

Using Gabmap

Therese Leinonen¹, Çağrı Çöltekin², and John Nerbonne²

¹University of Turku

²University of Groningen

February 2, 2015

Abstract

Gabmap is a freely available, open-source web application that analyzes the data of language variation, e.g. varying words for the same concepts, varying pronunciations for the same words, or varying frequencies of syntactic constructions in transcribed conversations. Gabmap is an integrated part of CLARIN (see e.g. <http://portal.clarin.nl>).¹. This article summarizes Gabmap's basic functionality, adding material on some new features and reporting on the range of uses to which Gabmap has been put. Gabmap is modestly successful, and its popularity underscores the fact that the study of language variation has crossed a watershed concerning the acceptability of automated language analysis. Automated analysis not only improves researchers' efficiency, it also improves the replicability of their analyses and allows them to focus on inferences to be drawn from analyses and other more abstract aspects of that study.

1 Introduction

Gabmap is a freely available, open-source web application that analyzes the data of language variation, e.g. varying words for the same concepts, varying pronunciations for the same words, or varying frequencies of syntactic constructions in transcribed conversations.

Other possibilities exist as well, but these are by far the most frequent uses to which Gabmap has been put. Nerbonne et al. (2011) reports on Gabmap's basic functionality and its implementation, so that this article can build on that, adding material on new functionality and reporting on the range of uses to which Gabmap has been put. Gabmap is modestly

¹We are grateful to CLARIN-NL and for their support of the project ADEPT (<http://www.clarin.nl/node/70#ADEPT>), which in turn produced Gabmap. CLARIN-NL was supported by the Netherlands Organization for Scientific Research (NWO)

26 successful, and its popularity underscores the fact that the study of language variation has
27 crossed a watershed concerning the acceptability of automated language analysis. Auto-
28 mated analysis not only improves researchers’ efficiency, it also improves the replicability of
29 their analyses and allows them to focus on inferences to be drawn from analyses and other
30 more abstract aspects of that study.

31 2 A Gabmap session

32 In this section, we show an example of a typical Gabmap session and the types of analyses
33 that can be conducted. For this purpose we use data from the Goeman-Taeldeman-Van
34 Reenen-project (GTRP; Goeman and Taeldeman 1996). The data consist of phonetic tran-
35 scriptions of Dutch dialects from the Netherlands and Belgium gathered during the period
36 1980—1995. These data are available as demo data on the Gabmap web site, which makes
37 it possible for users to try out the analyses described here directly in Gabmap.

38 2.1 Data

39 The dialect data can be prepared in a spreadsheet where rows represent sites and columns
40 represent linguistic variables. In the demo data, the columns are words and each cell in
41 the spreadsheet shows the pronunciation of a word in the International Phonetic Alphabet
42 (IPA) at one specific site:²

	boter	broden	zout
Aalsmeer	botər	brojə	zaut
Baardegem	botər	bruəs	zat
Coevorden	boetər	brodn	sɔlt _ɾ

44 Gabmap accepts tab-separated Unicode text files as input data, and most spreadsheet
45 software allow exporting data to text files with Unicode encoding.

46 Analysis in Gabmap is not restricted to transcribed pronunciation data; instead, any
47 kind of binary or numeric data can be used. When uploading data into Gabmap, the type
48 of data is specified, so that the data can be processed appropriately. For the phonetic
49 transcriptions in the example we choose *string data* as the type of data and *string edit*
50 *distance* as the type of processing (more about data processing in section 2.3).

51 In order to create dialect maps, the data file should be accompanied by a map file
52 with the geographical coordinates of the data sites and optionally borders of the country

²If there are several pronunciations available of a word from one site, these can be separated by “space slash space” in the data file.

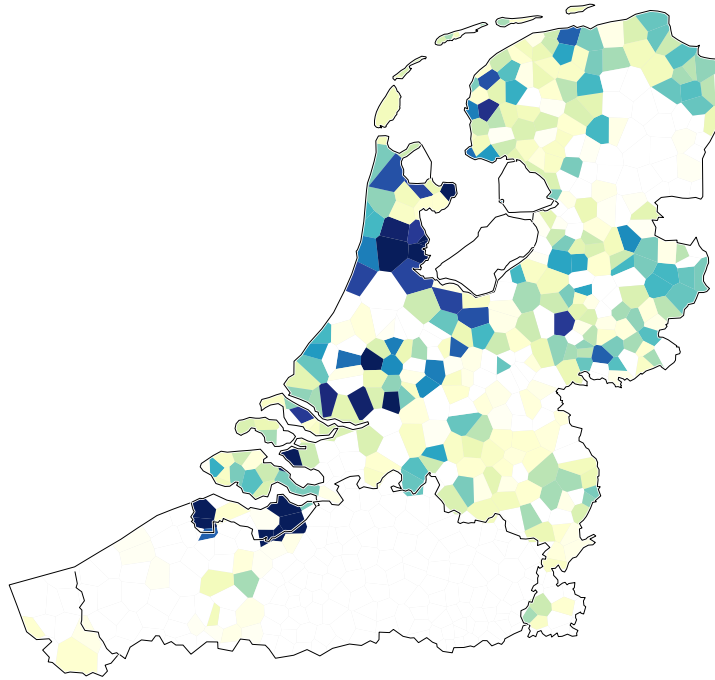


Figure 1. Distribution map of the character [ɫ].

53 or language area. The map file is a .kml or .kmz file that can be created in Google Earth
 54 or using the Google Maps service through any standard web browser. Using a map file is,
 55 however, not compulsory. Users might want to analyze language variation related to other
 56 factors than geography. The data rows might, for example, be individual speakers instead
 57 of sites. For analysis of this type of data, no map file is needed and Gabmap will create a
 58 pseudo map instead of real maps in the mapping functions. The statistical analyses, like
 59 cluster analysis and multidimensional scaling (see below), will, then, show how individual
 60 speakers group together based on their language use.

61 When a project is created, Gabmap offers several ways of inspecting the data. Summaries
 62 are created of the number of sites, number of words (or other linguistic variables), number
 63 of characters and number of tokens. In *Data overview* in Gabmap, we can, for example, see
 64 that the demo data has data from 613 places and that the number of different words (items)
 65 is 562. The total number of word transcriptions (instances) is 331,690, which is less than
 66 613×562 due to some missing data in the input table.

67 2.2 Distribution maps

68 Several types of distribution maps are offered in Gabmap. Figure 1 shows a map of one
 69 specific phonetic character in the data set. The character maps are part of the data overview
 70 function in Gabmap, where maps can be created of any character or token in the data set.
 71 Figure 1 shows the distribution of the velarized lateral approximant [ɫ]. White color means
 72 no instances at all of the character from a site, and the darker the color the higher the

Distribution map for RE "ə\$" in *dopen*

- *daupə* (1)
- *dəpə* (1)
- *dipə* (1)
- *doopə* (2)
- *dopə* (1)
- *dopə* (77)
- *doupə* (65)
- *dowpə* (1)
- *dœpə* (7)
- *duopə* (2)
- *dupə* (4)
- *duopə* (1)
- *duopə* (55)
- *duəpə* (7)
- *dypə* (1)
- *dyəpə* (7)
- *dyəpə* (7)
- *dyʌpə* (1)
- *døfə* (2)
- *døpə* (10)
- *døypə* (2)
- *dœjpə* (4)
- *dœjpə* (1)
- *dœpə* (2)
- *dœypə* (21)

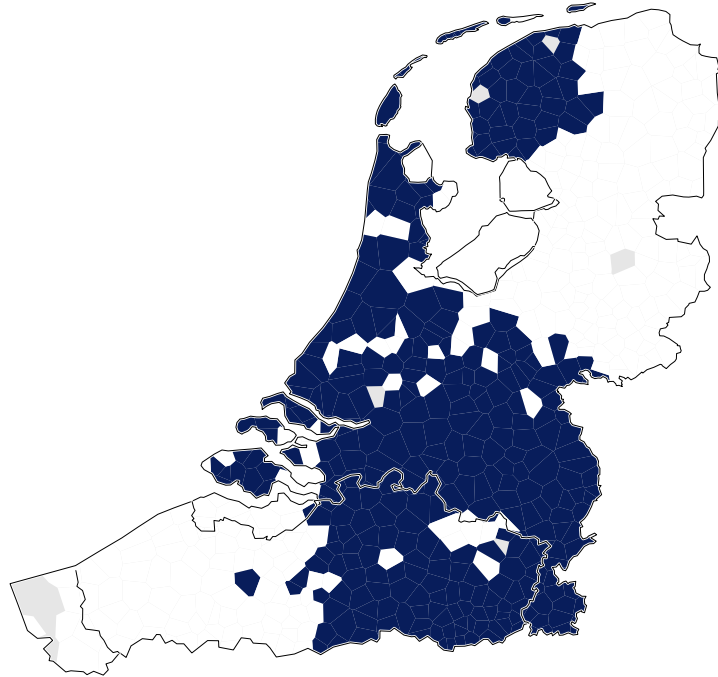


Figure 2. Map showing the distribution of pronunciations of the word *dopen* (‘to baptize’) ending in a schwa. To the left a part of the pronunciations selected by the regular expression is shown.

73 relative frequency of the character in the data at the given site. A map like this only gives
 74 a rough picture of the distribution of a speech sound, since the result depends on how
 75 well each data point has been sampled.³ Still, the map can give a rough overview of the
 76 distribution of a dialect feature and/or of the quality of the data. It is striking that the
 77 chosen phonetic symbol in Figure 1 is almost completely lacking in the data from Belgium.
 78 When a pattern like this is found, it could either mean that the distribution of the specific
 79 feature very closely follows the national border, or, it could mean that it was not transcribed
 80 with the same phonetic symbol by transcribers of the Flemish and Netherlandic Dutch data.
 81 In fact this is one of the indications that the Dutch and Flemish fieldworker-transcribers did
 82 not use the phonetic alphabet (Wieling, Heeringa, and Nerbonne 2007) in the same way; it
 83 turned out that the Flemish fieldworker-transcribers used many fewer symbols. See Wieling
 84 and Nerbonne (2011b) for a suggestion on how to correct for the differences in phonetic
 85 alphabet using dialectometric techniques.

86 Distribution maps of specific words can also be created in Gabmap. By first choosing a
 87 variable (word) and then a specific variant (pronunciation) a map is created which shows
 88 where the chosen variant can be found. Regular expressions can also be used to create
 89 distribution maps. Figure 2 was created by first choosing the word *dopen* (‘to baptize’) and
 90 subsequently using the regular expression ‘ə\$’ (‘\$’ to mark end-of-string) for selecting all

³Sites with a lot of missing data could by coincidence get too high or too low relative frequencies compared to other sites.

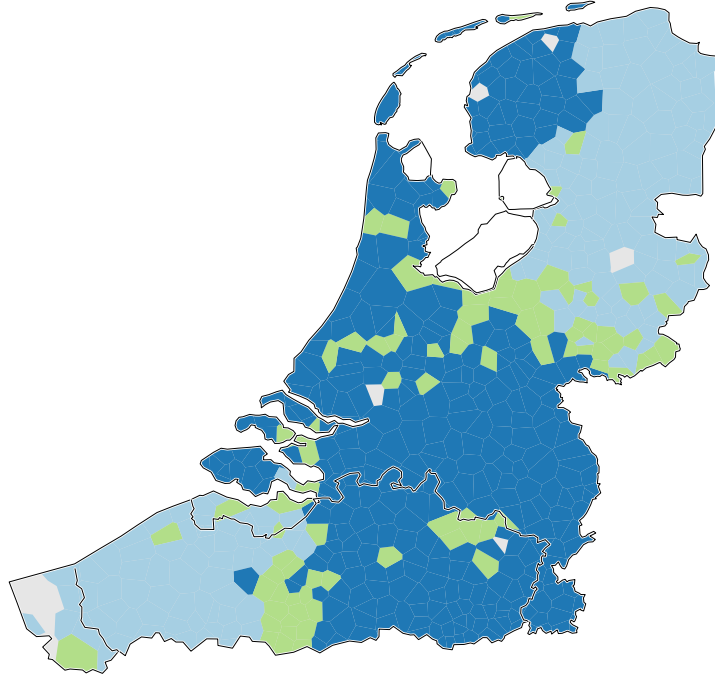


Figure 3. Map showing the distribution of three different types of endings for the word *dopen* in Dutch dialects: *-m* (light blue), *-n* (green) and *vowel* (dark blue). Gray spots are sites with missing data.

91 pronunciations ending with a schwa, illustrating one result of the weakening of unstressed
 92 syllables. In addition to creating the map, Gabmap shows a list of the chosen pronunciations.

93 The distribution maps in Gabmap can only show the presence or absence of a chosen
 94 feature. In traditional dialect maps, however, it is common to show the distribution of
 95 several different variants by using different symbols, patterns, or colors. For example, one
 96 might want to make a map of the word *dopen* showing the distribution of three different
 97 types of endings *-m* (e.g. [dopm]), *-n* (e.g. [dopən]) and ending in a *vowel* (e.g. [dopə]).
 98 This can be achieved in Gabmap by using a data file with a single variable (i.e. one data
 99 column):

	ending
Aalsmeer	vowel
Baardegem	-n
Coevorden	-m

100

101 When uploading the data, *categorical data* is used as data type and *binary comparison*
 102 as processing type. The map can be created as a *cluster map* in Gabmap. Since the clusters
 103 are coded in the uploaded data file, it does not matter which clustering algorithm is used,
 104 but the number of clusters should simply be the same as the number of different codes in
 105 the data file, which is 3 in the example case. The map is shown in Figure 3.

Aalsmeer NH — Aalten Gl

r	e	i	ɣ	ə		
r	e		x	ə	n	
		1	1		1	3

Figure 4. Example of computing of string edit distance.

106 2.3 Measuring linguistic distances

107 Dialectometric analyses are typically based on linguistic distances between pairs of sites
 108 in the data. The linguistic distances between sites are in turn calculated as the mean
 109 distances of the variables instantiated at both sites. Gabmap calculates these distances
 110 when a project is created. The distance measure used for string data is the *string edit*
 111 *distance* (or Levenshtein distance, Levenshtein 1966).

112 The string edit distance computes the minimal number of insertions, deletions and sub-
 113 stitutions needed to change one character string into another. Gabmap computes the dis-
 114 tance for all words and all pairs of sites and shows the alignments made (under *Measuring*
 115 *technique > alignments*). Figure 4 shows the alignment of the word *regen* (‘rain’) in the
 116 Aalsmeer dialect and the Aalten dialect. One deletion [i], one substitution [ɣ]~[x] and one
 117 insertion [n] is needed for the alignment, which results in a distance of 3. The linguistic
 118 distance between two sites is the average of the distances of the words available from both
 119 sites.⁴

120 For other types of data other distance measures can be chosen. For numeric dialect
 121 data the Euclidean distance is used, and for categorical data either binary comparison or
 122 the ‘Relative Identity Value’ (*Gewichteter Identitätswert*, Goebel 2006, p. 416), a weighted
 123 similarity index, can be used. Instead of uploading actual dialect data it is also possible to
 124 upload a matrix of any kind of distances into Gabmap.

125 The distances are displayed in Gabmap as beam maps or network maps (see Nerbonne
 126 et al. 2011, p. 79). Another possibility is to display the distances from one site to all other
 127 sites (*references point maps*), which shows how similar or different the dialects might sound
 128 to a speaker of a specific dialect. Figure 5 shows a reference point map where Coevorden in
 129 the north-east of the Netherlands is the reference point.

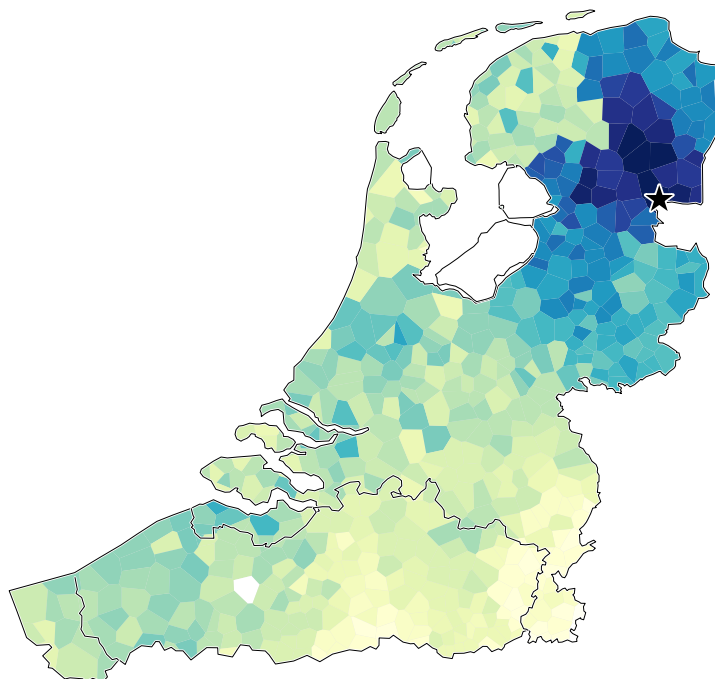


Figure 5. Reference point map. The lighter the color, the greater the linguistic distances from the starred reference site (Coevorden).

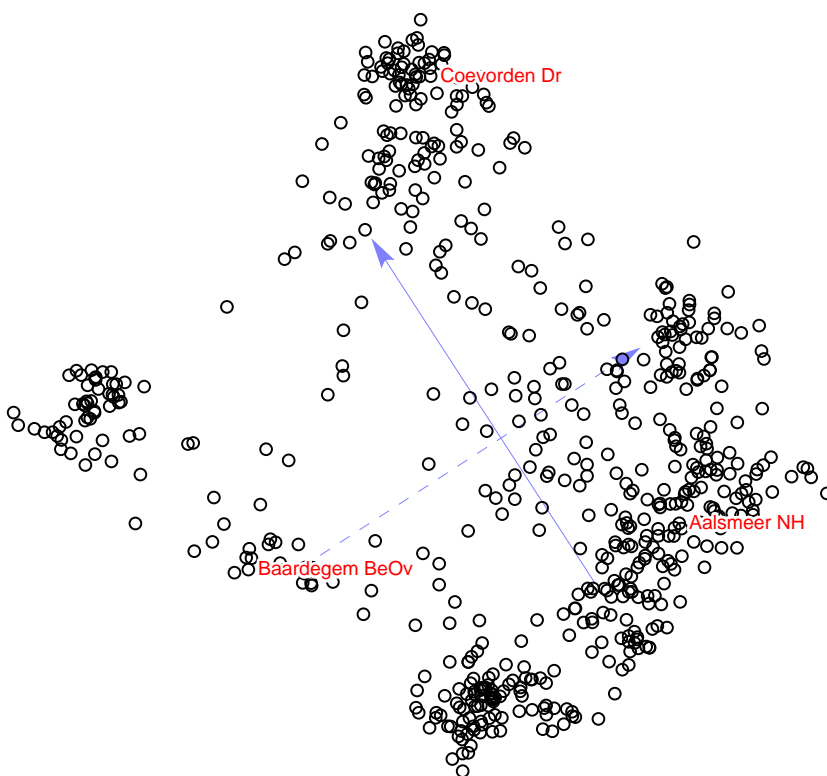


Figure 6. Plot of the result of multidimensional scaling in two dimensions. The labels of three reference sites are displayed.

2.4 Dialect continuum

Maps such as the reference point map in figure 5 only visualize the linguistic distances from one site to all other sites. In this map, there are very light areas to the north-east from the reference site Coevorden, as well as in the south of the language area. The map does not tell us whether these two areas are similar to or different from each other or not, only that both of them are linguistically very different from Coevorden. For an objective observer, a map that displays the linguistic relationships across all sites simultaneously might be more useful. This can be achieved by using multidimensional scaling (MDS).

MDS takes the full *sites* \times *sites* distance matrix as input and creates a representation in an n -dimensional space where the distances are approximations of the original linguistic distances.⁵

This can be compared to trying to create a map using only the distances between cities in kilometers as information. The results of MDS can be plotted in a Cartesian coordinate system (*mds plots* in Gabmap). Similar data points will be close to each other in the plot.

An example of this is found in Figure 6, where the labels of three example sites have been added. The first dimension of an MDS analysis always explains as much as possible of the variance in the data, and additional dimensions add maximally to the precision of the approximation of the distances, but each additional dimension explains less of the variance than the previous one. In Figure 6, the solid arrow represents the first dimension explaining 49% of the variance in the data (correlation between the original linguistic distances and the Euclidean distances between the MDS coordinates: $r = 0.70$) and the dashed arrow represents the second dimension explaining 23% of the variance ($r = 0.48$).⁶ Aalsmeer has the lowest value in the first dimension, while Baardegem has an intermediate value and Coevorden has a very high value. In the second dimension, on the other hand, Baardegem has a very low value, while Aalsmeer and Coevorden both have relatively high values. This means, that there are some linguistic features that Aalsmeer and Coevorden share (second dimension), but other features that are very different in these two dialects (first dimension). The plot clearly shows that there are some groups of dialects that cloud together, but also single sites which lie between those groups.

The results are easier to interpret if they are displayed on maps. The two first maps in Figure 7 show exactly the same results as Figure 6, but instead of displaying a coordinate

⁴If more than one pronunciation is available for a word from one site or both sites, an averaging procedure (ignoring identical pairs) is used (see Nerbonne and Kleiweg 2003, Sec. 3.2).

⁵On the use of multidimensional scaling in dialectology, see e.g. Embleton (1993), Heeringa (2004, pp. 156–161), Nerbonne (2011, pp. 487–489), and Embleton, Uritescu, and Wheeler (2013).

⁶If a map file is provided, the MDS plots produced by Gabmap are rotated using the Procrustes transformation (see, e.g., Peres-Neto and Jackson 2001), which has the effect that the sites presented in the MDS plot align with their geographic coordinates as closely as possible. The axes corresponding to the first two MDS dimensions are drawn on the graph.

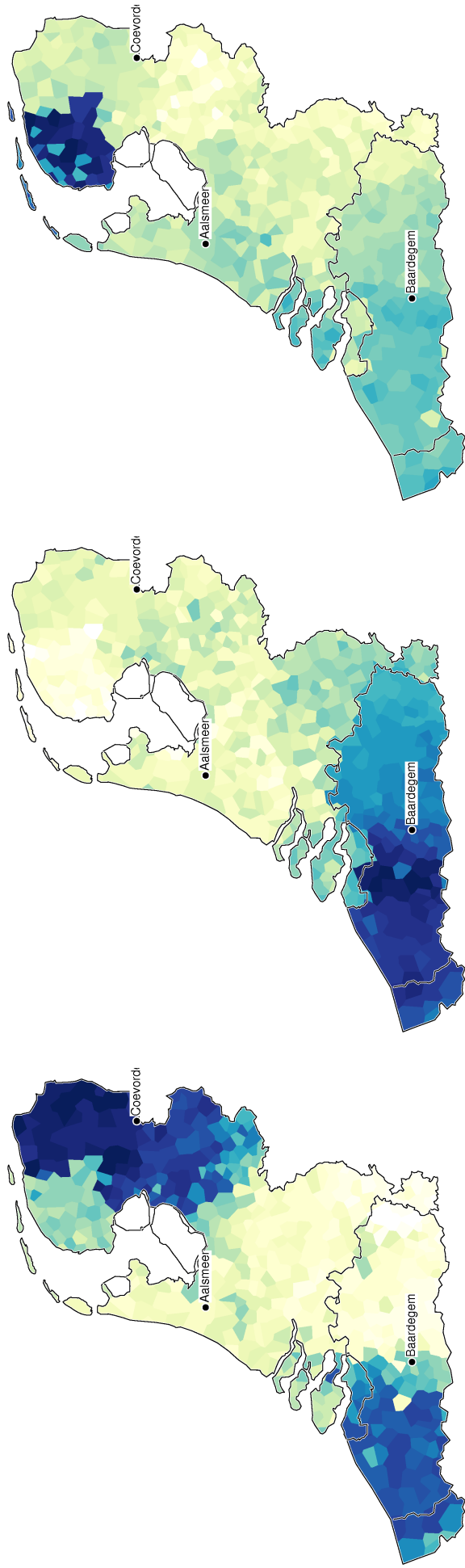


Figure 7. Maps of the first (left), second (center), and third (right) dimension of multidimensional scaling.

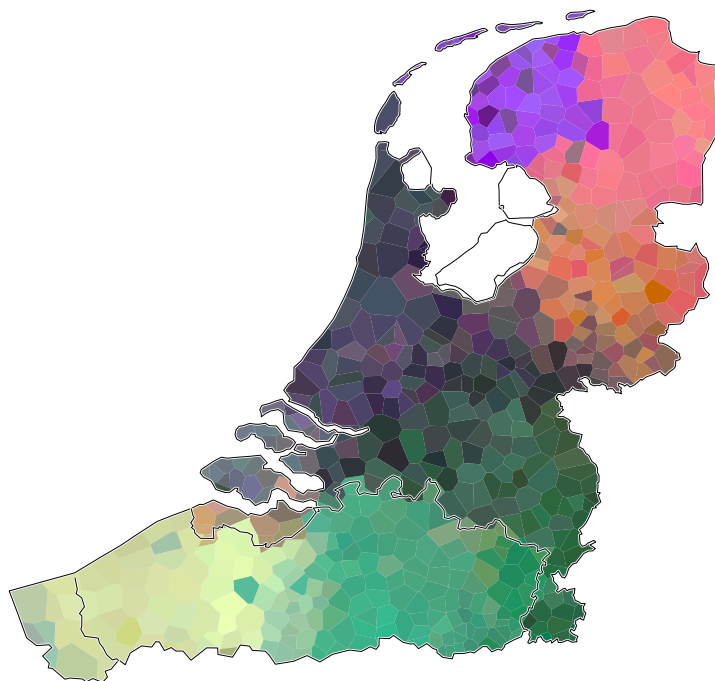


Figure 8. Map of multidimensional scaling applied to Dutch dialects ($r = 0.89$).

161 system, the area surrounding each site on the map has been colored according to the value
 162 of one dimension in the MDS analysis. The third map shows the third dimension ($r = 0.37$).
 163 Light color means high value, dark color low value. The maps show that the dimensions of
 164 the MDS represent different geographic distribution patterns: the first dimension shows a
 165 center–periphery effect, while the second dimension shows a northeast–southwest distribu-
 166 tion. The third dimension mainly distinguishes Frisian (dark area) from the Dutch dialects.
 167 Multidimensional scaling to three dimensions has almost always explained around 80 – 90%
 168 of the variance in the dialect data sets we have analysed, and it has been our experience that
 169 adding more than three dimensions to the analysis generally does not improve the solution
 170 much.

171 The maps in Figure 7 can be superimposed — or “put on top of each other” — using
 172 the red, green and blue (RGB) colors in order to show the aggregated dialectal differences,
 173 which gives the map in Figure 8.⁷ All the maps of MDS results are found in *mds maps* in
 174 Gabmap. Similar colors in these maps indicate that the dialects share many features. The
 175 sharpest dialect border in Figure 8 is found in the north where the Frisian dialects are very
 176 different from the Dutch dialects. Frisian is in fact officially recognized as a separate, but
 177 closely related, language with its own written standard. The rest of the map shows less
 178 crisp borders, reflecting instead rather continuous transitions from one dialect area into the
 179 other.

⁷For a detailed account of how this is achieved, see Heeringa 2004, pp. 161–163, Leinonen 2010, pp. 207–208, and Nerbonne 2011, pp. 489–491.

2.5 Identifying dialect groups

The MDS plot in Figure 6 shows that despite the continuous nature of the dialect data, the dialects also seem to cluster together to some extent forming dialect groups. Dialectologists often want to be able to identify these kinds of dialect groups and draw borders between dialect areas on maps. We can seek groups of sites and dialect areas using cluster analysis. Clustering algorithms aim at minimizing the differences within each group of data points, while maximizing the distances across groups. Several so-called hierarchical clustering methods are available in Gabmap. Cluster analysis is applied to the distance matrix consisting of the pair-wise aggregate linguistic distances between places, and groups are formed based on similarity.⁸

The results of cluster analysis are shown in maps in Gabmap, where each cluster is displayed by a unique color.⁹ Figure 9 shows the results of two different cluster algorithms: weighted average (left) and Ward’s method (right). The contrast in these maps highlights the fact that different clustering algorithms have different biases and can lead to very different results. Ward’s Method has a bias to favor equal size clusters, while weighted average is more faithful to the original linguistic distances. The figure shows that the map based on Ward’s method has seven quite large clusters of dialects, while the map of weighted average has five large clusters and two very small ones.

Not only do different clustering algorithms yield different results, each algorithm is also relatively unstable, meaning that small changes in the input data can lead to large changes in the cluster division. This is because each site is forced into a single cluster even in cases where the data might in fact be continuous. This can be compared to multidimensional scaling, which can show group structure in data, but also allows data points to float between groups or even show a truly continuous distribution (cf. Figure 6).

Because cluster analysis is a relatively instable method, *noisy clustering* (Nerbonne et al. 2008) has been implemented in Gabmap. In noisy clustering, cluster analysis is performed several times with different clustering methods and by contaminating the original distance matrix with (different) small amounts of random noise. The results of noisy clustering are displayed in a probabilistic dendrogram where percentages show how many times each cluster was encountered in the repeated clustering with noise. Clusters that appear in many runs of the analysis with added noise are particularly stable ones. For an example of noisy clustering in Gabmap, see Nerbonne et al. (2011, p. 83).

⁸For an introduction to cluster analysis and descriptions of differences between different cluster algorithms, see e.g. Jain and Dubes (1988), Manning and Schütze (1999, pp. 495–528), Heeringa (2004, pp. 146–156), and Prokić (2010, pp. 25–29).

⁹In contrast to MDS maps, the colors are arbitrary in the sense that similarity of colors does not indicate linguistic similarity. E.g. the light blue dialects are not necessarily any more similar to the dark blue dialects than to the red dialects.

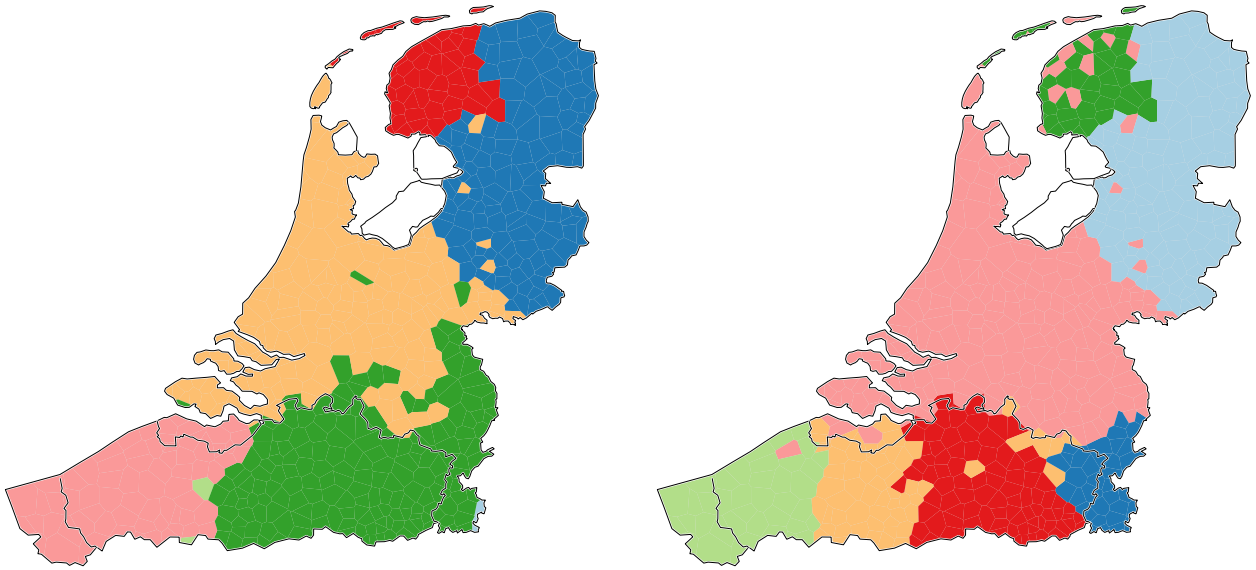


Figure 9. The results of two different cluster analyses: weighted average (left) and Ward's method (right). Seven clusters are displayed with seven distinct colors in both maps.

212 Another way of evaluating the results of cluster analysis is to compare the results of
 213 clustering to MDS (*cluster validation* in Gabmap). Figure 10 shows the two-dimensional
 214 MDS plot (Figure 6) colored according to the two different cluster analyses (Figure 9),
 215 respectively. Ward's method recognizes seven relatively large clusters, but at the cost of
 216 separating groups that are actually relatively similar according to the MDS (see, for example,
 217 the cloud of sites at the left side of the plot which belongs to one cluster according to
 218 weighted average but two different ones according to Ward's method). The two methods
 219 also disagree on how sites that fall between the clear clouds of sites are treated. Many of
 220 these are actually extreme points within a cluster, as indicated by the numbers added to
 221 the plots.

222 The comparison to MDS shows that, in this particular data set, the clusters might in
 223 fact not be as well separated on linguistic grounds as the cluster map might seem to suggest.
 224 Of course, the MDS plot only shows the first two dimensions of MDS which explain around
 225 72% of the variance, so some of the information used in the cluster analysis is not accounted
 226 for in the MDS solution. For example, the third dimension of MDS singles out Friesland (cf.
 227 Figure 7). which will make it a more distinct cluster than the two first dimensions of MDS
 228 suggest. Hence, the amount of variance explained by different dimensions of MDS should
 229 also be considered when using MDS for validating cluster analysis.

230 **2.6 Finding typical features or “shibboleths”**

231 The dialectometric methods we discussed so far aim to find and characterize dialect groups
 232 at an aggregate level. A large number of variables (e.g. words) are used for investigating

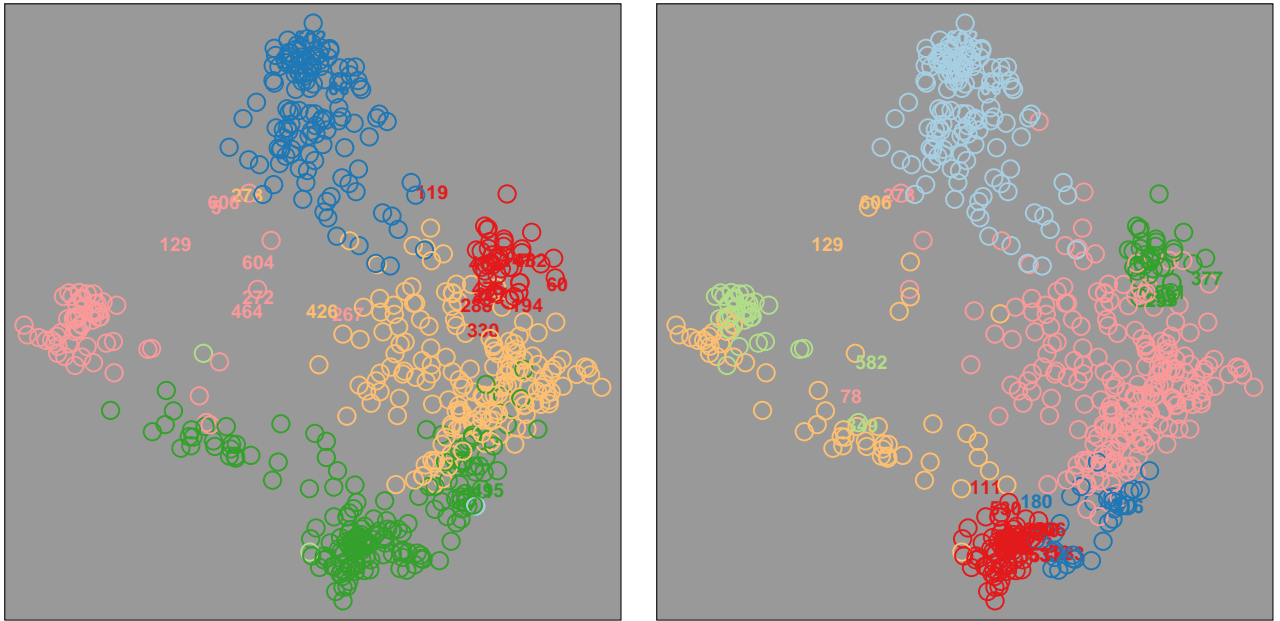


Figure 10. The results of weighted average (left) and Ward’s method (right) compared to multidimensional scaling. The colors above correspond to those in Figure 9 (left and right, respectively).

233 overall dialectal differences. Often, we want to know which variable or variables are most
 234 characteristic for a specific dialect area. Such variables, termed shibboleths, referring to a
 235 variant of speech that betrays where a speaker is from (*Judges 12:6*), can be identified with
 236 the ‘cluster determinants’ function of Gabmap.

237 The cluster determinants option in Gabmap implements the method described in Prokić,
 238 Coltekin, and Nerbonne (2012).¹⁰ The aim of the cluster determinants function is to find
 239 the items that are characteristic for a particular cluster, i.e. a set of sites. The method
 240 is related to the *Fisher’s linear discriminant* (Schalkoff 1992, 90ff) and the information re-
 241 trieval measures *precision* and *recall*. In essence, we would like to find items that distinguish
 242 sites in the target cluster from the sites outside it (possibly belonging to multiple clusters),
 243 but we also prefer the items that exhibit little variation within the target cluster. These
 244 two properties, *distinctiveness* and *representativeness*, together define how characteristic a
 245 particular item is for the target cluster.

246 Gabmap enables the investigation of typical linguistic elements (“cluster determinants”)
 247 in three steps. In a first step, the target cluster is determined. The user can obtain a
 248 clustering using any of the clustering options described in Section 2.5, selecting one of the
 249 clusters as the target cluster. Even if more than two clusters are determined by this process,
 250 the important distinction is between the (selected) target cluster and the rest of the sites.
 251 The structure outside the target cluster is not used. Alternatively, the sites in the target
 252 cluster can be defined manually, e.g. based on theoretical motivations. The procedure also
 253 allows automatic clustering at first step, and adjusting the result manually.

¹⁰An earlier method based loosely on Wieling and Nerbonne (2011a) is also available for categorical data.

Table 1. The top- and bottom-three ranked ‘shibboleths’ for the Frisian cluster. The scores in the column *between* represents the differences between the Frisian cluster and the rest of the Dutch speaking area in our data set with respect to each item. The higher the score, the more *distinctive* the item. The scores in the column *within* measures the variation of the item within the Frisian area. The lower the score (variation), the more *representative* the item. The overall score at the column labeled *score* is the difference *between* – *within*.

Item	between	within	score
vinden	0.03	-2.37	2.41
knieën	1.13	-1.20	2.34
zoet	1.17	-1.12	2.29
nog	0.22	0.28	-0.06
kaf	0.63	0.72	-0.09
elf	0.27	0.36	-0.09

254 In a second step, the user selects the target cluster, and generates a ranked list of items
 255 along with their representativeness and distinctiveness scores. The scores are presented after
 256 normalization, so that the average (randomly selected) item would get a score of zero. The
 257 items are ranked based on an (equally weighted) linear combination of the two scores (see
 258 Prokić, Coltekin, and Nerbonne 2012, for the details of normalization and combination of
 259 the scores). Depending on the application, one may prefer to select the items based on just
 260 representativeness or just distinctiveness, or possibly on a differently weighted combination
 261 of the two. Gabmap allows downloading the resulting table, which the user may then
 262 experiment further with.

263 Table 1 presents the top three and bottom three shibboleth candidates for the Frisian
 264 area we discussed in the previous sections. The first item ‘vinden’ scores high because of the
 265 fact that it is pronounced uniformly within the Frisian area (low within score). However, it is
 266 definitely not a distinctive item (between score close to zero). The pronunciation differences
 267 between the Frisian cluster and the rest of the Dutch speaking area is almost exactly what
 268 would be expected from the differences measured in the whole data set with respect this
 269 item, i.e. quite similar pronunciations are found in other areas even though the exact
 270 same pronunciation does not occur outside Friesland. The other two top items show more
 271 balanced representativeness and distinctiveness scores. The least likely candidates all show
 272 small scores of distinctiveness or representativeness, and their combination result in low
 273 scores (around zero). An expert would already get a sense of specific items for a particular
 274 cluster by eyeballing the ranked list. However, the next step in ‘Cluster determinants’
 275 function allows closer inspection of any item in the list.

276 In the last step, after determining the characteristic items, the user can select a particular
277 item and visualize the differences with respect to this particular item using beam maps (see
278 Nerbonne et al. 2011, p. 79), and list all forms (pronunciations) observed within or outside
279 the target area along with their frequencies. For example, looking at the item *vinden* we
280 identified as being representative of Frisian area, we observe that this item is pronounced
281 as [finə] in all 52 sites in this area. The exact pronunciation is not found elsewhere in our
282 larger area of interest. However, the distinctiveness score indicates that the pronunciation
283 differences (as measured by Levenshtein distance) between Frisian area and the rest of the
284 sites do not differ substantially. If we look at the second item in the list, *knieën*, we observe
285 that the item varies within the Frisian cluster, in total we observe 15 forms of the item, and
286 all except one of these forms are used exclusively within this cluster. The distinctiveness
287 score also indicates that the pronunciation difference between Frisian area and the rest is
288 over one standard deviation away from the typical pronunciation difference between two
289 sites with respect to this item. Further inspection of the forms recorded within the Frisian
290 area indicates that the pronunciation of *knieën* in this area almost always ends with an [s].
291 Similarly, all pronunciations of *zoet* (the third item in the list) in the Frisian area has a
292 initial [s], while this is rare in other sites in our data.

293 3 User experiences

294 3.1 Some statistics

295 It is difficult to characterize the users of Gabmap in detail, as we decided against requiring
296 users to identify themselves when developers of similar projects reported that mandatory
297 registration appears to depress the enthusiasm for web applications. We can report that
298 there were 45 users and 352 projects (excluding 10 guest users) as of late March, 2014.
299 This figure ignores those with completed projects whose accounts expire after two months
300 of no use (with one week's warning). The web server access for the last month indicates on
301 average 2795 hits and 71 visits per day.

302 We have also presented tutorials on Gabmap at the *Nordic Congress of Dialectologists*,
303 Uppsala, Aug. 20, 2010; at the *Tagung des Forums Sprachvariation*, Erlangen, Oct. 15,
304 2010; at the University of Potsdam, Dec. 7, 2010; in a poster at the 6th *International*
305 *Conference on Language Variation in Europe* (ICLaVE), Freiburg, June 30, 2011; at *Dig-*
306 *ital Humanities* 2011 (Stanford) with about 12 participants; at the conference *Methods in*
307 *Dialectology XIV* (London, Ontario, Aug. 2011) with 40 attending; at the conference *Com-*
308 *paring approaches to measuring linguistic differences* at the University of Gothenburg, Oct.
309 26, 2011; at the Society of Swedish Literature in Finland, Nov. 23-25, 2011; at the LOT
310 winter school of the Dutch National Research School for Linguistics (Tilburg, Jan, 2012);

311 at a *Digital Humanities* summer school in Leuven, Sept., 2012 with roughly 10 participants;
312 and at *Methods in Dialectology XV* conference (Groningen, Aug. 2014), with over twenty
313 participants. Users have been pleased at the ease with which analyses can be conducted.

314 **3.2 Examples of user work**

315 Gabmap has been used for various purposes in the three years since it was first launched;
316 these include not only linguistic and other research, but also the presentation of research to
317 professionals and to interested popular science audiences.¹¹ The recent *Methods in Dialectology XV*
318 conference included several talks which used Gabmap (and which are discussed
319 below) as well as talks which compared treatments to Gabmap (e.g., talks by Simon Pickl
320 and Fruzsina Vargha).

321 A number of users have especially exploited Gabmap's map-making facilities. Bouma
322 and Hermans (2012) use Gabmap to project the distribution of different syllable onsets
323 in medieval Dutch, and Wieling, Upton, and Thompson (2014) and Wieling (2013) use
324 Gabmap's facilities for analyzing numerical data (lexical frequency differences) to provide
325 analyses of the very large-scale BBC voice project. The work may be viewed in more detail
326 at <http://www.gabmap.nl/voices/> where users are encouraged to explore the lexical choices
327 of all the respondents, or to contrast men's and women's speech or the speech of the young
328 and old. Leinonen (in press) uses Gabmap's map-making facilities for analyzing data from
329 the dictionary of Swedish dialects in Finland. She uses the clustering facility for making
330 isogloss maps of single features with multiple variants as well as aggregating dialectometric
331 maps.

332 Castro (2011), on the other hand, uses Gabmap's clustering routines in his argument
333 that Southern Sui should be recognized as a separate dialect, distinct from Sandong Sui.
334 Coloma (2012) focuses on just ten features in modern Spanish and, like Castro, exploits
335 Gabmap's ability to process numeric data (differences in frequencies) and to invoke cluster-
336 ing and MDS. Scherrer (2012) introduces his own idea for measuring varietal distance based
337 on comparing the number of identical lexicalizations in Swiss German dialect corpora to
338 the number of cognates found there, and he uses Gabmap for MDS, clustering, and map-
339 ping even while examining the Cronbach's α score used in Gabmap to determine whether
340 samples are large enough to provide a geographical signal and using a Mantel test com-
341 paring distance matrices determined using different techniques. Moran and Prokić (2013)
342 investigated several endangered Dogon languages (spoken in Mali) emphasizing the need to
343 preserve what is possible in communities with few speakers. They made use of Gabmap's

¹¹Our thanks, too, to Erik R. Thomas, North Carolina, and Yonatan Belinkov, Tel Aviv, who referred us to their as yet unpublished work using Gabmap on Midwestern US varieties of English and on translations of the Hebrew Passover Haggadah, respectively.

344 probabilistic clustering routines as well as the mapping facilities. Reber (2013) focused not
345 on dialect speech, but rather on the range of place names found at different settlements, i.e.
346 the names of neighborhoods, fields, streets, paths, hills, peaks, rivers and other bodies of
347 water. The author uses Gabmap for clustering and mapping.

348 Uiboaed et al. (2013) investigated corpus-based morphosyntactic dialectometry by first
349 extracting corpus frequencies of various verbal “collostructions” (Stefanowitsch and Gries
350 2003) in Estonian and then examining the results for geographic cohesion using both corre-
351 spondence analysis and Gabmap’s clustering routines.

352 Mathussek (2013, pp. 248-251) uses Gabmap’s aggregating, dialectometric focus to an-
353 alyze middle Franconia (in northwest Bavaria) and to contrast the aggregate views with
354 perspectives from traditional research and from perceptual dialectology. The dialectometric
355 approach was crucial in identifying field worker boundaries in the data, which led her to
356 ignore phonetic details (diacritics) before proceeding, an issue which is the focus of Math-
357 ussek (submitted). Mathussek’s approach is emphatically pluralistic, and she notes that
358 it was the failure of initial dialectometric analyses to agree with traditional ones that led
359 her to pursue the possibility of field worker confounds. Mathussek (2014, Chap.4) discusses
360 Gabmap as means of returning to older data sets with new techniques — naturally, in order
361 to obtain new insights, or at least to examine older ideas from a fresh vantage point.

362 Šimičić et al. (2013) analyzed coastal Croatian dialects but also varieties from the Italian
363 provinces of Molise, attending to phonological and lexical variation. The two linguistic
364 levels correlated strongly ($r = 0.72$), and the authors interpret the differences to be due
365 to the stronger historical signal in phonology, and the greater volatility of the lexicon.
366 Due to the complicated history of Croatian migrations, one might have expected the usual
367 dialect areas and dialect continua not to emerge, and they indeed do not emerge from this
368 analysis. Instead the analysis uncovers a great many discontinuities, particularly on the
369 northern island of Istria, which the authors suggest ought to be attributed to migration.
370 The Štokavian and Čakavian varieties of the south were less diverse, and the varieties spoken
371 in Molise, Italy were very distinct from the others. The authors conclude methodologically
372 that the aggregating view inherent in Gabmap has advantages over the traditional analyses
373 based on isoglosses, in particular because it obviates the need to choose which isoglosses are
374 to be regarded as probative.

375 Mitterhofer (2013) used Gabmap to identify cognates and other related words in varieties
376 of Bena in Tanzania, comparing Gabmap’s edit-distance measures to the Summer Institute
377 of Linguistics’ “Survey on a Shoestring” (1990), and Bloem et al. (submitted) uses the
378 "cluster determinants" feature of Gabmap to identify characteristic mispronunciations in
379 foreign accents in English.

380 Snoek (2013) uses Gabmap to research lexical relations among Athapaskan languages in
381 order to improve the understanding of their historical relations, and Snoek (2014) provides

382 an article-length review of Gabmap targeted at researchers in language documentation.
383 The author analyzes phoneme strings denoting body-part terms in Northern Athapaskan
384 languages (in Canada and Alaska). The application of dialectometrical tools is appropriate
385 for these Athapaskan languages because their relations to one another are poorly established
386 in Amerindian scholarship. He adds to existing documentation by explaining how maps may
387 be produced for Gabmap using Google Earth, and he has some important warnings about
388 how Gabmap may handle transcriptions involving digraphs or trigraphs. Most intriguingly,
389 he shows how Gabmap's data examination facilities may be very useful even when researchers
390 do not aim at a quantitative analysis of their data. He concludes that "Gabmap is excellent
391 software that permits the mapping and comparison of linguistic data in a fast and generally
392 painless manner."

393 4 Conclusions

394 Gabmap offers a range of processing possibilities all geared to highlighting and tallying
395 linguistic differences. Nerbonne et al. (2011) sketched some of these, and the current paper
396 aims to supplement that one by describing other possibilities and also to review some of the
397 uses to which Gabmap has been put.

398 Gabmap would undoubtedly benefit from further use and also from the incorporation of
399 various advances in dialectometry since 2010, including more sensitive measures for pronun-
400 ciation differences that incorporate segment differences (Wieling, Margaretha, and Nerbonne
401 2012; List 2012), and we have noted that tutorial material as well as reference material at
402 various levels is invaluable. More material would be precious. The phylogenetic group-
403 ing procedures such as NeighborNet or Bayesian Monte Carlo Markov-Chain techniques
404 (Felsenstein 2004, Ch.16,18,35) are valuable historical perspectives for many of the ques-
405 tions dialectologists entertain. The maps Gabmap produces are not geo-referenced, and
406 this handicaps some interesting applications involving comparing the diffusion of linguistic
407 culture with other sorts of culture (Manni et al. 2008).

408 Gabmap is also open-source, and we would welcome proposals from others to incorporate
409 further processing possibilities into Gabmap, although we are also wary of the time that
410 might be needed to see this through successfully.

411 References

412 Bloem, J. et al. (submitted). "Automatically Identifying Characteristic Features of Non-
413 Native English Accents". In: *The Future of Dialects*. Ed. by M.-H. Côté, R. Knooihuizen,
414 and J. Nerbonne. Language Variation 1. Berlin: Language Science Press.

- 415 Bouma, G. and B. Hermans (2012). “Syllabification of Middle Dutch”. In: *The Second*
416 *Workshop on Annotation of Corpora for Research in the Humanities (ACRH-2)*. Ed. by
417 F. Mambrini, M. Passarotti, and C. Sporleder. Lisbon, pp. 27–38.
- 418 Castro, A. (2011). “Southern Sui: a Fourth Sui Dialect”. In: *Journal of the Southeast Asian*
419 *Linguistics Society* 4.2, pp. 1–31.
- 420 Coloma, G. (2012). “The Importance of Ten Phonetic Characteristics to Define Dialect
421 Areas in Spanish”. In: *Dialectologia: revista electrònica* 9, pp. 1–26.
- 422 Embleton, S. (1993). “Multidimensional Scaling as a Dialectometrical Technique. Outline
423 of a Research Project”. In: *Contributions to Quantitative Linguistics*. Ed. by R. Köhler
424 and B. Rieger. Language and Space. Dordrecht: Kluwer, pp. 267–276.
- 425 Embleton, S., D. Uritescu, and E. S. Wheeler (2013). “Defining Dialect Regions with Inter-
426 pretations: Advancing the MDS Approach”. In: *LLC: The Journal of Digital Scholarship*
427 *in the Humanities* 28.1, pp. 13–22.
- 428 Felsenstein, J. (2004). *Inferring Phylogenies*. 2nd ed. Sinauer Associates Sunderland.
- 429 Goebel, H. (2006). “Recent Advances in Salzburg Dialectometry”. In: *LLC: The Journal of*
430 *Digital Scholarship in the Humanities* 21.4, pp. 411–435.
- 431 Goeman, T. and J. Taeldeman (1996). “Fonologie en morfologie van de Nederlandse di-
432 alecten: Een nieuwe materiaalverzameling en twee nieuwe atlasprojecten”. In: *Taal en*
433 *Tongval* 48, pp. 38–59.
- 434 Heeringa, W. (2004). “Measuring Dialect Pronunciation Differences using Levenshtein Dis-
435 tance”. PhD thesis. University of Groningen.
- 436 Jain, A. K. and R. C. Dubes (1988). *Algorithms for Clustering Data*. Upper Saddle River,
437 New Jersey: Prentice-Hall, Inc.
- 438 Leinonen, T. (2010). “An Acoustic Analysis of Vowel Pronunciation in Swedish dialects”.
439 PhD thesis. University of Groningen.
- 440 — (in press). “Dialektgeografi och dialektometri”. In: *Sociolingvistik i praktiken*. Ed. by S.
441 Boyd and S. Ericsson. Lund: Studentlitteratur.
- 442 Levenshtein, V. I. (1966). “Binary Codes Capable of Correcting Deletions, Insertions, and
443 Reversals”. In: *Soviet Physics Doklady* 10, pp. 707–710.
- 444 List, J.-M. (2012). “Multiple Sequence Alignment in Historical Linguistics”. In: *Proceedings*
445 *of ConSOLE XIX*, pp. 241–260.
- 446 Manni, F. et al. (2008). “Do Surname Differences Mirror Dialect Variation?” In: *Human*
447 *Biology* 80.1, pp. 41–64.
- 448 Manning, C. D. and H. Schütze (1999). *Foundations of Statistical Natural Language Pro-*
449 *cessing*. Cambridge, Mass.: MIT press.
- 450 Mathussek, A. (2013). “Sprachräume im Mittelfranken”. In: *Handbuch zum Sprachatlas von*
451 *Mittelfranken. Dokumentaiton und Auswertung*. Ed. by H. H. Munske and A. Mathussek.
452 *Schriften zum Bayerischen Sprachatlas. Bnd 9. Heidelberg: Universitätsverlag Winter.*

- 453 Mathussek, A. (2014). *Sprachräume und Sprachgrenzen im Untersuchungsgebiet des Sprachat-*
454 *las von Mittelfranken: Traditionelle Dialektgeographie, Wahrnehmungsdialektologie, Di-*
455 *alektometrie*. Heidelberg: Universitätsverlag Winter.
- 456 — (submitted). “On the Problem of Field Worker Isoglosses”. In: *The Future of Dialects*.
457 Ed. by M.-H. Côté, R. Knooihuizen, and J. Nerbonne. Language Variation 1. Berlin:
458 Language Science Press.
- 459 Mitterhofer, B. (2013). *Lessons from a Dialect Survey of Bena: Analyzing Word Lists*. SIL
460 International Electronic Survey Report 2013-020. Summer Institute of Linguistics.
- 461 Moran, S. and J. Prokić (2013). “Investigating the Genealogical Relatedness of the Endan-
- 462 gered Dogon Languages”. In: *LLC: The Journal of Digital Humanities Scholarship* 28
463 (4), pp. 676–791.
- 464 Nerbonne, J. (2011). “Mapping Aggregate Variation”. In: *Language Mapping*. Ed. by A.
465 Lameli, R. Kehrein, and S. Rabanus. Language and Space 2. Berlin: De Gruyter, pp. 476–
466 501.
- 467 Nerbonne, J. and P. Kleiweg (2003). “Lexical Distance in LAMSAS”. In: *Computers and*
468 *the Humanities* 37.3, pp. 339–357.
- 469 Nerbonne, J. et al. (2008). “Projecting Dialect Differences to Geography: Bootstrap Cluster-
- 470 ing vs. Noisy Clustering”. In: *Data Analysis, Machine Learning, and Applications. Proc.*
471 *of the 31st Annual Meeting of the German Classification Society*. Ed. by C. Preisach
472 et al. Berlin: Springer, pp. 647–654.
- 473 Nerbonne, J. et al. (2011). “Gabmap – A Web Application for Dialectology”. In: *Dialectologia*
474 Special Issue II, pp. 65–89.
- 475 Peres-Neto, P. R. and D. A. Jackson (2001). “How well do Multivariate Data Sets Match?
476 The Advantages of a Procrustean Superimposition Approach over the Mantel Test”. In:
477 *Oecologia* 129.2, pp. 169–178.
- 478 Prokić, J., C. Coltekin, and J. Nerbonne (2012). “Detecting Shibboleths”. In: *Proceedings*
479 *of the EACL 2012 Joint Workshop of LINGVIS & UNCLH*, 72–80.
- 480 Prokić, J. (2010). “Family Resemblances”. PhD thesis. University of Groningen.
- 481 Reber, J. (2013). “Strukturen und Muster in der Namenwelt: Quantitative und qualita-
- 482 tive Untersuchungen zum Toponymenbestand der beiden Solothurner Amteien Dorneck-
- 483 Thierstein und Olten-Gösgen”. PhD thesis. Universität Basel.
- 484 Schalkoff, R. J. (1992). *Pattern Recognition*. New York: Wiley & Sons.
- 485 Scherrer, Y. (2012). “Recovering Dialect Geography from an Unaligned Comparable Cor-
- 486 pus”. In: *In proceedings of the EACL Workshop on Visualization of Language Patterns*
487 *and Uncovering Language History from Multilingual Resources (LINGVIS & UNCLH)*.
- 488 Šimičić, L. et al. (2013). “Diatopic Patterning of Croatian Varieties in the Adriatic Region”.
489 In: *Journal of Slavic Linguistics* 21.2, pp. 259–301.

- 490 Snoek, C. (2013). “Using Semantically Restricted Word-Lists to Investigate Relationships
491 among Athapaskan Languages”. In: *Approaches to measuring linguistic differences*. Ed.
492 by L. Borin and A. Saxena. Trends in Linguistics Studies and Monographs 265. Germany:
493 de Gruyter Mouton, pp. 231–248.
- 494 — (2014). “Review of Gabmap: Doing Dialect Analysis on the Web”. In: *Language Docu-
495 mentation and Preservation* 8, pp. 192–208.
- 496 Stefanowitsch, A. and S. T. Gries (2003). “Collostructions: Investigating the Interaction of
497 Words and Constructions”. In: *International Journal of Corpus Linguistics* 8.2, pp. 209–
498 243.
- 499 Uiboaed, K. et al. (2013). “Variation of Verbal Constructions in Estonian Dialects”. In:
500 *Literary and Linguistic Computing* 28.1, pp. 42–62.
- 501 Wieling, M. (2013). “Voices Dialectometry at the University of Groningen”. In: *Analyzing
502 21st-century British English: Conceptual and Methodological Aspects of the BBC ‘Voices’
503 Project*. Ed. by C. Upton and B. Davies. London: Routledge, pp. 208–218.
- 504 Wieling, M., W. Heeringa, and J. Nerbonne (2007). “An Aggregate Analysis of Pronunci-
505 ation in the Goeman-Taeldeman-van Reenen-Project Data”. In: *Taal en Tongval* 59.1,
506 pp. 84–116.
- 507 Wieling, M., E. Margaretha, and J. Nerbonne (2012). “Inducing a Measure of Phonetic
508 Similarity from Pronunciation Variation”. In: *Journal of Phonetics* 40.2, pp. 307–314.
- 509 Wieling, M. and J. Nerbonne (2011a). “Bipartite Spectral Graph Partitioning for Cluster-
510 ing Dialect Varieties and Detecting their Linguistic Features”. In: *Computer Speech &
511 Language* 25.3, pp. 700–715.
- 512 — (2011b). “Measuring Linguistic Variation Commensurably”. In: *Dialectologia* Special Is-
513 sue II. Spec. Iss. on Production, Perception and Attitude ed. by J. Nerbonne, S. Gron-
514 delaelers, D. Speelman & M.-P. Perea, pp. 141–162.
- 515 Wieling, M., C. Upton, and A. Thompson (2014). “Analyzing the BBC Voices data: Con-
516 temporary English Dialect Areas and their Characteristic Lexical Variants”. In: *LLC:
517 The Journal of Digital Scholarship in the Humanities* 29.1, pp. 107–117.