbonne 2009:193). So, the aggregate perspective – in Goebel’s parlance, the “the synthetic interpretation” of linguistic data (Goebel 2006:415) – is called for when the analyst’s attention is turned to the forest, not the trees. Aggregation mitigates the problem of feature-specific quirks, irrelevant statistical noise, and the problem of inherently subjective feature selection, and thus provides a more robust linguistic signal.

Second, compared to dialect atlas material (and we subsume here dialect dictionary material), corpora yield a more realistic linguistic signal. Atlas-based dialectometry typically aggregates observations such as ‘in the Yorkshire dialect, the lexeme *bus* is typically pronounced /bʊs/’, while corpus-based (that is to say, frequency-based) approaches seek generalizations along the lines of ‘in Nottinghamshire English, multiple negation is twice as frequent (6 occurrences per ten thousand words) in actual speech than in Yorkshire English (3 occurrences per ten thousand words)’. The atlas-based method has undeniable advantages. We emphasize, in particular, a fairly widespread availability of data sources and superb areal coverage. By contrast, dialect corpora are a rarer species, and their areal coverage is typically inferior to dialect atlases. Having said that, as a data source, corpora appear to have two major advantages over dialect atlases. First and foremost, the atlas signal is categorical, exhibits a high level of data reduction, and may hence be less accurate than the corpus signal, which can provide graded frequency information (cf. Wälchli 2009; Holman et al. 2007:413). This highlights the most crucial difference between atlas-based and corpus-based dialectometry: *corpus-based dialectometry is frequency-based dialectometry in its purest form* (which is why the approach outlined in this paper bears a certain similarity to the method of Hoppenbrouwers and Hoppenbrouwers 2001, discussed in some length in Heeringa 2004:16–20). The point is that although the exact cognitive status of

text frequencies is admittedly still unclear (for example, we do not currently know about the precise extent to which corpus frequencies correlate with psychological entrenchment; cf. Arppe et al. to appear), we do claim that text frequencies match better with perceptual salience than discrete atlas classifications; this is true even though some varieties of atlas-based dialectometry derive – with considerable computational effort – some form of commonness weighting, for instance at the phonetic segment level, from the atlas signal. Second, we note that the atlas signal is non-naturalistic and, basically, meta-linguistic in nature. It typically relies on elicitation and questionnaires, and is analytically twice removed (via fieldworkers *and* atlas compilers) from the analyst. By contrast, text corpora provide more direct access to language form and function, and may thus yield a more realistic and trustworthy picture (cf. Chafe 1992:84; Leech et al. 1994:58). The well-known major intrinsic drawback of the corpus-based method is that it is unable to deal with rare phenomena (cf. Penke and Rosenbach 2007:489; Haspelmath 2009:157–158) – but then again, it is arguable whether phenomena that are so infrequent that they cannot be described on the basis of a major text corpus should have a place in an aggregate analysis at all.

3. The data source: The *Freiburg Corpus of English Dialects* (FRED)

By way of a sample analysis, this paper will tap the *Freiburg Corpus of English Dialects* (henceforth: FRED) (see Hernández 2006; Szmrecsanyi and Hernández 2007 for details). The version of the corpus used in the present study contains 368 individual texts and spans approximately 2.44 million words of running text, consisting of samples (mainly transcribed so-called ‘oral history’ material) of dialectal speech from a variety of sources. Most of these samples were recorded between 1970 and 1990; in most cases, a fieldworker interviewed an informant about life, work, etc. in former days. The 431 informants sampled in the corpus are typically elderly people with a working-class background – so-called NORMs (*non-mobile old rural males*) (cf. Chambers and Trudgill 1998:29). The interviews were conducted in 156 different locations (that is, villages and towns) in 34 different pre-1974 counties in Great Britain including the Isle of Man and the Hebrides. The level of areal granularity investigated in the present study will be the county level. This leaves us with 34 objects (i.e. dialects or measuring points, which are listed in Table 1) that will be exemplarily subjected to dialectometrical analysis in the subsequent sections. Note that the corpus is annotated with longitude/latitude information for each of the locations sampled. From this annotation, county coordinates (mean longitude and latitude) can be calculated by computing the arithmetic mean of all the location coordinates associated with a particular county.

4. The empirical foundation: Defining the feature catalogue

The first step towards dialectometrical analysis is defining the FEATURE CATALOGUE as the empirical basis for the corpus-cum-aggregation endeavor. In keeping true to the spirit of dialectometrical analysis, the goal is to base the analysis on as many features as possible. In the case

map label	county	<i>a-priori</i> dialect area (Trudgill 1999)	mean longitude	mean latitude	no. words sampled in FRED
ANS	Angus	Sc Lowlands	-2.627	56.659	19,933
BAN	Banffshire	Sc Lowlands	-2.949	57.543	5,671
CON	Cornwall	Southwest of E	-5.502	50.175	107,290
DEN	Denbighshire	Wales	-3.743	53.146	5,789
DEV	Devon	Southwest of E	-3.681	50.378	97,229
DFS	Dumfriesshire	Sc Lowlands	-3.839	55.003	10,019
DUR	Durham	North of E	-1.703	54.890	28,086
ELN	East Lothian	Sc Lowlands	-2.954	55.945	40,193
GLA	Glamorganshire	Wales	-3.634	51.641	53,229
HEB	Hebrides	Hebrides	-7.038	57.502	73,209
MAN	Isle of Man	Isle of Man	-4.446	54.257	10,945
KCD	Kincardineshire	Sc Lowlands	-2.465	56.974	7,514
KEN	Kent	Southeast of E	0.835	51.246	177,055
LAN	Lancashire	North of E	-2.730	53.653	205,475
LEI	Leicestershire	E Midlands	-1.623	52.752	5,864
LND	London	Southeast of E	-0.068	51.504	110,878
MDX	Middlesex	Southeast of E	-0.382	51.594	31,794
MLN	Midlothian	Sc Lowlands	-3.265	55.918	32,040
NBL	Northumberland	North of E	-1.680	55.302	30,771
NTT	Nottinghamshire	E Midlands	-1.055	53.011	150,889
OXF	Oxfordshire	Southwest of E	-1.598	51.787	15,139
PEE	Peebleshire	Sc Lowlands	-3.377	55.721	14,975
PER	Perthshire	Sc Lowlands	-3.530	56.368	20,960
ROC	Ross and Cromarty	Sc Highlands	-4.776	57.808	10,495
SAL	Shropshire	E Midlands	-2.471	52.653	169,133
SEL	Selkirkshire	Sc Lowlands	-3.002	55.502	9,365
SFK	Suffolk	Southeast of E	1.699	52.555	312,600
SOM	Somerset	Southwest of E	-2.792	51.112	208,264
SUT	Sutherland	Sc Highlands	-4.676	58.144	11,025
WAR	Warwickshire	E Midlands	-1.968	52.574	8,271
WES	Westmorland	North of E	-2.962	54.428	157,590
WIL	Wiltshire	Southwest of E	-2.031	51.259	186,239
WLN	West Lothian	Sc Lowlands	-3.784	56.001	18,418
YKS	Yorkshire	North of E	-1.174	54.424	90,963

Table 1: $N = 34$ objects (i.e. FRED counties/dialects) considered in the present study: map labels, membership in *a-priori* dialect areas roughly following Trudgill's dialect division on pronunciation grounds (Trudgill 1999:Map 9), mean longitude, mean latitude, textual coverage (running words) in FRED.

study at hand, we surveyed the dialectological, variationist, and corpus-linguistic literature on morphosyntactic variability in varieties of English, and identified suitable dialect phenomena. This resulted in a list of $p = 57$ features, which overlaps with but is not identical to the comparative morphosyntax survey in Kortmann and Szmrecsanyi (2004) and the battery of morphosyntax features covered in the *Survey of English Dialects* (cf. Orton et al. 1978; Viereck et al. 1991). The features in the catalogue fall into eleven major grammatical domains: (i) pronouns and determiners (e.g. non-standard reflexives), (ii) the noun phrase (e.g. zero plural endings), (iii) primary verbs (e.g. the verb TO DO), (iv) tense & aspect (e.g. the present perfect with auxiliary BE), (v) modality (e.g. epistemic/deontic *must*), (vi) verb morphology (e.g. non-standard weak past tense and past participle forms), (vii) negation (e.g. *never* as a preverbal past tense negator), (viii) agreement (e.g. non-standard WAS), (ix) relativization (e.g. the relative particle *what*), (x) complementation (e.g. unsplit *for to*), and (xi) word order & discourse phenomena (e.g. lack of auxiliaries in *yes/no* questions). A detailed discussion of the features in the catalogue is beyond the scope of the present paper, but the Appendix provides the complete list of features.

A few comments on the case study's criteria for feature selection are in order, however. For a feature to be included in the catalogue, it did not matter whether the feature had previously been reported as having geographic variation or not. For instance, feature [31] (the negative suffix *-nae*) has a very clear and well-known regional distribution, but feature [10] (preposition stranding) does not according to the literature. Also note that the catalogue contains fairly categorical and thus somewhat salient non-standard features, which tend to be either largely present or absent – feature [31] (the negative suffix *-nae*) is again a good example – but also encompasses features whose variation is more statistical in nature, and thus arguably less salient (for example, features [8] and [9] on gradient genitive variation). The features included in the catalogue also differ in terms of their 'standardness' – feature [2] (standard reflexives), for instance, is examined with respect to the text frequency of perfect standard forms, while feature [28] (non-standard weak verb forms) is not really acceptable in Standard English. In short, the feature catalogue seeks to span as many features as possible, regardless of their geographic distribution, the scope of their variability, and their standardness. The rationale is that non-geographic and/or random variability will cancel out in the aggregate view. For practical purposes, however, two criteria had to be met for a candidate feature to be included in the catalogue. First, to ensure statistical robustness of text frequencies, the feature had to be relatively frequent. Specifically, the feature had to have a raw frequency of at least 100 raw hits in FRED. This criterion ruled out demonstrably infrequent phenomena such as resumptive relative pronouns, double modals, the relativizer *as*, and so on. Second, a candidate feature also had to be extractable – subject to a reasonable input of labor resources – by a human coder. This is why, for example, many hard-to-retrieve null phenomena (such as zero relativization) or features where semantics must be taken into consideration (such as gendered pronouns) are not considered in the catalogue.

5. Data mining: Extracting feature frequencies and creating a frequency matrix

The second step consists of extracting feature frequencies and creating a FREQUENCY MATRIX. In terms of the present study's feature catalogue, 31 sufficiently 'surfacy' features (e.g. the negator *ain't*) were extracted right away by software; 26 features in the catalogue (e.g. *don't* with 3rd person singular subjects) required more or less substantial manual disambiguation prior to extraction. Szmrecsanyi (2010b) provides detailed coding schemes and discusses the technicalities of the extraction process for all 57 features in the catalogue. Once feature frequencies are extracted, the analyst will normalize text frequencies – for example, to frequency per 10,000 words – if, as is the case with most relevant corpora, textual coverage of individual dialects varies. At this stage, we also recommend a *log*-transformation as a customary method to de-emphasize large frequency differentials and to alleviate the effect of frequency outliers (cf. Shackleton 2007:43), thus increasing reliability of the frequency matrix. To illustrate: In FRED, the county Cornwall has a textual coverage of 12 interviews totaling about 107,000 words of running text (interviewer utterances excluded). In this material, feature [34] (negative contraction, e.g. *they won't do anything*) occurs 326 times, which translates into a normalized text frequency of $326 \times 10,000/107,000 \approx 30$ occurrences per ten thousand words. A *log*-transformation of this frequency yields a value of $\log_{10}(30) \approx 1.5$.¹ This is the figure that characterizes this specific measuring point (Cornwall) in regard to feature [34].

The next step is to create an $N \times p$ frequency matrix in which the N objects (that is, dialects) are arranged in rows and the p features in columns, such that each cell in the matrix specifies a particular (normalized and *log*-transformed) feature frequency. Our case study thus yields a 34×57 frequency matrix: 34 British English dialects, each characterized by a vector of 57 text frequencies.

At this point, the analyst must assess the RELIABILITY of the frequency matrix: Are the features included in the catalogue a heterogeneous, mixed bag (for which an aggregate analysis would be meaningless), or is there a sufficient degree of consistency? Calculating a statistic known as Cronbach's α (cf. Cronbach 1951; Nunnally 1978) can address this issue.² Cronbach's α is, technically speaking, a coefficient measuring the average inter-item (in our case, inter-feature) correlation. Cronbach's α can take values between negative infinity and 1. An α value of 0 indicates that the features under consideration are not at all related, and a value of 1 means that all the features are perfectly correlated. Higher α values thus indicate a higher reliability of the frequency matrix. By convention – in dialectometry (cf. Heeringa 2004:173) and elsewhere (cf., for example, Bland and Altman 1997) – researchers aim for Cronbach's α values of .7 or higher. If a given frequency matrix yields a Cronbach's α value smaller than .7, there is a problem that should be addressed by expanding or altering the composition of the feature catalogue. Our case study's 34×57 frequency matrix yields a Cronbach's α value of .77, a score that comfortably passes the conventional threshold.

6. Aggregation: Obtaining a distance matrix

The task before us now is to convert the $N \times p$ frequency matrix into an $N \times N$ DISTANCE MATRIX. This transformation is an AGGREGATION step, in that the resulting distance matrix abstracts away from individual feature frequencies and specifies pairwise distances between the objects considered (similar to distance tables to be found in, e.g., road atlases). How do we calculate aggregate distances? Standard software packages offer a bewildering array of distance measures. Yet given the continuous nature of corpus-derived frequency vectors, we advocate usage of the well-known and fairly straightforward EUCLIDEAN DISTANCE MEASURE (see, for instance, Aldenderfer and Blashfield 1984:25) unless there is a good reason not to use it. Drawing on the Pythagorean theorem (cf. Nishisato 2007:77), the Euclidean distance measure defines the distance between two objects a and b as the square root of the sum of all p squared frequency differentials:

$$d(a, b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_p - b_p)^2} = \sqrt{\sum_{i=1}^p (a_i - b_i)^2} \quad (1)$$

where p is the number of features, a_1 is the frequency of feature 1 in object a , b_1 is the frequency of feature 1 in object b , a_2 is the frequency of feature 2 in object a , and so on. The Euclidean distance measure is, for one thing, interpretationally convenient: In two-dimensional space, it yields the distance between two points that one would measure with a ruler, which is why the measure is also sometimes referred to as “ruler distance” (Giles 2002:139). Furthermore, the Euclidean distance measure is theory-neutral in that all features receive the same weight in the distance calculation. Having said that, we stress that bigger frequency *differentials* receive proportionally more weight than smaller frequency differentials, which must appear as a desirable property to all those who believe that corpus frequencies mirror some sort of psychological and perceptual reality.

The chart in Figure 1 illustrates the aggregation process. In step ①, we start out with a fictional 3×2 frequency matrix, which has 6 cells specifying frequencies of 2 features in 3 dialects. In step ②, we calculate three distances: the distance between dialects a and b (which we commonsensically define as identical to the distance between dialects b and a), the distance between dialects a and c , and the distance between dialects b and c . In step ③, we enter these distances into a 3×3 distance matrix, which has $3 \times \frac{3-1}{2} = 3$ unique cells, i.e. dialect/dialect pairings. The other cells are redundant in that the distance between a given dialect and itself is always zero, and the distances in the upper right half of the matrix would mirror the distances in the lower left half of the matrix.

7. Visualization, analysis, and interpretation

Distance matrices can be analyzed in a myriad ways, not all of which may make sense for a particular set of research questions. This section sketches some ways to visually represent, analyze,

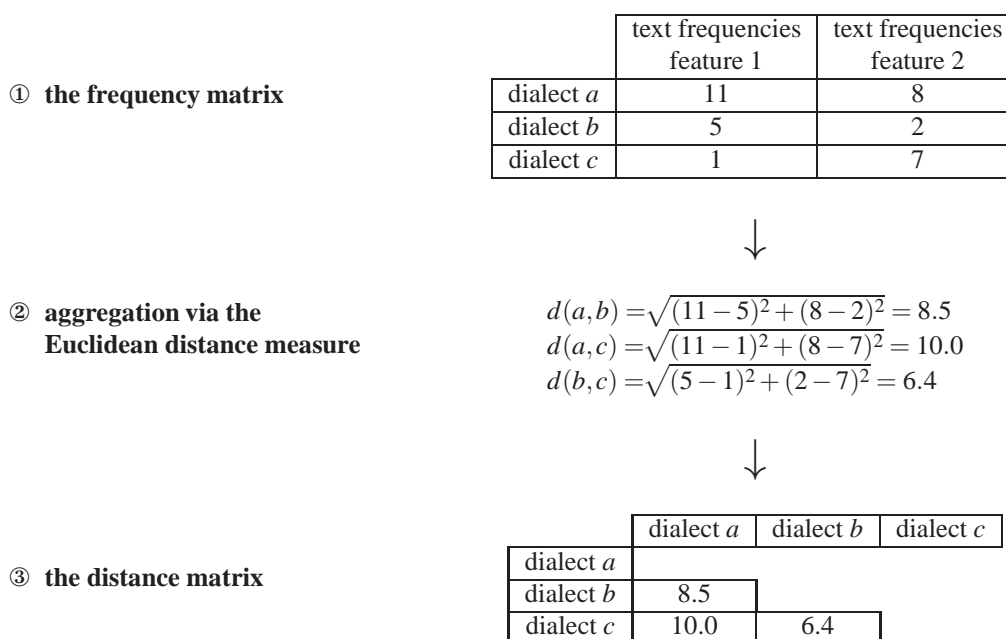


Figure 1: Converting a fictional 3×2 frequency matrix into a 3×3 distance matrix utilizing Euclidean distance as an aggregation measure.

and interpret distance matrices (geo)linguistically. Section 7.1 utilizes so-called ‘network maps’ to project aggregate dialect distances and similarities to geography. In section 7.2, we rely on ‘continuum maps’ to probe the extent to which joint dialectal variability is structured in terms of a dialect continuum. Section 7.3 draws on ‘cluster maps’ to explore the existence of dialect areas. Section 7.4 marshals correlative statistical analysis techniques to gauge the explanatory power of a number of language-external distance measures.³

7.1. Projecting aggregate distances and similarities to geography

In this section, we will discuss how to investigate the distribution of aggregate dialect distances. In this spirit, the table in Figure 2 provides a number of SUMMARY STATISTICS which describe the distribution of morphosyntactic dialect distances in Great Britain. Our distance matrix spanning $N = 34$ FRED dialects yields $34 \times 33/2 = 561$ pairwise distances. Mean morphosyntactic distance is 5.41 Euclidean distance points. This distance roughly corresponds to the distance between two hypothetical dialects *a* and *b* where dialect *a* attests a normalized text frequency of 2 hits per 10,000 words for each of the 57 features, while dialect *b* attests a normalized text frequency of approximately 10 hits per 10,000 words for each of the 57 features. As for the dataset-internal dispersion around the mean, we are dealing with a standard deviation of 1.11. Given that the distances are, as we shall see shortly, normally distributed, this is another way of saying that roughly two thirds of the 561 county/county pairings score a distance within 1.11 points of the mean, and that 95% of all pairwise distances do not deviate more than 2.22 points

measuring point pairings	561
mean	5.41
standard deviation	1.11
minimum	2.32
median	5.40
maximum	8.14
skewness	-.06
kurtosis	-.37

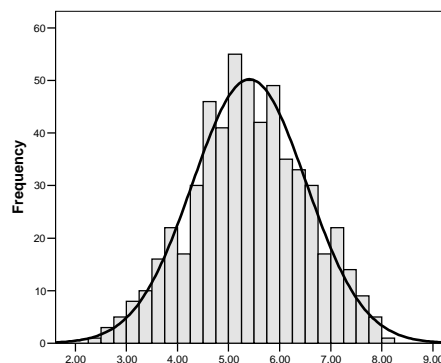


Figure 2: Aggregate morphosyntactic distances: summary statistics (left) and histogram (right).

from the mean. The minimum observable distance in the dataset is 2.32 points (this happens to be the morphosyntactic distance between the dialects spoken in the county of Somerset and the county of Wiltshire, two neighboring counties located in the Southwest of England). The median distance is 5.40 Euclidean distance points. This is the distance that separates the higher half of the distance sample from the lower half. The maximum observable distance in the dataset is 8.14 points, which is the distance between the dialects spoken in the county of Denbighshire in Wales and the county of Kincardineshire in the Scottish Lowlands.

Are pairwise morphosyntactic distances between the counties normally distributed? The histogram in Figure 2, which plots the frequency of a number of distance brackets, suggests that they roughly are. Numerically speaking, the skewness value of $-.06$ suggests that there is only a very slight negative skew, such that there is a greater number of larger distances than smaller distances. As for ‘peakedness’, the kurtosis value of $-.37$ indicates that the distribution of distances is a bit flatter than it would be in a perfectly normal distribution. Having said that, skewness and kurtosis values of ± 1.0 are by convention (cf. Meyers et al. 2006:90) taken to be indicative of a normal – albeit not perfectly normal – distribution.

The maps in Figure 3 are Groningen-style NETWORK MAPS (cf. Nerbonne and Heeringa 1997) that project dialect distances to geography without much statistical ado.⁴ The simple idea behind the left map in Figure 3 is that dialects that are close linguistically are linked by darker, more blueish lines, while linguistically more distant dialects are linked by proportionally lighter, more yellowish lines. Visual inspection of the map reveals that we are dealing with a network of comparatively strong morphosyntactic links in England, and with a somewhat looser network structure in Scotland. Within England, we observe particular strong link bundles in the South and in the North. Northumberland seems to link well to some Scottish measuring points, and turning back to the literature we note that both Ellis (1889) and Trudgill (1999) actually consider the traditional dialect spoken in Northern Northumberland a Scots variety. The Hebrides have strong morphosyntactic ties to measuring points all over Great Britain; notice in this connection that as a Scottish Highlands variety, Hebridean English is a relatively young dialect which Trudgill, for example, does not in fact categorize as a traditional dialect on account of the fact that it has “become English-speaking only relatively recently” (Trudgill 1999:5). It is therefore not

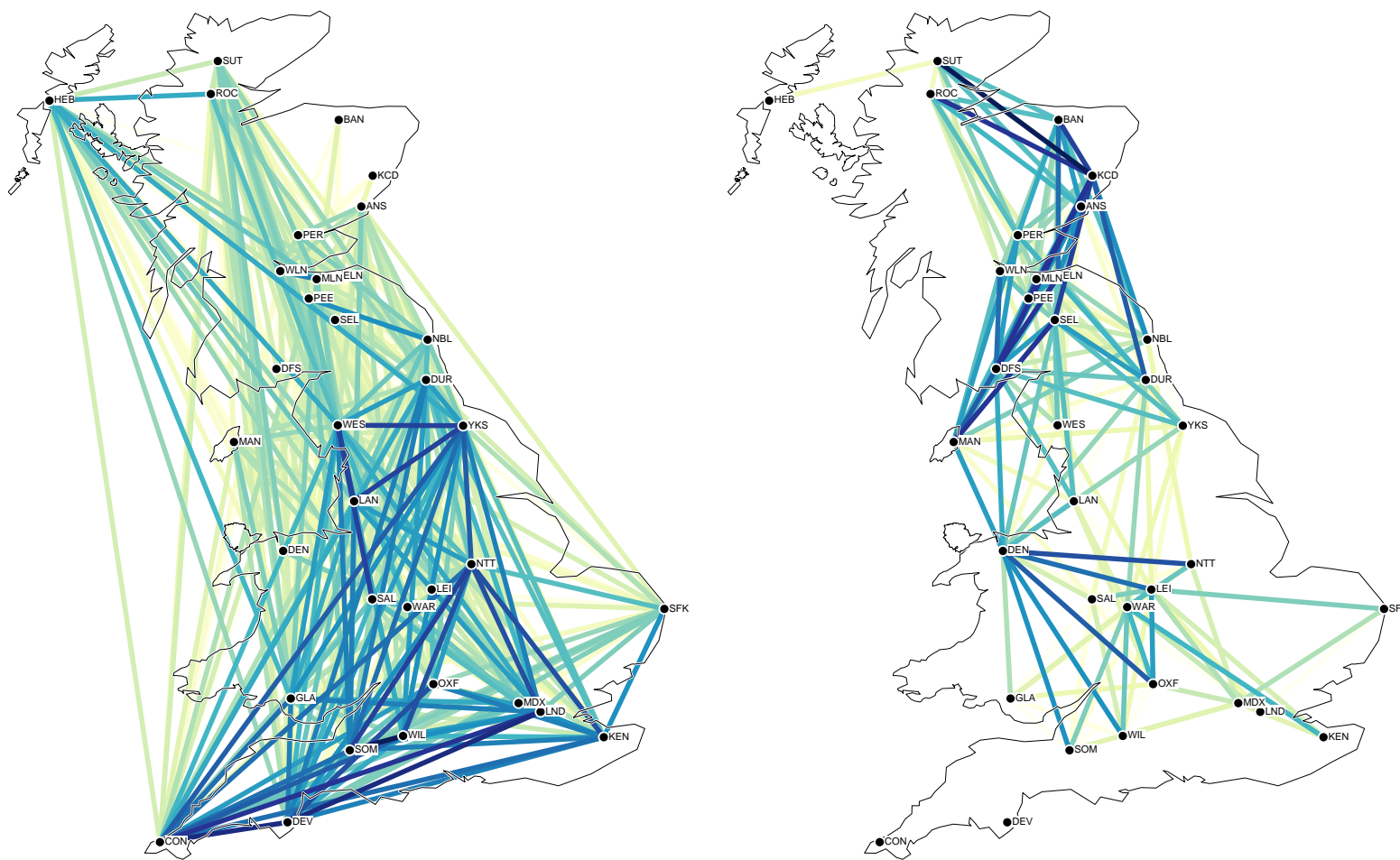


Figure 3: Projecting aggregate morphosyntactic dialect relationships to geography: network maps. Link blueness is directly proportional to dialectal similarity (left) or dialectal distance (right).

particularly surprising that Hebridean English lacks a clear-cut morphosyntactic profile of its own and bears similarities to dialects all over the place.

As a mirror image of sorts to the previous discussion, the right map in Figure 3 – a REVERSE NETWORK MAP – highlights morphosyntactic *dissimilarities* in the dataset. This particular map omits links between dialects that are more than 250 km apart, which mainly serves presentational purposes by enhancing readability without swamping the reader with an abundance of dissimilarity links. In Figure 3's reverse network map we observe, first, the striking tangle of dissimilarities covering much of Northern England and Scotland. In this connection, it is especially Banffshire and Kincardineshire that radiate strong beams of dissimilarity to other measuring points in the FRED network. In England, there is a web of modest dissimilarities involving measuring points in the Midlands (Shropshire, Warwickshire, Leicestershire, and Nottinghamshire). Warwickshire is additionally dissimilar to a number of sites in the Southwest of England (Somerset, Wiltshire, and Oxfordshire). As for Wales, observe the strong signal of dissimilarity emanating from the county of Denbighshire in Northern Wales; by contrast, Glamorganshire in Southern Wales is fairly inconspicuous. This difference between the two measuring points in Wales is not surprising. We know that there are robust lexical, phonological, and grammatical differences between dialects in the North and in the South of Wales (Penhallurick 1993), many of which can be traced back to an influx of Southwestern English speakers to the South of Wales beginning as early as the end of the eleventh century AD (Penhallurick 2004:98). By comparison, Denbighshire English in the North of Wales is a younger dialect with less well-established historical links to English English dialects.

This section has provided a first impression of the overall picture by considering summary statistics as well as network maps to represent aggregate distances in geographic map space. In what follows, we push deeper into the geographic structure of linguistic variation, subjecting the distance matrix to a good deal of statistical processing and subsequently projecting the output to geography.

7.2. Exploring dialect continua

Many dialectologists and geolinguists assume that geographic proximity predicts dialectal similarity. Nerbonne and Kleiweg (2007:154) refer to this axiom as the “Fundamental Dialectology Principle”. This section presents ways to depict the extent to which linguistic distance is directly proportional to geographic distance such that there are “no real boundaries, but only gradual transitions” (Bloomfield 1984 [1933]: 341).

To approach this issue cartographically, we turn to so-called CONTINUUM MAPS, a signature visualization technique developed in Groningen (cf. Nerbonne et al. 1999; Heeringa 2004). On the cartographic side, we set the scene by utilizing customary Voronoi tessellation (Voronoi 1907; Goebel 1984) to assign each dialect site on the map a convex polygon such that each point within the polygon is closer to the generating dialect site than to any other dialect site. Notice that when areal coverage is very fine-grained (as is usually the case in dialect atlases), it makes sense to exhaustively tessellate map space into Voronoi polygons. However, our case study covers Great Britain with $N = 34$ sampling points, which is why we prefer to limit the radius of the Voronoi

polygons to approximately 50km in order to do visual justice to the areal coverage of the dialect corpus. The next step is a computational one and subjects the data to MULTIDIMENSIONAL SCALING (MDS) (see Kruskal and Wish 1978; Embleton 1993). MDS is an exploratory statistical technique used to reduce a higher-dimensional dataset to a lower-dimensional representation which is more amenable to visualization. The task here is to scale down a $N - 1$ dimensional distance matrix (in which each object is characterized by its distance to the other $N - 1$ objects in the matrix) to a three-dimensional representation, in which each object has a coordinate in three artificial MDS dimensions. These coordinates are then mapped to the red–green–blue color scheme, giving each of the Voronoi polygons a distinct hue. On the interpretational plane, then, smooth color transitions between dialect polygons emphasize the continuum-like nature of the dialect landscape; abrupt color transitions point to the necessity of alternative explanations.⁵

Turning back to our case study, in Figure 4 we find two continuum maps that explore and, in fact, correlate (cf. Goebel 2005) the dialect landscape to the geographic landscape in Great Britain. The left-hand map is based on scaling a distance matrix detailing not linguistic distances but as-the-crow-flies geographic distances between dialect sites, thus depicting, for reference purposes, a perfect continuum. The right-hand map in Figure 4 visually represents an MDS solution that scales actual morphosyntactic distances. Statistically speaking, the MDS distances underlying the left visualization capture 100% ($r = 1.0$) of the variance in the original as-the-crow-flies distance matrix. The MDS solution depicted in the right-hand map captures about 89.5% ($r = .95$) of the distance variance in the original linguistic dataset, which is a fairly good score.

In all, the mosaic pattern in the morphosyntax continuum map suggests that the morphosyntactic dialect landscape in Great Britain is less continuum-like than it could be. It is true that there are some fairly nice micro-continua, especially so in the Southwest of England and in the Central and Northern Scottish Lowlands. Note also that dialects spoken in the North of England fade rather smoothly into Southern Scottish Lowlands dialects. But we also observe rather abrupt transitions between the Central Scottish Lowlands – comprising dialects spoken in West Lothian, Midlothian, and East Lothian – and Southern Scottish dialects (Peebleshire and Selkirkshire). In England, the dialects spoken in Middlesex and Warwickshire are outliers. In Wales, it is Denbighshire that does not fit into the picture.

At this point, it is instructive to abandon the aggregate perspective for a moment and to reconsider the actual features on which the analysis is based. In this spirit, to aid interpretation of continuum maps the analyst can correlate frequency vectors with MDS dimensions to identify those features that are most robustly implicated in the overall dimensionality (cf. Heeringa 2004:266-271). In the case study at hand, this produces the following top correlations between color shades and feature frequencies in the right-hand map displayed in Figure 4:

- Increased text frequencies of feature [33] (multiple negation, as in *you didn't want no beer*) correlate best with MDS dimension 1, which yields reddish tones ($r = .82, p < .001$).
- Increased text frequencies of feature [31] (the negative suffix *-nae*, as in *I cannae mind of ever being laid off*) correlate best with MDS dimension 2, which yields greenish tones ($r = .75, p < .001$).
- Increased text frequencies of feature [28] (non-standard weak past tense and past participle

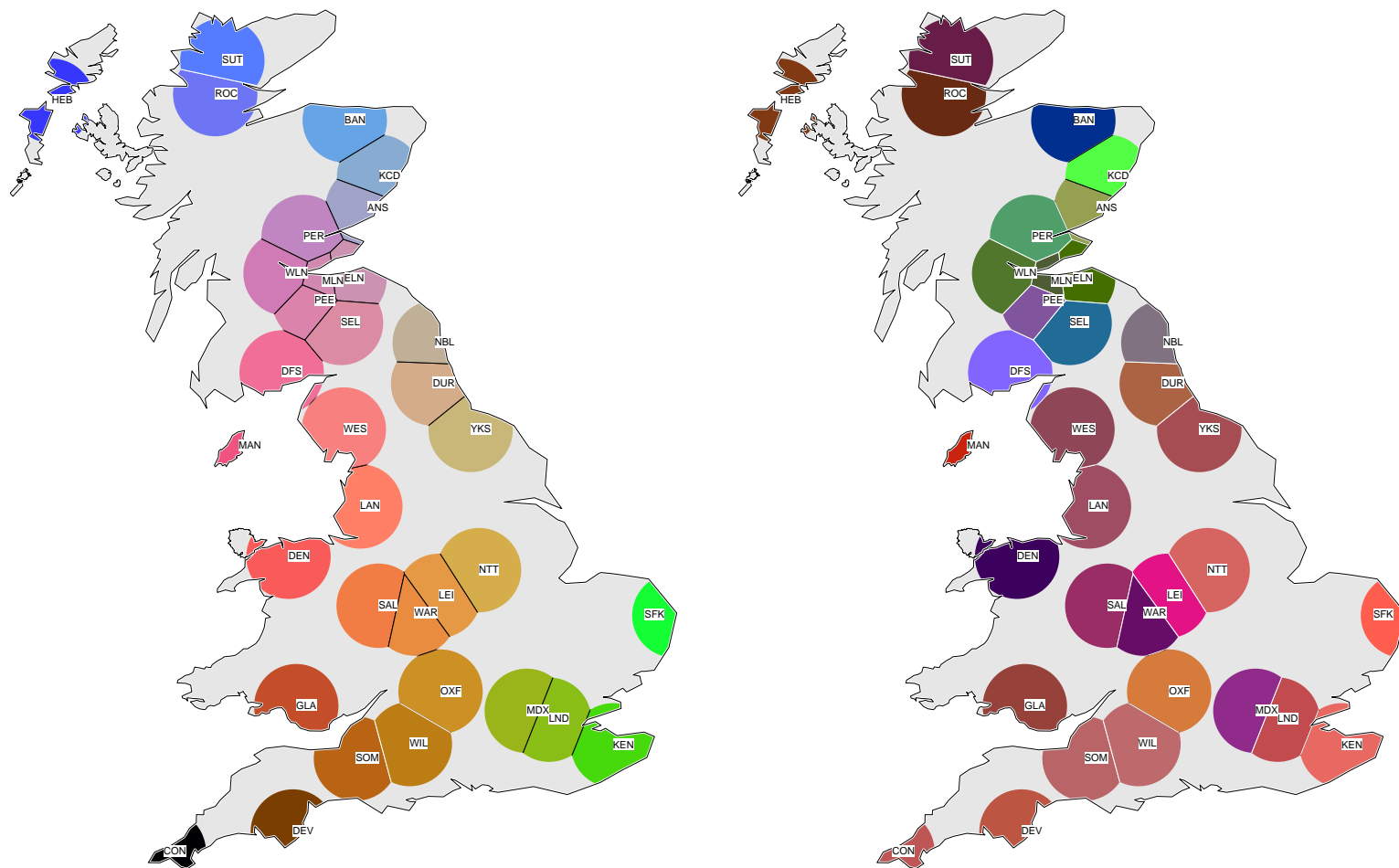


Figure 4: Continuum maps: a perfect continuum for reference purposes (left) versus the actual linguistic situation (right). Similar color hues indicate morphosyntactic similarity. Input left projection: as-the-crow-flies distances (correlation with distances in the original distance matrix: $r = 1.00$). Input right projection: morphosyntactic distances (correlation with distances in the original distance matrix: $r = .95$)

forms, as in *we runned up a bill*) correlate best with MDS dimension 3, which yields blueish tones ($r = .72, p < .001$).

By way of an interim summary, we have seen in this section that viewing aggregate morphosyntactic variability in a dialect continuum perspective can be instructive. Yet this approach does not necessarily tell the whole story: The existence of abrupt dialect transitions, as in Figure 4, suggests that linguistic variability may be organized in terms of dialect areas rather than continua. The next section is dedicated to investigating this hypothesis more closely.

7.3. The dialect area scenario

The tacit assumption guiding the foregoing discussion was that linguistic similarity between dialects is inversely proportional to geographic distance between dialects. There is, however, an alternative view, according to which dialect landscapes may be geographically organized along the lines of geographically coherent and linguistically homogeneous “areas within which similar varieties are spoken” (Heeringa and Nerbonne 2001:375). In this view, we should find linguistic boundaries between rather than within dialect areas. In this section, we discuss methods to explore this view.

The dialect area scenario lends itself to visualization via CLUSTER MAPS, a cartographic technique that is common in all strands of dialectometry and which projects the outcome of cluster analysis to geography (cf., for example, Goebel 2007:Map 18; Heeringa 2004:Figure 9.6). Exactly as with continuum maps, the starting point is a Voronoi tessellation of map space. Subsequently, the $N \times N$ distance matrix is subjected not to MDS, but to HIERARCHICAL AGGLOMERATIVE CLUSTER ANALYSIS (cf. Jain et al. 1999), a statistical technique used to group a number of objects (in this study, dialects) into a smaller number of discrete clusters. While there are many different clustering algorithms, we prefer ‘Ward’s Minimum Variance Method’ (Ward 1963), an algorithm that tends to create small and even-sized clusters and which is popular both in corpus linguistics (for example, Gries and Wulff 2005) and in dialectometry (cf., for instance, Goebel 2008).⁶ Cluster analysis initially yields a so-called DENDROGRAM (cf. Figures 5 and 6), which depicts cophenetic distances between the clustered objects. The optimal number of clusters is determined by, e.g. diagramming the number of clusters against the fusion coefficient and spotting the ‘elbow’ in the resulting graph (cf. Aldenderfer and Blashfield 1984:54). Finally, each of the clusters is assigned a distinct color and the Voronoi polygons are colorized accordingly.⁷

Applying these steps to our dataset on dialect variability in Great Britain yields Figures 5 and 6. Figure 5 – which is based not on morphosyntactic but on as-the-crow-flies geographic distances – will serve as the non-linguistic reference point for our discussion. The map suggests that on strictly geographic grounds and according to Ward’s method, Great Britain can be partitioned into three coherent areas: a green region comprising the South of England plus the county of Glamorganshire in Southern Wales; a red region containing the North of England plus the county of Denbighshire in Northern Wales plus the county of Dumfriesshire in Southern Scotland; and a blue region encompassing Scotland minus the county of Dumfriesshire.

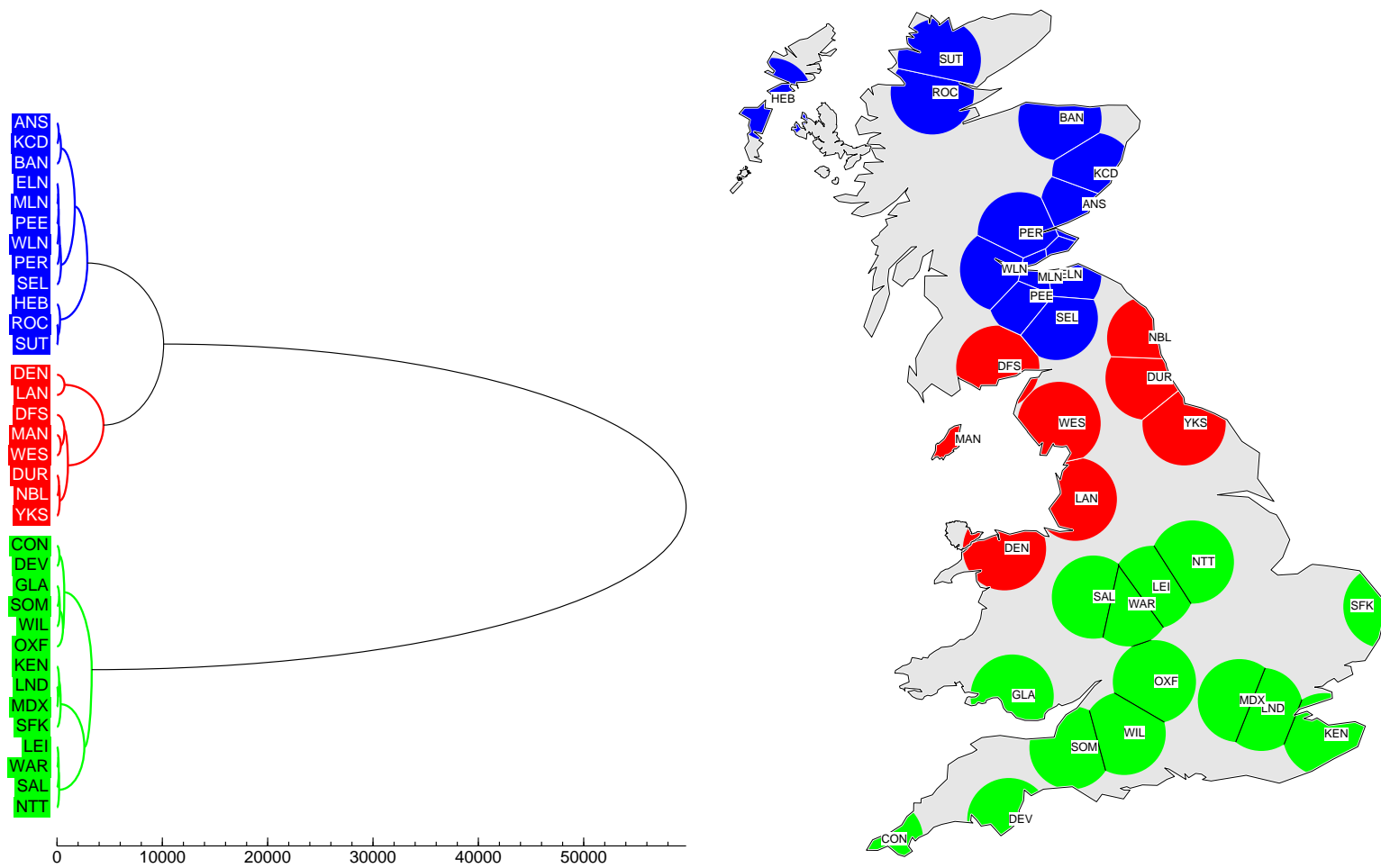


Figure 5: A perfect regionalization for reference purposes: Clustering as-the-crow-flies distances (hierarchical agglomerative cluster algorithm: WARD). Displayed: 3-cluster solution. Left: dendrogram. Right: cluster map. Colors indicate regional membership.

Compare this landscape to Figure 6, which visually depicts a 5-cluster regionalization on morphosyntactic grounds. There is clearly some similarity between the geographic and linguistic partitioning, although we note that there is also a good deal of geographic incoherence in the morphosyntax division. A more detailed account would highlight the following differences between the maps in Figures 5 and 6:

- The yellow dialect grouping in the morphosyntax map encompasses some isolated outliers in England (Middlesex and Warwickshire), Wales (Denbighshire), and a geographically coherent sub-cluster of Scottish Highland dialects (Ross and Cromarty and Sutherland) plus the Hebrides.
- The morphosyntax visualization has a small yet geographically coherent Central Scottish Lowlands dialect area (in light blue), comprising the counties of East Lothian, Midlothian, and West Lothian.
- Also in contrast to the geographic division, the dark blue Scottish cluster in the morphosyntax map includes both Dumfriesshire as well as Northumberland, a dialect site that is actually located in political England.
- In the morphosyntax partitioning, the red Northern England area additionally comprises Shropshire and Leicestershire in what is often referred to as the English Midlands, as well as Glamorganshire in Southern Wales.
- Compared to the geographic map, the greenish Southern English area is smaller in the morphosyntax partitioning: Linguistically speaking, Shropshire, Leicestershire, and Glamorganshire are – as we have seen – red dialects while Middlesex and Warwickshire are yellow outliers. Durham, a county that is geographically located in Northern England, is grouped with the Southern English English dialects.

The morphosyntax dendrogram in Figure 6 demonstrates that the by far most fundamental split in the dataset occurs between English English dialects (red and green) and other dialects (including yellow outliers). The second most crucial split is between Northern English English dialects (red) and Southern English English dialects (green). The least important split is the one between Central Scottish Lowlands dialects (light blue) and other Scottish Lowlands dialects (dark blue).

The name of the game in this section was classification. In this spirit, we have discussed ways to categorize dialects into discrete clusters. Our case study suggests that despite some geographic incoherence and the presence of outliers, Great Britain can be divided into three major morphosyntactic dialect areas: the Scottish Lowlands versus the North of England versus the South of England.

7.4. Quantifying the explanatory power of language-external predictors

Dialectometry is intrinsically quantitative, yet the foregoing discussion has relied heavily on interpreting cartographic projections to geography. In this section, we introduce techniques used

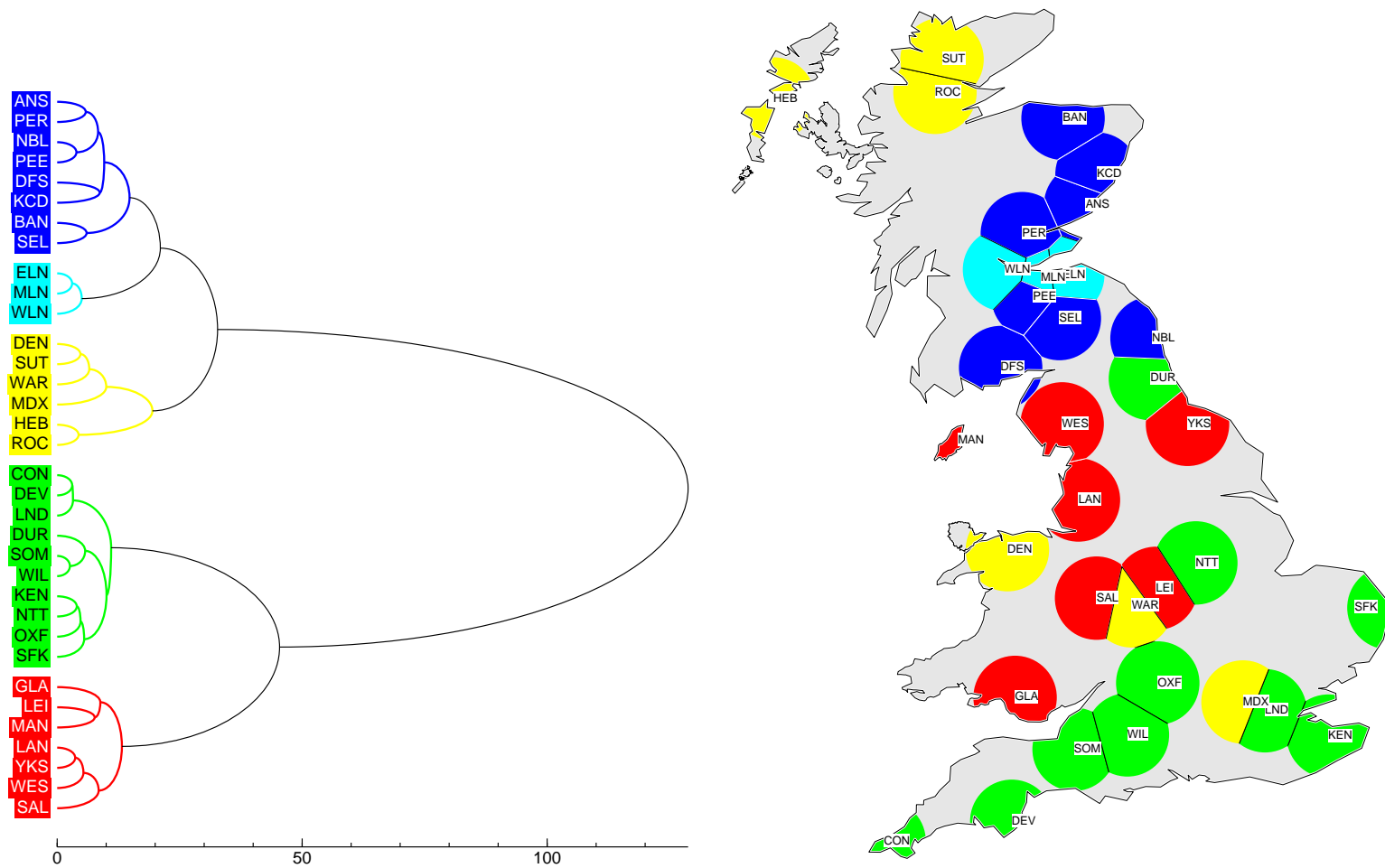


Figure 6: Actual morphosyntax clusters: Clustering morphosyntactic distances (hierarchical agglomerative cluster algorithm: WARD). Displayed: 5-cluster solution. Left: dendrogram. Right: cluster map. Colors indicate dialect area or dialect grouping membership.

to correlate language-external parameters with linguistic distances, for the sake of precisely quantifying the extent to which dialect distances are predictable from language-external factors.

To this end, the analyst typically starts out with an $N \times N$ linguistic distance matrix and creates parallel language-external distance matrices – one for each predictor to be tested. In the simplest case, each of these language-external distance matrices is then correlated with the linguistic distance matrix by calculating, e.g., a Pearson product-moment correlation coefficient. The language-external predictor that scores the highest coefficient is the best predictor of linguistic distances.

To exemplify, we revisit our dataset on dialect variability in Great Britain. We begin by correlating our 34×34 morphosyntactic distance matrix with three language external distance matrices:

1. AS-THE-CROW-FLIES DISTANCES. Using a trigonometry formula on the FRED county coordinates, it is computationally trivial to calculate pair wise as-the-crow-flies distances (these actually underlie the left-hand projection in Figure 4 and the cluster map in Figure 5).⁸ A proxy for the likelihood of social contact, as-the-crow-flies distance is the most common geographic distance measure in the literature (for example, Goebel 2001; Gooskens and Heeringa 2004; Nerbonne et al. 1996; Shackleton 2007).
2. LEAST-COST TRAVEL TIMES. Speakers do not have wings, so it is reasonable to assume that what really matters for dialect distances is how much time it would take a human traveler to get from point A to point B (cf. Gooskens 2005; Szmercsanyi 2010a). To calculate this measure, we turned to Google Maps (<http://maps.google.co.uk/>), which has a route finder tool that allows the user to enter longitude/latitude pairings for two locations to obtain a least-cost travel route and, crucially, an estimate of the total travel time. We queried Google Maps for all $34 \times 33/2 = 561$ dialect pairings, thus obtaining pair wise least-cost-travel time estimates.⁹
3. LINGUISTIC GRAVITY INDICES. Trudgill (1974) suggested a gravity model to account for geographic diffusion, claiming that “the interaction (M) of a centre i and a centre j can be expressed as the population of i multiplied by the population of j divided by the square of the distance between them” (1974: 233). Using Trudgill’s formula, we calculated linguistic gravity values for each of the 561 dialect pairings in our database, feeding in least-cost travel time as geographic distance measure and early twentieth century population figures¹⁰ (in thousands) as a proxy for speaker community size.

Figure 7 provides three scatterplots that graph morphosyntactic distances against the language-external distance measures. In all three cases, there is a highly significant relationship, and the direction of the effect is the theoretically expected one throughout: Increasing as-the-crow-flies distance and increasing least-cost travel time predict increasing morphosyntactic distance; conversely, increasing linguistic gravity indices predict decreasing morphosyntactic distance. The R^2 values reported in Figure 7 suggest that as-the-crow flies distance accounts for 4.4% of the morphosyntactic variance, least-cost travel time for 7.4%, and linguistic gravity for 24.1%. Hence, by factoring in speaker community size in addition to geographic distance, we can explain up to a quarter of the variance in morphosyntactic dialect distances. That Trudgill’s linguistic gravity model turns out to be the most successful predictor of morphosyntactic distances

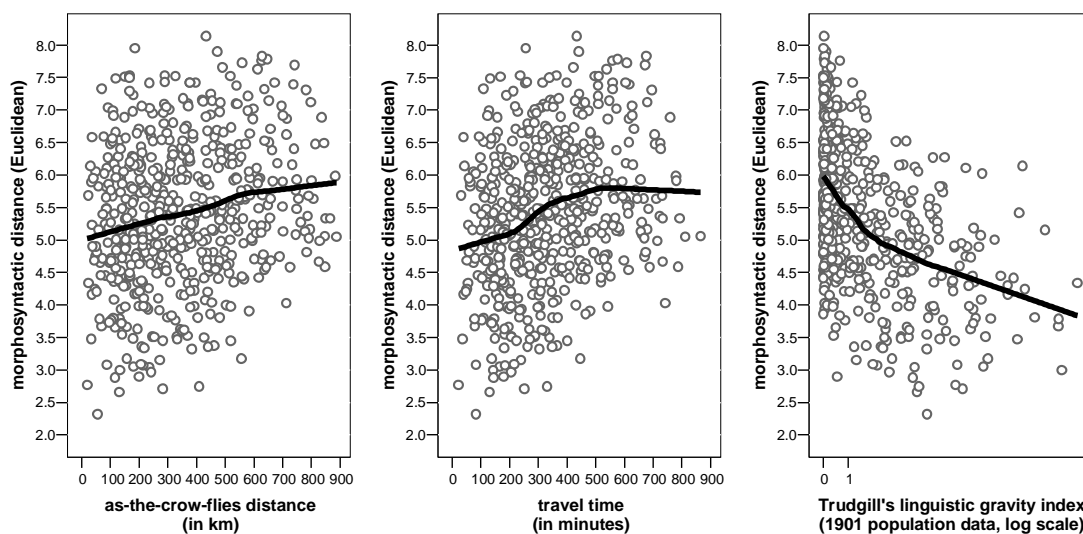


Figure 7: Correlating distance matrices: morphosyntactic distances versus (i) as-the-crow-flies distances (left) ($r = .21$, $p < .001$, $R^2 = 4.4\%$), (ii) least-cost travel times (middle) ($r = .27$, $p < .001$, $R^2 = 7.4\%$), and (iii) Trudgill's linguistic gravity indices (right; log scale) ($r = -.49$, $p < .001$, $R^2 = 24.1\%$; logarithmic estimate). Each dot represents one of $N = 561$ unique dialect pairings. Solid lines are LOESS curves estimating the overall nature of the relationship.

in the dataset is interesting, given that some previous research has failed to detect a significant effect of linguistic gravity in, e.g., Dutch dialects (cf. Heeringa et al. 2007; Nerbonne and Heeringa 2007). We conclude that unlike in Dutch dialectology, the model works rather well for morphosyntactic variation in traditional British English dialects.

Having said that, however, we emphasize that in comparison to previous research the R^2 values reported here are rather low anyway. For example, Shackleton (2007), in his study of phonetic variation in the *Survey of English Dialects*, reports R^2 values of up to 66% for the relationship between phonetic and geographic distances in England; Spruit et al. (2009), in an atlas-based study on aggregate syntactic distances in Dutch dialects, calculate an R^2 value of 45% for the relation between syntax and geography. So, given that the continuous distance measures in Figure 7 somewhat fail us, could it be that the dialect partition presented in Section 7.3 is a more potent predictor of morphosyntactic dialect distances? The problem is that binary dialect area (or: dialect grouping) memberships do not elegantly lend themselves to a straightforward continuous quantification with which one may furnish a full-blown $N \times N$ distance matrix that could be correlated with the original dialect distance matrix and depicted in a scatterplot. This is why we resort to a technique known as Permutational Multivariate Analysis of Variance Using Distance Matrices (PERMANOVA; cf. Anderson 2001)¹¹, which is analogous to MANOVA (Multivariate Analysis of Variance) but designed specifically to analyze distance matrices. The goal was to test how successful dialect area membership is in predicting continuous dialect distances between measuring points. It turns out that the 5-cluster partition depicted in Figure 6 explains

about one third of the variance in morphosyntactic distances ($R^2 = 35.9\%$, $p = .001$), which is a good deal more than we can explain by regressing as-the-crow-flies distance, travel time, or linguistic gravity against morphosyntactic distances.

This section was an exercise in number crunching, and we have seen that the relationship between dialect distances and properties of geographic space is amenable to fairly precise quantification. In traditional British English dialects, as-the-crow-flies distance turns out to be a fairly poor predictor of aggregate morphosyntactic distances, explaining no more than about 4% of the overall variability. By factoring in parameters such as size of speaker community and travel distance, the analyst can boost the share of the linguistic variability that is accounted for to about 25%. At the same time, we have also seen that the cluster analytic dialect partition presented in the previous section explains more than a third of the morphosyntactic variance in the dataset. In other words, the dialect area scenario appears to be more appropriate for our dataset than the dialect continuum scenario.

8. Conclusion

This paper has presented methodologies which can be used to combine corpus-based variation studies with aggregative-dialectometrical analysis and visualization methods. We have argued that this synthesis is desirable for two principal reasons. First, multidimensional objects, such as dialects, call for aggregate analysis techniques; second, vis-à-vis linguistic atlas material, corpora yield an arguably more trustworthy frequency signal. To exemplify the empirical potential of corpus-based dialectometry, we have drawn on a major dialect corpus to study aggregate relations between 34 traditional British English dialects, on the basis of joint variability in text frequencies of 57 morphosyntactic features. The analysis has demonstrated that linguistic variability between British English dialects demonstrably provides a geographic signal, and that this signal has a number of interesting facets.

Needless to say, the line of analysis sketched in this paper is extendable in many ways. First and foremost, due to the case study character of the investigation, we often stopped where real interpretation would start. More in-depth scholarship would seek to interpret the findings uncovered in this paper in a wider analytical and theoretical context, considering, among other things, the literature on dialect genesis, dialect formation, and historical dialect variability; theories on the status of, and constraints on, dialect grammar and grammatical variability between dialects from a typological perspective; and previous research on the issue of the contact-induced diffusibility versus universality of grammatical features.

Lastly, we should add that the methodology outlined in this paper is of course not limited to morphosyntactic phenomena. Phonology, lexis, and even pragmatics are all in principle amenable to dialectometrical analysis using a corpus-based approach. There is even the intriguing possibility of aggregating not ‘surfacy’ feature frequencies but ‘deep’ feature conditionings (e.g. via probabilistic regression weights), which is something that simply cannot be done on the basis of decontextualized atlas or dictionary data. As for suitable databases, corpus-based dialectometry can be applied to *any* corpus in which we find geographic variability. This includes not

only dialect corpora in the traditional sense (such as the *Freiburg Corpus of English Dialects* analyzed here), but also corpora sampling geographically non-contiguous regional language varieties (such as the *International Corpus of English*; cf. Greenbaum 1996) or corpora concerned with variation in written, not spoken, language (such as the letters-to-the-editor corpus presented in Grieve 2009). In short, there are a great many research opportunities just waiting to be tapped.

Appendix: The feature catalogue

A. Pronouns and determiners

- [1] non-standard reflexives (e.g. *they didn't go theirself*)
- [2] standard reflexives (e.g. *they didn't go themselves*)
- [3] archaic *thee/thou/thy* (e.g. *I tell thee a bit more*)
- [4] archaic *ye* (e.g. *ye'd dancing every week*)
- [5] *us* (e.g. *us couldn't get back, there was no train*)
- [6] *them* (e.g. *I wonder if they'd do any of them things today*)

B. The noun phrase

- [7] synthetic adjective comparison (e.g. *he was always keener on farming*)
- [8] the *of*-genitive (e.g. *the presence of my father*)
- [9] the *s*-genitive (e.g. *my father's presence*)
- [10] preposition stranding (e.g. *the very house which it was in*)
- [11] cardinal number + *years* (e.g. *I was there about three years*)
- [12] cardinal number + *year-Ø* (e.g. *she were three year old*)

C. Primary verbs

- [13] the primary verb TO DO (e.g. *why did you not wait?*)
- [14] the primary verb TO BE (e.g. *I was took straight into this pitting job*)
- [15] the primary verb TO HAVE (e.g. *we thought somebody had brought them*)
- [16] marking of possession – HAVE GOT (e.g. *I have got the photographs*)

D. Tense and aspect

- [17] the future marker BE GOING TO (e.g. *I'm going to let you into a secret*)
- [18] the future markers WILL/SHALL (e.g. *I will let you into a secret*)
- [19] WOULD as marker of habitual past (e.g. *he would go around killing pigs*)
- [20] *used to* as marker of habitual past (e.g. *he used to go around killing pigs*)
- [21] progressive verb forms (e.g. *the rest are going to Portree School*)
- [22] the present perfect with auxiliary BE (e.g. *I'm come down to pay the rent*)
- [23] the present perfect with auxiliary HAVE (e.g. *they've killed the skipper*)

E. Modality

- [24] marking of epistemic and deontic modality: MUST (e.g. *I must pick up the book*)
- [25] marking of epistemic and deontic modality: HAVE TO (e.g. *I have to pick up the book*)
- [26] marking of epistemic and deontic modality: GOT TO (e.g. *I gotta pick up the book*)

F. Verb morphology

- [27] *a*-prefixing on *-ing*-forms (e.g. *he was a-waiting*)
- [28] non-standard weak past tense and past participle forms
(e.g. *they knowed all about these things*)
- [29] non-standard past tense *done* (e.g. *you came home and done the home fishing*)
- [30] non-standard past tense *come* (e.g. *he come down the road one day*)

G. Negation

- [31] the negative suffix *-nae* (e.g. *I cannae do it*)
- [32] the negator *ain't* (e.g. *people ain't got no money*)
- [33] multiple negation (e.g. *don't you make no damn mistake*)
- [34] negative contraction (e.g. *they won't do anything*)
- [35] auxiliary contraction (e.g. *they'll not do anything*)
- [36] *never* as past tense negator (e.g. *and they never moved no more*)
- [37] WASN'T (e.g. *they wasn't hungry*)
- [38] WEREN'T (e.g. *they weren't hungry*)

H. Agreement

- [39] non-standard verbal *-s* (e.g. *so I says, What have you to do?*)
- [40] *don't* with 3rd person singular subjects (e.g. *if this man don't come up to it*)
- [41] standard *doesn't* with 3rd person singular subjects
(e.g. *if this man doesn't come up to it*)
- [42] existential/presentational *there is/was* with plural subjects
(e.g. *there was children involved*)
- [43] absence of auxiliary BE in progressive constructions
(e.g. *I said, How Ø you doing?*)
- [44] non-standard WAS (e.g. *three of them was killed*)
- [45] non-standard WERE (e.g. *he were a young lad*)

I. Relativization

- [46] *wh*-relativization (e.g. *the man who read the book*)
- [47] the relative particle *what* (e.g. *the man what read the book*)
- [48] the relative particle *that* (e.g. *the man that read the book*)

J. Complementation

- [49] *as what* or *than what* in comparative clauses
(e.g. *we done no more than what other kids used to do*)
- [50] unsplit *for to* (e.g. *it was ready for to go away with the order*)
- [51] infinitival complementation after BEGIN, START, CONTINUE, HATE, and LOVE
(e.g. *I began to take an interest*)
- [52] gerundial complementation after BEGIN, START, CONTINUE, HATE, and LOVE
(e.g. *I began taking an interest*)
- [53] zero complementation after THINK, SAY, and KNOW
(e.g. *they just thought [Ø it isn't for girls]*)
- [54] *that* complementation after THINK, SAY, and KNOW
(e.g. *they just thought [that it isn't for girls]*)

K. Word order and discourse phenomena

- [55] lack of inversion and/or of auxiliaries in *wh*-questions and in main clause *yes/no*-questions (e.g. *where Ø you put the shovel?*)
- [56] the prepositional dative after the verb GIVE (e.g. *she gave [a job] [to my brother]*)
- [57] double object structures after the verb GIVE (e.g. *she gave [my brother] [a job]*)

Notes

¹Text frequencies of 0 (for instance, feature [4] [archaic *ye*, as in *ye'd dancing every week*] does not occur in material from Cornwall) were rendered as 0.1, which yields a *log*-transformed value of $\log_{10}(0.1) = -1$.

²Standard statistical software packages, such as R and SPSS, can easily calculate the Cronbach's α statistic.

³On a technical note, all cartographic projections and most of the non-trivial statistical analyses (such as multidimensional scaling and cluster analysis) presented in this section were created using Peter Kleiweg's *RuG/L04* dialectometry software package (available for free at <http://www.let.rug.nl/~kleiweg/L04/>). The input required by *RuG/L04* is (i) the longitude/latitude coordinates provided in Table 1, (ii) a (linguistic) distance matrix, and (iii) a polygon map – which can be created using Google Earth – that defines the boundaries of a land mass and/or political borders. Note, along these lines, that there is another major dialectometry package: the *Visual Dialectometry* (VDM) software developed in Salzburg (Haimerl 2006). VDM is also free and comes, as an added bonus, with a graphical user interface.

⁴The maps were created using the *RuG/L04* package's `maplink` module.

⁵The visualization and statistical analysis techniques in this section draw on the *RuG/L04* package's `mDS` (method: Kruskal's MDS) and `maprgb` modules.

⁶Simple clustering can be unstable, so the analysis in this section utilizes a procedure known as “clustering with noise” (Nerbonne et al. 2008): The original distance matrix is clustered repeatedly, adding some random amount of noise ($c = \sigma/2$) in each run. This exercise yields a cophenetic distance matrix which provides consensus (and thus more stable) cophenetic distances between dialects.

⁷The visualization and statistical analysis techniques in this section draw on the *RuG/L04* package's `cluster` module (which implements clustering with noise; cf. fn 6), as well as on the `mapclust` and `den` modules.

⁸The *RuG/L04* dialectometry software package provides a module (`l12dst`) that can do this job automatically.

⁹These estimates assume travel by car (Google Maps' ‘walking’ option yields a matrix with a substantially lower correlation with linguistic distances). We fully acknowledge that matching linguistic data sourced from speakers born around the beginning of the twentieth century with travel estimates based on twenty-first century transportation infrastructure is convenient but clearly suboptimal; what is really needed are *historic* travel estimates, which, alas, are in short supply. Still, we submit that the procedure is not fatally flawed, as modern infrastructure can be argued to actually follow, to a large extent, historical travel routes, trade patterns, and avenues of social contact.

¹⁰Specifically, we used 1901 population figures, as published in the *Census of England and Wales, 1921* and the *Census of Scotland, 1921*. These documents are available online at <http://histpop.org/>.

¹¹To conduct the analysis, we utilized the statistical software package R: library `vegan`, function `adonis` (<http://vegan.r-forge.r-project.org/>).

References

- Aldenderfer, M. S. and R. K. Blashfield (1984). *Cluster Analysis*. Quantitative Applications in the Social Sciences. Newbury Park, London, New Delhi: Sage Publications.
- Anderson, M. J. (2001). A new method for non-parametric multivariate analysis of variance. *Austral Ecology* 26, 32–46.
- Anderwald, L. and B. Szmrecsanyi (2009). Corpus linguistics and dialectology. In A. Lüdeling and M. Kytö (Eds.), *Corpus Linguistics. An International Handbook*, Handbücher zur Sprache und Kommunikationswissenschaft/Handbooks of Linguistics and Communication Science, pp. 1126–1139. Berlin, New York: Mouton de Gruyter.
- Arppe, A., G. Gilquin, D. Glynn, M. Hilpert, and A. Zeschel (to appear 2011). Cognitive Corpus Linguistics: Five points of debate on current theory and methodology. *Corpora* 5(2).
- Auer, P., P. Baumann, and C. Schwarz (to appear). Vertical vs horizontal change in the traditional dialects of southwest Germany: a quantitative approach. *Taal en Tongval*.
- Bland, J. M. and D. G. Altman (1997). Statistics notes: Cronbach's alpha. *British Medical Journal* 314, 572.
- Bloomfield, L. (1984 [1933]). *Language*. Chicago: University of Chicago Press.
- Chafe, W. (1992). The importance of corpus linguistics to understanding the nature of language. In J. Svartvik (Ed.), *Directions in Corpus Linguistics*, pp. 79–97. Berlin, New York: Mouton de Gruyter.
- Chambers, J. K. and P. Trudgill (1998). *Dialectology* (2nd ed.). Cambridge, New York: Cambridge University Press.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika* 16(3), 297–334.
- Ellis, A. J. (1889). *The Existing Phonology Of English Dialects Compared With That Of West Saxon Speech*, Volume V of *On Early English Pronunciation*. London: Trübner & co.
- Embleton, S. (1993). Multidimensional scaling as a dialectometrical technique: Outline of a research project. In R. Köhler and B. Rieger (Eds.), *Contributions to quantitative linguistics*, pp. 267–276. Dordrecht: Kluwer.
- Giles, D. (2002). *Advanced research methods in psychology*. Hove, New York: Routledge.
- Goebel, H. (1982). *Dialektometrie: Prinzipien und Methoden des Einsatzes der Numerischen Taxonomie im Bereich der Dialektgeographie*. Wien: Österreichische Akademie der Wissenschaften.
- Goebel, H. (1984). *Dialektometrische Studien: Anhand italoromanischer, rätoromanischer und galloromanischer Sprachmaterialien aus AIS und ALF*. Tübingen: Niemeyer.
- Goebel, H. (2001). Arealtypologie und Dialektologie. In M. Haspelmath, E. König, W. Oesterreicher, and W. Raible (Eds.), *Language Typology and Language Universals / La typologie des langues et les universaux linguistiques / Sprachtypologie und sprachliche Universalien: An International Handbook / Manuel international / Ein internationales Handbuch*, Volume 2, pp. 1471–1491. Berlin, New York: Walter de Gruyter.

- Goebel, H. (2005). La dialectométrie corrélatrice. Un nouvel outil pour l'étude de l'aménagement dialectal de l'espace par l'homme. *Revue de linguistique romane* 69, 321–367.
- Goebel, H. (2006). Recent Advances in Salzburg Dialectometry. *Literary and Linguistic Computing* 21(4), 411–435.
- Goebel, H. (2007). A bunch of dialectometric flowers: a brief introduction to dialectometry. In U. Smit, S. Dollinger, J. Hüttner, G. Kaltenböck, and U. Lutzky (Eds.), *Tracing English through time: Explorations in language variation*, pp. 133–172. Wien: Braumüller.
- Goebel, H. (2008). Le Laboratoire de dialectométrie de l'Université de Salzbourg. *Zeitschrift für französische Sprache und Literatur* 118(1), 35–55.
- Goebel, H. and G. Schiltz (1997). A dialectometrical compilation of CLAE 1 and CLAE 2: Isoglosses and dialect integration. In W. Viereck and H. Ramisch (Eds.), *Computer developed linguistic atlas of England (CLAE)*, Volume 2, pp. 13–21. Tübingen: Max Niemeyer Verlag.
- Gooskens, C. (2005). Traveling time as a predictor of linguistic distance. *Dialectologia et Geolinguistica* 13, 38–62.
- Gooskens, C. and W. Heeringa (2004). Perceptive evaluation of Levenshtein dialect distance measurements using Norwegian dialect data. *Language Variation and Change* 16(3), 189–207.
- Greenbaum, S. (1996). *Comparing English worldwide: the International Corpus of English*. Oxford, New York: Clarendon Press.
- Gries, S. T. and S. Wulff (2005). Do foreign language learners also have constructions? Evidence from priming, sorting, and corpora. *Annual Review of Cognitive Linguistics* 3, 182–200.
- Grieve, J. (2009). *A Corpus-Based Regional Dialect Survey of Grammatical Variation in Written Standard American English*. Phd dissertation, Northern Arizona University.
- Haimerl, E. (2006). Database Design and Technical Solutions for the Management, Calculation, and Visualization of Dialect Mass Data. *Literary and Linguistic Computing* 21(4), 437–444.
- Haspelmath, M. (2009). Welche Fragen können wir mit herkömmlichen Daten beantworten? *Zeitschrift für Sprachwissenschaft* 28, 157–162.
- Heeringa, W. (2004). *Measuring dialect pronunciation differences using Levenshtein distance*. Phd dissertation, University of Groningen.
- Heeringa, W., K. Johnson, and C. Gooskens (2009). Measuring Norwegian dialect distances using acoustic features. *Speech Communication* 51(2), 167–183.
- Heeringa, W. and J. Nerbonne (2001). Dialect areas and dialect continua. *Language Variation and Change* 13(3), 375–400.
- Heeringa, W., J. Nerbonne, R. v. Bezooijen, and M. R. Spruit (2007). Geografie en inwoneraantallen als verklarende factoren voor variatie in het Nederlandse dialectgebied. *Nederlandse Taal- en Letterkunde* 123(1), 70–82.
- Hernández, N. (2006). *User's Guide to FRED*. URN: urn:nbn:de:bsz:25-opus-24895, URL: <http://www.freidok.uni-freiburg.de/volltexte/2489/>. Freiburg: University of Freiburg.

- Holman, E. W., C. Schulze, D. Stauffer, and S. Wichmann (2007). On the relation between structural diversity and geographical distance among languages: observations and computer simulations. *Linguistic Typology* 11, 393–421.
- Hoppenbrouwers, C. and G. Hoppenbrouwers (2001). *De indeling van de Nederlandse streektaalen. Dialecten van 156 steden en dorpen geklasseerd volgens de FFM*. Assen: Koninklijke Van Gorcum.
- Jain, A. K., M. N. Murty, and P. J. Flynn (1999). Data clustering: a review. *ACM Computing Surveys* 31(3), 264–323.
- Kortmann, B. and B. Szmrecsanyi (2004). Global synopsis: morphological and syntactic variation in English. In B. Kortmann, E. Schneider, K. Burridge, R. Mesthrie, and C. Upton (Eds.), *A Handbook of Varieties of English*, Volume 2, pp. 1142–1202. Berlin/New York: Mouton de Gruyter.
- Kruskal, J. B. and M. Wish (1978). *Multidimensional Scaling*, Volume 11 of *Quantitative Applications in the Social Sciences*. Newbury Park, London, New Delhi: Sage Publications.
- Leech, G. N., B. Francis, and X. Xu (1994). The use of computer corpora in the textual demonstrability of gradience in linguistic categories. In C. Fuchs and B. Victorri (Eds.), *Continuity in Linguistic Semantics*, pp. 57–76. Amsterdam, Philadelphia: Benjamins.
- Leinonen, T. (2008). Factor Analysis of Vowel Pronunciation in Swedish Dialects. *International Journal of Humanities and Arts Computing* 2(1-2), 189–204.
- Meyers, L. S., G. Gamst, and A. J. Guarino (2006). *Applied multivariate research: design and interpretation*. Thousand Oaks: Sage Publications.
- Nerbonne, J. (2006). Identifying Linguistic Structure in Aggregate Comparison. *Literary and Linguistic Computing* 21(4), 463–475.
- Nerbonne, J. (2009). Data-driven dialectology. *Language and Linguistics Compass* 3(1), 175–198.
- Nerbonne, J. and W. Heeringa (1997). Measuring Dialect Distance Phonetically. In J. Coleman (Ed.), *Workshop on Computational Phonology, Special Interest Group of the Association for Computational Linguistics*, pp. 11–18. Madrid.
- Nerbonne, J. and W. Heeringa (2007). Geographic distributions of linguistic variation reflect dynamics of differentiation. In S. Featherston and W. Sternefeld (Eds.), *Roots: Linguistics in Search of its Evidential Base*, pp. 267–297. Berlin, New York: Mouton de Gruyter.
- Nerbonne, J., W. Heeringa, and P. Kleiweg (1999). Edit Distance and Dialect Proximity. In D. Sankoff and J. Kruskal (Eds.), *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*, pp. v–xv. Stanford: CSLI Press.
- Nerbonne, J., W. Heeringa, E. van den Hout, P. van de Kooi, S. Otten, and W. van de Vis (1996). Phonetic Distance between Dutch Dialects. In G. Durieux, W. Daelemans, and S. Gillis (Eds.), *CLIN VI: Proc. of the Sixth CLIN Meeting*, pp. 185–202. Antwerp: Centre for Dutch Language and Speech (UIA).
- Nerbonne, J. and P. Kleiweg (2007). Toward a Dialectological Yardstick. *Journal of Quantitative Linguistics* 14(2), 148–166.
- Nerbonne, J., P. Kleiweg, and F. Manni (2008). Projecting dialect differences to geography: bootstrapping clustering vs. clustering with noise. In C. Preisach, L. Schmidt-Thieme, H. Burkhardt, and R. Decker (Eds.), *Data Analysis, Machine Learning, and Applications. Proceedings of the 31st Annual Meeting of the German Classification Society*, pp. 647–654. Berlin: Springer.

- Nishisato, S. (2007). *Multidimensional nonlinear descriptive analysis*. Boca Raton: Chapman & Hall/CRC.
- Nunnally, J. C. (1978). *Psychometric Theory*, Volume New York. McGraw-Hill.
- Orton, H., S. Sanderson, and J. D. A. Widdowson (1978). *The Linguistic Atlas of England*. London, Atlantic Highlands, N.J.: Croom Helm.
- Penhallurick, R. (1993). Welsh English: a national language? *Dialectologia et Geolinguistica 1*, 28–46.
- Penhallurick, R. (2004). Welsh English: phonology. In B. Kortmann, E. Schneider, K. Burridge, R. Mesthrie, and C. Upton (Eds.), *A Handbook of Varieties of English*, Volume 1, pp. 98–112. Berlin/New York: Mouton de Gruyter.
- Penke, M. and A. Rosenbach (2007). What counts as evidence in linguistics? An introduction. *Studies in Language 28*(3), 480–526.
- Séguy, J. (1971). La relation entre la distance spatiale et la distance lexicale. *Revue de Linguistique Romane 35*, 335–357.
- Shackleton, R. G. J. (2007). Phonetic variation in the traditional English dialects: a computational analysis. *Journal of English Linguistics 35*(1), 30–102.
- Speelman, D. and D. Geeraerts (2008). The role of concept characteristics in lexical dialectometry. *International Journal of Humanities and Arts Computing 2*(1–2), 221–242.
- Spruit, M. R. (2005). Classifying Dutch dialects using a syntactic measure: the perceptual Daan and Blok dialect map revisited. *Linguistics in the Netherlands 22*(1), 179–190.
- Spruit, M. R. (2006). Measuring syntactic variation in Dutch dialects. *Literary and Linguistic Computing 21*(4), 493–506.
- Spruit, M. R., W. Heeringa, and J. Nerbonne (2009). Associations among Linguistic Levels. *Lingua 119*(11), 1624–1642.
- Szmrecsanyi, B. (2008). Corpus-based dialectometry: aggregate morphosyntactic variability in British English dialects. *International Journal of Humanities and Arts Computing 2*(1–2), 279–296.
- Szmrecsanyi, B. (2010a). Geography is overrated. In S. Hansen, C. Schwarz, P. Stoeckle, and T. Streck (Eds.), *Dialectological and folk dialectological concepts of space*. Berlin: Walter de Gruyter.
- Szmrecsanyi, B. (2010b). *The morphosyntax of BrE dialects in a corpus-based dialectometrical perspective: feature extraction, coding protocols, projections to geography, summary statistics*. URN: urn:nbn:de:bsz:25-opus-73209, URL: <http://www.freidok.uni-freiburg.de/volltexte/7320/>. Freiburg: University of Freiburg.
- Szmrecsanyi, B. and N. Hernández (2007). *Manual of Information to accompany the Freiburg Corpus of English Dialects Sampler ("FRED-S")*. URN: urn:nbn:de:bsz:25-opus28598, URL: <http://www.freidok.uni-freiburg.de/volltexte/2859/>. Freiburg: University of Freiburg.
- Trudgill, P. (1974). Linguistic change and diffusion: description and explanation in sociolinguistic dialect geography. *Language in Society 2*, 215–246.

- Trudgill, P. (1999). *The dialects of England* (2nd ed.). Cambridge, MA, Oxford: Blackwell.
- Viereck, W. (1985). Linguistic atlases and dialectometry: The survey of English dialects. In J. M. Kirk, S. Sanderson, and J. D. A. Widdowson (Eds.), *Studies in linguistic geography: The dialects of English in Britain and Ireland*, pp. 94–112. London: Croom Helm.
- Viereck, W., H. Ramisch, H. Händler, P. Hoffmann, and W. Putschke (1991). *The computer developed linguistic atlas of England*. Tübingen: Niemeyer.
- Voronoi, G. (1907). Nouvelles applications des paramètres continus à la théorie des formes quadratiques. *Journal für die Reine und Angewandte Mathematik* 133, 97–178.
- Wälchli, B. (2009). Data reduction typology and the bimodal distribution bias. *Linguistic Typology* 13(1), 77–94.
- Ward, J. H. J. (1963). Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association* 58, 236–244.