

# Probabilistic corpus-based dialectometry

Q1 4 Christoph Wolk,<sup>1\*</sup> and Benedikt Szmrecsanyi<sup>2</sup>

Q2 5 <sup>1</sup> University of Giessen

6 <sup>2</sup> KU Leuven

7 Researchers in dialectometry have begun to explore measurements based on fundamentally quantitative metrics, often  
8 sourced from dialect corpora, as an alternative to the traditional signals derived from dialect atlases. This change of data  
9 type amplifies an existing issue in the classical paradigm, namely that locations may vary in coverage and that this affects  
10 the distance measurements: pairs involving a location with lower coverage suffer from greater noise and therefore  
11 imprecision. We propose a method for increasing robustness using generalized additive modeling, a statistical technique  
12 that allows leveraging the spatial arrangement of the data. The technique is applied to data from the British English dialect  
13 corpus FRED; the results are evaluated regarding their interpretability and according to several quantitative metrics. We  
14 conclude that data availability is an influential covariate in corpus-based dialectometry and beyond, and recommend that  
15 researchers be aware of this issue and of methods to alleviate it.

## 16 1. INTRODUCTION

17 In this paper, we discuss new technologies to measure  
18 dialect distances on the basis of dialect corpora that  
19 sample authentic, naturalistic usage data. More speci-  
20 fically, we show how probabilistic modeling techniques  
21 can be used to address a thorny challenge in corpus-  
22 based dialectometry: the amount of text and speech that  
23 is available to cover particular dialect locations is typi-  
24 cally not constant but variable, and this variability can  
25 seriously confound unless it is neutralized by using the  
26 mathematics of uncertainty.

27 Let us fix some basic concepts and methodological  
28 preliminaries at the outset. Practitioners of traditional  
29 dialectology study “interesting” dialect phenomena,  
30 one feature at a time, often only in a handful of dialects;  
31 cross-feature comparisons remain rather impres-  
32 sionistic through the bundling and/or grading of iso-  
33 glosses and the resulting areal classifications. In  
34 contrast, DIALECTOMETRY is the branch of geolinguistics  
35 dedicated to measuring, visualizing, and analyzing  
36 aggregate dialect similarities or distances as a function  
37 of properties of geographic space (see Goebel 1984, 2007;  
38 Heeringa, 2004; Nerbonne, 2009; Nerbonne, Heeringa &  
39 Kleiweg, 1999; Séguy, 1971 for foundational work).

40 As for DATA SOURCES, traditional dialectometry draws  
41 on dialect atlases as its data source. Take, for example,  
42 Goebel (1982), who investigates similarities between  
43 Italian dialects on the basis of 696 linguistic features that  
44 are mapped in the *Sprach- und Sachatlas Italiens und der  
Südschweiz* (AIS), an atlas that covers Italy and

southern Switzerland. It is important to bear in mind  
that the data that (most) linguistic atlases provide speak  
primarily to the issue of what informants KNOW about  
their dialects. To study LANGUAGE USAGE analysts typi-  
cally turn not to surveys and dialect atlases but to  
DIALECT CORPORA. Linguistic corpora are principled and  
broadly representative collections of naturalistic texts or  
speech that sample usage data—a data type that is  
increasingly popular in dialectology (Anderwald &  
Szmrecsanyi, 2009; Grieve, 2016) and beyond (see the  
papers in Szmrecsanyi & Wälchli, 2014).

CORPUS-BASED DIALECTOMETRY (henceforth: CBDM),  
then, combines the study of dialectometric research  
questions with corpus-linguistic methodologies. CBDM  
utilizes aggregation methodologies to explore quanti-  
tative and distributional usage patterns extracted from  
dialect corpora (see Szmrecsanyi, 2008, 2011, 2013;  
Szmrecsanyi & Wolk, 2011; Wolk, 2014; Wolk &  
Szmrecsanyi, 2016). Turning to corpora enables analysts  
to address questions about usage versus knowledge,  
production/comprehension versus intuition, chaos  
versus orderliness, and so on.

In this contribution, we discuss a recent methodolo-  
gical advance in CBDM. To do justice to the fact that  
textual coverage of individual dialects may (and typi-  
cally does) vary in dialect corpora (e.g. dialect A may be  
represented by 20 interviews, but dialect B by only 5  
interviews), the first wave of CBDM approaches merely  
normalized text frequencies prior to aggregation: so  
instead of saying that a particular linguistic variant  
occurred, e.g., 100 times in total in material from some  
particular dialect, a normalized measure (“the linguistic  
variant occurs 30 times per 10,000 words of running  
text”) was used to calculate dialect distances. The

\*Address for correspondence: Otto-Behaghel-Straße 10 B 403, 35394  
Giessen, Germany, +49 641 - 99 301 53, christoph.b.wolk@anglistik.  
uni-giessen.de

79 innovation we discuss in the present paper introduces  
 80 statistical modeling into the CBDM pipeline: subcorpus  
 81 size turns out to be a crucial covariate of aggregate  
 82 distances, and so probabilistic CBDM draws on sto-  
 83 chastic reasoning to take this covariate more seriously  
 84 than normalization-based CBDM does. The outcome, as  
 85 we shall see, is a less noisy and arguably more accurate  
 86 portrayal of aggregate dialect relationships.

87 This paper is structured as follows: Section 2 moti-  
 Q3 88 vates our approach by discussing the role of data  
 89 availability in dialectometric analyses, with a particular  
 90 focus on frequency- and corpus-based dialectometry.  
 91 We will also introduce some of the technicalities of our  
 92 particular solution to the challenges this factor provides.  
 93 The remainder of the paper will work in detail through a  
 94 case study concerning British dialects. Section 3 will  
 95 introduce the data set, while sections 4 and 5 will pre-  
 96 sent two related techniques, straightforward frequency-  
 97 based CBDM and the probabilistically enhanced ver-  
 98 sion, as well as their results on the present dataset. The  
 99 final section will conclude by reviewing and discussing  
 100 these outcomes in light of the discussion in section 2.

## 101 2. ON THE INFLUENCE OF DATA AVAILABILITY 102 IN DIALECTOMETRY

103 The principal question underlying the CBDM approach  
 104 is the following: How can we derive accurate repre-  
 105 sentations of linguistic divergence from naturally  
 Q4 106 occurring discourse?<sup>2</sup> This is, in principle, not too dif-  
 107 ferent from questions regularly asked in atlas-based  
 108 dialectometry, which has greatly benefited from a long  
 109 methodological discussion and a varied set of techni-  
 110 ques (e.g., Heeringa, 2004). In general, a dialectometric  
 111 analysis proceeds as follows:

- 112 1. Establish a feature set (lexical items, pronunciations  
 113 etc) by which the dialects are to be compared.
- 114 2. Determine the realizations of these features in each  
 115 location.
- 116 3. Compare all features in all locations and derive a  
 117 numerical value indicating the degree of  
 118 dissimilarity.
- 119 4. Aggregate over all features to yield a composite  
 120 score, or distance, for each pair of locations.
- 121 5. Analyze the resulting scores using exploratory and/  
 122 or confirmatory data analysis.

123 Steps 3–5 in particular have seen considerable  
 124 advancement and extension. Steps 1 and 2 typically rely  
 125 on the results of large survey projects, compiled into  
 126 dialect atlases. Atlas data are in many ways well-suited  
 127 for such analyses; for instance, the network of locations  
 128 tends to have a rather fine mesh, which allows a geo-  
 129 graphically high resolution. Particularly crucial for  
 130 present purposes is that the amount of data per location

tends to be quite homogeneous: sites usually have 131  
 similar amounts of informants, and ideally each infor- 132  
 mant has a complete set of responses to the survey 133  
 items. In practice, atlas data are not always complete. 134  
 The problems that missing entries can cause, and how 135  
 to mitigate them, have been recognized as important 136  
 issues almost from the inception of dialectometry. 137  
 Goebel (1977: 46) discusses the issue and employs a 138  
 method that has since become the standard treatment 139  
 (see also, e.g., Nerbonne & Kleiweg, 2007: 159): features 140  
 which are missing for one or both locations in a pair are 141  
 to be completely left out of the analysis for that pair, 142  
 continuing as if they were never in the feature set. 143  
 While this introduces noise into the measurements— 144  
 Goebel (1993: 286) terms this the *missing data effect*—in 145  
 general, it has not led to major problems for dialecto- 146  
 metry. The well-documented increase in robustness 147  
 that dialectometry achieves through aggregation seems 148  
 sufficient to cancel out minor noise resulting from this 149  
 (Nerbonne, 2009). Nevertheless, other approaches have 150  
 been suggested by Viereck (1988: 546): missing entries 151  
 could be included, estimated based on the closest 152  
 neighbor, or an aggregate of neighbors. Taking geo- 153  
 graphic information into account like this can reduce 154  
 noise, and therefore diminish the missing data effect. 155  
 The cost of this is, of course, that the additional 156  
 assumption may not be warranted—that a location may 157  
 well be unlike its neighbors with regard to the missing 158  
 feature, and the resulting analysis is biased. This mat- 159  
 ches a central trade-off in statistical learning, the “bias- 160  
 variance trade-off” (James, Witten, Hastie & Tibshirani, 161  
 2013). Statistical learning—such as learning the dis- 162  
 tances between dialects from samples of individual 163  
 informants’ judgments and productions—depends 164  
 both on the data set itself and on the specific char- 165  
 acteristics and implicit assumptions inherent in the 166  
 method. As one increases the flexibility of a method, 167  
 incidental properties of the data (such as missing entries 168  
 at certain locations) will have greater weight; lower 169  
 flexibility, however, leads to greater dominance of the 170  
 assumptions that the method makes. Both yield the 171  
 danger of distorted and poorly generalizable results. 172  
 For categorical atlas-based signals, increased flexibility 173  
 may often be more desirable, as incidental distortions 174  
 should even out in the long run. 175

176 One recent extension of the dialectometric tool box  
 177 involves adapting the methods to data of a different  
 178 nature: instead of relying on (typically questionnaire-  
 179 based) categorical surveys filtered through dialect  
 180 atlases, authentic speech is tapped directly. This is made  
 181 possible through the emergence of specialized dialect  
 182 corpora, i.e., naturally occurring linguistic material,  
 183 typically whole interviews, collected from dialectologi-  
 184 cally appropriate speakers. We believe this to be a cen-  
 185 tral locus for the advancement of the dialectometric

186 project: it opens new possibilities for frequency-based  
 187 analyses that integrate well with the kind of usage-based  
 188 approaches that are gaining ground in many branches of  
 189 linguistics, and allows tackling research questions that  
 190 crucially depend on frequency. Nevertheless, this  
 191 change of data source necessitates a retooling of the  
 192 usual approaches in dialectometry. As we shall demon-  
 193 strate, for typical dialect corpora, the importance of the  
 194 *data availability* factor greatly increases, and relying only  
 195 on aggregation may lead to wrong inferences.

196 How can one establish comparability of frequency  
 197 measurements between corpora, and therefore loca-  
 198 tions? In a realistic scenario, it is highly likely that the  
 199 two corpora are unequal in size, be it due to availability  
 200 of raw material (because, e.g., more interviews happen  
 201 to have been conducted in location A than in location B)  
 202 or corpus design. But when one corpus is substantially  
 203 larger than the other, similar counts do not imply simi-  
 204 lar usage rates. The usual solution in corpus linguistics  
 205 involves *normalization*, i.e., transforming raw counts into  
 206 occurrence rates. The process is straightforward:

$$(\text{raw counts} / \text{text size}) \times \text{normalization constant}$$

207 Hence, the total number of occurrences in a (sub)corpus  
 208 is divided by the number of words, and, to make the  
 209 numbers more easily interpretable, the resulting figure  
 210 is scaled by a fixed number, yielding the number of  
 211 occurrences per, say, ten thousand words (*pttw*). As  
 212 these normalized values have a common basis, they can  
 213 be compared with one another and their difference can  
 214 be quantified.  
 215

216 There are, as we shall show, cases in which the nor-  
 217 malization procedure may lead to biased results, and  
 218 these bear a strong resemblance to the missing data  
 219 effect. While the process will always yield a numerical  
 220 value of the difference between the corpora, the accu-  
 221 racy of this value crucially depends on the text size prior  
 222 to normalization, and in particular on that of the smaller  
 223 corpus. To illustrate this, consider a hypothetical feature  
 224 with a (population) frequency value of 1 *pttw* in two  
 225 communities; from the first, we sample a corpus of one  
 226 hundred thousand words (C1), from the second a cor-  
 227 pus of only ten thousand words (C2). We should expect  
 228 the normalized value for C1 to approximate the popu-  
 229 lation value relatively well. C2, however, may well be  
 230 far from the true value—it would not be surprising to  
 231 see the feature completely absent, or have a normalized  
 232 frequency two to three times as high as in C1. The dis-  
 233 tance between the corpora is quite likely to be sub-  
 234 stantially higher than the difference between the  
 235 populations—exactly zero. On the other hand, C1 will  
 236 likely be more similar to a population where there are  
 237 actual frequency differences; higher or lower depend-  
 238 ing on the accident of chance. Note that this is a prop-  
 239 erty of the corpora and their sizes; when moving

240 toward combining such individual measurements into  
 241 a multi-feature dialectometric analysis, this will apply  
 242 to all of them individually. It follows that the distances  
 243 between similar points are likely to be too high, and the  
 244 distances between dissimilar points may be too high or  
 245 even too low. Goebel (1993) reports that the missing data  
 246 effect for categorical data, with missing entries left out  
 247 of the analysis, entails “measurement results which are  
 248 too high [i.e., similar, as Goebel uses similarity instead of  
 249 distance] in comparison to the general trend” (286). In  
 250 frequency-based analyses, absence corresponds to a  
 251 frequency of zero and the feature is not left out of the  
 252 analysis. This suggests that measurements affected by  
 253 such issues would be prone to be more *dissimilar* than  
 254 the trend suggests at close locations, but may be too  
 255 similar at more distant locations.

256 As we shall show, this is not just a hypothetical issue,  
 257 but has direct implications for dialectometric practice.  
 258 While we will be focusing on frequencies, similar con-  
 259 siderations apply to proportions of categorical alter-  
 260 nants and possibly even categorical atlas realizations.  
 261 Streck (2014), for example, reports that his corpus-  
 262 oriented analysis of phonological variation in south-  
 263 western German yielded a substantially stronger rela-  
 264 tionship between geographic distance and linguistic  
 265 similarity after removing the locations with the least  
 266 amount of data. In principle, only datasets that are  
 267 complete and large are fully safe, although minor  
 268 amounts of size discrepancies or empty cells need not  
 269 display any adverse effect.

270 How can the effect that the amount of data has on the  
 271 result be mitigated in corpus-based analyses? Similar  
 272 solutions apply as for the missing data approaches  
 273 described above, but with a crucial difference: instead of  
 274 a categorical presence signal we have a gradient quality  
 275 signal applying to all features at the same time, which  
 276 makes the consensus solution for atlas-based dialecto-  
 277 metry unavailable—at least without dropping the loca-  
 278 tion completely. One way around the problem is to just  
 279 accept it, to not “fill in missing data artificially” (Goebel,  
 280 1991: 283). The benefit here is that every step of the  
 281 analysis is purely based on actually observed data, and  
 282 few additional assumptions are necessary. This comes at  
 283 a cost: the usual visualizations and maps hide the fact  
 284 that some measurements are more imprecise than oth-  
 285 ers, and various statistical results become difficult to  
 286 interpret. A second possibility involves restricting the  
 287 analysis to those locations where there is ample data.  
 288 This seems statistically reliable, but will typically result  
 289 in substantial reductions in geographic coverage. The  
 290 third possibility is to use a larger corpus. Our data sug-  
 291 gest that effects of differences in the amount of text is still  
 292 detectable even as the number of words in both corpora  
 293 increases (see section 6 below), but it is clear that having  
 294 more data leads to more precise measurements even

295 with simple normalization. This is clearly the best solu-  
 296 tion, but is in general not feasible for spoken dialect  
 297 data, due to difficulties in acquiring material and the  
 298 costs of transcription. For studying geographic varia-  
 299 tion in written material, such as letters to the editor in  
 300 regional newspapers (Grieve, 2016) or Twitter data  
 301 (Eisenstein, 2018; Huang, Guo, Kasakoff & Grieve,  
 302 2016), where corpus sizes can reach dozens of millions  
 303 or even billions of words, normalized values alone may  
 304 suffice fully. The final possibility, and one that is  
 305 applicable to the relatively small spoken dialect cor-  
 306 pora, is to use a method that is less sensitive to the  
 307 incidental variance in the data through the use of  
 308 additional assumptions. Given the geographic nature of  
 309 dialectal data, the best candidates for such analyses are  
 310 those that build on the fundamental dialectological  
 311 postulate that “geographically proximate varieties tend  
 312 to be more similar than distant ones” (Nerbonne &  
 313 Kleiweg, 2007; see also Tobler’s (1970) first law of geo-  
 314 graphy: “Everything is related to everything else, but  
 315 near things are more related than distant things.”).<sup>3</sup> The  
 316 observed values should be considered in their spatial  
 317 context, for example by discounting large differences  
 318 between close locations if they are based on little data.  
 319 Of course, this correction should not be too strong, and  
 320 the method should be able to still find differences  
 321 between close locations if this is warranted.

322 Within the frequency- and corpus-based dialectolo-  
 323 gical literature, several methods have been proposed  
 324 that can achieve this (e.g., Grieve, 2009; Pickl, Spettl,  
 325 Pröll, Elspaß, König & Schmidt, 2014). We believe that  
 326 *generalized additive models* (GAMs) have particularly nice  
 327 properties for this purpose (Wood, 2006). GAMs share  
 328 conceptual similarities with generalized linear models,  
 329 which are familiar to many linguists (e.g., in the form of  
 330 VARBRUL models, cf. Sankoff, 1987). GAMs extend  
 331 this by including smooth terms, smooth functions  
 332 whose shape is determined from the data. They can be  
 333 two-dimensional, and are therefore suitable to repre-  
 334 sent dialectal “frequency landscapes” that represent the  
 335 geographic distribution of features as mountains and  
 336 valleys of high and low usage. GAMs are not new in  
 337 geolinguistics; they were previously used successfully  
 338 for dialectometric purposes in Wieling, Nerbonne &  
 339 Baayen (2011) and Wieling, Montemagni, Nerbonne &  
 340 Baayen (2013). Our use differs quite substantially from  
 341 theirs: instead of building a single model based on the  
 342 linguistic distances of all points to a reference variety,  
 343 we build individual models for each feature. This model  
 344 represents an improved estimate of how frequent that  
 345 feature is in each subcorpus, and is itself suitable for  
 346 visualization and interpretation. An example can be  
 347 seen in Figure 1, which displays the smooth term for  
 348 one of the features that we include in our analysis:  
 349 multiple negation, as in (1)

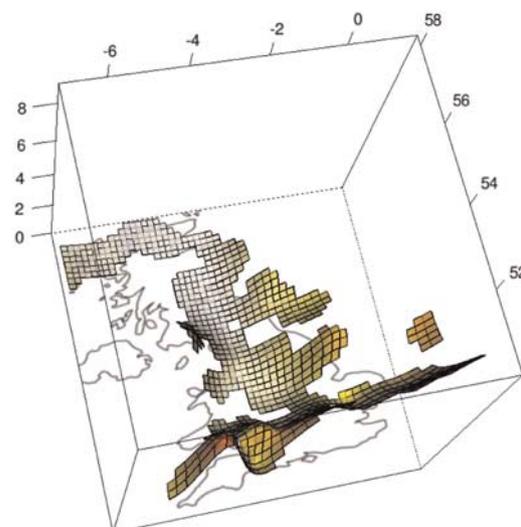


Figure 1. Frequency landscape for feature [33], multiple negation.

(1) ‘cause you dare not say nothing ... <LND\_001 > 350

The GAM estimates a frequency of around 8 *pttw* in 351  
 southern England, which drops as one moves north. 352  
 The Isle of Man also shows a higher rate of usage. 353

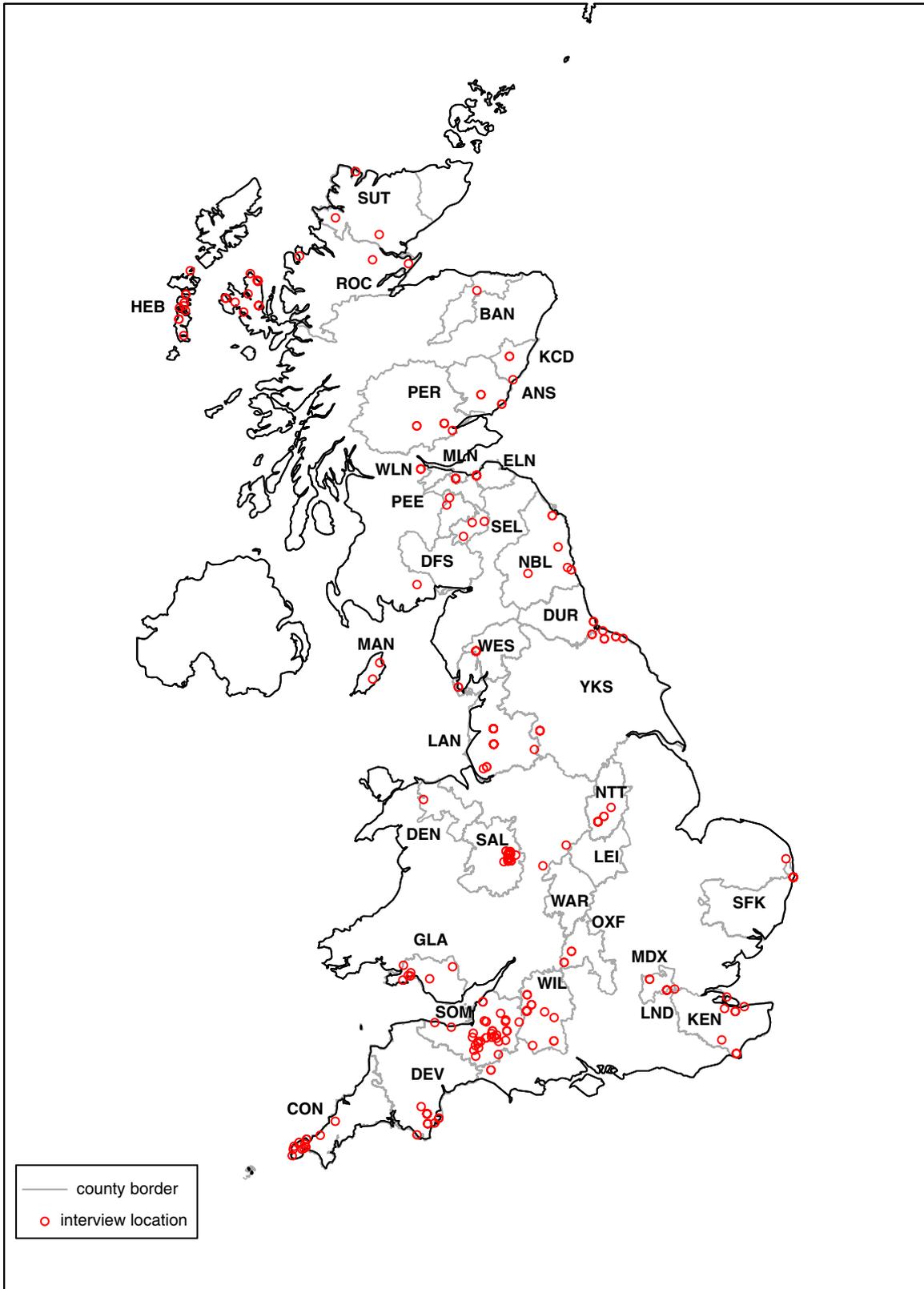
As is apparent from the plot, such landscapes can be 354  
 quite complex. Nevertheless, the GAM implementation 355  
 used here limits this complexity by means of general- 356  
 ized cross-validation: the effect of leaving out indivi- 357  
 dual data points is determined by analyzing subsets of 358  
 the data, and the final result is chosen such that single 359  
 points do not have undue influence. This yields a 360  
 landscape that can be steeply sloped when this is war- 361  
 ranted based on the data, but where the pattern tends to 362  
 be flat when it is not. 363

It is important to note that such models have the 364  
 capabilities of generalized linear models, and can 365  
 therefore include speaker- and/or text-oriented covari- 366  
 ates, such as speaker age or some measure of con- 367  
 versational interactivity. We make only limited use of 368  
 this capability in the present paper, but believe that it is 369  
 a central advantage of this method. What is crucial, 370  
 however, is that both normalization and modeling both 371  
 yield numerical values representing the same thing: 372  
 estimates of frequency. Therefore, both can be analyzed 373  
 in exactly the same way, and are easily compared and 374  
 contrasted. It is to this task that we turn next. 375

### 3. DATA 376

#### 3.1. Corpus Database 377

This case study taps into FRED, the *Freiburg Corpus of* 378  
*English Dialects* (Hernández, 2006; Szmrecsanyi & Her- 379  
 nández, 2007), a major dialect corpus that covers tradi- 380  
 tional dialect speech (mainly transcribed so-called “oral 381



Map 1. Counties and interview locations that contributed data to the corpus as used in the analyses to follow.

**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), textual coverage (running words) in FRED.

map label	county	a-priori dialect area	no. words sampled in FRED
ANS	Angus	Sc Lowlands	19,899
BAN	Banffshire	Sc Lowlands	5,655
CON	Cornwall	Southwest of E	107,072
DEN	Denbighshire	Wales	5,794
DEV	Devon	Southwest of E	97,080
DFS	Dumfriesshire	Sc Lowlands	9997
DUR	Durham	North of E	28,069
ELN	East Lothian	Sc Lowlands	40,190
GLA	Glamorganshire	Wales	53,110
HEB	Hebrides	Hebrides	72,761
MAN	Isle of Man	Isle of Man	10,930
KCD	Kincardineshire	Sc Lowlands	7,509
KEN	Kent	Southeast of E	176,908
LAN	Lancashire	North of E	205,342
LEI	Leicestershire	E Midlands	5,864
LND	London	Southeast of E	110,802
MDX	Middlesex	Southeast of E	31,794
MLN	Midlothian	Sc Lowlands	32,040
NBL	Northumberland	North of E	30,644
NTT	Nottinghamshire	E Midlands	150,810
OXF	Oxfordshire	Southwest of E	15,109
PEE	Peebleshire	Sc Lowlands	14,955
PER	Perthshire	Sc Lowlands	20,896
ROC	Ross and Cromarty	Sc Highlands	10,475
SAL	Shropshire	E Midlands	168,572
SEL	Selkirkshire	Sc Lowlands	9,325
SFK	Suffolk	Southeast of E	312,560
SOM	Somerset	Southwest of E	207,502
SUT	Sutherland	Sc Highlands	10,967
WAR	Warwickshire	E Midlands	8,269
WES	Westmorland	North of E	157,562
WIL	Wiltshire	Southwest of E	185,928
WLN	West Lothian	Sc Lowlands	18,418
YKS	Yorkshire	North of E	90,816

382 history" material) all over England, Scotland, and  
 383 Wales.<sup>5</sup> The vast majority of the interviews were recor-  
 384 ded between 1970 and 1990. In most cases, a fieldwor-  
 385 ker interviewed an informant about life, work, etc., in  
 386 former days. The informants sampled in the corpus are  
 387 typically elderly people with a working-class back-  
 388 ground, so-called 'non-mobile old rural males' (Cham-  
 389 bers & Trudgill, 1998, 29).

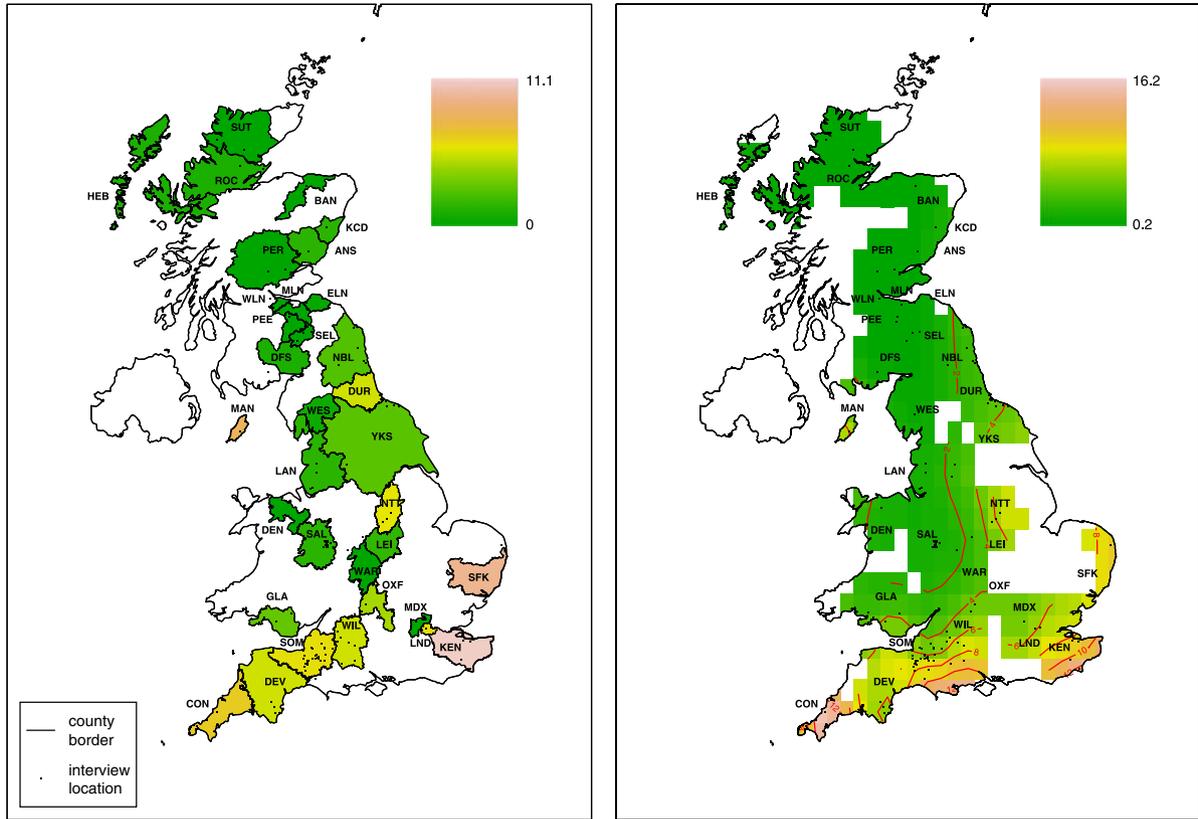
390 The version of FRED we use here consists of inter-  
 391 views with 376 informants and spans approximately 2.4  
 392 million words of running text. The interviews were  
 393 conducted in 34 pre-1974 counties in Great Britain,

including the Isle of Man and the Hebrides. To mitigate  
 data sparsity, the level of areal granularity investigated  
 in the present study will be the level of individual  
 counties. Map 1 displays the county boundaries and  
 interview locations. See Wolk (2014, chapter 3) for an in-  
 depth discussion of the dataset. (Table 1)

### 3.2. Features

The analysis in this paper is based on the feature set  
 used in Szmrecsanyi (2013). The set comprises 57 fea-  
 tures, and overlaps with, but is not identical to, the  
 comparative morphosyntax survey in Kortmann &  
 Szmrecsanyi (2004) and the battery of morphosyntax  
 features covered in the Survey of English Dialects  
 (Orton, Sanderson & Widdowson 1978; Viereck,  
 Ramisch, Händler, Hoffmann & Putschke, 1991). The  
 features in the catalog fall into eleven major gramma-  
 tical domains: (i) pronouns and determiners (e.g., non-  
 standard reflexives, as in (2)), (ii) the noun phrase (e.g.,  
 the *s*-genitive, as in (3)), (iii) primary verbs (e.g., the verb  
*to do*, as in (4)), (iv) tense & aspect (e.g., the present  
 perfect with auxiliary *be*, as in (5)), (v) modality (e.g.,  
 epistemic/deontic *must*, as in (6)), (vi) verb morphology  
 (e.g., non-standard weak past tense and past participle  
 forms, as in (7)), (vii) negation (e.g., *never* as a preverbal  
 past tense negator, as in (8)), (viii) agreement (e.g., non-  
 standard *was* as in (9)), (ix) relativization (e.g., the rela-  
 tive particle *what*, as in (10)), (x) complementation (e.g.,  
 unsplit *for to*, as in (11)), and (xi) word order and dis-  
 course phenomena (e.g., lack of auxiliaries in *yes/no*  
 questions, as in (12)). We cannot discuss the features  
 in much detail here, but the Appendix provides the  
 complete list of features. See Szmrecsanyi (2013, chap-  
 ter 3) for guidelines regarding feature selection and  
 Szmrecsanyi (2010) for a detailed description of the  
 feature extraction procedure.

- (2) But old Silvain, he used to look after hisself, really  
<CON\_003 >
- (3) But his wife is dead my brother Quentin's wife is  
dead two years ago <GLA\_002 >
- (4) I don't know <CON\_001 >
- (5) Joe, if you weigh them up, and you 're got an odd  
britch, I could do with a pair o' them <SFK\_038 >
- (6) [...] we must have been tough nuts you know,  
really. <LAN\_009 >
- (7) Oh, but you wouldnae be telled the wages.  
<BAN\_001 >
- (8) [...] and they never moved no more, neither one of  
them, never tried to <CON\_005 >
- (9) they was half gypsies you see? <OXF\_001 >
- (10) See that up on the top there, the stamp what you  
hammer in... <WIL\_024 >
- (11) For to screw down the cover on the churn  
<CON\_002 >



Map 2. Normalized frequency (left) and model-based frequency prediction (right) for multiple negation [33]. Yellow colors indicate areas where high frequencies are observed or predicted, green colors indicates low (observed or predicted) frequencies.



447 (12) They didn't want him prosecuted? <DEV\_001 >

448 **4. THE NORMALIZATION-BASED CBDM**  
 449 **APPROACH**

450 How does the normalization-based CBDM approach  
 451 without probabilistic enhancements following Szmrec-  
 452 sanyi (2013) study dialect relationships as a function of  
 453 geographic space? We first determine the text frequency  
 454 of the features in the corpus material: how often do we  
 455 find particular features—say, multiple negation—in  
 456 interviews from particular locations. Next, we normal-  
 457 ize text frequencies to frequency per 10,000 words  
 458 because textual coverage of individual dialects varies,  
 459 and round the result to whole numbers. At this stage,  
 460 we also perform a log-transformation, which is a cus-  
 461 tomary method to de-emphasize large frequency dif-  
 462 ferentials and to alleviate the effect of frequency outliers  
 463 (Shackleton, 2007, 43), thus increasing reliability of the  
 464 measurements. For features that are absent from a cor-  
 465 pus, the value is set to -1, corresponding to a frequency  
 466 of 0.1 *pttw*, as the logarithmic transformation requires  
 467 its input to be larger than 0. Let us illustrate the proce-  
 468 dure: in FRED, the county Cornwall has a textual cover-  
 469 age of 12 interviews totaling about 107,000 words of

running text (interviewer utterances excluded). In this  
 material, feature [34] (negative contraction, as in (13))  
 occurs 326 times, which translates into a normalized  
 text frequency of  $326 \times 10,000 / 107,000 \approx 30$  occurrences  
 per ten thousand words.

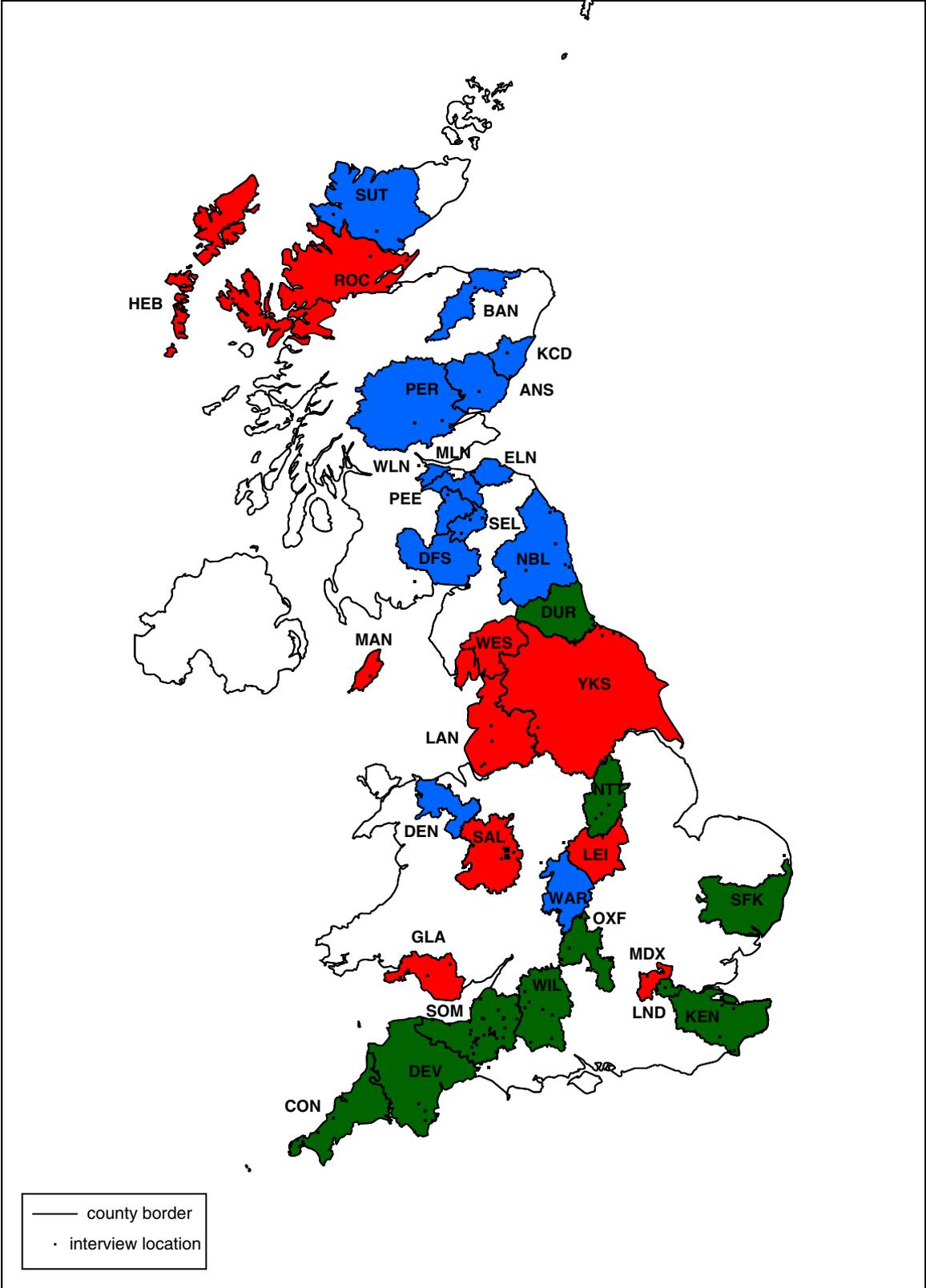
(13) They won't do anything. <WES\_011 >

A *log*-transformation of this frequency yields a value of  
 $\log_{10}(30) \approx 1.5$ . This is the measurement that characterizes  
 this specific measuring point (Cornwall) in regard to  
 feature [34].

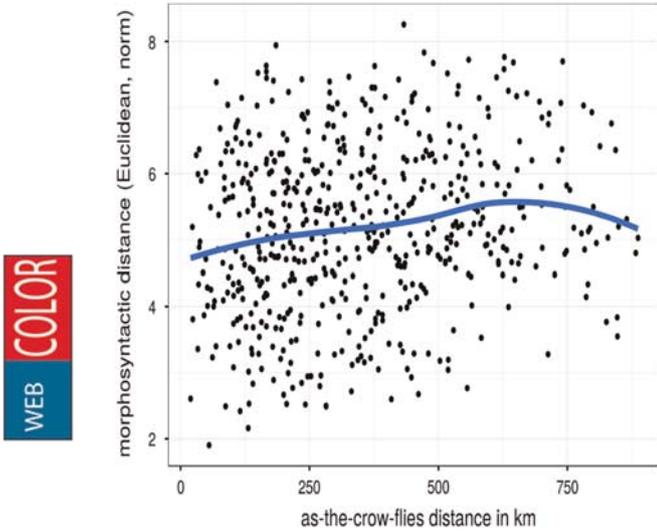
In the next step, we create an  $N \times p$  frequency matrix,  
 in which the  $N = 34$  objects (that is, dialects) are arran-  
 ged in rows and the  $p = 57$  features in columns, such  
 that each cell in the matrix specifies a particular (nor-  
 malized and log-transformed) feature frequency. Our  
 case study thus yields a  $34 \times 57$  frequency matrix: 34  
 British English dialects, each characterized by a vector  
 of 57 text frequencies. To illustrate, Map 2 (left) projects  
 feature frequencies of feature [33] (multiple negation) to  
 geography. (Note: parallel maps for the other 56 fea-  
 tures in the catalog are available in the online  
 appendix.)

Frequency matrices can serve as input to a number of  
 multivariate analysis techniques, such as Principal  
 Component analysis or Factor Analysis (see, e.g.,

470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494



Map 3. Cluster map based on the normalization-based CBDM approach. Displayed: 3-cluster solution.



**Figure 2.** Correlating normalization-based morphosyntactic distances with as-the-crow-flies distances ( $r = .19$ ,  $p < .001$ , logarithmic  $R^2 = 3.6\%$ ).

495 Grieve, 2014). Most aggregational procedures cus-  
 496 tomary in dialectometry, however, are empirically  
 497 based on so-called distance matrices, which are  
 498 obtained by transforming an  $N \times p$  frequency matrix  
 499 into an  $N \times N$  distance matrix. This transformation  
 500 abstracts away from individual feature frequencies and  
 501 instead provides pairwise distances between the dialect  
 502 objects considered. To create a distance matrix, we  
 503 relied on the Euclidean distance measure (Aldenderfer  
 504 & Blashfield, 1984, 25), which defines the distance  
 505 between two dialect objects as the square root of the  
 506 sum of all  $p$  squared frequency differentials.

507 In this paper, we analyze distance matrices in two  
 508 ways: via cluster maps<sup>1</sup> and by correlating linguistic  
 509 Q distances with geographic distances<sup>4</sup>. Cluster maps are  
 510 a staple analysis technique in dialectometry (see, e.g.,  
 511 Goebel, 2007, Map 18; Heeringa, 2004, Figure 9.6)—the  
 512  $N \times N$  distance matrix is subjected to hierarchical  
 513 agglomerative cluster analysis (Jain, Murty & Flynn,  
 514 1999), a statistical technique used to group a number of  
 515 objects (in this study, dialects) into a smaller number of  
 516 discrete clusters.<sup>6</sup> Each of the clusters is assigned a dis-  
 517 tinct color, and the clusters are subsequently visually  
 518 depicted on a map.

519 Thus Map 3 projects a 3-cluster categorization based  
 520 on the normalization-based distances to geography. We  
 521 can see that there clearly is a geographic signal in the  
 522 dataset: Scottish counties are colored in blue, Northern  
 523 English dialects tend to belong to the red cluster, and  
 524 Southern English dialects tend to be assigned to the  
 525 green cluster. That being said, it is clear that there is also  
 526 a good deal of geographic incoherence and noise: there  
 527 are blue counties in Wales and England, red spots in  
 Southern Wales, the Scottish Highlands, and the

Hebrides, and Durham in the North of England is  
 mysteriously green. 528

529  
 530 We move on to correlating linguistic distances with  
 531 geographic distances, for the sake of precisely quanti-  
 532 fying the extent to which normalization based dialect  
 533 distances are predictable from geographic distance  
 534 (specifically: pair wise as-the-crow-flies distances,  
 535 which can be easily calculated from longitude/latitude  
 536 information) between dialect locations. The relationship  
 537 is visually depicted in Figure 2. There is a significant  
 538 relationship, but as-the-crow flies distance accounts for  
 539 only 3.4% of the normalization-based morphosyntactic  
 540 variance; a sub-linear logarithmic relationship only fares  
 541 marginally better at 3.6%. This is not a big share: in the  
 542 realm of syntax-focused atlas-based dialectometry,  
 543 analysts have reported  $R^2$  values of up to 45% (Spruit,  
 544 Heeringa & Nerbonne, 2009). Compared to that, the  
 545 geolinguistic signal in our normalization-based dataset  
 546 is quite weak.

**5. THE MODEL-BASED APPROACH** 547

548 As we have argued in section 2, the results of the  
 549 method discussed in the previous section may be  
 550 influenced by imprecise measurements in some features  
 551 and/or locations. We now move on to a method that  
 552 can alleviate this, namely regression modeling  
 553 using GAMs.

554 Regression models are essentially statistical models,  
 555 in which the values of one variable are represented as  
 556 combinations of the effects of other variables, the so-  
 557 called predictors. The effects of the individual pre-  
 558 dictors are determined from the data during the fitting  
 559 process. Generalized additive models, in particular,  
 560 allow the estimation of complex non-linear predictor  
 561 behavior in one or more dimensions, and thus allow the  
 562 representation in maps. When building regression  
 563 models for linguistic phenomena, the analyst faces a  
 564 bewildering amount of choices, ranging from the basic  
 565 representation of the data over the precise model spe-  
 566 cification to the details of the fitting process. The prin-  
 567 ciples guiding our selection processes for present  
 568 purposes are the following: First, the models and their  
 569 results should be as straightforwardly comparable to  
 570 the normalization based results as possible; the fewer  
 571 deviations from the process outlined in the previous  
 572 section, the better. This enables comparative analysis  
 573 that clearly shows where the methods differ. Second,  
 574 where possible, we should choose the methods such  
 575 that the result is still responsive to local conditions and  
 576 the frequency patterns at individual locations; after all,  
 577 simply parroting geography for its own sake would be  
 578 dialectologically meaningless. The aggregation process  
 579 can alleviate the impact of overfitting to a degree, but  
 580 may struggle on severely underfit data.

581 The first choice to be made pertains to the repre-  
 582 sentation of the outcome. Many linguistic phenomena  
 583 can be analyzed and represented in several ways; con-  
 584 sider, for example, feature [2], non-standard reflexive  
 585 pronouns such as *hissself* in (2) where Standard British  
 586 English would use *himself*. This feature could be studied  
 587 in terms of its frequency (e.g., how often are such non-  
 588 standard forms used per ten thousand words) or its share  
 589 of all constructions fulfilling similar functions (e.g., what  
 590 is the share of non-standard forms among all reflexives?).

591 Different operationalizations allow different aspects  
 592 of the data to shine through: proportions are more  
 593 robust toward variation in the base frequency of the  
 594 feature—a generally higher frequency of reflexive use in  
 595 some areas may or may not be linguistically relevant.  
 596 Even where it is relevant, frequency appears to mix two  
 597 different aspects of the phenomenon. To give a hypo-  
 598 theoretical extreme example, a county with ten reflexive  
 599 pronouns *pttw*, all of which are non-standard, is intu-  
 600 itively different with regard to non-standardness from  
 601 one where ten of 100 reflexives *pttw* are non-standard,  
 602 even if the normalized frequency of the non-standard  
 603 variant is the same. Nevertheless, there is considerable  
 604 debate whether pure frequencies or relative metrics are  
 605 what is ultimately of relevance to cognition and lin-  
 606 guistic theory (Bybee, 2010; Gries, 2012). For present  
 607 purposes, we decided to model only frequencies, as this  
 608 minimizes the differences between both normalization-  
 609 based and probabilistic analyses, and removes the  
 610 selection of relevant contexts as a source of errors. There  
 611 are several model types that allow modeling such fre-  
 612 quencies, of which perhaps the best-known is the Pois-  
 613 son regression. One assumption of this method,  
 614 however, is that mean and variance of the dependent  
 615 variable are the same. Linguistic material often violates  
 616 these assumptions, especially content words and other  
 617 “bursty” features, i.e., those that stray from even dis-  
 618 persion throughout a text/corpus (Manning & Schütze,  
 619 2000: 547). This is particularly troubling as grammatical  
 620 features have been shown to reliably occur more often  
 621 after they have already appeared (Branigan, 2007;  
 622 Szmrecsanyi, 2006) and are therefore likely to be bursty.  
 623 An alternative is negative binomial regression, which  
 624 includes an additional parameter (*theta*), allowing the  
 625 shape (and therefore variance) of the distribution to  
 626 vary. This parameter can be pre-specified, or deter-  
 627 mined from the data. This makes the negative binomial  
 628 distribution more appropriate for word and/or gram-  
 629 matical feature distributions. Note, however, that there  
 630 are still potential issues—these models may still suffer  
 631 from overdispersion, and especially from zero inflation,  
 632 i.e., more observations of zero than the distribution  
 633 allows (Hilbe, 2007). Models specifically designed for  
 634 such situations exist, but are difficult to operate, are not  
 available in standard tools, and may lead to model

fitting issues. For these reasons, overdispersed models  
 are sometimes recommended against,<sup>7</sup> and we will  
 employ regular negative binomial regression. The soft-  
 ware package we use, *mgcv* version 1.8–10 (Wood, 2006),  
 allows two major ways of determining the additional  
 coefficient for the negative binomial distribution. One  
 makes use of restricted maximum likelihood (REML),  
 the other of generalized cross-validation (GCV); both  
 also affect the other estimations in the model, and in  
 particular the general shapes of the frequency land-  
 scapes (as in Figure 1) that are our primary interest. In  
 our sample, the GCV approach seems to lead to more  
 varied, hilly landscapes, whereas the REML-based  
 models are flatter, and may remove too much of the  
 geographic specificity in the data. We therefore use  
 GCV where possible, with the search space for the theta  
 parameter rather wide (ranging from 0.01 to 50). How-  
 ever, there is a small number of features where the  
 GCV-based model either does not converge ([4], [9],  
 [29], and [39]), or leads to degenerate estimates for the  
 family parameter ([27], [31], [40] and [43]). In those  
 cases,<sup>8</sup> the REML-based model was used, which is more  
 robust in convergence (*mgcv* documentation: *negbin*). To  
 include geographic location in the models, we follow  
 the practice recommended by Wieling et al. (2014: 679)  
 and use *thin plate regression splines*, which are a “highly  
 suitable approach to model the influence of geography  
 in dialectology.” Each interview is coded with the geo-  
 location of the informant’s village or town, and we use  
 this information directly for modeling instead of first  
 aggregating on the county level.

One of the big advantages of the GAM-based  
 approach is that the models can easily be combined  
 with further information, whether it pertains to the  
 sociolinguistic situation (such as speaker character-  
 istics) or the immediate linguistic context (which may  
 often be feature-specific, such as subject type). Taking  
 such information into account can increase the reli-  
 ability of the analysis by reducing the effect of non-  
 geographic variation in the data, and can serve as cor-  
 roborating evidence for the analysis—if non-geographic  
 variables have the expected effects, this should raise our  
 confidence in both data and models—and is interesting  
 from the single-feature perspective. Nevertheless, there  
 are also significant costs to this. Feature-based annota-  
 tion can require tremendous effort, which is not feasible  
 for holistic analyses that cover a large number of fea-  
 tures. Information that is trivial to code, such as a  
 speaker’s age, may not be available in the corpus  
 metadata. This applies in our case: while most speakers  
 in the dialect corpus we are using have metadata con-  
 taining the sociolinguistically relevant predictors gen-  
 der and age, or have information that allows relatively  
 accurate estimation (such as birth and interview dec-  
 ades), some do not. This forces us to choose between

635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689

690 either excluding these factors in the models or excluding  
 691 ing the speakers where the information is not available  
 692 (which would remove some counties completely).  
 693 Finally, adding such predictors would also add another  
 694 difference to the normalization-based approach in the  
 695 previous section. We therefore restrict ourselves mostly  
 696 to geography-only models, but do present a brief sum-  
 697 mary of the sociolinguistic patterns observed in the next  
 698 section. Previous research suggests that, at least for this  
 699 dataset, age and gender are not substantially significant  
 700 factors, as their effect on the aggregate level is very  
 701 restricted (Wolk, 2014). For the models included in this  
 702 brief summary, age, gender, and their interaction are  
 703 always included. To make the models more easily  
 704 interpretable, the age variable was centered around the  
 705 mean, so that frequency differences between genders  
 706 are evaluated toward the center of the observed data,  
 707 and not at a hypothetical value for speakers that are  
 708 zero years old.

709 **5.1. Sociolinguistic Predictors**

710 In this section, we report on the models that contain  
 711 the sociolinguistic predictors age and gender. The  
 712 purpose of this is twofold: First, to demonstrate that  
 713 our approach can integrate sociolinguistic aspects and  
 714 thus link dialectometry proper to social dialectology.  
 715 Second, from sociolinguistic and dialectometric  
 716 research, we have clear ideas on the kinds of patterns  
 717 that should be expected. If our modeling results vio-  
 718 late these expectations, it would be curious and  
 719 potentially troubling; if, on the other hand, they lar-  
 720 gely match the expectations, our confidence in the  
 721 results can be raised. To make our expectations more  
 722 concrete, a large body of literature (e.g., Labov’s (2003:  
 723 266) principle 2 and the evidence presented in support  
 724 of it) suggests that women tend to use more standard  
 725 forms, particularly in linguistically stable conditions.  
 726 Similarly, many traditional dialectal forms are  
 727 believed to be receding; we therefore assume, as per  
 728 the apparent time hypothesis, that older speakers will  
 729 tend to have higher frequencies for these features. We  
 730 apply statistical significance at the customary thresh-  
 731 old of .05 as a crude filter to keep out particularly  
 732 unreliable signals.

733 For female speakers, we find that there are four clearly  
 734 non-standard features where the usage frequency is  
 735 clearly lower than for male speakers: [43] auxiliary  
 736 deletion, as in (14), [47] relativizer *what*, as in (10), [49] *as*  
 737 *what/than what*, as in (15), and [50] *unsplit for to*, as in (11).

- 738 (14) They gettin’ too much. <DEV\_007 >  
 739 (15) [...] but years ago they were a lot harder than what  
 740 they are today [SAL\_013]  
 741  
 742

In contrast, there is only one feature where  
 women use a non-standard feature more often: [36]  
*never* as a past tense negator, as in (8) above. This  
 feature is an atypical case: Cheshire, Edwards & Whittle  
 (1989) note that although a range of authors consider it  
 non-standard, they also attest its widespread use even  
 in formal written English. It is not implausible that this  
 feature may behave differently from more clearly  
 stigmatized features. There are some features that cover  
 both standard and non-standard usages. Feature [37],  
*wasn’t* as in (16), is one such case. The frequency of this  
 form is clearly going to depend on the prevalence of the  
*was - weren’t* split: if the preferred negated form is  
 always *weren’t*, we should expect a lower frequency of  
*wasn’t*. Female speakers, however, use *wasn’t* more  
 often than men. The remaining features with significant  
 differences are features that also exist in Standard  
 English and either have non-standard extensions or are  
 features undergoing language change. These include  
 pronominal forms, primary and modal verbs, and  
 features involved in classic grammatical alternations  
 such as *that*/zero complementizers.

- (16) There wasn’t a great deal. <LND\_004 >

For speaker age, we find a much clearer picture.  
 We find that, as expected, older speakers have a clear  
 tendency to use more archaic and non-standard  
 features. These include [28] non-standard weak verb  
 forms, as in (7) above, [33] multiple negation, as in (1)  
 above, and [50] *unsplit for to*, as in (11) above, as well as  
 the following:

- (17) [27] *a*-prefixing: And the other week she was  
 a-telling me, she said, [...] <KEN\_003 >  
 (18) [30] non-standard *come*: [...] he come home on a  
 Saturday afternoon [...] <LND\_006 >  
 (19) [32] *ain’t*: He says, You ain’t had your rotten teeth  
 out. <NTT\_012 >  
 (20) [39] non-standard verbal *-s*: [...]so I goes round  
 see, and hits the belt like that with mi hand  
 <WIL\_001 >

Again, we also find standard grammatical features,  
 often ones that are undergoing language change. Our  
 results here often match those reported in the literature;  
 to exemplify, [16] possessive *have got*, as in (21), and [26]  
*(have) got to* as a marker of epistemic or deontic modality,  
 as in (22), are both used less often by older speakers. Both  
 match previous results by Tagliamonte (2004), where,  
 for the birth dates covered in our corpus, older speakers  
 showed a dispreference for the variants involving *got*  
 while younger speakers used them more often. Interac-  
 tions turn out not to matter too much: only four features  
 exhibit a significant interaction between age and gender.

- 801 (21) Each class have got their own form captain [...] 853  
 802 <HEB\_023> 854  
 803 (22) Ehr, you 've got to bit them up proper [...] 855  
 804 <SOM\_005> 856

805 Overall, the results of the models including socio- 857  
 806 linguistic information often seem to match our expecta- 858  
 807 tions. There are, however, some effects that would be 859  
 808 expected based on the literature but fail to show up.  
 809 Feature [24], *must* as a marker of epistemic or deontic  
 810 modality as in (6), for instance, is a feature that is widely  
 811 considered an “obsolescing form” (Jankowski, 2004:  
 812 101). However, there is no general trend for lower fre-  
 813 quencies with younger speakers, only for female speak-  
 814 ers as an interaction. Feature [17], the *going to* future,  
 815 which has been shown to still grammaticalize and  
 816 expand in British dialects (Tagliamonte, Durham &  
 817 Smith, 2014), exhibits no pattern based on speaker age.  
 818 Nevertheless, we consider these results to be reassuring  
 819 —what we do find is quite plausible, and largely in line  
 820 with our expectations.

## 821 5.2. Geography

822 We now turn to the models that contain all speakers,  
 823 but do not include any predictors beyond geography  
 824 and, as an offset, text size. We first report summary  
 825 information on the individual feature models, then take  
 826 the aggregational perspective. To generate distances  
 827 from the individual models, we follow the method  
 828 outlined earlier for the normalization-based CBDM  
 829 approach as closely as possible. The major difference is  
 830 that, as input to the distance calculation, we use the  
 831 model predictions *pttw* instead of the normalized  
 832 counts. There is another minor technicality, which  
 833 involves exceedingly rare phenomena. For the normal-  
 834 ized counts, absences were coded as  $-1$ , corresponding  
 835 to  $0.1$  *pttw*, instead of taking the logarithm, as the  
 836 logarithm of zero is undefined. The models will in  
 837 general not predict exactly zero tokens, but numbers  
 838 that are arbitrarily close to zero, and therefore without  
 839 lower bound under logarithmic transformation. This  
 840 means that, in contrast to the process on normalized  
 841 values, rare features would have an undue influence.  
 842 Instead, we enforce a lower bound of again  $-1$  for the  
 843 logarithmically transformed frequency, keeping the  
 844 resulting values in a similar range for both processes.

845 Looking at the models individually, we find that for the  
 846 majority of features geographic information is significant.  
 847 Only eleven of the 57 features have a geographic  
 848 smoother with a non-significant *p*-value.<sup>10</sup> The GAM  
 849 solution used here also reports the proportion of the total  
 850 deviance in the data that the model explains; in some  
 851 cases, a non-significant (or marginally significant one) can  
 852 still have considerable explanatory power. The major

example here is [9] (the *s*-genitive), where the geographic 853  
 smoother is only marginally significant, but the model 854  
 nevertheless explains almost 46.6 percent of the deviance. 855  
 Similar cases are [43] auxiliary deletion, as in (14) above, at 856  
 37.3 percent and, to a lower degree, [49] *as what/than what*, 857  
 as in (15) above, and [52] gerundial complementation, as 858  
 in (23), where the model explains about 10 percent. 859

- (23) Oh, it was before I started working [...] 860  
 <GLA\_001> 861  
 862  
 863

Most of the other non-significant features hover 864  
 between 1 and about 7 percent of the deviance. This 865  
 pattern holds across the full dataset: the *p*-value 866  
 correlates negatively (at around  $r = -0.41$ ) with explana- 867  
 tory power. In other words, the lower the *p*-value, the 868  
 better geography can explain the observed distribution 869  
 of features. Restricting our analysis to those cases where 870  
 the smoother is significant, we find a left-skewed 871  
 distribution, with the peak at around 20 percent of 872  
 explained deviance. The median is at around 28 percent, 873  
 the mean slightly higher at 32 ( $SD = 20$ ). There is a small 874  
 number of outlier features where the simple model 875  
 accounts for over 80 percent, namely [22] the present 876  
 perfect with *be*, as in (5) above, [27] *a-prefixing*, as in (24), 877  
 and the negator *nae* [31], as in (25). 878  
 879

- (24) [...] What 's he 're a-jumpin' at? <SFK\_006> 880  
 (25) Ach, it s- s- sounds good, but it wasnae really. 881  
 <MLN\_006> 882  
 883  
 884

All of these are features with a very marked 885  
 geographic distribution, with high peaks in individual 886  
 counties or regions (the Scottish Lowlands for [31], 887  
 Suffolk for the others), and almost complete absence 888  
 elsewhere. On the other side of the spectrum, there are a 889  
 few features where geography is relevant, but not very 890  
 informative. Looking at the sixteen cases where geogra- 891  
 phy accounts for less than 20 percent of the deviance, we 892  
 find mostly features of Standard English: the pronouns 893  
*us* and *them*, relativizers, the *of*-genitive, zero comple- 894  
 mentation and so on. There are only three clearly non- 895  
 standard features in this list: [1] non-standard reflexives, 896  
 as in (2) above, [43] *there is/was* with plural subjects, as in 897  
 (26), and [44] non-standard *was*, as in (9) above. All of 898  
 these are among the most widespread features in Britain 899  
 or even worldwide, as the relevant surveys attest 900  
 (Britain, 2010; Kortmann, 2004, 2013). In short, the 901  
 features that are particularly weakly influenced by 902  
 geography are those that are available in most locations 903  
 —a quite plausible result. The remaining features, about 904  
 half of the total feature set, lie in the region ranging from 905  
 20 to about 60 percent. 906  
 907

- (26) There was all kinds of bits of quirks to keep the job 908  
 as quiet as ever they possible could. <YKS\_008> 909

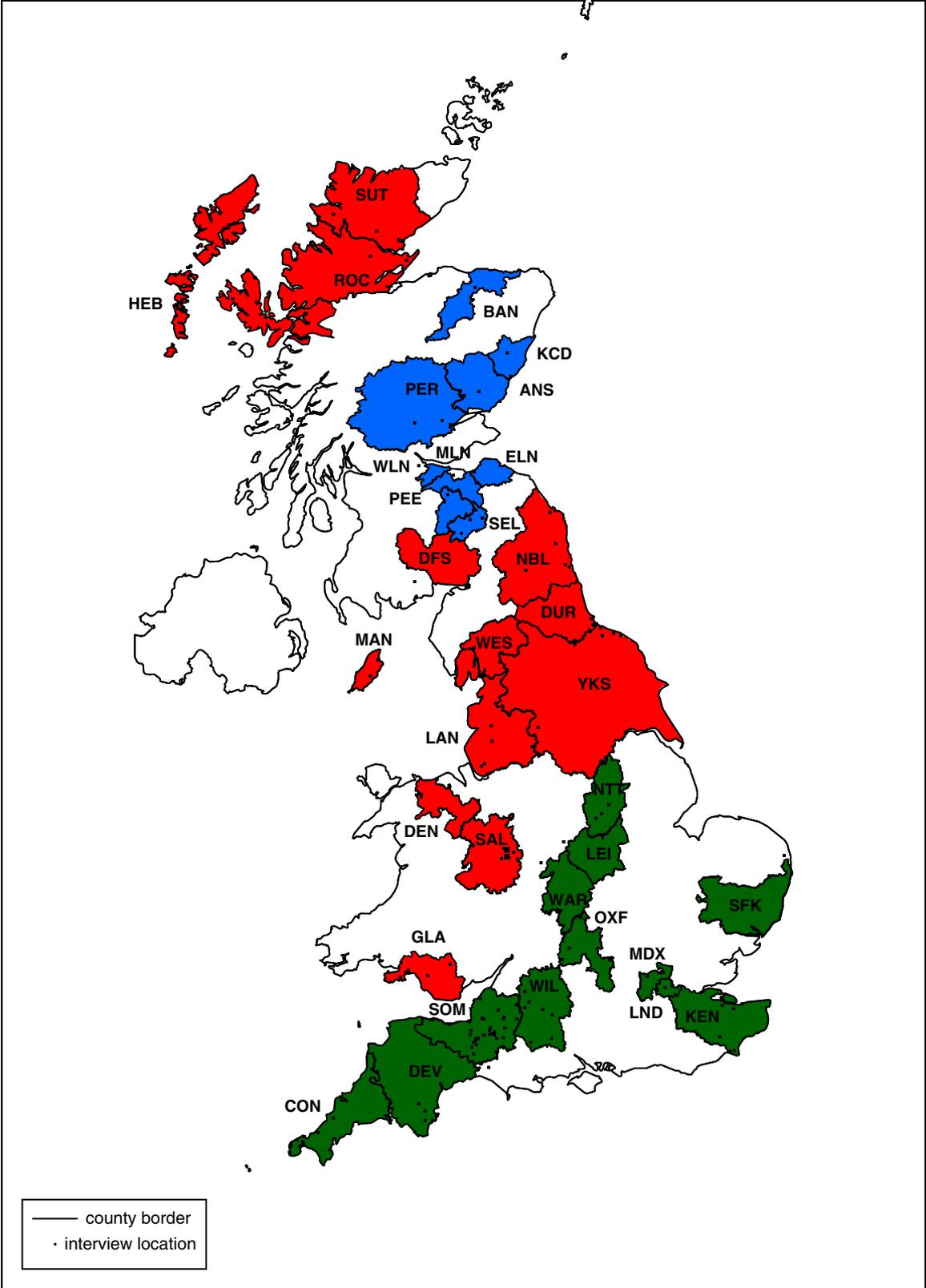
910 Map 2 (right) shows the resulting geographic pattern  
 911 for one feature, namely multiple negation [33], as in (1).  
 912 The geographic pattern displayed in the background  
 913 (i.e., the probabilistic pattern) is highly significant  
 914 ( $p < 0.01$ ) and very informative at almost 48 percent of  
 915 explained deviance. Lighter yellow/orange tones indi-  
 916 cate higher frequencies, while the dark green indicates  
 917 lower frequencies, in this case virtual absence. The small  
 918 black dots in the background show the interview loca-  
 919 tions; the display of the smoother is constrained to the  
 920 area surrounding these points. The left-hand side of  
 921 Map 2 displays the normalization-based values for the  
 922 same feature. The legends in the top right corners  
 923 illustrate the color gradients used, and highlight the  
 924 endpoints of the scale, in observations *pttw*. To give an  
 925 example, in the county with the most tokens, Kent, we  
 926 have 11 tokens *pttw*, and the GAM predicts a higher  
 927 value of 16.2 as the highest value for any individual  
 928 location. Observe that, in this case, the highest predic-  
 929 tion is even higher than the observed value; this results  
 930 from within-county variation patterns.<sup>6</sup> The lowest  
 931 observed county value is 0, complete absence, while the  
 932 lowest predicted value is 0.2. For these areas, there is  
 933 not enough variation to assert a lower value, even if no  
 934 token was observed. Even in Scotland, which seems  
 935 uniformly green in the normalized values, there are  
 936 seven observations in around 200,000 words, which  
 937 suggests an overall frequency of 0.34. Of course, there is  
 938 variation within Scotland, and Angus alone accounts  
 939 for almost half of the Scottish tokens. Nevertheless, it is  
 940 plausible that our best guess even for the lowest-  
 941 frequency areas is a value that is very close to zero, but  
 942 slightly higher. The red lines indicate the overall shape  
 943 of the frequency landscape: the feature is most pre-  
 944 valent along the southern and, to a lesser degree, east-  
 945 ern coasts of England, and decreases as one moves  
 946 north or west from there.

947 Using these models, we can calculate the predictions  
 948 for the counties, more precisely their centers, and pro-  
 949 ceed as previously described. To recapitulate briefly,  
 950 the resulting per-county frequencies *pttw* are collected  
 951 into an  $N \times p$  frequency matrix. This frequency matrix is  
 952 then used to calculate an  $N \times N$  distance matrix using  
 953 the Euclidean distance measure. The distance measure  
 954 we then subject to hierarchical clustering with noise  
 955 using Ward's method. The result is displayed in Map 4.  
 956 Despite the substantial differences between the meth-  
 957 ods—one working with straightforward normalized  
 958 frequencies, the other with an elaborate post-processing  
 959 of the raw data using generalized additive modeling—  
 960 the large-scale distribution is a remarkably similar tri-  
 961 partite division into a southern English area, a northern  
 962 English area, and a Scottish Lowlands area. Gone,  
 963 however, are the many outliers that were present in  
 964 Map 2; all clusters are now geographically contiguous,

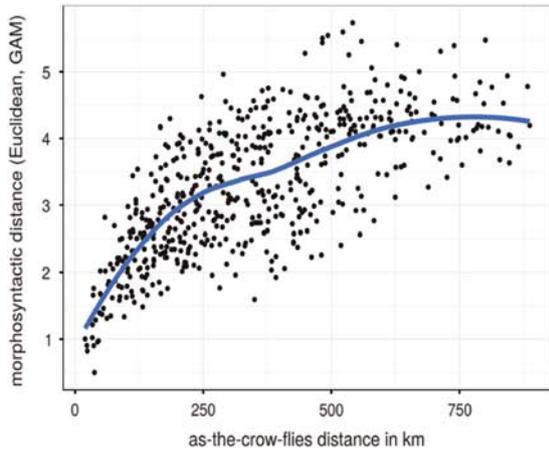
965 with the exception of the Scottish Highlands and the  
 966 Hebrides, which again show the largest similarity to the  
 967 northern English cluster. Looking at the results in  
 968 slightly greater detail, we find that there are smaller  
 969 differences at the borders of these areas: Dumfriesshire  
 970 and Northumberland have moved from the Scottish  
 971 Lowlands cluster to the northern English cluster, and  
 972 the Midlands counties east of Shropshire now exclu-  
 973 sively group with the south. This increase in areal  
 974 cohesion is also reflected in the relationship between  
 975 geographic and GAM-derived linguistic distances, dis-  
 976 played in Figure 3. It bears noting that the curve has the  
 977 sublinear curve predicted by Seguy's law (Nerbonne,  
 978 2010), in contrast to the one resulting from the  
 979 normalization-based procedure. Furthermore, the  
 980 explanatory power of (logarithmically transformed)  
 981 geography increases greatly, from less than 4 percent to  
 982 58.1 percent. This is hardly surprising—the assumption  
 983 of geographic coherence is central to the model for-  
 984 mulation. Nevertheless, we can interpret this as an  
 985 upper boundary: it is quite likely that having more data  
 986 for particularly sparse counties, where the model  
 987 assumptions carry greater weight, would increase the  
 988 variability somewhat, but less likely to lead to a reduc-  
 989 tion. This also means that there is significant linguistic  
 990 information left in the aggregated model predictions, as  
 991 58 percent is very high, but much closer to other dia-  
 992 lectometrical estimates (e.g., the 45 percent reported in  
 993 Spruit et al., 2009, for syntactic distances in Dutch dia-  
 994 lects) than to (a dialectologically almost meaningless)  
 995 100 percent.

## 996 6. DISCUSSION AND CONCLUSION

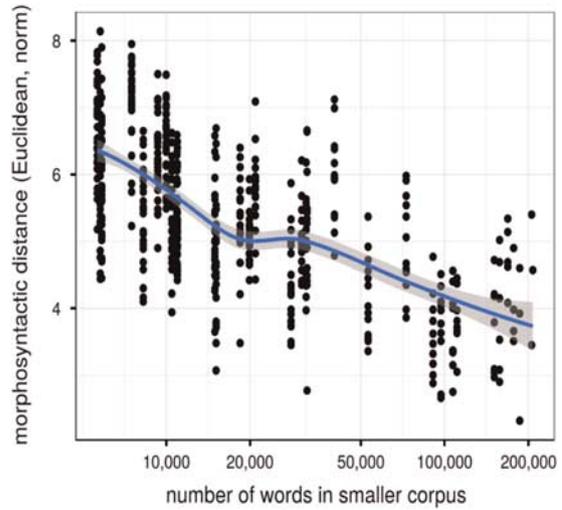
997 Let us briefly recapitulate our approach and main  
 998 results. We began by discussing data availability in  
 999 dialectometry, the *missing data effect*, and why it may be  
 1000 particularly troubling for frequency-based dialectome-  
 1001 try, including corpus-based variants. After outlining  
 1002 the data used in this paper, we presented the  
 1003 normalization-based approach to CBDM and showed  
 1004 that on the present dataset, it yields plausible results,  
 1005 with some peculiarities: the normalization-based solu-  
 1006 tion suffered from outliers that were hard to explain,  
 1007 and the relationship between geographic and linguistic  
 1008 distances had both a much lower explanatory power  
 1009 than one would expect based on the dialectometrical  
 1010 literature and a shape that is linear rather than sub-  
 1011 linear. In the previous section, we showed that a tech-  
 1012 nically more sophisticated solution based on  
 1013 generalized additive models yields large-scale dialect  
 1014 areas consistent with the first method. There were,  
 1015 however, major differences with regard to the peculia-  
 1016 rities, as we had hypothesized. The presence of outlying  
 1017 locations is a characteristic that Goebel (1991) reports for



Map 4. Cluster map based on the model-based CBDM approach. Displayed: 3-cluster solution.



**Figure 3.** Correlating model-based morphosyntactic distances with as-the-crow-flies distances (logarithmic  $r = .76$ ,  $p < .001$ ,  $R^2 = 58.1\%$ ).



**Figure 4.** Number of words in smaller corpus of county pair plotted against their morphosyntactic normalization-based distance.

WEB  
COLOR

WEB  
COLOR

1018 locations in atlas-based dialectometry where data is  
 1019 sparse. The second issue has been linked to such spar-  
 1020 sity as well: Streck (2014) reports that removing the 40  
 1021 percent of location pairs that involve those locations  
 1022 with the least amount of data almost doubled the per-  
 1023 centage of variance that as-the-crow-flies distance can  
 1024 explain. We argued earlier that processing the raw data  
 1025 with GAMs can alleviate these symptoms, and indeed  
 1026 this has turned out to be the case. The model-derived  
 1027 distances correlate more strongly, and sublinearly, with  
 1028 the spatial configuration of their locations, and the  
 1029 outliers have vanished. We will now discuss the two  
 1030 methods by comparing some of their properties  
 1031 directly, and argue that such modeling is an appro-  
 1032 priate choice for dialectometric purposes.

1033 If our hypothesis that the outliers and low fit result  
 1034 from sparsity is correct, we would expect the counties  
 1035 that have reduced coverage to behave systematically  
 1036 differently from those with ample text. The distances,  
 1037 however, apply to pairs of locations, while the number  
 1038 of words is a property of an individual location. As our  
 1039 hypothesis predicts that smaller sizes should have the  
 1040 greatest impact, it is sensible to associate each distance  
 1041 with the minimum size of either subcorpus involved.  
 1042 This way, a pairing that is assigned the value of 50,000  
 1043 ensures that both corpora at least reach that level of  
 1044 coverage, and the higher this number is, the more con-  
 1045 fidence we can place in the distance measurement of  
 1046 this pair. A small downside to this is that the right-hand  
 1047 side of the scale thins out—the county with the lowest  
 1048 amount of running text (Banffshire), contributes 33  
 1049 individual points to the analysis, as it is the smallest  
 1050 corpus in all its 33 pairings. This makes this county  
 1051 particularly prominent visually. As the number of  
 1052 words for a county increases, the number of points on  
 1053 the corresponding spot on the x-axis decreases, as there

are fewer and fewer counties that have at least as much  
 text. Figure 4 displays the minimum number of words  
 on the x-axis on a logarithmic scale, and the y-axis  
 shows the normalization-based distance. The blue line  
 is a LOESS smoother that indicates the overall trend.  
 The data follow an almost perfectly straight downward-  
 sloping line, with a small plateau of stability after 20,000  
 words. In other words, the relationship is a logarithmic  
 decay: as the number of words in a county subcorpus  
 increases, its distance to larger subcorpora decreases,  
 but the rate at which it does so decreases as well. This  
 relationship is visually strong, and this is confirmed by  
 examining the correlation numerically: the log-  
 transformed minimum number of words can account  
 for 44 percent of the variance in the normalization-  
 based distances. Note also that there is no discernible  
 relationship between minimum size and geographic  
 distance ( $r = -0.08$ , and no visible pattern when plotted)  
 —any such relationship is therefore not due to the spa-  
 tial distribution of locations. In the model-derived dis-  
 tances, this pattern disappears almost completely, with  
 only 2 percent of the variance in linguistic distances  
 attributable to minimum text size. All of this is con-  
 sistent with the hypothesis. While we do expect the  
 logarithmic decay to cease at some point (after all, it is  
 unlikely that the actual difference between dialects is  
 arbitrarily close to zero), none is apparent so far—the  
 influence of corpus size on the normalized distances  
 diminishes, but does not vanish.

How good, then is our solution to this problem—  
 geographic smoothing using generalized additive mod-  
 els? It seems to pass a gauntlet of consistency checks—  
 strength and shape of the relationship of linguistic and  
 geographic distances, influence of geography, and  
 coherence and interpretability of the resulting areal

1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088

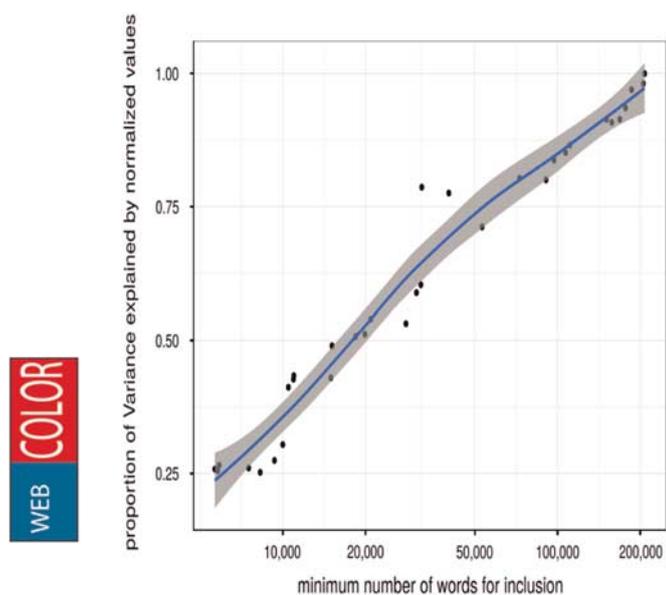


Figure 5. Percent of explained variance between normalization-based and model-based distance as a function of the minimum corpus size required for inclusion.

1089 classification. Nevertheless, this comes at a cost: the  
 1090 Fundamental Dialectological Postulate has to be accep-  
 1091 ted beforehand, limiting the inferences one can make  
 1092 from the data. The smoothing may also be too aggres-  
 1093 sive, presenting a limited picture of the actually exist-  
 1094 ing variation. Given that there is no external measure for  
 1095 morphosyntactic differences between these locations, we  
 1096 cannot directly evaluate exactly how accurate this strat-  
 1097 egy is. What we can do, however, is compare it to  
 1098 another strategy—a tactical retreat to the counties that  
 1099 are most plentiful in terms of the textual coverage, where  
 1100 normalization should work best. Taking all the counties,  
 1101 we find that there is a clear but limited relationship  
 1102 between normalization- and model-based metrics, with  
 1103 a linear  $R^2$  of 0.26. We can now successively drop the  
 1104 county with the lowest amount of text and see how this  
 1105 changes the result (without recalculating the models).  
 1106 The development is displayed in Figure 5, with the  
 1107 minimum number of words required for inclusion on the  
 1108  $x$ -axis (again on a log scale) and the linear  $R^2$  on the  
 1109  $y$ -axis. There is a clear relationship: the more one re-  
 1110 stricts attention to well-covered areas, the more the re-  
 1111 sults of the two methods approximate one another. The  
 1112  $R^2$  for the logarithmic relationship is 0.85. In words, the  
 1113 two analyses increasingly resemble one another. This  
 1114 means that there are few downsides to modeling com-  
 1115 pared to exclusion; including the low-data counties  
 1116 yields largely the same results for the high-data coun-  
 1117 ties, but modeling also leads to plausible results for  
 1118 the rest, and allows them to contribute to the analysis.

1119 We wish to make a final point: Dialectometry has  
 1120 heralded the beneficial properties of aggregation for

noise reduction and pattern identification (see e.g.,  
 Nerbonne, 2009: 129, who considers it “at the heart of  
 the benefits of dialectometry”). Our analysis confirms  
 this, but also shows important limitations: there are  
 biases that are impervious to aggregation; for noise  
 that persists across features, or limits which features  
 can be compared for individual location pairs, aggre-  
 gation may simply not be enough. Our example was  
 related to frequency-based dialectometry, but similar  
 concerns should apply for categorical data. We urge  
 scholars to be mindful of this, and include analyses  
 based on data availability in cases where the data  
 basis is not exactly equivalent for all locations.

### Acknowledgements

The research presented in this article is partially  
 based on work done for the first author’s PhD dis-  
 sertation as well as the second author’s book *Gram-  
 matical variation in British English dialects*, but has  
 been substantially reworked and expanded. We are  
 grateful to Bernd Kortmann, Guido Seiler, Peter Auer,  
 John Nerbonne, Joan Bresnan, and to the audience at  
 Methods in Dialectology XV where an earlier version  
 of this research was presented. We use cartographic  
 material provided by Natural Earth, GADM, the Scot-  
 tish Government Spatial Data Infrastructure, and the  
 Great Britain Historical GIS Project; licensing terms  
 can be found in the appendix. Funding from the  
 Freiburg Institute for Advanced Studies (FRIAS) is  
 gratefully acknowledged.

### Supplementary material

To view supplementary material for this article, please  
 visit <http://dx.doi.org/10.1017/jlg.2018.6>

### Notes

- <sup>1</sup> For the purposes of this paper, we consider frequency  
 differences of linguistic phenomena as observed in  
 naturalistic corpora a (potentially noisy) proxy of the  
 underlying linguistic differences between locations. Some  
 readers may prefer methods that are explicitly based on  
 linguistic variables, showing likelihood of use as op-  
 posed to surface frequency. Our method can also in-  
 corporate this approach; examples can be found in  
 Wolk (2014). The advantage of surface frequencies  
 for present purposes is that they can be directly ap-  
 plied to all the features under study equally.
- <sup>2</sup> The somewhat coarsely meshed network of interview  
 locations makes it hard to read continuum maps, as  
 subtle color comparisons over areas separated by white  
 space are visually challenging. Interested readers can  
 find the continuum maps for both approaches in the  
 online appendix.
- <sup>3</sup> An anonymous reviewer points out that the Fundamen-  
 tal Dialectological Postulate need not hold, as linguis-  
 tic

1171 behavior may be influenced by factors other than distance,  
 1172 including migration or political borders. We agree in gen-  
 1173 eral, but consider it applicable in situations with a long  
 1174 settlement history; also note that both the principle and  
 1175 our application of it are intended as empirical general-  
 1176 izations and allow for deviations given sufficient evidence.  
 1177 <sup>4</sup> The use of measures of spatial autocorrelation, namely  
 1178 Moran's I, instead of correlation coefficients was recom-  
 1179 mended to us by an anonymous reviewer. We agree that  
 1180 there are clear advantages to using Moran's I to investigate  
 1181 the distribution of individual features, but see no clear way  
 1182 to extend this method to patterns of aggregate distances.  
 1183 <sup>5</sup> The interview locations are unfortunately not balanced  
 1184 throughout the regions, and may be more concentrated in  
 1185 certain counties and spread out elsewhere. We will not  
 1186 make claims about places not covered in our data.  
 1187 <sup>6</sup> We specifically used Ward's minimum variance method  
 1188 (Ward, 1963), an algorithm that tends to create small and  
 1189 even-sized clusters. Note that simple clustering can be  
 1190 unstable, which is why we clustered with noise (Ner-  
 1191 bonne, Kleiweg, Manni & Heeringa, 2008): The original  
 1192 distance matrix was clustered repeatedly, adding some  
 1193 random amount of noise ( $c = \sigma/2$ ) in each run.  
 1194 <sup>7</sup> For example, Paul Allison advises against overdispersed  
 1195 models at [http://statisticalhorizons.com/zero-inflated-](http://statisticalhorizons.com/zero-inflated-models)  
 1196 [models](http://statisticalhorizons.com/zero-inflated-models).  
 1197 <sup>8</sup> For the set of models that include social information on the  
 1198 speakers, the list of features fit using REML because of  
 1199 convergence problems is [19], [39], [55], while the set with  
 1200 too large family parameters remains the same.  
 1201 <sup>9</sup> The full list is: [7] synthetic comparison, [9] the *g*-genitive,  
 1202 [10] preposition stranding, [17] the *going to* future, [36]  
 1203 preverbal *never*, [43] auxiliary deletion, [49] *as what/than*  
 1204 *what*, [51] infinitival and [52] gerundial complementation,  
 1205 [55] lack of inversion, and [57] the prepositional dative.  
 1206 <sup>10</sup> In other cases, especially those involving exceptionally  
 1207 high values in low-data counties, the GAM prediction may  
 1208 be much lower. See for example feature [31], *-nae*, where  
 1209 the highest observed value is 144 *pttw*, but the highest  
 1210 prediction only 79.

## 1211 References

1212 Aldenderfer, Mark S. & Roger K. Blashfield. 1984. *Cluster*  
 1213 *analysis*. London: Sage Publications.  
 1214 Anderwald, Lieselotte & Benedikt Szmrecsanyi. 2009. Corpus  
 1215 linguistics and dialectology. In Lüdeling, Anke & Merja Kytö  
 1216 (eds.), *Corpus linguistics: An international handbook*,  
 1217 Handbücher Zur Sprach- Und  
 1218 Kommunikationswissenschaft / Handbooks of Linguistics  
 1219 and Communication Science 29(1), 1126–39. Berlin: Mouton  
 1220 de Gruyter.  
 1221 Branigan, Holly. 2007. Syntactic priming. *Language and*  
 1222 *Linguistics Compass* 1(1/2). 1-16. doi: 10.1111/j.1749-  
 1223 818X.2006.00001.x.  
 1224 Britain, David. 2010. Grammatical variation in the  
 1225 contemporary spoken English of England. In Kirkpatrick,  
 1226 Andy (ed.), *The handbook of world Englishes*, 37-58. London:  
 1227 Routledge.

Bybee, Joan L. 2010. *Language, usage and cognition*. Cambridge: 1228  
 Cambridge University Press. 1229  
 Chambers, Jack K. & Peter Trudgill. 1998. *Dialectology*, 2nd edn. 1230  
 Cambridge: Cambridge University Press. 1231  
 Cheshire, Jenny, Viv Edwards & Pamela Whittle. 1989. Urban 1232  
 British dialect grammar: the question of dialect levelling. 1233  
*English World-Wide* 10. 185-225. 1234  
 Eisenstein, Jacob. 2018. Identifying regional dialects in on-line 1235  
 social media. In Boberg, Charles, John Nerbonne & Dominic 1236  
 Watts (eds.), *The handbook of dialectology*. New York: Wiley- 1237  
 Blackwell. 1238  
 Goebel, Hans. Rätoromanisch versus Hochitalienisch versus 1239  
 Oberitalienisch. *Ladinia* 1. 39-71. 1240  
 Goebel, Hans. 1982. *Dialektometrie: Prinzipien und methoden des*  
*einsatzes der numerischen taxonomie im bereich der*  
*dialektgeographie*. Wien: Österreichische Akademie der 1241  
 Wissenschaften. 1242  
 Goebel, Hans. 1984. *Dialektometrische studien: Anhand*  
*Italoromanischer, Rätoromanischer Und Galloromanischer*  
*sprachmaterialien aus AIS und ALF*. Tübingen: Niemeyer. 1243  
 Goebel, Hans. 1993. Dialectometry: A short overview of the 1244  
 principles and practice of quantitative classification of 1245  
 linguistic atlas data. In Köhler, Reinhard & Burghard Rieger 1246  
 (eds.), *Contributions to Quantitative Linguistics: Proceedings of*  
*the First International Conference on Quantitative Linguistics,*  
*QUALICO, Trier, 1991, 277–315*. Dordrecht: Kluwer. 1247  
 Goebel, Hans. 2007. A bunch of dialectometric flowers: A brief 1248  
 introduction to dialectometry. In Smit, Ute, Stefan Dollinger, 1249  
 Julia Hüttner, Gunter Kaltenböck & Ursula Lutzky (eds.), 1250  
*Tracing English through time: Explorations in language variation,*  
 133-172. Wien: Braumüller. 1251  
 Great Britain Historical GIS Project. 2004. Great Britain 1252  
 Historical GIS. University of Portsmouth. 1253  
 Gries, Stefan Th. 2012. Frequencies, probabilities, and 1254  
 association measures in usage-/exemplar-based linguistics: 1255  
 Some necessary clarifications. *Studies in Language* 36(3). 1256  
 477-510. 1257  
 Grieve, Jack. 2014. A comparison of statistical methods for the 1258  
 aggregation of regional linguistic variation. In Szmrecsanyi, 1259  
 Benedikt & Bernhard Wälchli (eds.), *Aggregating dialectology,*  
*typology, and register analysis: Linguistic variation in text and*  
*speech*, *Lingua & Litterae* 2853-2888. Berlin: Walter de 1260  
 Gruyter. 1261  
 Grieve, Jack. 2016. *Regional variation in written American English*. 1262  
 Cambridge: Cambridge University Press. 1263  
 Heeringa, Wilbert. 2004. *Measuring dialect pronunciation*  
*differences using Levenshtein distance*. Groningen,  
 Netherlands: University of Groningen dissertation. 1264  
 Hilbe, Joseph M. 2007. *Negative binomial regression*. Cambridge: 1265  
 Cambridge University Press. 1266  
 Hernández, Nuria. 2006. *User's guide to FRED*. Freiburg: 1267  
 University of Freiburg. 1268  
 Huang, Yuan, Diansheng Guo, Alice Kasakoff & Jack Grieve. 1269  
 2016. Understanding US regional linguistic variation with 1270  
 Twitter data analysis. *Computers, Environment and Urban*  
*Systems* 59. 244-255. 1271  
 Jain, Anil K., M. Narasimha Murty & Patrick J. Flynn. 1999. 1272  
 Data clustering: A review. *ACM Computing Surveys* 31(3). 1273  
 264-323. 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286

- 1287 James, Gareth, Daniela Witten, Trevor Hastie & Robert  
1288 Tibshirani. 2013. *An introduction to statistical learning: With*  
1289 *applications in R*. New York: Springer. 1348
- 1290 Jankowski, Bridget. 2004. A transatlantic perspective of  
1291 variation and change in English deontic modality. *Toronto*  
1292 *Working Papers in Linguistics* 23(2). 85-113. 1349
- 1293 Kortmann, Bernd & Benedikt Szmrecsanyi. 2004. Global  
1294 synopsis: Morphological and syntactic variation in  
1295 English. In Bernd Kortmann, Edgar Schneider, Kate  
1296 Burridge, Rajend Mesthrie & Clive Upton (eds.), *A Handbook*  
1297 *of varieties of English* vol. 2. 1142-1202. Berlin: Mouton de  
1298 Gruyter. 1350
- 1299 Manning, Chris & Hinrich Schütze. 1999. *Foundations of*  
1300 *statistical natural language processing*. Cambridge, MA: MIT  
1301 Press. 1351
- 1302 Nerbonne, John. 2009. Data-driven dialectology. *Language and*  
1303 *Linguistics Compass* 3(1). 175-198. 1352
- 1304 Nerbonne, John. 2010. Measuring the diffusion of  
1305 linguistic change. *Philosophical Transactions of the Royal*  
1306 *Society B: Biological Sciences* 365. 3821-3828. 1353
- 1307 Nerbonne, John & Peter Kleiweg. 2007. Toward a  
1308 dialectological yardstick. *Journal of Quantitative Linguistics* 14  
1309 (2). 148-166. 1354
- 1310 Nerbonne, John, Wilbert Heeringa & Peter Kleiweg. 1999. Edit  
1311 distance and dialect proximity. In David Sankoff & Joseph B.  
1312 Kruskal (eds.), *Time Warps, String Edits and Macromolecules:*  
1313 *The Theory and Practice of Sequence Comparison*. Stanford: CSLI  
1314 Press. 1355
- 1315 Nerbonne, John, Peter Kleiweg, Franz Manni & Wilbert  
1316 Heeringa. 2008. Projecting dialect differences to geography:  
1317 Bootstrapping clustering vs. clustering with noise. In  
1318 Christine Preisach, Lars Schmidt-Thieme, Hans Burkhardt &  
1319 Reinhold Decker (eds.), *Data analysis, machine learning, and*  
1320 *applications. Proceedings of the 31st Annual Meeting of the*  
1321 *German Classification Society*, 647-654. Berlin: Springer. 1356
- 1322 Orton, Harold, Stewart Sanderson & J. D. A. Widdowson. 1978.  
1323 *The Linguistic Atlas of England*. London: Croom Helm. 1357
- 1324 Pickl, Simon, Aaron Spettl, Simon Pröll, Stephan Elspaß,  
1325 Werner König & Volker Schmidt. 2014. Linguistic distances  
1326 in dialectometric intensity estimation. *Journal of Linguistic*  
1327 *Geography* 2(1). 25-40. 1358
- 1328 Sankoff, David. 1987. Variable Rules. In Ulrich Ammon,  
1329 Norbert Dittmar & Klaus J. Mattheier (eds.), *Sociolinguistics /*  
1330 *Soziolinguistik: An International Handbook / Ein internationales*  
1331 *Handbuch, Vol.2*, 984-997. Berlin: de Gruyter. 1359
- 1332 Séguy, Jean. 1971. La relation entre la distance spatiale et la  
1333 distance lexicale. *Revue de Linguistique Romane* 35. 335-357. 1360
- 1334 Shackleton, R. G. 2007. Phonetic variation in the traditional  
1335 English dialects: A computational analysis. *Journal of English*  
1336 *Linguistics* 35(1). 30-102. doi: 10.1177/0075424206297857. 1361
- 1337 Spruit, Marco René. 2006. Measuring syntactic variation in  
1338 Dutch dialects. *Literary and Linguistic Computing* 21(4).  
1339 493-506. 1362
- 1340 Spruit, Marco René, Wilbert Heeringa & John Nerbonne. 2009.  
1341 Associations among linguistic levels. *Lingua* 119(11). 1624-  
1342 1642. doi: 10.1016/j.lingua.2009.02.001. 1363
- 1343 Streck, Tobias. 2014. Alemannisch quantitativ. Zur  
1344 Erklärungskraft der Geografie für aggregierte  
1345 Dialektunterschiede. In Pia Bergmann, Karin Birkner, Peter  
1346 Gilles, Helmut Spiekermann & Tobias Streck (eds.),  
1347 *Sprache im Gebrauch: Räumlich, zeitlich, interaktional.*  
1348 *Festschrift für Peter Auer*, 157-173. Heidelberg: Winter  
1349 (OraLingua 9). 1364
- Szmrecsanyi, Benedikt. 2006. *Morphosyntactic persistence in*  
1350 *spoken English: A corpus study at the intersection of variationist*  
1351 *sociolinguistics, psycholinguistics, and discourse analysis*. Berlin:  
1352 Mouton de Gruyter. 1353
- Szmrecsanyi, Benedikt. 2008. Corpus-based dialectometry:  
1354 Aggregate morphosyntactic variability in British English  
1355 dialects. *International Journal of Humanities and Arts Comput-*  
1356 *ing* 2(1/2). 279-296. 1357
- Szmrecsanyi, Benedikt. 2010. *The morphosyntax of BrE dialects in*  
1358 *a corpus-based dialectometrical perspective: Feature extraction,*  
1359 *coding protocols, projections to geography, summary statistics*.  
1360 Freiburg: University of Freiburg. URN: urn:nbn:de:bsz:25-  
1361 Opus-73209, URL: [http://www.freidok.uni-freiburg.de/](http://www.freidok.uni-freiburg.de/volltexte/7320/)  
1362 [volltexte/7320/](http://www.freidok.uni-freiburg.de/volltexte/7320/). 1363
- Szmrecsanyi, Benedikt. 2011. Corpus-based dialectometry: A  
1364 methodological sketch. *Corpora* 6(1). 45-76. 1365
- Szmrecsanyi, Benedikt. 2013. *Grammatical variation in British*  
1366 *English dialects: A study in corpus-based dialectometry*.  
1367 Cambridge: Cambridge University Press. 1368
- Szmrecsanyi, Benedikt & Nuria Hernández. 2007. *Manual of*  
1369 *information to accompany the Freiburg Corpus of English Dialects*  
1370 *Sampler ("FRED-S")*. Freiburg: University of Freiburg. 1371
- <http://www.freidok.uni-freiburg.de/volltexte/2859/>.  
1372 1373
- Szmrecsanyi, Benedikt & Bernhard Wälchli. eds. 2014.  
1374 *Aggregating dialectology, typology, and register analysis:*  
1375 *Linguistic variation in text and speech. Lingua & Litterae.*  
1376 28 Berlin: Walter de Gruyter. 1377
- Szmrecsanyi, Benedikt & Christoph Wolk. 2011. Holistic  
1378 corpus-based dialectology. *Brazilian Journal of Applied*  
1379 *Linguistics/Revista Brasileira de Linguística Aplicada* 11(2).  
1380 561-592. 1378
- Tagliamonte, Sali. 2004. *Have to, gotta, must:*  
1381 *Grammaticalisation, variation and specialization in English*  
1382 *deontic modality*. In Christian Mair & Hans Lindquist (eds.),  
1383 *Corpus approaches to grammaticalization in English*, 33-55.  
1384 Amsterdam: Benjamins. 1385
- Tagliamonte, Sali A., Mercedes Durham & Jennifer Smith.  
1386 2014. Grammaticalization at an early stage: future *be going to*  
1387 *in conservative British dialects. English Language and*  
1388 *Linguistics* 18(1). 75-108. 1389
- Tobler, Waldo. 1970. A computer movie simulating urban  
1390 growth in the Detroit region. *Economic Geography* 46(2).  
1391 234-240. 1392
- Viereck, Wolfgang. The Computerisation and quantification of  
1393 linguistic data: Dialectometrical methods. In Alan R. Thomas  
1394 (ed.), *Methods in Dialectology : Proceedings of the Sixth*  
1395 *International Conference Held At the University College of North*  
1396 *Wales, 3rd-7th August 1987*. 523-550. Clevedon: Multilingual  
1397 Matters. 1398
- Viereck, Wolfgang, Heinrich Ramisch, Harald Händler, Petra  
1399 Hoffmann & Wolfgang Putschke. 1991. *The Computer*  
1400 *Developed Linguistic Atlas of England*. Tübingen: Niemeyer. 1401
- Ward, Joe H. Jr. 1963. Hierarchical grouping to optimize an  
1402 objective function. *Journal of the American Statistical*  
1403 *Association* 58. 236-244. 1404

1405	Wieling, Martijn, Simonetta Montemagni, John Nerbonne & R. Harald Baayen. 2014. Lexical differences between Tuscan dialects and standard Italian: Accounting for geographic and sociodemographic variation using generalized additive mixed modeling. <i>Language</i> 90(3). 669-692.		
1406			
1407			
1408			
1409			
1410	Wieling, Martijn, John Nerbonne & R. Harald Baayen. 2011. Quantitative social dialectology: Explaining linguistic variation geographically and socially. <i>PLoS ONE</i> 6 (9).		
1411			
1412			
1413			
1414	Wolk, Christoph. 2014. <i>Integrating aggregational and probabilistic approaches to language variation</i> . Freiburg: University of Freiburg dissertation.		
1415			
1416			
1417	Wolk, Christoph & Benedikt Szmrecsanyi. 2016. Top-down and bottom-up advances in corpus-based dialectometry. In Marie-Hélène Côté, Remco Knooihuizen & John Nerbonne (eds.), <i>The future of dialects</i> , 225-244. Berlin: Language Science Press.		
1418			
1419			
1420			
1421			
1422	Wood, Simon N. 2006. <i>Generalized additive models: an introduction with R</i> . Boca Raton, FL: Chapman & Hall/CRC.		
1423			
1424	<b>Appendix: the feature catalogue</b>		
1425	A. Pronouns and determiners		
1426	[1] non-standard reflexives (e.g. <u>they didn't go</u>		
1427	<u>themselves</u> )		
1428	[2] standard reflexives (e.g. <u>they didn't go themselves</u> )		
1429	[3] archaic <u>thee/thou/thy</u> (e.g. <u>I tell thee a bit more</u> )		
1430	[4] archaic <u>ye</u> (e.g. <u>ye'd dancing every week</u> )		
1431	[5] <u>us</u> (e.g. <u>us couldn't get back, there was no train</u> )		
1432	[6] <u>them</u> (e.g. <u>I wonder if they'd do any of them things</u>		
1433	<u>today</u> )		
1434			
1435			
1436	B. The noun phrase		
1437			
1438	[7] synthetic adjective comparison (e.g. <u>he was always</u>		
1439	<u>keener on farming</u> )		
1440	[8] the of-genitive (e.g. <u>the presence of my father</u> )		
1441	[9] the s-genitive (e.g. <u>my father's presence</u> )		
1442	[10] preposition stranding (e.g. <u>the very house which it</u>		
1443	<u>was in</u> )		
1444	[11] cardinal number + <u>years</u> (e.g. <u>I was there about</u>		
1445	<u>three years</u> )		
1446	[12] cardinal number + <u>year-Ø</u> (e.g. <u>she were three</u>		
1447	<u>year old</u> )		
1448			
1449			
1450	C. Primary verbs		
1451			
1452	[13] the primary verb TO DO (e.g. <u>why did you not</u>		
1453	<u>wait?</u> )		
1454	[14] the primary verb TO BE (e.g. <u>I was took straight into</u>		
1455	<u>this pitting job</u> )		
1456	[15] the primary verb TO HAVE (e.g. <u>we thought some-</u>		
1457	<u>body had brought them</u> )		
1458	[16] marking of possession – HAVE GOT (e.g. <u>I have got</u>		
1459	<u>the photographs</u> )		
1460			
1461			
		D. Tense and aspect	1462
			1463
		[17] the future marker BE GOING TO (e.g. <u>I'm going to let</u>	1464
		<u>you into a secret</u> )	1465
		[18] the future markers WILL/SHALL (e.g. <u>I will let you</u>	1466
		<u>into a secret</u> )	1467
		[19] WOULD as marker of habitual past (e.g. <u>he would</u>	1468
		<u>go around killing pigs</u> )	1469
		[20] <u>used to</u> as marker of habitual past (e.g. <u>he used to</u>	1470
		<u>go around killing pigs</u> )	1471
		[21] progressive verb forms (e.g. <u>the rest are going to</u>	1472
		<u>Portree School</u> )	1473
		[22] the present perfect with auxiliary BE (e.g. <u>I'm come</u>	1474
		<u>down to pay the rent</u> )	1475
		[23] the present perfect with auxiliary HAVE (e.g. <u>they've</u>	1476
		<u>killed the skipper</u> )	1477
			1479
		E. Modality	1480
			1481
		[24] marking of epistemic and deontic modality: MUST	1482
		(e.g. <u>I must pick up the book</u> )	1483
		[25] marking of epistemic and deontic modality: HAVE	1484
		TO (e.g. <u>I have to pick up the book</u> )	1485
		[26] marking of epistemic and deontic modality: GOT	1486
		TO (e.g. <u>I gotta pick up the book</u> )	1487
			1488
			1489
		F. Verb morphology	1490
			1491
		[27] a-prefixing on -ing-forms (e.g. <u>he was a-waiting</u> )	1492
		[28] non-standard weak past tense and past participle	1493
		forms (e.g. <u>they knewed all about these things</u> )	1494
		[29] non-standard past tense <u>done</u> (e.g. <u>you came</u>	1495
		<u>home and done the home fishing</u> )	1496
		[30] non-standard past tense <u>come</u> (e.g. <u>he come down</u>	1497
		<u>the road one day</u> )	1498
			1500
		G. Negation	1501
			1502
		[31] the negative suffix -nae (e.g. <u>I cannae do it</u> )	1503
		[32] the negator ain't (e.g. <u>people ain't got no money</u> )	1504
		[33] multiple negation (e.g. <u>don't you make no damn</u>	1505
		<u>mistake</u> )	1506
		[34] negative contraction (e.g. <u>they won't do anything</u> )	1507
		[35] auxiliary contraction (e.g. <u>they'll not do anything</u> )	1508
		[36] <u>never</u> as past tense negator (e.g. <u>and they never</u>	1509
		<u>moved no more</u> )	1510
		[37] WASN'T (e.g. <u>they wasn't hungry</u> )	1511
		[38] WEREN'T (e.g. <u>they weren't hungry</u> )	1512
			1514
		H. Agreement	1515
			1516
		[39] non-standard verbal -s (e.g. <u>so I says, What have</u>	1517
		<u>you to do?</u> )	1518
		[40] <u>don't</u> with 3 <sup>rd</sup> person singular subjects (e.g. <u>if this</u>	1519
		<u>man don't come up to it</u> )	1520

1521	[41]	standard <u>doesn't</u> with 3 <sup>rd</sup> person singular subjects (e.g. <u>if this man doesn't come up to it</u> )	[50]	unsplit <u>for to</u> (e.g. <u>it was ready for to go away with the order</u> )	1545
1522					1546
1523	[42]	existential/presentational <u>there is/was</u> with plural subjects (e.g. <u>there was children involved</u> )	[51]	infinitival complementation after BEGIN, START, CONTINUE, HATE, and LOVE (e.g. <u>I began to take an interest</u> )	1547
1524					1548
1525	[43]	absence of auxiliary BE in progressive constructions (e.g. <u>I said, How <math>\emptyset</math> you doing?</u> )	[52]	gerundial complementation after BEGIN, START, CONTINUE, HATE, and LOVE (e.g. <u>I began taking an interest</u> )	1549
1526					1550
1527	[44]	non-standard WAS (e.g. <u>three of them was killed</u> )			1551
1528	[45]	non-standard WERE (e.g. <u>he were a young lad</u> )	[53]	zero complementation after THINK, SAY, and KNOW (e.g. <u>they just thought <math>\emptyset</math> it isn't for girls</u> )	1552
1530					1553
1531	I. Relativization		[54]	that complementation after THINK, SAY, and KNOW (e.g. <u>they just thought that it isn't for girls</u> )	1554
1532					1555
1533	[46]	<u>wh</u> -relativization (e.g. <u>the man who read the book</u> )			1558
1534			K. Word order and discourse phenomena		1559
1535	[47]	the relative particle <u>what</u> (e.g. <u>the man what read the book</u> )			1560
1536			[55]	lack of inversion and/or of auxiliaries in <u>wh</u> -questions and in main clause <u>yes/no</u> -questions (e.g. <u>where you put the shovel?</u> )	1561
1537	[48]	the relative particle <u>that</u> (e.g. <u>the man that read the book</u> )			1562
1538			[56]	the prepositional dative after the verb GIVE (e.g. <u>she gave [a job] to [my brother]</u> )	1563
1540					1565
1541	J. Complementation		[57]	double object structures after the verb GIVE (e.g. <u>she gave [my brother] [a job]</u> )	1566
1542					1567
1543	[49]	<u>as what</u> or <u>than what</u> in comparative clauses (e.g. <u>we done no more than what other kids used to do</u> )			
1544					
1568					