

# A multivariate spatial analysis of vowel formants in American English

Jack Grieve,<sup>1\*</sup> Dirk Speelman,<sup>2</sup> and Dirk Geeraerts<sup>2</sup>

<sup>1</sup> Centre for Forensic Linguistics, School of Languages and Social Science, Aston University, Birmingham, UK

<sup>2</sup> Quantitative Lexicology and Variational Linguistics Research Unit, Department of Linguistics, University of Leuven, Leuven, Belgium

This paper presents the results of a multivariate spatial analysis of thirty-eight vowel formant variables measured in 236 cities from across the contiguous United States, based on the acoustic data from the *Atlas of North American English*. The results of the analysis both confirm and challenge the results of the *Atlas*. Most notably, while the analysis identifies similar patterns as the *Atlas* in the West and the Southeast, the analysis finds that the Midwest and the Northeast are distinct dialect regions that are considerably stronger than the traditional Midland dialect region identified in the *Atlas*. The analysis also finds evidence that a vowel shift is actively shaping the language of the Western United States.

## 1. Introduction

Most research on regional linguistic variation in American English has been based on the subjective analysis of linguistic survey data (e.g. Kurath, 1949; Carver, 1987; Labov et al., 2006). The traditional and standard approach to data analysis in American dialectology involves manually analyzing maps that plot the values of numerous linguistic variables across a region to identify individual and common patterns of regional variation. Usually isoglosses are drawn to divide each map into the regions where the different values of the linguistic variable predominate. Common patterns of regional linguistic variation are then identified by searching for linguistic variables with isoglosses that follow similar paths. Finally, dialect regions are identified based on how these bundles of isoglosses divide the region.

Although the traditional approach to the analysis of regional linguistic variation follows a logical series of steps, each stage of the analysis ultimately relies on the judgment of the dialectologist. Descriptive statistics and replicable procedures are sometimes used to help guide the analysis, but key decisions such as the selection of the linguistic variables that define a particular region (e.g. Carver, 1987) or the design of the algorithm that generates isoglosses (e.g. Labov et al., 2006) are still based on the judgment of the dialectologist, making it difficult to replicate analyses and to choose between competing theories of dialect regions. A statistical approach to the analysis of regional linguistic variation can help to resolve issues such as these. Despite the

advantages of such an approach, statistical analysis is uncommon in American dialectology, perhaps because the statistical methods commonly used to analyze language variation and change do not allow for regional variation to be analyzed following the same series of steps as the traditional approach.

For example, the types of statistical methods commonly used in variationist sociolinguistics to analyze the relationships between linguistic variables and social variables are unsuitable to analyze the relationships between linguistic variables and regional variables, such as longitude and latitude, because these spatial relationships are often nonlinear. Linear patterns such as a change from the Northeast to Southwest can be identified using standard variationist methods, but other types of patterns, such as a single central cluster, cannot. While variationist methods are of limited use in regional dialectology, a quantitative approach to the analysis of regional linguistic variation known as *dialectometry* is common in Europe (see Séguéy, 1971, 1973; Goebel, 1982, 2006; Heeringa, 2004; Nerbonne, 2006). In dialectometry, patterns of aggregated regional linguistic variation are identified using replicable and statistically justified methods; however, because dialectometry does not follow the same steps as a traditional analysis, it does not allow for regional patterns to be identified in a way that is satisfactory to many dialectologists. In particular, dialectometry studies do not generally analyze individual linguistic variables or identify subsets of linguistic variables that exhibit similar patterns of regional variation.

Although the standard statistical methods used in variationist sociolinguistics and dialectometry cannot replace the traditional approach to the analysis of regional linguistic variation, a statistical approach known as a *multivariate spatial analysis* identifies patterns of regional linguistic variation following the same

---

\*Address for correspondence: Jack Grieve, Lecturer in Forensic Linguistics, Centre for Forensic Linguistics, School of Languages and Social Sciences, Aston University, Aston Triangle B4 7ET, Birmingham, UK.  
(Email: j.grieve1@aston.ac.uk)

series of steps as a traditional analysis (Grieve, 2009; Grieve et al., 2011). This new approach to the analysis of regional linguistic variation is based on spatial autocorrelation statistics (Grieve, 2009, 2011, 2012), which allow for significant patterns of spatial clustering to be identified in the values of individual linguistic variables in a manner that is similar to plotting isoglosses. The results of the spatial autocorrelation analysis are then subjected to a factor analysis to identify common patterns of regional variation in a manner that is similar to identifying bundles of isoglosses. A multivariate spatial analysis therefore identifies both individual and common patterns of regional linguistic variation in a way that is similar to the traditional approach to the analysis of regional linguistic variation.

This paper describes a multivariate spatial analysis of the acoustic vowel formant data from the *Atlas of North American English* (Labov et al., 2006). As one would expect, given the results of the *Atlas*, this analysis finds that most vowel formant variables are regionally patterned in American English. However, although the analysis identifies similar regional patterns as those presented in the *Atlas*, the method also identifies additional patterns that had previously gone unnoticed, which in some cases challenge traditional taxonomies and theories of American dialect regions.

## 2. Data

The multivariate spatial analysis reported in this paper was based on the acoustic vowel data gathered for the *Atlas of North American English* (Labov et al., 2006). Data collection for the *Atlas* took place during the 1990s through linguistic interviews conducted over the telephone with informants in cities from across English-speaking North America. On average two to four informants were selected per city, with each informant being interviewed for approximately thirty to forty-five minutes. In total, 762 informants were interviewed for the *Atlas*. The recordings for 439 of these informants were then subjected to an acoustic analysis in order to measure the values of forty-eight vowel formant variables (see Chapter 10 of the *Atlas* for the raw maps), which consist of the average formant 1 and formant 2 values for twenty-four distinct vowel measures, including a number of vowels measured in different phonological contexts.

The analysis reported here was restricted to thirty-eight of the forty-eight acoustic vowel formant variables from the *Atlas* (i.e. average formant 1 and formant 2 values for nineteen vowel measures). Ten vowel formant variables (formant 1 and formant 2 for /oy/, /aeh/, /eyr/, /uwr/, /ah/) were excluded from this analysis because they were missing data for over 10 percent of the locations in the dataset. Twelve of

the remaining vowel formant variables (formant 1 and formant 2 for /uh/, /u/, /ay0/, /iw/, /uwc/, /ohr/) were also missing data, but because this missing data accounted for less than 5 percent of the locations in the dataset, these variables were retained and the missing data were replaced by the mean value for that variable across all locations.<sup>1</sup> The nineteen vowel measures analyzed in this study are listed in Table 1, including the phonetic symbol from the *Atlas*, the IPA equivalent, and an example of the vowel in context. In addition, Table 1 is organized based on the system of vowel categorization used in the *Atlas*. This system distinguishes between three levels of height and between front and back and round and unrounded vowels, as well as between short and long vowels, with the long vowels being further divided into ingliding vowels and both front and back upgliding vowels. The nineteen vowel measures are also plotted in Figure 1 based on their average formant 1 and formant 2 values across the entire dataset.

The analysis reported here was also restricted to 402 of the 439 informants whose recorded interviews were subjected to an acoustic analysis. In particular, Canadian informants were excluded from the analysis to control for the influence of national linguistic variation and Alaskan informants were excluded from the analysis because as extreme geographical outliers their inclusion would confound the spatial analysis. In addition, the one speaker from Bloomington, Illinois was also excluded from the analysis because he is an extreme outlier on formant 2 for all vowels. As well as removing these informants, the dataset analyzed here was also pooled across all informants from the same city, reducing the number of cases in the dataset from 402 informants to 236 cities. Pooling the data by location is required for a multivariate spatial analysis but has several additional advantages. Most important, maps based on pooled data are easier to interpret, especially when there is considerable variation in the number of informants per location, as is the case here.

The final dataset analyzed in this study therefore consists of thirty-eight vowel formant variables measured across 236 cities from across the contiguous United States. Before subjecting this regional linguistic data matrix to a multivariate spatial analysis, the raw values of the thirty-eight vowel formant variables were mapped across the 236 locations. Examples for four of the variables are presented in Maps 1–4. Map 1 shows that /eyc/ tends to be raised in the North and lowered in the South and the West. Map 2 shows that /ae/ tends to be raised in the Midwest and lowered in the Northeast and the West. Map 3 shows that /oh/ tends to be backed in the East and fronted in the West. Finally, Map 4 shows that /ayv/ tends to be backed in the Midland and fronted across most of the rest of the

Table 1. Vowel measures

Vowel measure	IPA vowel	Context restrictions	Example	Length	Position	Height	Glide type
/i/	[i]		<i>bit</i>	Short	Front	High	
/e/	[ɛ]		<i>bet</i>	Short	Front	Mid	
/æ/	[æ]		<i>bat</i>	Short	Front	Low	
/u/	[u]		<i>book</i>	Short	Back	High	
/uh/	[ʌ]		<i>but</i>	Short	Back	Mid	
/o/	[ɑ]		<i>cot</i>	Short	Back	Low	
/iyc/	[i]	Word internally	<i>beat</i>	Long	Front	High	Front upglide
/eyc/	[eɪ]	Word internally	<i>bait</i>	Long	Front	Mid	Front upglide
/ayv/	[aɪ]	Before voiced consonants	<i>bide</i>	Long	Back	Low	Front upglide
/ay0/	[aɪ]	Before voiceless consonants	<i>bite</i>	Long	Back	Low	Front upglide
/iw/	[ju]		<i>suit</i>	Long	Front	High	Back upglide
/uwc/	[u]	Word internally	<i>boot</i>	Long	Back	High	Back upglide
/uwf/	[u]	Word finally	<i>boo</i>	Long	Back	High	Back upglide
/owr/	[ou]	Before /r/	<i>boar</i>	Long	Back	Mid	Back upglide
/owc/	[ou]	Before other consonants	<i>boat</i>	Long	Back	Mid	Back upglide
/aw/	[aʊ]		<i>bout</i>	Long	Back	Low	Back upglide
/ohr/	[ɔ]	Before /r/	<i>north</i>	Long	Back	Mid	Rounded inglide
/oh/	[ɔ]	Before other consonants	<i>caught</i>	Long	Back	Mid	Rounded inglide
/ahr/	[ɑ]	Before /r/	<i>start</i>	Long	Back	Low	Unrounded inglide

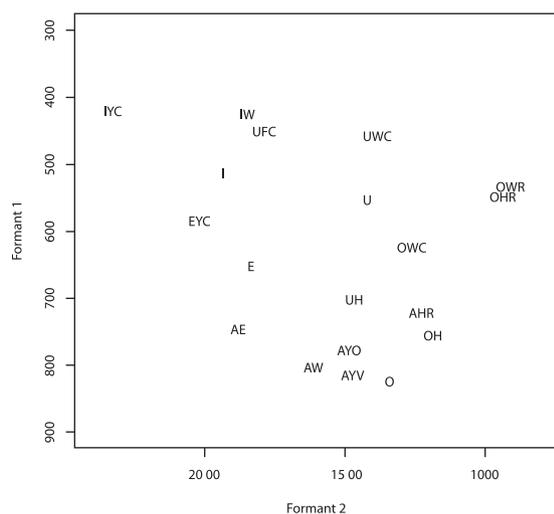


Figure 1. Average formant 1 and formant 2 values for all vowel measures.

United States. In each of these cases a regional pattern is discernible in the raw maps; however, because these patterns are not absolute, their significance is unclear.

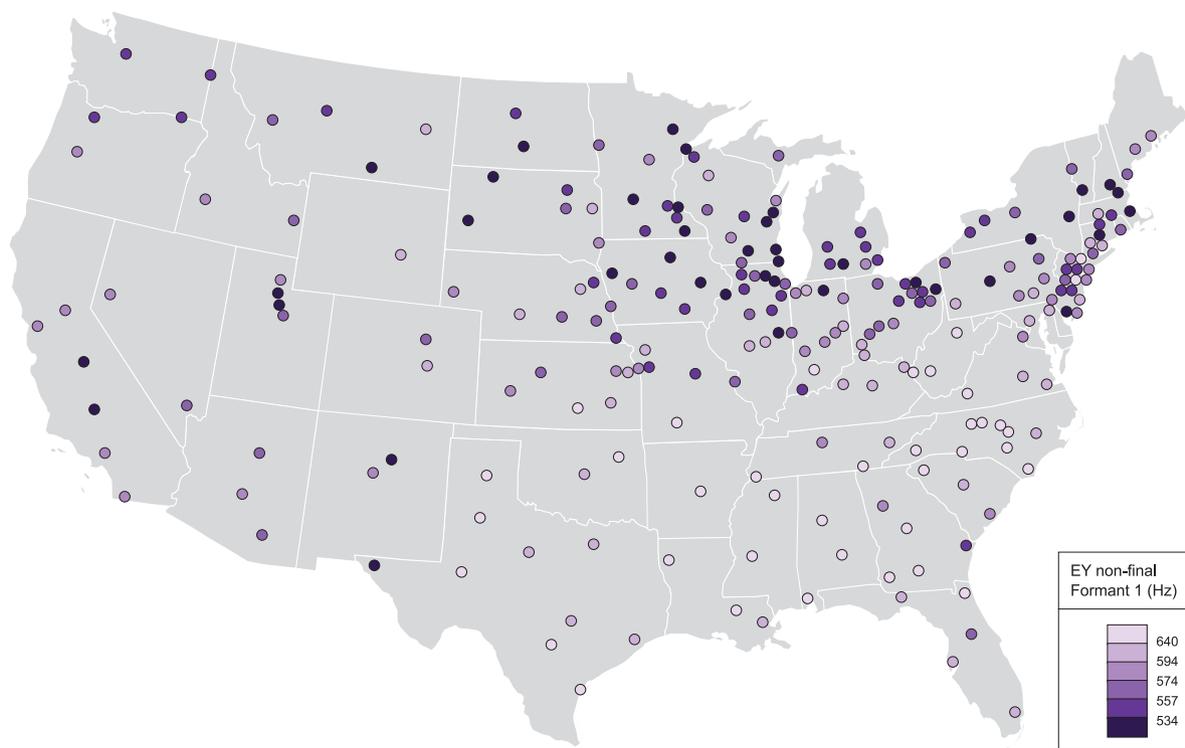
### 3. Spatial autocorrelation analysis

In order to identify statistically significant patterns of regional variation in the values of the thirty-eight individual vowel formant variables, each variable was subjected to a spatial autocorrelation analysis (Grieve, 2011, 2012; Grieve et al., 2011).<sup>2</sup> First, each variable

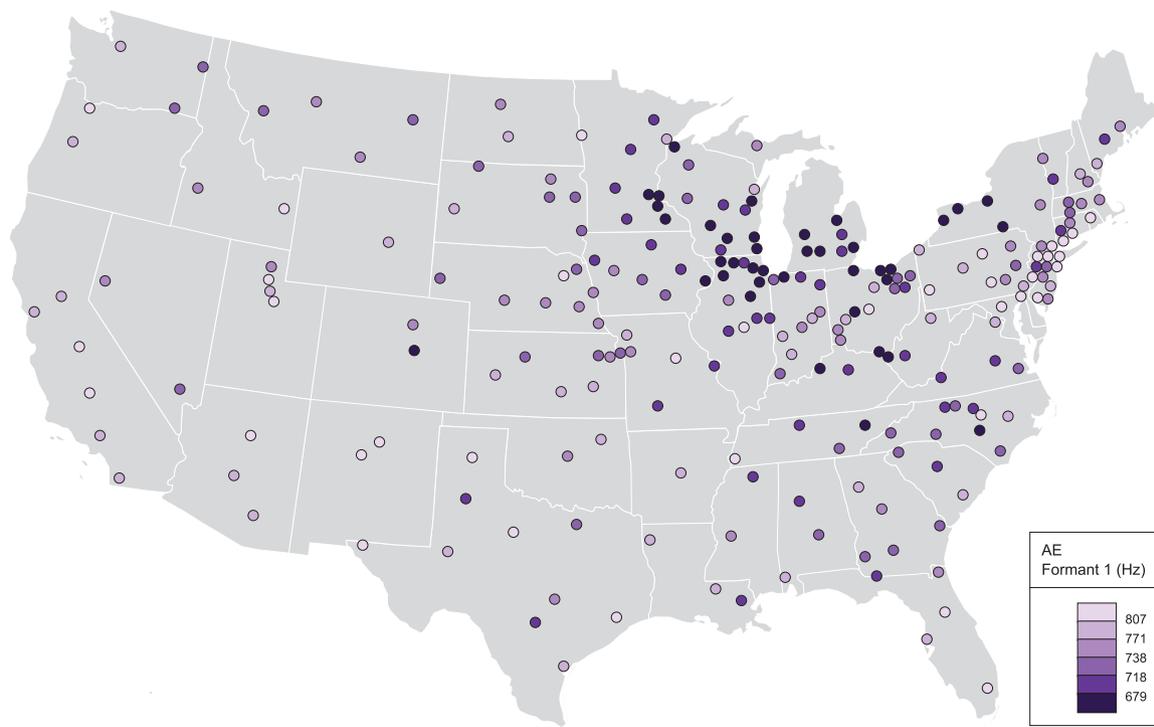
was subjected to an analysis of global spatial autocorrelation using global Moran's *I* (Moran, 1948; Odland, 1988) to identify variables exhibiting significant levels of spatial clustering. Second, each variable was subjected to an analysis of local spatial autocorrelation using local Getis-Ord *G<sub>i</sub>* (Ord and Getis, 1995) to identify the locations of any high- and low-value clusters.

To calculate both spatial autocorrelation measures, it is necessary to define a *spatial weighting* function, which is a set of rules that assigns a weight to the comparison of every pair of locations so that comparisons between locations that are close together are given greater weight than comparisons between locations that are far apart (Odland, 1988; Grieve, 2011, 2012; Grieve et al., 2011). In this study, a reciprocal spatial weighting function was used, which assigns a weight to each pair of locations based on the reciprocal of the distance between the locations so that the weight decreases with distance. A reciprocal weighting function was selected because it allows for the measures of spatial autocorrelation to be based primarily on the closest locations. This is important because some of the dialect regions identified by the *Atlas* are relatively narrow, such as the Midland dialect region or the extension of the Northern Cities region around the Great Lakes. A weighting function that focuses primarily on nearby locations is most suitable for replicating these types of results.

In order to interpret the significance of Moran's *I*, a standardized z-score was obtained under the assumption



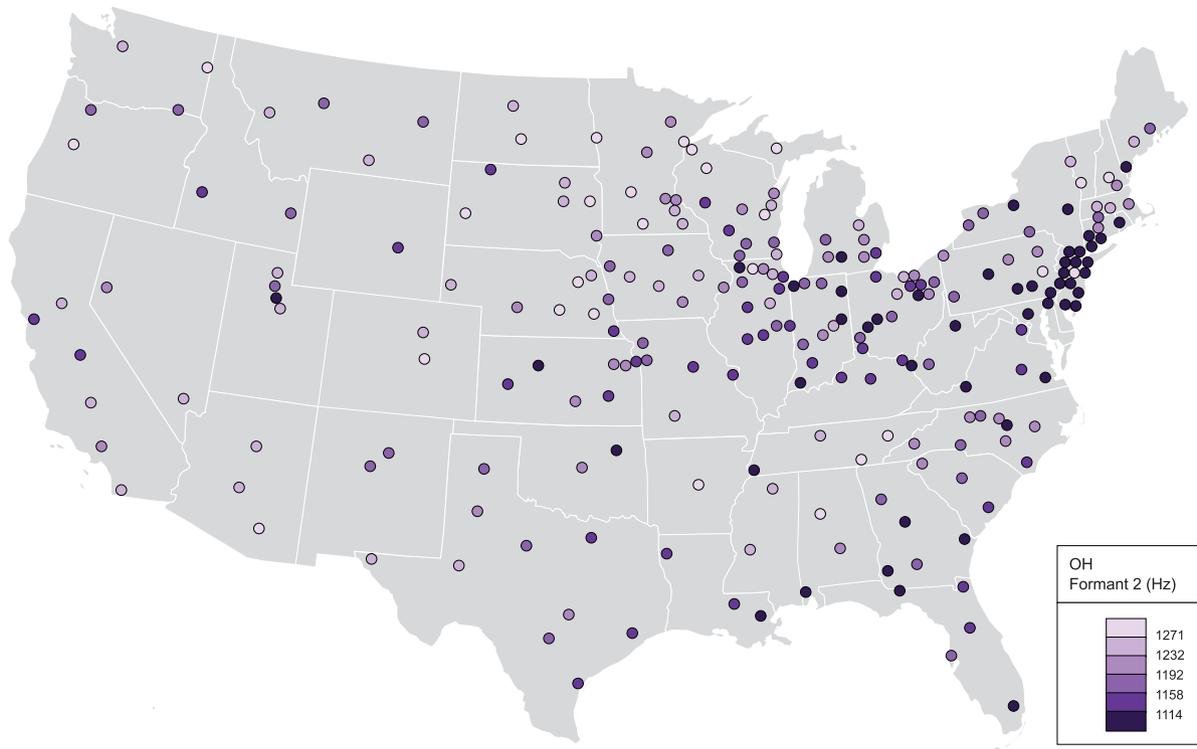
**Map 1.** Raw value map for non-word-final /ey/ (e.g. *bait*) on formant 1.



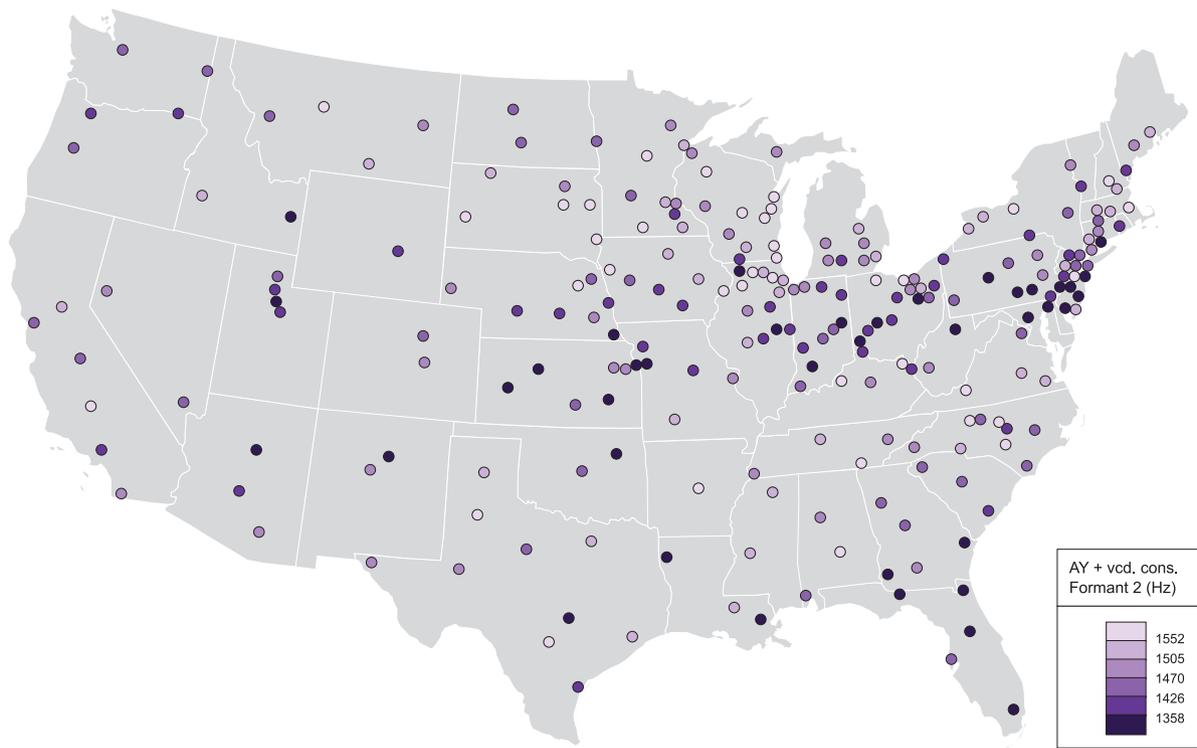
**Map 2.** Raw value map for /ae/ (e.g. *bat*) on formant 1.

of randomization (Odland, 1988). This z-score was then interpreted based on a one-tailed test of significance, because the goal of the analysis was to identify positive

spatial autocorrelation. A variable was deemed to exhibit significant global autocorrelation if the computed z-score was larger than or equal to 3.01—corresponding to a



Map 3. Raw value map for /oh/ (e.g. *caught*) on formant 2.



Map 4. Raw value map for /ayv/ before voiced consonants (e.g. *bide*) on formant 2.

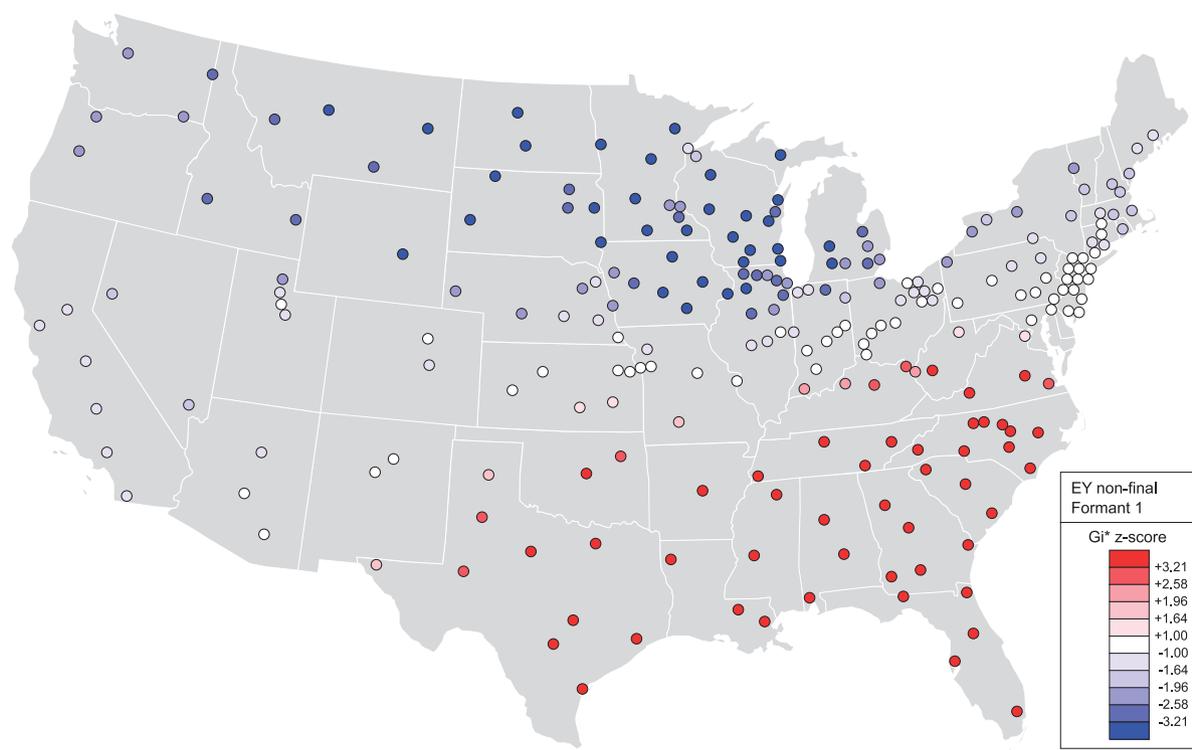
**Table 2.** Global autocorrelation results (reciprocal weighting function,  $n = 236$ )

Vowel	Formant	Mean (Hz)	Moran's I	z-score	$p$ (one-tailed)
/oh/	1	755	0.446	30.56	< 0.0001
/oh/	2	1177	0.303	20.82	< 0.0001
/aw/	2	1600	0.251	17.3	< 0.0001
/uwc/	2	1373	0.227	15.62	< 0.0001
/uwf/	2	1787	0.202	13.95	< 0.0001
/ae/	1	744	0.199	13.79	< 0.0001
/owc/	2	1267	0.195	13.47	< 0.0001
/ahr/	2	1227	0.169	11.76	< 0.0001
/o/	2	1334	0.169	11.74	< 0.0001
/iw/	2	1843	0.143	9.98	< 0.0001
/u/	2	1425	0.134	9.35	< 0.0001
/eyc/	1	584	0.134	9.34	< 0.0001
/uh/	2	1447	0.116	8.14	< 0.0001
/eyc/	2	2017	0.105	7.4	< 0.0001
/ae/	2	1869	0.098	6.9	< 0.0001
/uwf/	1	452	0.075	5.35	< 0.0001
/e/	2	1826	0.067	4.85	< 0.0001
/i/	2	1933	0.067	4.8	< 0.0001
/ay0/	1	777	0.066	4.75	< 0.0001
/ay0/	2	1481	0.054	3.96	< 0.0001
/ahr/	1	721	0.054	3.94	< 0.0001
/ayv/	2	1462	0.052	3.81	< 0.0001
/e/	1	653	0.047	3.48	0.0003
/iyc/	1	422	0.047	3.46	0.0003
/ohr/	2	925	0.041	3.08	0.001
<hr/>					
/o/	1	822	0.039	2.95	0.0016
/owr/	1	533	0.037	2.8	0.0026
/owr/	2	906	0.032	2.47	0.0068
/uh/	1	702	0.032	2.45	0.0071
/owc/	1	625	0.03	2.33	0.0099
/aw/	1	804	0.026	2.08	0.0188
/iyc/	2	2322	0.021	1.71	0.0436
/uwc/	1	456	0.02	1.63	0.0516
/i/	1	517	0.012	1.09	0.1379
/iw/	1	425	0.011	1.04	0.1492
/ohr/	1	548	0.01	0.95	0.1711
/u/	1	552	0.002	0.42	0.3372
/ayv/	1	813	0	0.27	0.3936

one-tailed 0.0013 significance level, which was selected based on a Bonferroni correction for thirty-eight variables ( $0.05/38 = 0.0013$ ). A Bonferroni correction controls for the fact that, every time a variable is added to the analysis, the likelihood that a significant pattern will be found by chance increases. A significant positive value for Moran's  $I$  indicates that the values of the variable exhibit spatial clustering, where nearby locations tend to have similar values at a greater degree than would be expected by chance. The results of the global spatial autocorrelation analysis are presented in Table 2, which for each vowel formant

variable lists its mean value, Moran's  $I$ , the associated z-score, and the associated one-tailed  $p$ -value. Based on the global spatial autocorrelation analysis, twenty-five out of the thirty-eight variables were found to exhibit significant levels of spatial clustering at the adjusted 0.0013 significance level.<sup>3</sup>

In addition to conducting a global spatial autocorrelation analysis, a local spatial autocorrelation analysis was conducted using local Getis-Ord  $G_i$  to identify the locations of spatial clusters in the values of each vowel formant variable. For each location, Getis-Ord  $G_i$  returns a z-score indicating the degree to which that location is



**Map 5.** Local autocorrelation map for non-word-final /ey/ on formant 1.

surrounded by locations with similar values. A significant negative Getis-Ord *Gi* z-score indicates that the location is part of a low-value cluster, whereas a significant positive Getis-Ord *Gi* z-score indicates that the location is part of a high-value cluster. A Getis-Ord *Gi* z-score was interpreted as significant if it was larger than or equal to  $\pm 3.21$ , which corresponds to a two-tailed 0.0013  $\alpha$  level, based on the Bonferroni correction described above, although in this case a two-tailed test of significance was used instead of a one-tailed test of significance, because the goal of the analysis was to identify both high- and low-value clusters.

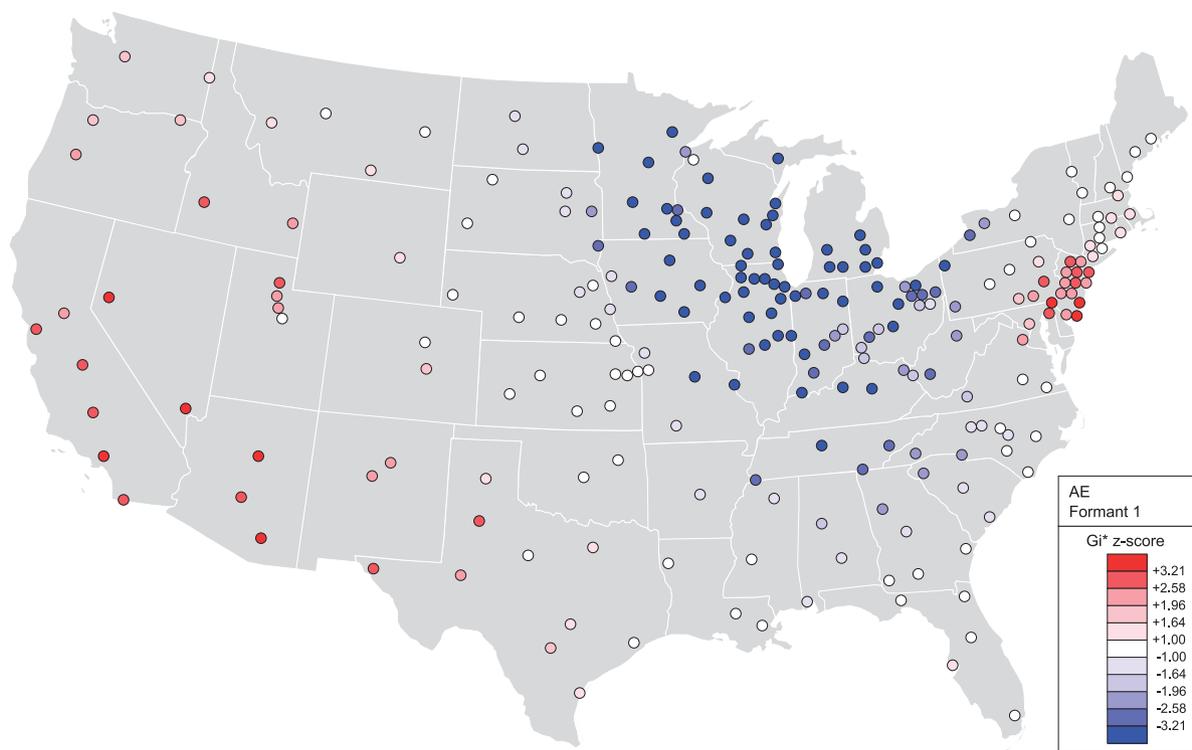
In order to visualize the patterns of spatial clustering identified by the local spatial autocorrelation analysis, the Getis-Ord *Gi* z-scores for each variable were plotted over the cities in the dataset. This is essentially a statistical method for plotting isoglosses. These local spatial autocorrelation maps identify clear high- and low-value clusters in the majority of the vowel formant variables. Local autocorrelation maps are provided for four variables in Maps 5–8, corresponding to the raw maps presented in Maps 1–4. These maps support the interpretation of the raw maps presented above: /eyc/ on formant 1 contrasts the South and West with the North (Map 5), /ae/ on formant 1 contrasts the Midwest with the Northeast and the West (Map 6), /oh/ on formant 2 contrasts the Northeast with the West (Map 7), and /ay/ on

formant 2 contrasts the Midland with the rest of the United States (Map 8).

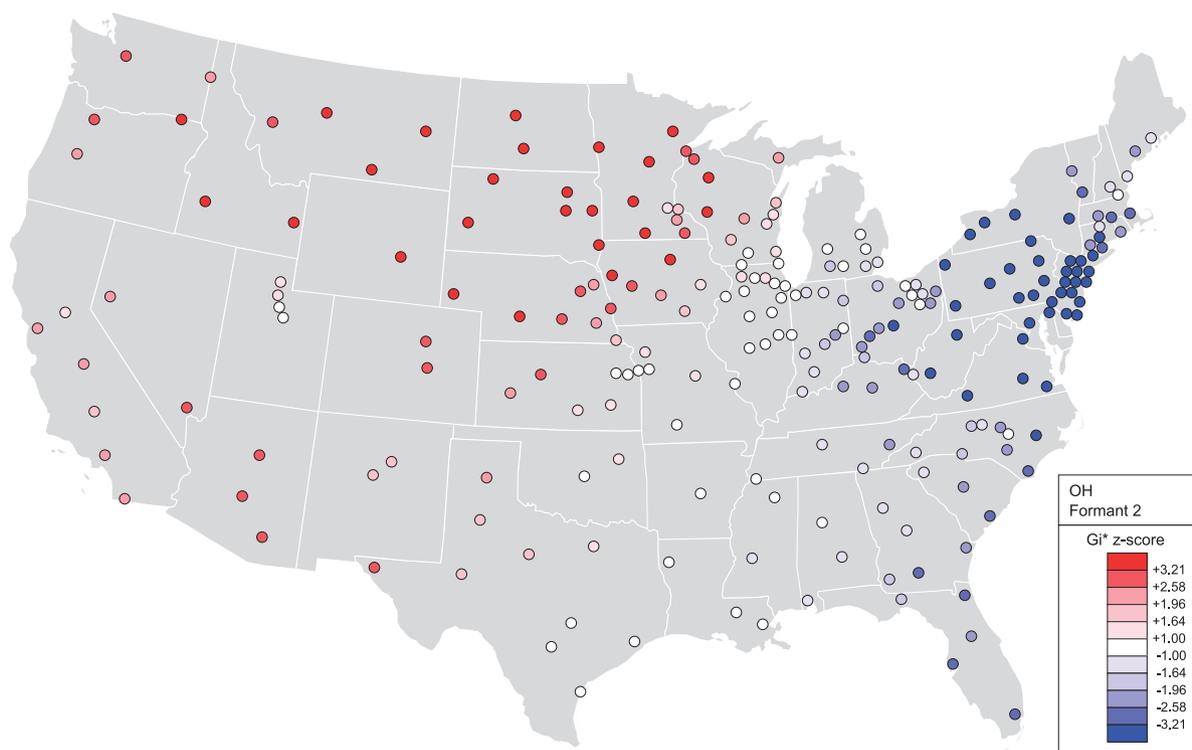
Before moving on, it is important to emphasize that conducting a local spatial autocorrelation analysis results in the loss of fine details present in the raw maps. For example, the value of a location that differs from the values of surrounding locations will be smoothed over by the local spatial autocorrelation analysis. That outlier location could be noise introduced through data collection or it could represent a true regional pattern, but in either case this variation will be lost. This is not a limitation of the local spatial autocorrelation analysis, which cannot identify spatial clusters if the dataset does not include a sufficient number of locations in that region for clusters to be formed. A regional linguistic dataset has a certain level of resolution and the local spatial autocorrelation maps must be interpreted accordingly.

#### 4. Factor analysis

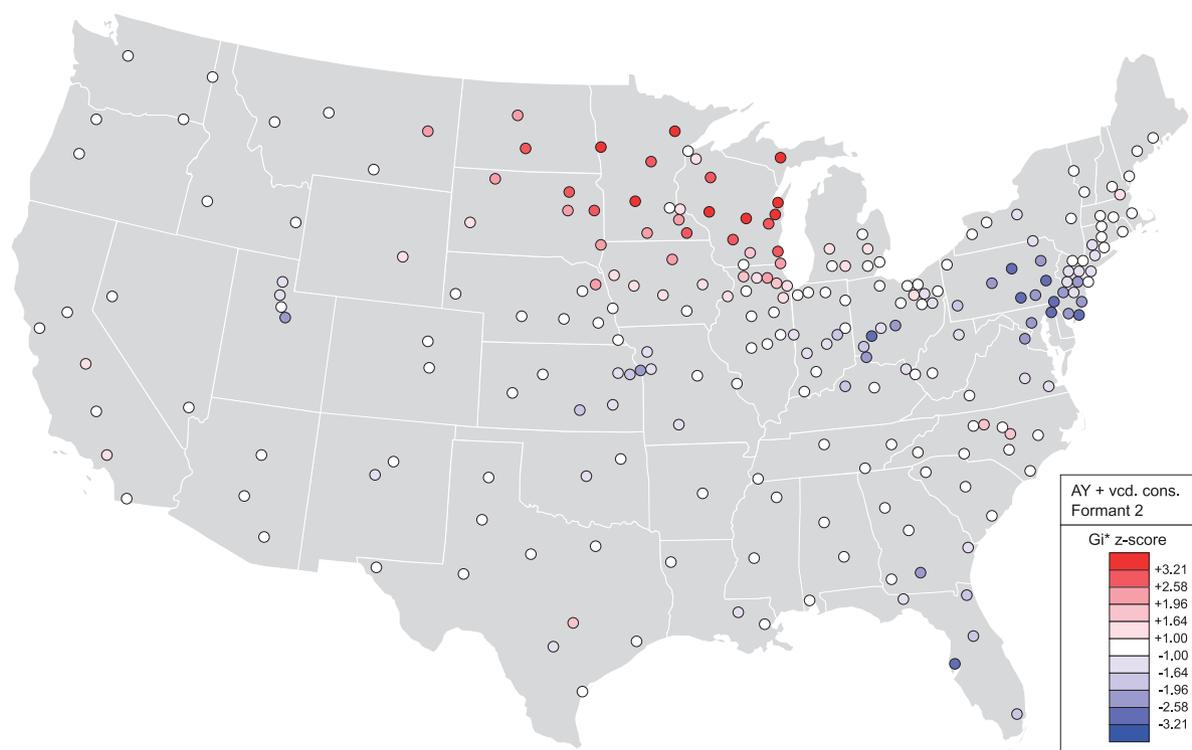
In order to identify common patterns of regional variation in the values of the thirty-eight vowel formant variables, the Getis-Ord *Gi* z-scores for the complete set of vowel formant variables were subjected to a factor analysis (Grieve et al., 2011). Given a set of variables measured over a set of cases, a factor analysis extracts a series of factors that each represent



Map 6. Local autocorrelation map for /ae/ on formant 1.



Map 7. Local autocorrelation map for /oh/ on formant 2.



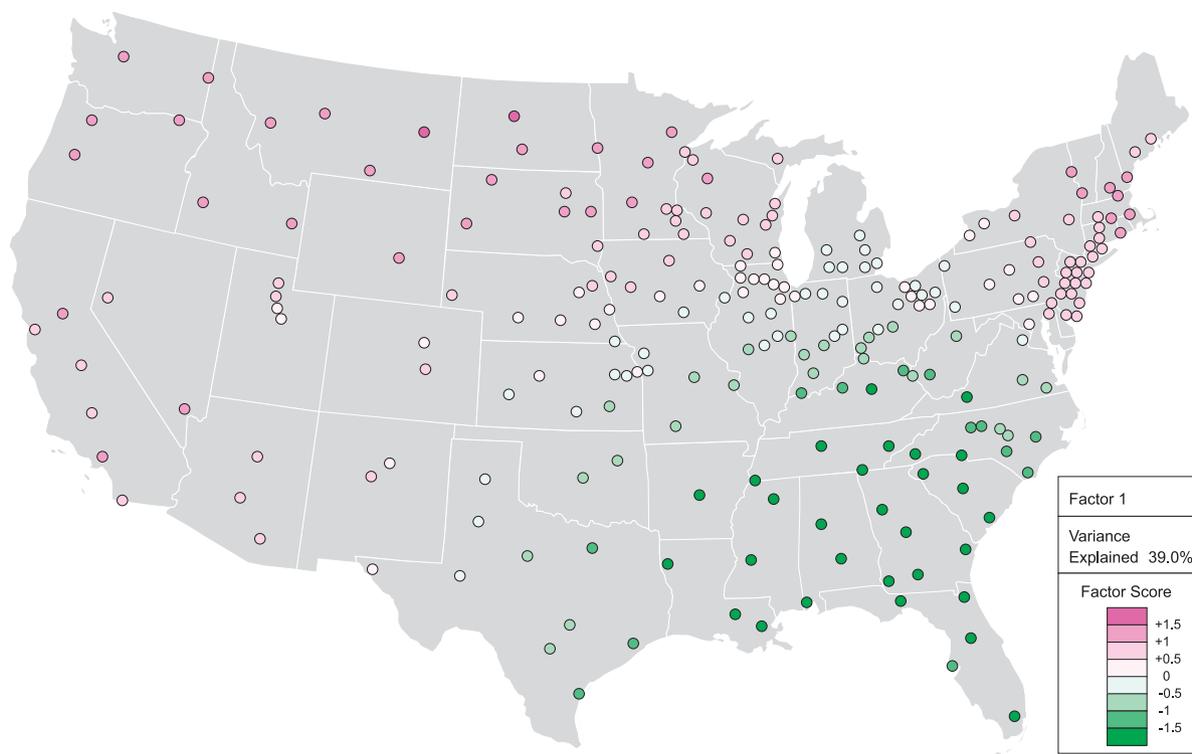
**Map 8.** Local autocorrelation map for /ay/ before voiced consonants on formant 2.

a common pattern of variation in that dataset, as well as the variables associated with each of these patterns (Hair et al., 2006). Because the local Getis-Ord  $G_i^*$  z-scores represent the location of spatial clusters in the values of the vowel formant variables, subjecting this dataset to a factor analysis identifies common patterns of spatial clustering, as well as the specific vowel formant variables that are associated with each of these patterns.<sup>4</sup> Subjecting the results of the local spatial autocorrelation analysis to a factor analysis is therefore essentially a statistical method for identifying bundles of isoglosses.

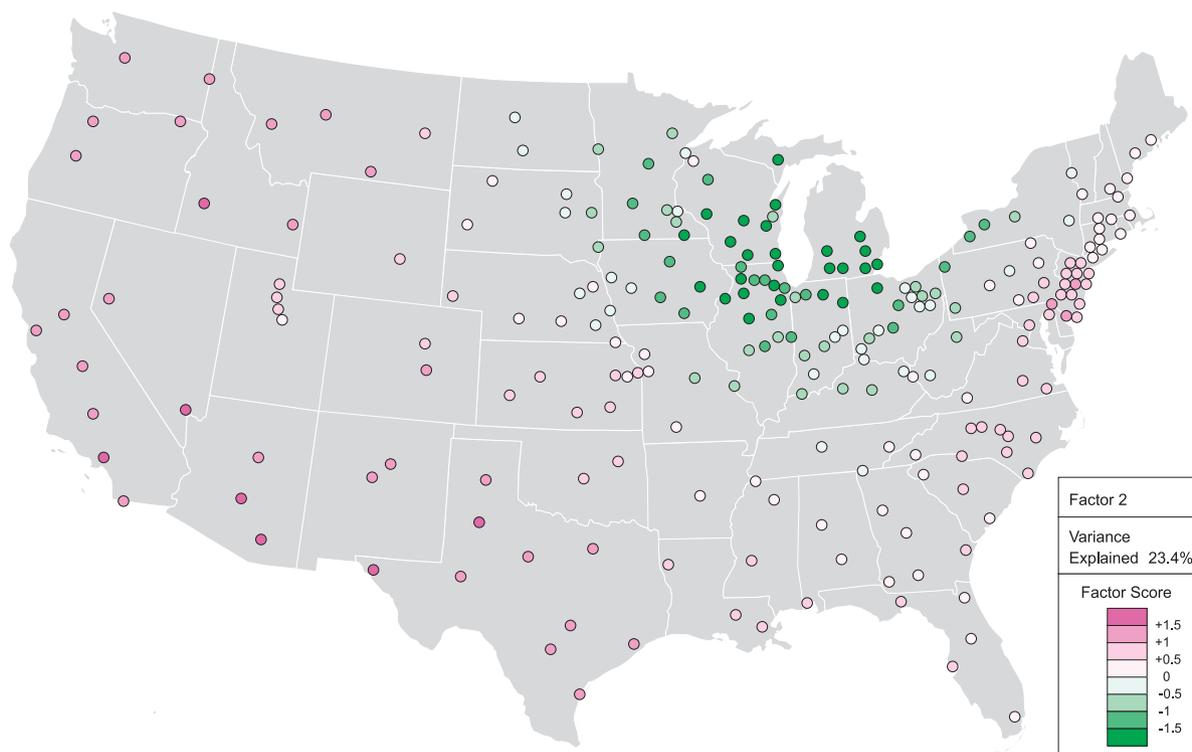
Before conducting a factor analysis, it is necessary to select the number of factors to be extracted, which can be determined by identifying the point where retaining further factors would explain relatively little additional variance. In this case four factors were selected because together these factors accounted for 86 percent of the variance in the values of the thirty-eight vowel formant variables, with additional factors explaining relatively little additional variance, indicating that the regional patterns exhibited by the thirty-eight vowel formant variables can largely be accounted for by these four basic patterns. The final factor analysis was thus run to extract four factors. In addition, the four factors were rotated using varimax rotation to limit the number of factors onto which each of the variables load, causing the factors to more clearly reflect the

spatial patterns visible in the local autocorrelation maps for the individual linguistic variables.

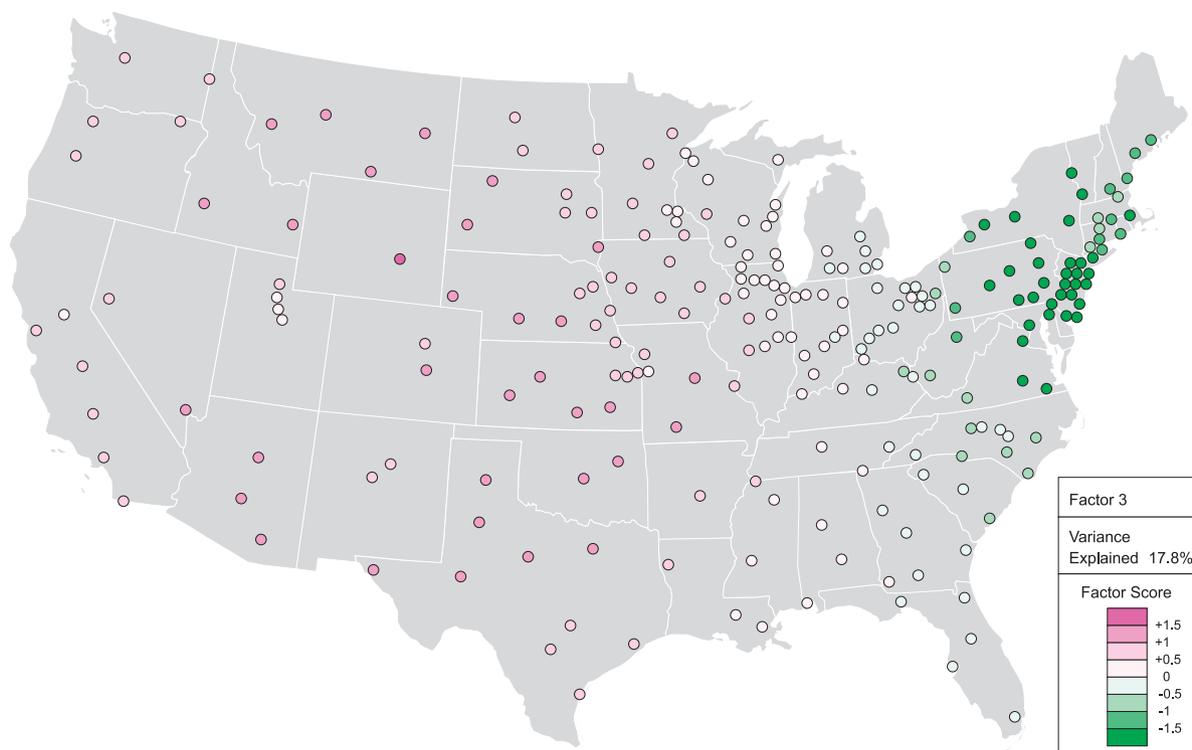
Each factor was analyzed in three ways. First, the factor scores for each factor were mapped across the 236 cities in the dataset in order to visualize the common patterns of spatial clustering represented by each factor. These factor maps, which are presented in Maps 9–12, contrast regions with positive factor scores (in magenta) to regions with negative factor scores (in green). Second, the factor loadings, which are presented in Table 3, were inspected. A high positive or negative factor loading indicates that the vowel formant variable exhibits the basic pattern of spatial clustering represented by that factor. In addition, for a formant 1 variable, a positive factor loading indicates that the vowel measure is lowered in regions with positive factor scores and raised in regions with negative factor scores, whereas a negative factor loading indicates that the vowel measure is raised in regions with positive factor scores and lowered in regions with negative factor scores. Similarly, for a formant 2 variable, a positive factor loading indicates that the vowel measure is fronted in regions with positive factor scores and backed in regions with negative factor scores, whereas a negative factor loading indicates that the vowel measure is backed in regions with positive factor scores and fronted in regions with negative factor scores. Table 3 also lists the uniqueness values for each vowel formant



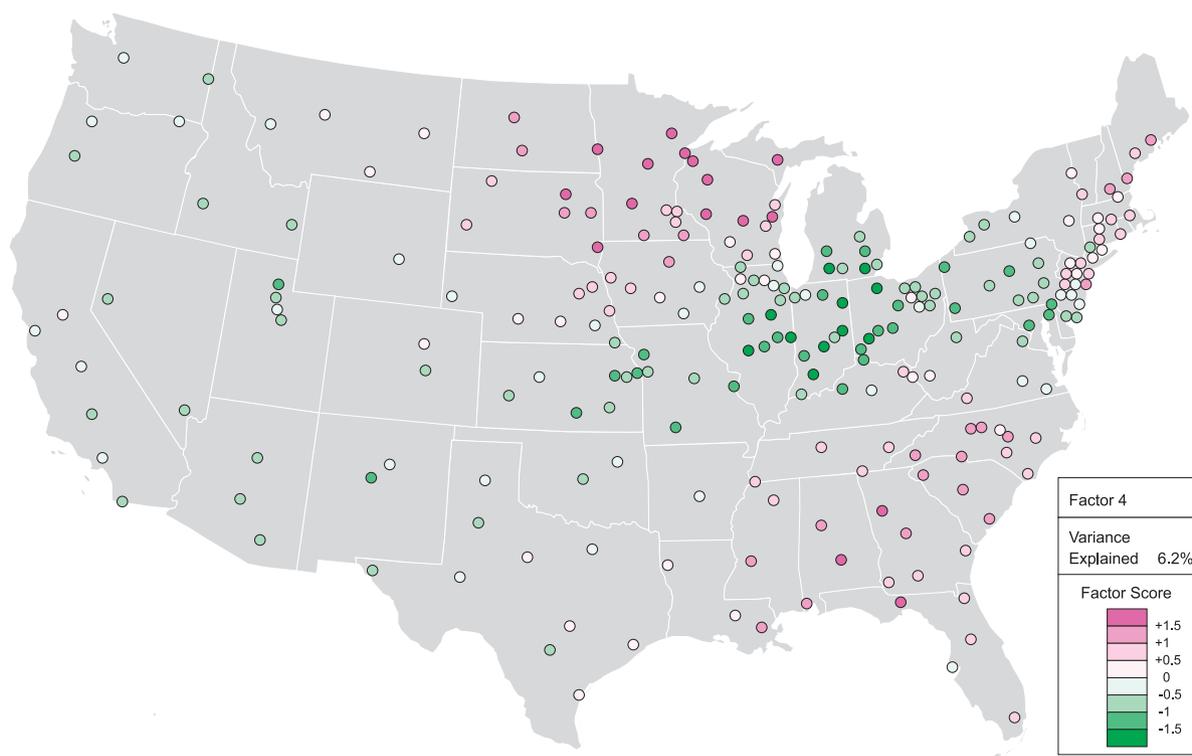
Map 9. Factor 1.



Map 10. Factor 2.



Map 11. Factor 3.



Map 12. Factor 4.

**Table 3.** Factor loadings (> 0.300) and uniqueness values

Vowel	Formant	Uniqueness	Factor 1	Factor 2	Factor 3	Factor 4
/i/	1	0.214	0.828			
/i/	2	0.115		0.898		
/e/	1	0.175	0.691		-0.413	-0.370
/e/	2	0.097		0.871		
/ae/	1	0.051		0.927		
/ae/	2	0.206		-0.852		
/o/	1	0.406	0.382	-0.595		
/o/	2	0.140	0.352	-0.839		
/uh/	1	0.156	0.915			
/uh/	2	0.068		0.627	0.714	
/u/	1	0.267	0.827			
/u/	2	0.080		0.496	0.812	
/iyc/	1	0.069	-0.689		-0.522	0.325
/iyc/	2	0.061	0.839			-0.420
/eyc/	1	0.014	-0.866	0.448		
/eyc/	2	0.012	0.980			
/ayv/	1	0.478				0.650
/ayv/	2	0.214	0.317	-0.428	0.455	0.544
/ay0/	1	0.098	-0.524	0.595	0.498	
/ay0/	2	0.067	0.950			
/iw/	1	0.269			-0.762	0.324
/iw/	2	0.033	-0.758	0.521	0.333	
/uwc/	1	0.276	-0.388	0.736		
/uwc/	2	0.019	-0.786	0.517	0.308	
/uwf/	1	0.180		0.574	-0.642	
/uwf/	2	0.023	-0.618	0.484	0.520	-0.302
/owc/	1	0.147	-0.857			
/owc/	2	0.025	-0.822	0.503		
/aw/	1	0.300	0.644		0.372	-0.348
/aw/	2	0.022	-0.668	0.621	-0.370	
/oh/	1	0.064			0.947	
/oh/	2	0.024	0.315		0.901	
/ahr/	1	0.172	0.753		0.396	
/ahr/	2	0.093	0.761	-0.554		
/ohr/	1	0.148	0.490		0.634	0.395
/ohr/	2	0.183	0.774		0.362	
/owr/	1	0.116	0.676	-0.388	0.483	
/owr/	2	0.092	0.933			

variable, which in all cases are relatively low (far below 0.800) indicating that the four factors account relatively well for the regional patterns exhibited by all thirty-eight variables. Third, the positions of the nineteen vowel measures were plotted based on the average formant 1 and formant 2 values for the vowel measures in cities with the highest (>1.00) and lowest (<-1.00) factor scores for each factor. In essence, these plots, which are presented in Figure 2, show the vowel spaces for the average informants in each of the opposing regions identified by the four factors.<sup>5</sup>

Factor 1 accounts for 39.0 percent of the variance in the dataset and contrasts the Southeast with the North and the West, with the approximate area of transition

between these two regions running along the northern borders of the Virginias, through northern Ohio, Indiana and Illinois, and central Iowa (Map 9). Numerous vowel formant variables distinguish between these two regions, as indicated by the large number of variables that load strongly on Factor 1. Most notably, nine long vowels load strongly on Factor 1 for both formant 1 and formant 2, with /eyc/, /iyc/ and /ay0/ tending to be lowered and backed in the Southeast, /ahr/, /owr/ and /ohr/ tending to be raised and backed in the Southeast, /uwc/ and /owc/ tending to be lowered and fronted in the Southeast, and /aw/ tending to be raised and fronted in the Southeast. Formant 2 values for two additional long

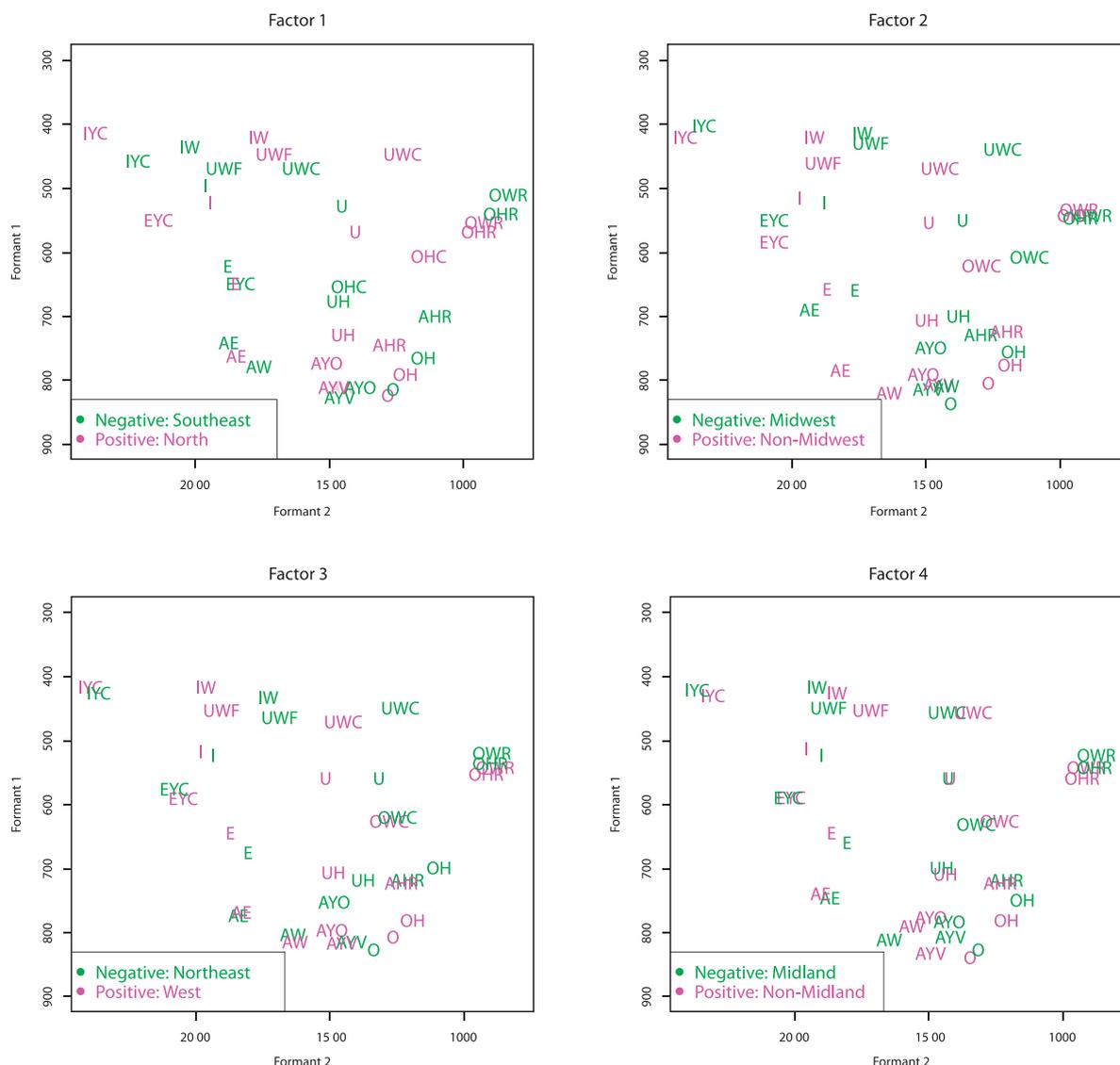
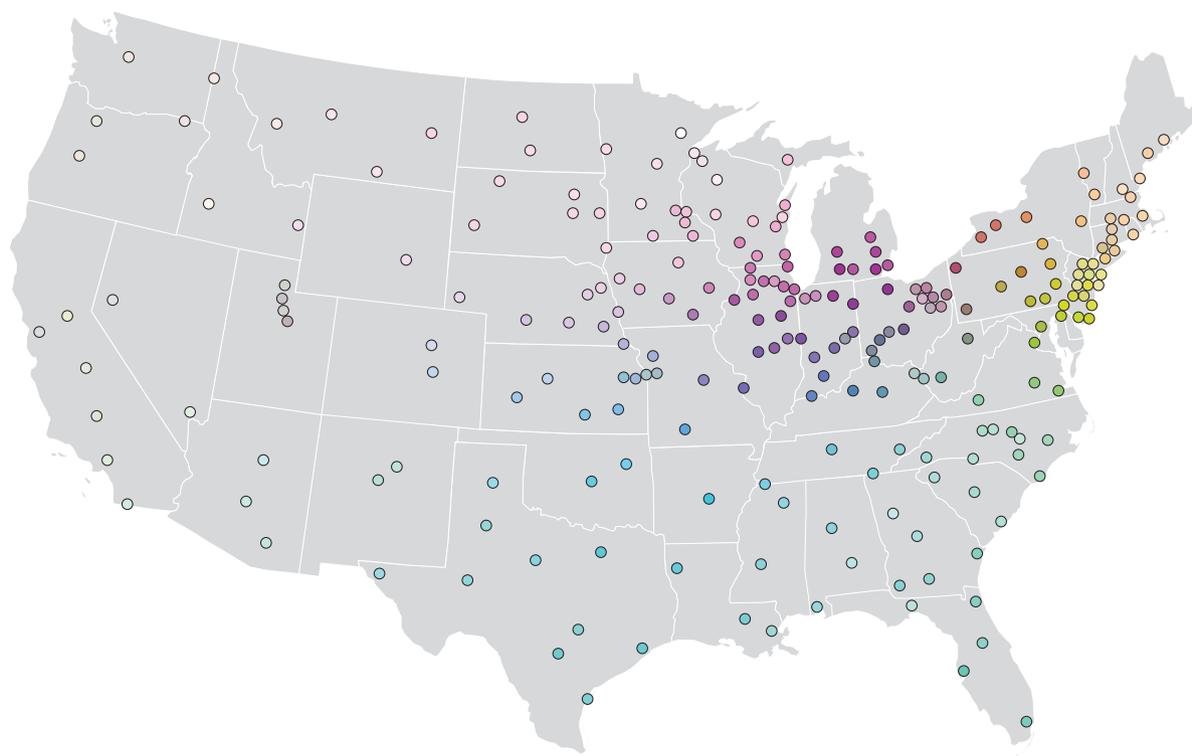


Figure 2. Average vowel positions for locations with high- and low-factor scores.

vowels also load strongly on Factor 1, with /iw/ and /uwf/ tending to be fronted in the Southeast. Finally, formant 1 values for four short vowels also load strongly on Factor 1, with /i/, /e/, /uh/ and /u/ all tending to be raised in the Southeast. In addition, other variables also load weakly on Factor 1 and shift slightly in the opposing vowel spaces for this factor, but these variables all load and shift to greater degrees on other factors.

Factor 2 accounts for 23.4 percent of the variance in the dataset and contrasts the Midwest with the rest of the United States, especially the West and the Mid-Atlantic (Map 10). The Midwest region identified here encompasses the core Midwestern states, but also stretches into New York, Pennsylvania, West Virginia, Kentucky, Tennessee, Missouri, Iowa, Nebraska, and the Dakotas. The factor loadings identify numerous

vowel formant variables that distinguish between the Midwest and the rest of the United States. Most notably, all six short vowels load on Factor 2. In particular, two short vowels load strongly on Factor 2 for both formant 1 and formant 2, with /ae/ tending to be raised and fronted in the Midwest and with /o/ tending to be lowered and fronted in the Midwest. Four additional short vowels also load strongly for formant 2, with /i/, /e/, /uh/ and /u/ all tending to be backed in the Midwest. Seven long vowel formant variables also load strongly on Factor 2, with /iw/, /owc/ and /aw/ tending to be backed in the Midwest, /uwc/ and /uwf/ tending to be raised and backed in the Midwest, /ahr/ tending to be fronted in the Midwest, and /ayo/ tending to be raised in the Midwest. In addition, other variables also load weakly on Factor 2 and shift slightly in the opposing vowel spaces for this



**Map 13.** CMYK map for Factors 1, 2, 3, and 4.

factor, but these variables all load and shift to greater degrees on other factors.

Factor 3 accounts for 17.8 percent of the variance in the dataset and contrasts the Northeast, which extends into eastern Ohio and Michigan as well as the Virginias and the Carolinas, with the West, especially the Great Plains and the Mountain States (Map 11). Ten long vowel measures load strongly on Factor 3. Most notably, /oh/ for both formant 1 and formant 2 loads strongly on Factor 3, with /oh/ tending to be lowered and fronted in the West. Similarly, /ohr/, /owr/ and /ay0/ also load strongly on Factor 3, with /ohr/ tending to be lowered and fronted in the West, and with /owr/ and /ay0/ tending to be lowered in the West. In addition, /iw/, /uwf/, /uwc/ and /iyc/ load strongly on Factor 3, with /iw/ and /uwf/ tending to be fronted and raised slightly in the West, with /uwc/ tending to be fronted in the West, and with /iyc/ tending to be very slightly raised in the West. Finally, three short vowels also load strongly on Factor 3, with both /u/ and /uh/ tending to be fronted in the West, and with /e/ tending to be raised in the West. In addition, other variables also load weakly on Factor 3 and shift slightly in the opposing vowel spaces for this factor, but these variables all load and shift to greater degrees on other factors.

Finally, Factor 4 accounts for 6.2 percent of the variance in the dataset and contrasts the Midland with

the rest of the country, especially the Southeast and the Upper Midwest (Map 12). The Midland as identified by Factor 4 stretches from Philadelphia, Baltimore and southern New Jersey through all of Pennsylvania and western New York State, and into northern West Virginia and Kentucky, southern Michigan, and all of Ohio, Indiana, Illinois, Missouri, and Kansas, and to a lesser extent across the West. Only two vowel formant variables load strongly on Factor 4, with /ayv/ tending to be both raised and backed in the Midland. In addition, other variables also load weakly on Factor 4 and shift slightly in the opposing vowel spaces for this factor, but these variables all load and shift to greater degrees on other factors. Overall, the position of all of the vowels have changed far less in the average vowel spaces for the Midland and the non-Midland regions identified by Factor 4, than in the other opposing vowel spaces identified by the factor analysis. This is not surprising given the relatively small amount of variance explained by this factor.

In addition to the individual factor maps, the four sets of factor scores were mapped simultaneously using CMYK (cyan, magenta, yellow, black) mapping, which is presented in Map 13. A hue was defined for each location representing the scores of all four factors at that location by associating each factor with one of the four CMYK color parameters. These hues were then mapped across the 236 cities to produce a single

overall picture of continuous regional linguistic variation in the dataset. Map 13 shows a clear regional pattern, plainly derived from the four factor maps reproduced in Maps 9–12. This aggregated factor map identifies at least four clear clusters, consisting of the Northeast, the Midwest, the Southeast, and the West. Texas also stands out as a region of transition between the West and the South. In addition, within the Northeast there is a clear distinction between New England and the Mid-Atlantic States, while within the Midwest there is a clear distinction between the Lower and Upper Midwest.

### 5. Cluster analysis

In addition to aggregating the four sets of factor scores using CMYK mapping to produce an overall map of continuous regional linguistic variation, it is also possible to identify dialect regions by clustering the locations based on their factor scores using an agglomerative hierarchical cluster analysis (HCA) (Hair et al., 2006), as is common in dialectometry (e.g. see Prokić & Nerbonne, 2008; Nerbonne & Heeringa, 2009; Wieling & Nerbonne, 2010; Grieve et al., 2011). In particular, Ward's method for hierarchical clustering (Ward, 1963) was used because it tends to identify the clearest dialect regions and because it is one of the most common methods for hierarchical clustering in dialectometry.<sup>6</sup> The results of the cluster analysis are represented by a tree diagram called a *dendrogram*, which shows the order in which the clusters were formed, and which can be used to identify clusters and subclusters of observations in the dataset. These clusters can then be mapped as dialect regions. Subjecting the results of the factor analysis to a cluster analysis is therefore essentially a statistical method for identifying dialect regions based on how bundles of isoglosses divide a region.

Based on the dendrogram reproduced in Map 13, the HCA identified five clear dialect regions: the Northeast, the Southeast, the West, the Lower Midwest, and the Upper Midwest, with the Upper and Lower Midwest further combining to form a Midwestern super-region. These five dialect regions are mapped in Figure 3 and the average vowel spaces for these five dialect regions are plotted in Map 14. Although the cluster analysis also groups the Northeast with the Southeast, and the Midwest with the West, because these two super-regions are formed so late in the cluster analysis they are not particularly meaningful. On the other hand, the most distinct internal clusters within the major clusters identified above are important, although as one descends further down the dendrogram the clusters begin to lose spatial consistency. First, Texas and the South Central States are separated from the rest

of the West. Second the Lower Midwest is divided into northern and southern subregions. Third, the Northeast is divided into New England and the Mid-Atlantic States. Fourth, the Upper Midwest is divided into northern and southern subregions.

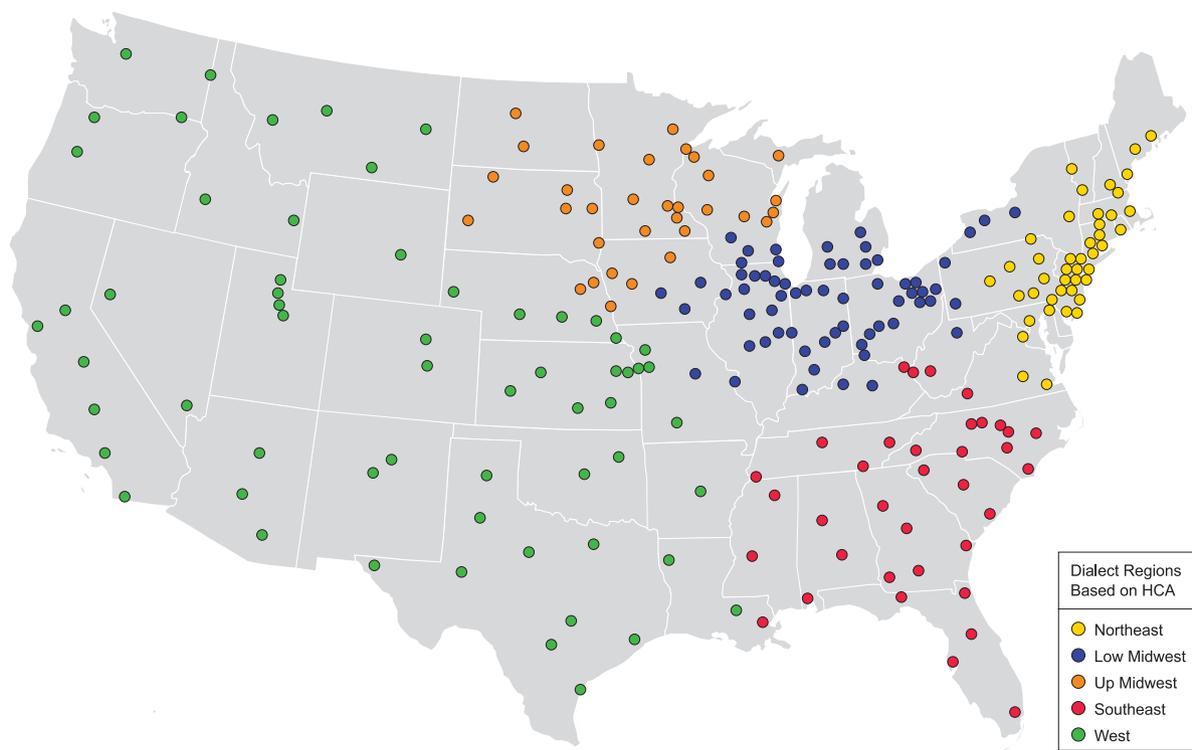
### 6. Discussion

The multivariate spatial analysis of the thirty-eight vowel formant variables identified clear and consistent patterns of regional variation in American English. This basic result was to be expected—the *Atlas* has already shown that vowel formant variables are regionally patterned in American English—but the analysis presented here has identified a somewhat different picture of American dialect regions than is presented in the *Atlas*. For comparison, the approximate dialect regions identified in the *Atlas* are presented in Map 15 based on map 11.5 and figure 11.9 from the *Atlas*.

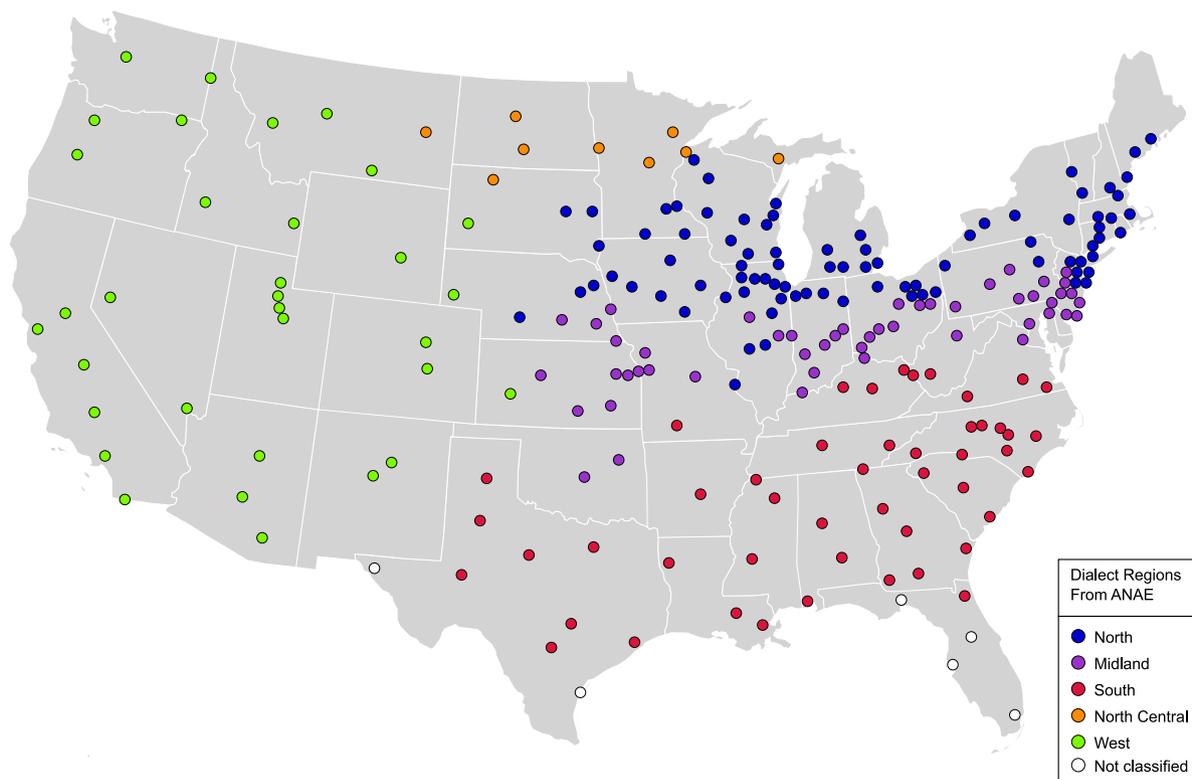
The multivariate spatial analysis identified four major patterns of regional variation. Factor 1 contrasts the Southeast with the North (see Map 9). The Southeastern region is similar to the Southern region identified in the *Atlas*, although the region identified here extends further north. The Northern region is also similar to the Northern region identified in the *Atlas*, although the region identified here is larger, stretching from the East Coast to the Southwest. Factor 1 also identifies an inland zone of transition that is similar to the Midland as identified in the *Atlas*. Factor 2 contrasts the Midwest with the East Coast and the West (see Map 10). Unlike the Northern region identified by Factor 1, the Midwest stretches further south and does not extend to either coast. The *Atlas* does not identify a distinct Midwestern region, but the Inland North region identified in the *Atlas* is at the core of the Midwestern region identified here. Factor 3 contrasts the Northeast with the West, especially the Great Plains and the Mountain States (see Map 11). The Western region is very similar to the Western region identified in the *Atlas*; however, no North-eastern region is recognized by the *Atlas*, which divides that part of the United States between the North and the Midland. Finally, Factor 4 contrasts the Midland with the rest of the United States (see Map 12). The Midland region is similar to the Midland region identified in *Atlas*, although the region identified here is somewhat larger.

Based on these four common patterns of spatial clustering, two maps of American dialect regions were generated. First, a continuous picture of American dialect regions was generated through a CMYK mapping of the four sets of factor scores (Map 13). Second, a categorical picture of American dialect





Map 14. Hierarchical cluster analysis 5 cluster map based on Factors 1, 2, 3, and 4.



Map 15. Dialect regions in the *Atlas*.

regions was generated through a cluster analysis based on the four sets of factor scores (Map 14), which is easier to interpret and is especially useful for comparing the results of this analysis to the dialect regions identified in the *Atlas* (see Map 15). Overall, there are numerous similarities between these maps. Most notably, the West and the South are identified as dialect regions in both analyses and have very similar dimensions. The main difference between the two analyses lies in the Northeastern quarter of the country. The *Atlas* first divides the region north-to-south, identifying Northern and Midland dialect regions, whereas the multivariate spatial analysis first divides the region east-to-west.

The multivariate spatial analysis does not identify the Midland as a strong dialect region because three of the four factors, which together account for 80 percent of the regional variance in the dataset, separate the Mid-Atlantic from the Lower Midwest. Only Factor 4, which accounts for 6 percent of the regional variance in the dataset, combines the Mid-Atlantic with the Lower Midwest, which is the defining characteristic of the traditional Midland dialect region (e.g. see Kurath, 1949). Because this analysis is not based on the complete dataset analyzed in the *Atlas*,<sup>7</sup> the strength of the Midland may be underestimated here; however, most English vowels have been included in this analysis and it therefore appears that the distinction between the Northeast and the Midwest is stronger than the distinction between the North and the Midland.<sup>8</sup> Nevertheless, a weak Midland signal is present in this dataset. Given the consistent identification of a strong Midland dialect region in previous American dialect surveys, this result suggests that the traditional Midland dialect region is in the process of being replaced by Northeastern and Midwestern dialect regions.

While the patterns of regional linguistic variation identified by the multivariate spatial analysis differ in some ways from the patterns identified in the *Atlas*, it appears that these results can largely be explained by chain shifts (Gordon, 2002) as proposed in the *Atlas* (Labov, 2004; Labov et al., 2006). In particular, the first two common patterns of regional variation identified by the factor analysis, which account for a majority of the regional variance in the dataset, identify both the regions and the vowel formant variables associated with the Southern Shift and the Northern Cities Shift, respectively.

The Southern Shift explains the majority of vowel formant variables that load on Factor 1 (see Table 3). As described in the *Atlas*, the Southern Shift begins with the fronting of /ay/, followed by the lowering and backing of /ey/ and /iy/, the raising and fronting of /i/, /e/ and /ae/, and the raising of /ahr/. In

addition, the fronting of /uw/ and /ow/ and the raising of /ohr/ are sometimes associated with this shift. Aside from /ae/, all of these vowel formant measures load on Factor 1, although sometimes only on one of the predicted formants. The Southeastern region identified by Factor 1 is also characterized by a vowel space where all of these vowels have shifted as predicted by the Southern Shift, with the exception of /ay/ (see Figure 2). While three of the four /ay/ vowel formant variables do load on Factor 1, /ay/ is not fronted in the Southeastern vowel space as predicted by the Southern Shift, but rather is lowered and backed. This is particularly surprising given that /ay/ fronting is the first step of the Southern Shift. In addition to these predicted vowel measures, /aw/ also loads on Factor 1 and is fronted in the Southeast. Although this vowel is not associated with the Southern Shift, its movement appears to be in line with the shift nonetheless. In addition, /u/ and /uh/ both load on Factor 1, with both vowels being higher in the Southeast, perhaps filling the space left behind by the fronting of /uw/.

Similarly, the Northern Cities Shift explains most vowel formant variables that load on Factor 2 (see Table 3). As described in the *Atlas*, the Northern Cities Shift begins with the fronting and raising of /ae/, followed by the fronting of /o/ and the lowering of /oh/, and the backing of /e/, /uh/, and /i/. Aside from /oh/, all of the vowel measures involved in the Northern Cities Shift load on Factor 2, although sometimes only on one of the predicted formants. The Midwestern region identified by Factor 2, which is centered around the Northern Cities, is also characterized by a vowel space where all of these vowels have shifted as predicted by the Northern Cities Shift, with the exception of /oh/, which is slightly higher in the Midwest (see Figure 2). Several additional vowel measures also load on Factor 2. In particular, /uwc/, /uwf/, /iw/, /u/, and /owc/ are all backed in the Midwestern vowel system. Although none of these variables are associated with the Northern Cities Shift, these movements appear to be related to the backing of /e/, /uh/, and /i/. In addition, /aw/ is backed, /ay0/ is raised, and /ahr/ is fronted in the Midwestern vowel system; however, the relationship, if any, between these changes and Northern Cities Shift is unclear.

While the Southern and Northern Cities Shifts described in the *Atlas* explain the majority of vowel formant variables that load on the first two factors, no chain shift discussed in the *Atlas* can explain the variables that load on Factor 3 (see Table 3), which contrasts the West, especially the Great Plains and the Mountain States, with the Northeast. Most notably, the West is associated primarily with the lowering and

fronting of /oh/ (resulting in the low-back merger with /o/) as well as the lowering and fronting of /ohr/, the fronting and slight raising of /iw/ and /uwf/, the lowering of /ay0/, the raising of /e/, and the fronting of /uh/, /u/ and /uwc/. All of these movements can be seen by comparing the average vowel space for the Western region identified by Factor 3 to the average vowel space for the Northeastern region identified by Factor 3 (see Figure 2). Some of these features are identified by the *Atlas* as being characteristic of western speech, specifically the fronting of /uw/ (as well as the lack of fronting of /ow/ and /aw/, which accompany /uw/ fronting in the Southeast) and the fronting and lowering of /oh/ (resulting in the low back merger with /o/). However, the other vowel formant variables loading on Factor 3 are not listed as western features in the *Atlas* nor is a western vowel shift identified in the *Atlas*.<sup>9</sup>

A western shift has been discussed in other studies, which have analyzed the language spoken in several western states including California (Hinton et al., 1987; Fought, 1999), Utah (Di Paolo, 1988), Nevada (Fridland, 2008), Arizona (Hall-Lew, 2004, 2005), Oregon (Conn, 2000; Ward, 2003), and Texas (Koops, 2010). This Western Shift is defined somewhat differently in these various studies but primarily involves the fronting of /uw/ and /ow/ (Hinton et al., 1987; Ward, 2003; Hall-Lew, 2004, 2005; Fridland, 2008), as well as the fronting of /u/ (Fought, 1999; Hall-Lew, 2004, 2005; Fridland, 2008; Koops, 2010), and occasionally the fronting of /o/ (Ward, 2003) and the raising of /ae/ (Hall-Lew, 2004, 2005; although cf. Conn, 2000).

The acoustic data from the *Atlas*, however, tells a somewhat different story. The fronting of /uw/ and /u/ are both identified by Factor 3 as being strong western features, but /ae/, /o/ and /ow/ do not load on Factor 3. In fact, /ae/ is at its lowest average position in the West, and /o/ is near its backest average position in the West (see Figure 4). Furthermore, as discussed in the *Atlas*, the relative stability of /ow/ in the West is a defining feature of the region, compared with the Southeast, for example, where /uw/ fronting is accompanied by /ow/ fronting (as identified by Factor 1). However, as described above, numerous other vowel formant variables also load on Factor 3, which have not been identified as part of a Western Shift in previous research. In particular, the fronting of /uh/, /iw/ and /uwf/, the raising of /e/, and the lowering of /ay0/, all are identified as Western features by Factor 3 and all appear to be related to the fronting of /uw/ and /u/. Furthermore, it is possible that all of these changes are triggered by the low back merger (i.e. the fronting and lowering of /oh/), which is also identified by Factor 3 as

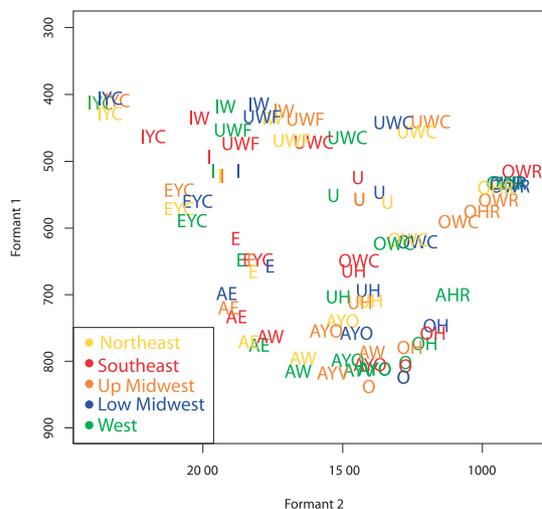


Figure 4. Average vowel spaces for the five clusters identified by the cluster analysis.

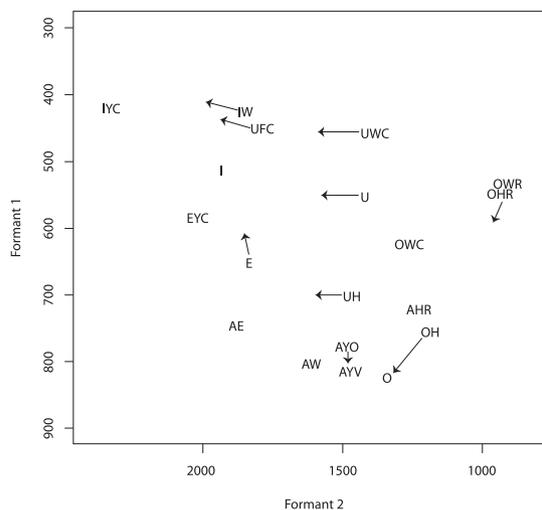


Figure 5. The Western Shift.

a Western feature. Given that these vowel shifts appear to form a chain of interrelated changes, it appears that Factor 3 may have identified a distinct Western Shift that involves a large number of vowels that has only been partially observed in previous research. This putative Western Shift is diagrammed in Figure 5.

It is important to note, however, that although the Western Shift identified here spans the West, it appears to be strongest in the Mountain States and the Great Plains, and weakest in the regions where most previous studies have been conducted, including California, Utah, Nevada, Arizona, and Oregon. It is therefore possible that there are two related but distinct shifts underway in the West, both perhaps a reaction to the low back merger, with a Californian Shift operating on the West Coast and in the Southwest, and with a Central Shift, as

identified by Factor 3, operating in the Mountain States and the Great Plains.

Finally, while chain shifts provide an internal explanation for the patterns identified by this analysis, the standard theory that American dialect regions correspond to historical settlement patterns does not provide a satisfactory external explanation for these results. Settlement patterns cannot account for the distinct Midwestern and Northeastern dialect regions identified in this study, because both regions were settled by people from the New England and Midland settlement hearths. On the other hand, the four major dialect regions identified here, including the Northeast and the Midwest, clearly correspond closely to the four major modern American cultural regions. This finding supports a cultural theory of American dialect regions, where American dialect regions correspond to contemporary cultural regions (see Grieve, 2009). This cultural theory of American dialect regions also explains the decline of the Midland dialect region and the emergence of Midwestern and Northeastern dialect regions identified in this study, which follow the apparent decline of the traditional Midland cultural region and emergence of the modern Midwestern and Northern cultural regions.

### Acknowledgments

The authors thank William Labov for his comments on this analysis and for providing the data upon which this analysis is based, along with Sharon Ash and Charles Boberg. They also thank the members of the QLVL Research Unit at the University of Leuven, Emily Waibel, and four anonymous reviewers for their comments on this paper.

### Notes

- <sup>1</sup> The imputation of missing data has very little effect on the results of the analysis because the missing data is minimal (0.02 percent of the total data) and because the variables with missing data follow the same basic patterns as the other vowel formant variables.
- <sup>2</sup> Each formant for a vowel is analyzed individually because vowels can pattern differently on the two formants.
- <sup>3</sup> It is notable that formant 2 variables tend to exhibit higher levels of global spatial autocorrelation than formant 1 variables.
- <sup>4</sup> A similar technique, known as a principal component analysis, was used in the *Atlas* to cluster informants; however, the analysis was based on the raw variables, the component loadings were not reported, and the component scores were not mapped. A factor analysis was used instead of principal component analysis in this study to focus on the identification of common patterns of regional variation (see Nerbonne, 2006; Grieve et al.,

2011), although a principal component analysis of this dataset would have produced similar results.

- <sup>5</sup> In all cases vowel measures that load strongly on a particular factor change position in the pair of opposing vowel spaces identified by that factor (see Figure 2). Nevertheless, vowel measures can also change position in the opposing vowel spaces identified by factors upon which they do not load if that change is smaller than the change in the opposing vowel spaces identified by factors upon which they do load.
- <sup>6</sup> There are several other possible hierarchical clustering algorithms that could have been applied (e.g. see Heeringa, 2004). The main reason why Ward's method is used here (and in many dialectometry studies) is because it tends to identify contiguous dialect regions, which is the goal of the cluster analysis, whereas most other clustering algorithms tend to identify clusters that include relatively large numbers of geographic outlier locations. In this case, it is acceptable to select the clustering algorithm that gives the best results. This is because a cluster analysis does not test if there are dialect regions in a regional linguistic dataset; dialect regions are assumed to exist (in this case based on the results of the spatial autocorrelation analyses) and the cluster analysis is used to identify their location. Consequently, clustering algorithms that tend to identify contiguous clusters of locations should generally be preferred to clustering algorithms that do not.
- <sup>7</sup> The analysis presented here is based only on the most complete acoustic vowel data available, whereas the *Atlas* is based on vowel formant variables that were excluded from this analysis due to missing data (see section 2), as well as additional phonetic and phonological measures.
- <sup>8</sup> The analysis presented here also aligns closely with the analysis of lexico-grammatical variation in a modern corpus of written American English (see Grieve et al., 2011; Grieve, 2013).
- <sup>9</sup> The vowel formant variables loading on Factor 3 were interpreted as identifying a Western Shift as opposed to a Northeastern Shift because the low back merger, which is identified by Factor 3 as characteristic of the Western region, is known to be a change in progress. However, it is also possible that Factor 3 has identified distinct Western and Northeastern features.

### References

- Carver, Craig M. 1987. *American regional dialects*. Ann Arbor: University of Michigan Press.
- Conn, Jeffrey C. 2000. *The story of /ae/ in Portland*. Portland, OR: Portland State University thesis.
- Di Paolo, Marianna. 1988. Pronunciation and categorization in sound change. In K. Ferrara, B. Brown, K. Walters & J. Baugh (eds), *Linguistic Change and Contact: Sixteenth Annual Conference on New Ways of Analyzing Variation*, 84–92. Austin, TX: Department of Linguistics, University of Texas.
- Fought, Carmen. 1999. A majority sound change in a minority community: /u/-fronting in Chicano English. *Journal of Sociolinguistics* 3: 5–23.

- Fridland, Valerie. 2008. Patterns of /uw/, /upsilon/, and /ow/ fronting in Reno, Nevada. *American Speech* 83: 432–454.
- Goebel, Hans. 1982. *Dialektometrie; Prinzipien und Methoden des Einsatzes der numerischen Taxonomie im Bereich der Dialektgeographie*. Vienna: Verlag der Osterreichischen Akademie der Wissenschaften.
- Goebel, Hans. 2006. Recent advances in Salzburg dialectometry. *Literary and Linguistic Computing* 21: 411–435.
- Gordon, Matthew J. 2002. Investigating chain shifts and mergers. In J. Chambers, P. Trudgill & N. Schilling-Estes (eds), *The handbook of language variation and change*, 245–266. London: Blackwell.
- Grieve, Jack. 2009. *A corpus-based regional dialect survey of grammatical variation in written Standard American English*. San Francisco, AZ: Northern Arizona University dissertation.
- Grieve, Jack. 2011. A regional analysis of contraction rate in written standard American English. *International Journal of Corpus Linguistics* 16: 514–546.
- Grieve, Jack. 2012. A statistical analysis of regional variation in adverb position in a corpus of written standard American English. *Corpus Linguistics and Linguistic Theory* 8: 39–72.
- Grieve, Jack. 2013. A statistical comparison of regional phonetic and lexical variation in American English. *Literary and Linguistic Computing* 8: 82–107.
- Grieve, Jack, Dirk Speelman & Dirk Geeraerts. 2011. A statistical method for the identification and aggregation of regional linguistic variation. *Language Variation and Change* 23: 193–221.
- Hair, Joseph F. Jr., William C. Black, Barry J. Babin, Rolph E. Anderson & Ronald L. Tatham. 2006. *Multivariate data analysis*, 6th edn. Englewood Cliffs, NJ: Prentice-Hall.
- Hall-Lew, Lauren. 2004. The western vowel shift in Northern Arizona. Unpublished Manuscript. <http://www.lcl.ed.ac.uk/~lhlew/index.html>, Accessed July 8, 2013.
- Hall-Lew, Lauren. 2005. One shift, two groups: When fronting alone is not enough. In M. Baranowski, U. Horesh, K. Evans & G. Nguyen (eds), *Selected papers from NWAV 32*: 105–116. Philadelphia PA: University of Pennsylvania Working Papers in Linguistics 10.2.
- Heeringa, Wilbert J. 2004. *Measuring dialect pronunciation differences using Levenshtein distance*. Groningen: University of Groningen dissertation.
- Hinton, Leanne, Birch Moonwomon, Sue Bremner, Herb Luthin, Mary Van Clay, Jean Lerner & Hazel Corcoran. 1987. It's not just valley girls: A study of California English. In J. Aske, N. Beery, L. Michaelis & H. Filip (eds), *Thirteenth Annual Meeting of the Berkeley Linguistics Society*, 117–127. Berkeley, CA: Berkeley Linguistics Society.
- Koops, Christian. 2010. /u/-fronting is not Monolithic: Two types of fronted /u/ in Houston anglos. University of Pennsylvania Working Papers in Linguistics 16. Article 14. <http://repository.upenn.edu/pwpl/vol16/iss2/14>.
- Kurath, Hans. 1949. *Word geography of the eastern United States*. Ann Arbor, MI: University of Michigan Press.
- Labov, William. 2004. *Principles of linguistic change, internal factors*. Hoboken, NJ: Wiley-Blackwell.
- Labov, William, Sharon Ash & Charles Boberg. 2006. *Atlas of North American English: Phonetics, phonology, and sound change*. New York: Mouton de Gruyter.
- Moran, P. A. P. 1948. The interpretation of statistical maps. *Journal of the Royal Statistical Society, Series B* 37: 243–251.
- Nerbonne, John. 2006. Identifying linguistic structure in aggregate comparison. *Literary and Linguistic Computing* 21: 463–476.
- Nerbonne, John & Wilbert J. Heeringa. 2009. Measuring dialect differences. In P. Auer & J. E. Schmidt (eds), *Language and space: An international handbook of linguistic variation, Vol. 1: Theories and Methods* (Handbooks of Linguistics and Communication Science 30.1), 550–567. Berlin/New York: De Gruyter Mouton.
- Odland, John D. 1988. *Spatial autocorrelation*. Beverly Hills, CA: Sage Publications.
- Ord, J. K. & Arthur Getis. 1995. Local spatial autocorrelation statistics: Distributional issues and an application. *Geographical Analysis* 27: 286–306.
- Prokić, Jelena & John Nerbonne. 2008. Recognizing groups among dialects. *International Journal of Humanities and Arts Computing* 1: 153–172.
- Séguy, Jean. 1971. La relation entre la distance spatiale et la distance lexicale. *Revue de linguistique romane* 35: 335–357.
- Séguy, Jean. 1973. La dialectométrie dans l'Atlas linguistique de la Gascogne. *Revue de linguistique romane* 37: 1–24.
- Ward, Joe H. Jr. 1963. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association* 58: 236–244.
- Ward, Michael. 2003. *Portland Dialect Study: The Fronting of /ow, u, uw/ in Portland, Oregon*. Portland, OR: Portland State University thesis.
- Wieling, Martijn & John Nerbonne. 2010. Hierarchical bipartite spectral graph partitioning to cluster dialect varieties and determine their most important linguistic features. TextGraphs-5 Workshop on Graph-Based Methods for NLP 16: 33–41.