

# Grouping Spanish-speaking countries by dialect: An exploratory corpus dialectometric approach

David Ellingson Eddington  
Brigham Young University  
eddington@byu.edu



Received: 12-02-2022  
Accepted: 28-05-2022  
Published: 04-06-2022

How to cite: Eddington, David E. 2022. Grouping Spanish-speaking countries by dialect: An exploratory corpus dialectometric approach. *Isogloss. Open Journal of Romance Linguistics* 8(1)/9, 1-30.  
DOI: <https://doi.org/10.5565/rev/isogloss.207>

## Abstract

---

The present study attempts to cluster Spanish-speaking countries into dialect regions by computational means. The frequencies of 592 lexical and grammatical features for 21 countries were obtained from the Corpus del Español-Web Dialects. Principal components analysis and hierarchical clustering analyses used the resulting data to group countries into dialect regions. A number of algorithms were used to rank features in terms of how much they aided in dialect classification, which allowed grouping based on a smaller set of features.

Six dialect zones were identified: European (Spain), Southern Cone (Uruguay, Argentina), Southern Central America (Costa Rica, Panama), Caribbean (Puerto Rico, Dominican Republic), Northern Central America (Nicaragua, El Salvador, Guatemala, Honduras), Andean South America (Bolivia, Paraguay, Chile, Peru). However, different subsets of features, and different clustering algorithms produced groupings that varied somewhat. The bulk of the variation dealt with where Cuba, Ecuador, Mexico, Venezuela, Colombia, and the US fit into the dialect regions.

---

The difficulties of the computational approach to dialect classification are discussed. Allowing computer algorithms to determine dialect boundaries appears objective. However, interpreting a principal components analysis entails a degree of subjectivity. Furthermore, the plethora of different classification algorithms allows the researcher to choose the one that produces the desired outcome.

**Keywords:** dialectometry, Spanish dialects, corpus approach, statistical analysis.

---

## 1. Introduction

Early studies in the field of dialectology were carried out by interviewing speakers, extracting features from their speech, and then placing isoglosses on a map to delineate where the boundaries between features existed spatially. Since isoglosses for different features are notorious for not coinciding with each other, except where topographic features such as oceans and mountain ranges are found, determining the exact boundary between dialects was difficult. The use of isoglosses is also problematic in another way since linguistic features are rarely binary, as isoglosses suggest, but scalar in nature. What is more, sociolinguistic studies have uncovered the wealth of variation that exists within the bounds of what may be considered a single dialect.

Delineating dialect areas in the Spanish-speaking world has been the focus of many linguistic investigations (see Rodríguez Vázquez 2019 for a review). The early studies carried out by Armas y Céspedes (1882) and Wagner (1920) were followed by more substantial investigations that divide the Spanish-speaking world into different dialect areas in a number of different ways. Canfield (1962) makes a tripartite division, while Henríquez Ureña (1921) posits five regions. Zamora & Guitart's (1988) division includes nine and Rona (1964) suggests 16. The differences between the dialect boundaries that have been proposed is principally the result of different criteria employed by each researcher. For example, Resnick (1975) bases his on eight phonetic differences. On the one hand, we hope that precise boundaries will be found once enough features have been considered. On the other hand, reducing the complexities of language and language variation to a series of lines on a map can sometimes seem like a futile endeavor (Alba 1992).

In any event, the invention of the internet, the widespread availability of powerful computers, and the existence of large corpora have led to innovative approaches to dialect studies. The most notable characteristic of contemporary approaches is that they do not depend on small numbers of features, but follow the advice of researchers who argue that dialectology must aggregate large number of features to obtain maximally reliable results (e.g. Nerbonne 2009, Séguy 1971). Among these is the use of data from Twitter. Every second 6,000 short messages are broadcast to the world as tweets, and many of these contain geotags that allows their authors to be mapped in space. The sheer amount of language data produced in tweets makes many a linguist feel like a kid in a candy store. Some have taken advantage of the data to delineate dialect boundaries in the US (Huang et al 2016), while others have examined dialect boundaries within a single Spanish-speaking country such as Columbia (Rodríguez-Díaz et al. 2018) and Spain (Aliaga Jiménez 2003, Donoso & Sánchez 2017, García Mouton 1991, Moreno Fernández 1991). More germane to the

present paper are studies of tweets in the Spanish-speaking world (e.g. Brown 2016, Gonçalves & Sánchez 2014).

The copious amount of data produced by tweeters requires a systematic way to examine them. One approach is that of Tellez et al. (2021) who take an include-almost-everything approach to their analysis of 800 million tweets, in which they only exclude the 100 most frequent words, and very infrequent words, but retain everything else. They give their results in terms of scalar similarities rather than setting firm dialect boundaries. Gonçalves & Sánchez (2016) take a more manageable sample of 331 words that represent 46 different concepts taken from the VARILEX project (Tinoco & Ueda 2007). For example, a ‘merry-go-round’ is known as a *caballitos*, *calesita*, *carrusel*, *tiovivo*, or *machina* in different regions. Their analysis combines countries into three groups: 1) Spain, 2) Uruguay, Argentina, Paraguay, 3) all other countries (see also Moreno Fernández & Ueda 2018).

Another approach to determining dialect boundaries involves data from surveys. For example, Burrige et al. (2019) used the results of the Cambridge Online Survey of World Englishes to map dialect areas. That survey was principally based on vocabulary differences such as different words for traffic circle, tennis shoes, and pill bugs. In a similar vein, the VARILEX database (Tinoco & Ueda 2007) contains 2382 words representing 206 different concepts in Spanish. Among these are words for closet, ring, and suspenders. Speakers from 47 Spanish-speaking cities, principally capital cities, were asked to choose which word they use for each concept. Based on the results of the survey, Ueda (2009) suggests six major dialect areas: 1) Spain, 2) Caribbean: Puerto Rico, Cuba, Dominican Republic 3) Mexico, 4) Central America: Guatemala, Costa Rica, Panama, Honduras, Nicaragua, Colombia, Venezuela, 5) Andean countries: Peru, Bolivia, Ecuador, 6) Southern Cone: Chile, Uruguay, Paraguay, Argentina.

An innovative approach to grouping countries is that of Quesada Pacheco (2014), which involved asking speakers from each country which countries sounded most similar to their own. According to speakers’ perceptions eight divisions exist: 1) Cuba, Puerto Rico, Panama, Northern Venezuela, Northern Colombia, 2) Ecuador, Peru, Bolivia, Southern Venezuela, Southern Colombia, 3) Uruguay, Chile, Paraguay, Argentina, 4) Mexico, 5) Guatemala, 6) Honduras, 7) Costa Rica, 8) El Salvador. What makes his findings unusual is the fine-grained differentiation between Central America countries, rather than their conglomeration.

All of these approaches fall into what has been called dialectometry which can be defined as using computational means to study dialectal differences (see Wieling & Nerbonne 2015 for an overview). One subdomain of dialectometry is corpus-based dialectometry, which involves computational analyses of corpora (Szmrecsanyi 2011). One example of using corpora to delineate dialect boundaries is Grieve’s study (2012) of letters written to the editor in a number of major cities in the US. In those letters he examined variation in adverb placement (e.g. *often repeated* versus *repeated often*). In another corpus Grieve (2011) studied the variable use of contractions (e.g. *don’t* versus *do not*, 2011) in American English. In both studies, and used the results to computationally determine dialect areas.

The first attempts at grouping Spanish-speaking countries according to their linguistic similarities relied on the author’s personal experience, reports of dialect features reported by other researchers, and small dialect surveys. However, recent advances in technology have made it possible to expand on previous work by examining much more extensive data sets. For example, a number of large-scale

surveys have been carried out and linguistic features have been examined in large collections of tweets. One method that has yet to be applied to the task of associating Spanish-speaking countries into dialect groups is using extant corpora. In the present study, the Corpus del Español-Web Dialects (Davies 2017) is used to this end. In addition to establishing country clusterings another goal of the study is highlighting the features that are most helpful in making those clusterings. In the following sections, a number of statistical and computational methods are described which are applied to the task. The novel use of these methods to investigate the question at hand means that the results can be viewed as exploratory in nature. The challenges that this corpus-based dialectometry approach presents are discussed as well.

## 2. Data Set

One way subjectivity may creep into an analysis is in the choice of the features used. In order to address this issue a large number of features should be included computational algorithms should be used to choose the most important ones. The analysis described below is based on 592 features and their corpus frequencies (DOI 10.17605/OSF.IO/892MW). These features were chosen since they have been used in previous variationist studies. Of these, 45 come from Eddington (2021). These include six grammatical differences such as the frequency of the use of present perfect versus preterite (e.g. *Esta mañana he comido huevos* ‘This morning I’ve eaten eggs’), the use of present subjunctive in embedded clauses with present tense matrix clauses (e.g. *Le pedí que no lo haga/hiciera* ‘I asked him not to do it’), five nominal gender variations (*el/la sarten* ‘the pan’), and 34 vocabulary differences (*valija / maleta* ‘suitcase’). Also included were the items from VARILEX used by Gonçalves & Sánchez (2014). This consisted of 454 lexical items for 43 concepts (e.g. ‘sidewalk’ *acera, andén, badén, calzada, contén, escarpa, vereda*). The per million frequencies of these lexical items in each country was obtained from the Corpus de Español-Web Dialects corpus (Davies 2017, extracted data: DOI 10.17605/OSF.IO/892MW). This 2 billion word corpus comprises from 24 to 440 million words per country, 78% of which are from Latin American sources. Some of the lexical variants from VARILEX had a frequency of zero in the Corpus del Español. This meant that they were not helpful for the purposes of the present study and were not included the analysis. In addition to these words, a handful of other lexical items were added as well (e.g. *zumos, coges, damascos, gizes* ‘juice, to get, apricot, chalk’) whose per million frequency also came from Davies (2017).

In many ways using the Corpus de Español-Web Dialects serves the purposes of the study well. It contains a sizable amount of data from each country. What is more, it was designed to represent less formal levels of language since 60% of it derives from blogs. There are, however, serious limitations. For example, information about the writers (e.g. age, gender, social class) is not available. The possibility also exists that a writer from one country may be included in the corpus from another country. These are issues that apply to the whole of corpus linguistics. More importantly, level of granularity is based on country boundaries. That means that different varieties in a single country are lumped together. It is hoped that the by country results reported herein may serve as a guide for future research that makes use of subtler geographic distinctions.

### 3.1. Evaluating different feature subsets

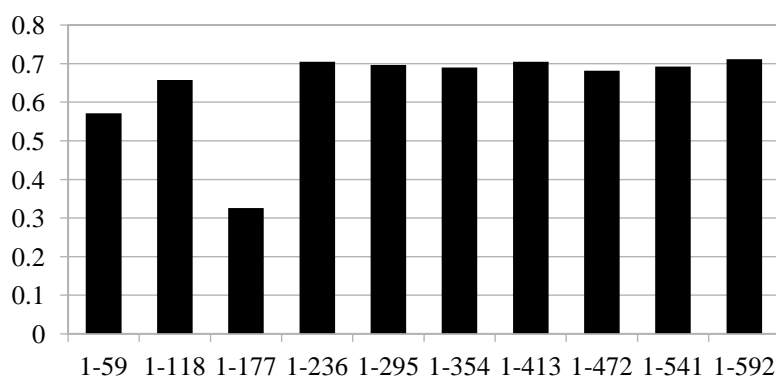
The idea behind agglomerating large numbers of features is to eliminate the influence that a handful of features may have, and also to keep the researcher from picking and choosing only those features that may align with his or her preconceived notions of what countries should be clustered together. However, an important aspect of clustering is determining how many features provide optimal results, which features those are, and if systematically eliminating features helps the task of classification.

One method of reducing the number of variables to a more manageable size is principal components analysis (PCA). PCA is a commonly used procedure in dialect studies of the kind presented here (e.g. Huang et al. 2016, Manni et al. 2008, Moreno Fernández & Ueda 2018). For a discussion of how PCA compares to other methods see Leino et al. 2008. PCA is an exploratory procedure designed to reduce a large number of variables to a more manageable and easily interpretable format. Rather than eliminating variables altogether, PCA creates new variables from the original variables called principal components. These components retain as much of the original variation in the data as possible. The principal components are ordered so that the first one accounts for more of the variation than the second, and so on. In general, the first two components are the most important.

In order to compare analyses with different subsets of variables, there is a need for a standard of comparison. Geography itself was used in the initial analyses. This was done by calculating the distance in kilometers between each country's capital city as a point of reference. The Euclidean distance between the first two dimensions of each of the principal components analyses (PCA) described below was calculated, and a correlation between the distance between each capital city and the euclidean distance between the two PCA dimensions was performed. The resulting  $r^2$  provides a point of comparison. Analyses with similar  $r^2$  values place the dialects in a similar geographical position in contrast to analyses with different  $r^2$  values.

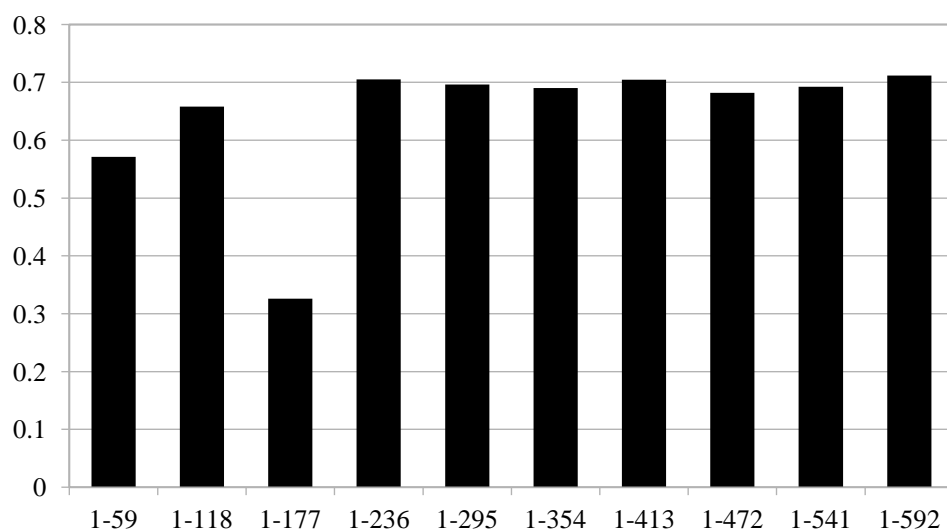
The first question to be examined is how many features are needed to make good dialect groupings. To answer this, the 592 features were ordered randomly and divided into ten groups of about 59 features each. A PCA was run on each group and the resulting  $r^2$  values appear in Figure 1. The wide discrepancies between subsets is a clear indicator that some groups of features situate countries differently with respect to each other. Perhaps this problem is the result of agglomerating too few features.

**Figure 1.**  $R^2$  between PCA and capital city distances for each subset of 59 features



To address the issue of how many features are needed, the 10 subsets were recombined in the following manner. The first subset comprises the first 1-59 features, the second adds another 59 features and comprises features 1-118. Each subset grows in this manner until the last group includes all 592 features. The results of these analyses appear in Figure 2. What is clear is that once 236 features are included an overlap between the results of the PCA and the distance between capital cities reaches about 70%, a level no subset of 59 features alone reaches. However, increasing the number of features beyond that point does little to change the spatial grouping of the countries in relation to their dialectal features. The clear takeaway is that good predictions may be made with only a subset of the features considered.

**Figure 2.**  $R^2$  Between PCA and capital city distances for subsets of increasing numbers of features

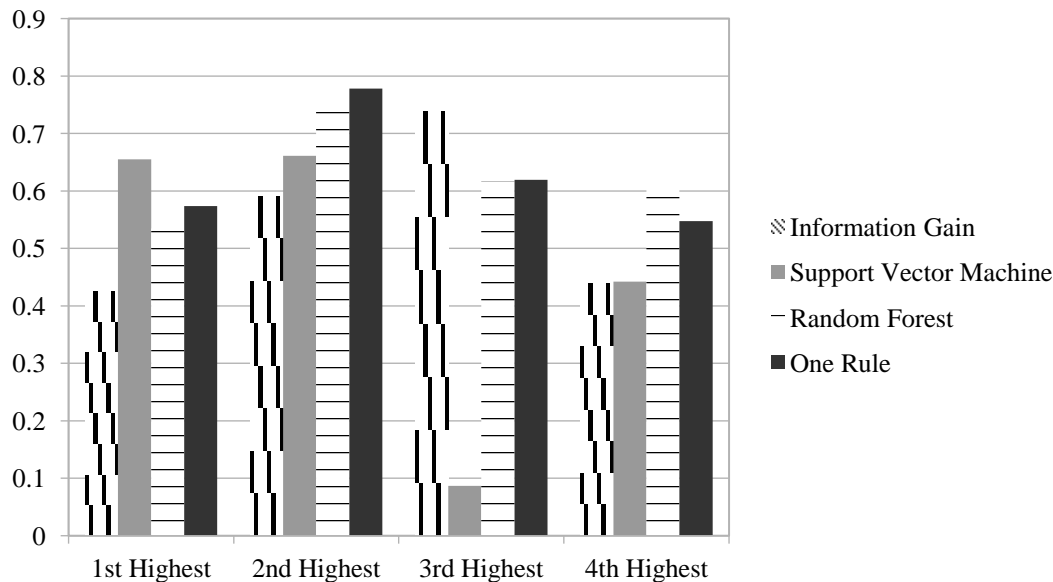


The question now is how to reduce the set of features used by finding the most relevant ones. Given the large numbers of algorithms that one may choose to eliminate variables, it could be tempting to sort through them until one is found that supports the researcher's hypothesis. In order to avoid this, four algorithms were applied to the data, all of which ranked the features in terms of how well they classify the countries. The first was a random forest analysis carried out in R (R Core Team 2020). The remaining three were carried out using WEKA (Frank et al. 2016). Holte's (1993) one rule attribute evaluator determines the worth of each feature in terms of how much it contributes to classifying countries, and ranks the features accordingly. The information gain attribute evaluator calculates the value of each feature by measuring its information gain in respect to each country. The SVM classifier (Guyon et al. 2002) uses a support vector machine algorithm to determine the value of each feature in classifying countries. The resulting rankings from each algorithm were ordered from best to worst and divided into quartiles of 148 features. The rankings produced by each algorithm appear in the appendix.

As Figure 3 illustrates, the feature ranking algorithms produce different results. What is surprising is that the highest ranked 148 features chosen by three of the algorithms place the country's dialects farther from their capitals than do the second highest set of 148 features. However, this may merely suggest that the features that best group the countries do not correlate as highly with the geographical location of the country's capital cities, since the  $r^2$  is not a measure of goodness of linguistic fit.

Nevertheless, we can now begin to use the data to determine dialect areas, as well as to hone in on the features that best define those dialect areas.

**Figure 3.**  $R^2$  Between PCA and capital city distances for subsets of ranked features



An initial foray into dialect grouping was done with a hierarchical clustering dendrogram (Seol 2020) using the Jamovi statistical package (Jamovi 2021) which is a graphical user interface for R (R Core Team 2020). Hierarchical clustering is a method commonly used in dialect studies (e.g. Leino, Antti, & Saara Hyvönen, Nagy et al. 2006, Moreno Fernández & Ueda 2018, Sato & Hefernan 2018). It clusters countries according to their similarities based on the features. Although we have seen that some features are more important than other in making dialectal groupings, as a point of comparison the dendrogram in Figure 4 was built using all 592 features. The algorithm groups the southern cone countries of Uruguay, Paraguay, and Argentina together, and places Spain on its own branch, both of which seem intuitive. However, considering Cuba and Chile as isolated dialects, and combining countries in North, Central, and South America into a single dialect group runs counter to all previous classifications.

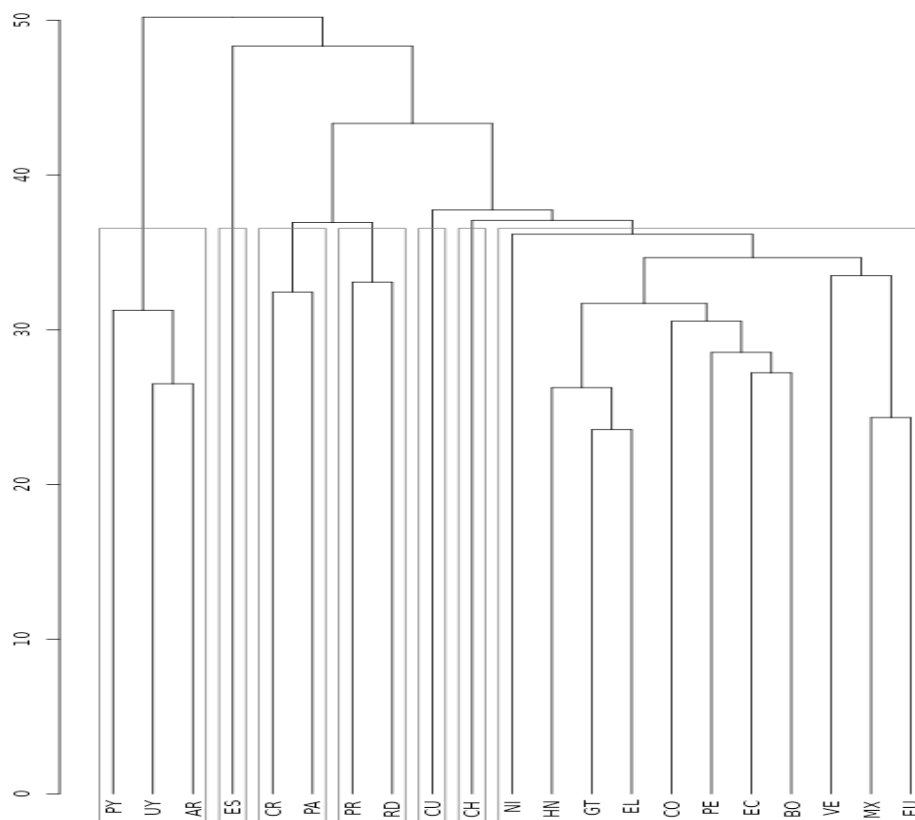
### 3.2. Evaluating the most highly ranked subset of features

Four algorithms were used to evaluate the worth of the features in terms of their ability to find similarities between countries. The 148 highest ranked features chosen by the support vector machine achieved the greatest overlap (66%) with the geographical location of the capital cities (Figure 3) and is used in this first analysis. When these features are considered a smaller, but tentative, grouping of countries into seven dialect areas emerges (Figure 5):

1. European (Spain)
2. Southern Cone (Uruguay, Argentina)
3. Southern Central America (Costa Rica, Panama)
4. Caribbean (Puerto Rico, Dominican Republic)
5. North America (Cuba, United States, Mexico)
6. Northern Central America (Nicaragua, El Salvador, Guatemala, Honduras)

## 7. Andean and Northern South America (Bolivia, Paraguay, Chile, Venezuela, Colombia, Ecuador, Peru)

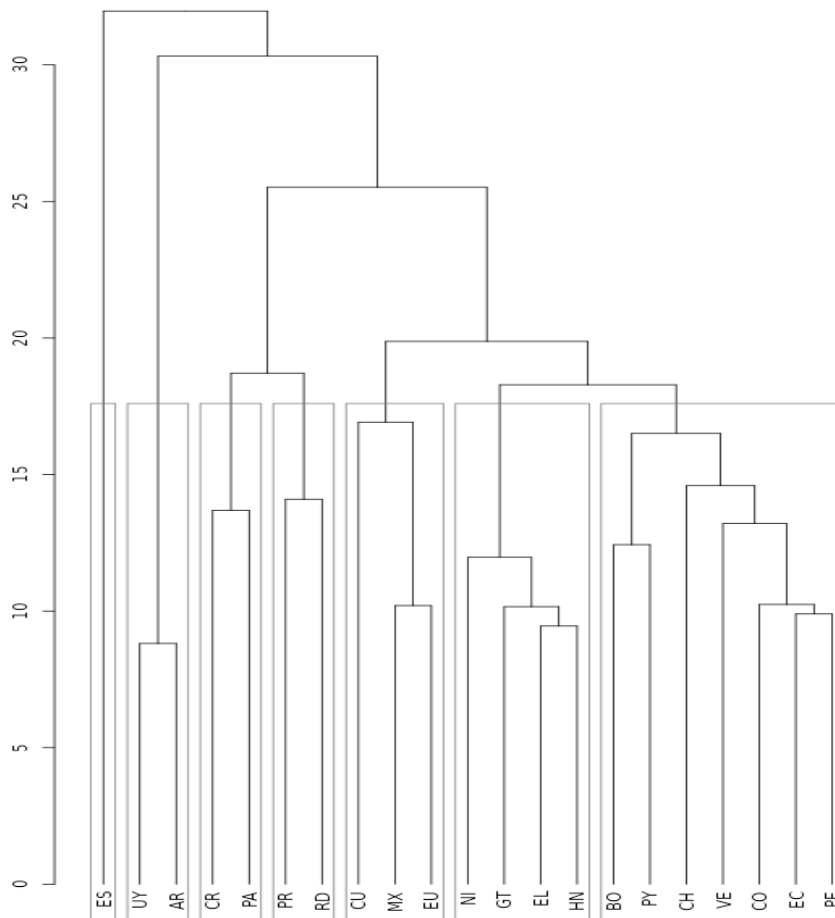
**Figure 4.** Hierarchical clustering dendrogram using all 592 features



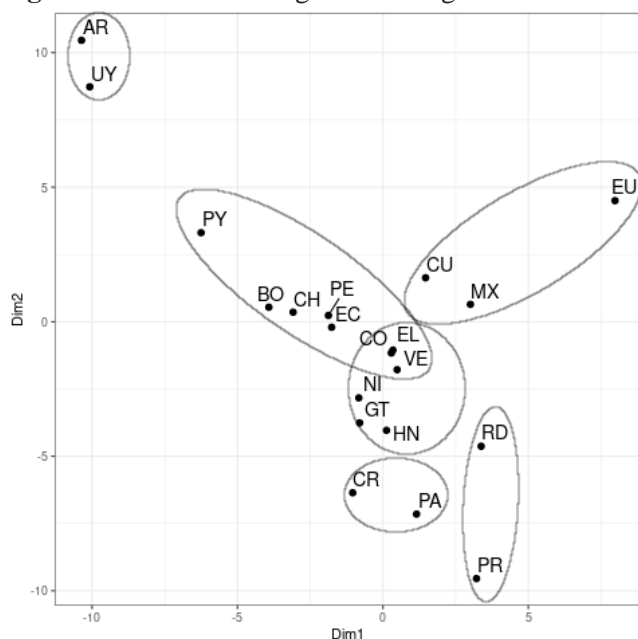
PY=Paraguay, UY=Uruguay, AR=Argentina, CR=Costa Rica, PA=Panama, PR=Puerto Rico, RD=Dominican Republic, CU=Cuba, NI=Nicaragua, GT=Guatemala, EL=El Salvador, CO=Colombia, PE=Peru, EC=Ecuador, BO=Bolivia, VE=Venezuela, MX=Mexico, EU=United States

Another way of visualizing these dialect groups is by plotting the first two dimensions of the PCA of the same 148 highest ranked support vector machine features (Figure 6). There are three reasons for using both PCAs and cluster dendograms to evaluate groupings. The first is that when different methods produce similar outcomes, the results are considered more robust, and less algorithm specific. Secondly, while dendograms cluster countries, they do not give a sense of distance between countries and clusters than the two dimensional representation in a PCA provides. Finally, presenting the results of both analyses reduces the chances that a specific method is chosen simply because it yields the expected outcome. In any event, Spain is an outlier falling far from the Latin American countries. For this reason it was excluded from Figure 6 so that the remaining countries would be better separated.



**Figure 5.** Hierarchical Clustering Dendrogram Using the 148 Highest Ranked SVM Features

The dendrogram clusterings are indicated in the two dimensional space in Figure 6 with ellipses. It should be apparent that there are significant differences between the dendrogram and the PCA plot. For example, the dendrogram places Cuba and El Salvador in different groups, while in the PCA plot the two countries nearly overlap. The dendrogram's placement of Cuba along with Mexico and the United States, rather than with other Caribbean countries, is unusual. While the three countries are grouped in the hierarchical dendrogram, the PCA plot places Cuba closer to other Latin American countries than to the US. Most dialectologists include Cuba alongside the Dominican Republic and Puerto Rico (Canfield, 1962, Henríquez Ureña 1921, Rona 1964, Quesada Pacheco 2014, Zamora & Guitart 1988, Wagner 1920). On the other hand, the similarities between the countries capture the fact that the largest dialectal influences on the Spanish of the US are arguably Mexico and Cuba. The dialectal separation of Uruguay and Argentina from other South American nations has been noted since the 19th century (Armas Céspedes 1882). Paraguay falls closest to these countries, and some dialectologists include Paraguay, or parts of Paraguay, in southern cone varieties (Cahuzac 1980, Canfield, 1962, Henríquez Ureña 1921, Rona 1964, Quesada Pacheco 2014, Zamora & Guitart 1988).

**Figure 6.** PCA Plot Using the 148 Highest Ranked SVM Features

The dialectal placement of Central American countries is debated. Some cluster them with Mexico (Cahuzac 1980, Henríquez Ureña 1921), while others consider Central America a separate, but united dialect area (Rona 1964, Zamora & Guitart 1988). Still others separate Central American countries into different dialects (Canfield 1962, Quesada Pacheco 2014). The data from the present study support the existence of two dialects in Central America.

What countries Venezuela and Colombia should be aligned with is debated as well. This is principally because both countries are divided between Caribbean and non-Caribbean dialect zones. The similarities between Colombia and the Andean countries (Peru, Ecuador, Chile, Bolivia) that fall out in the present study may be explained by the fact that the majority of Colombians do not live on the Caribbean. As a result, this may skew the data compiled in the Corpus del Español toward non-Caribbean Colombian speech. The same argument, however, does not hold up for Venezuela, where the bulk of the population is concentrated on the Caribbean coast. In Figure 6, Venezuela and Colombia appear in overlapping area of the ellipses representing Northern Central American and Andean and Caribbean South American dialects, falling particularly close to El Salvador. In contrast, the dendrogram does not capture the similarity between Venezuela and Colombia and Central American Countries.

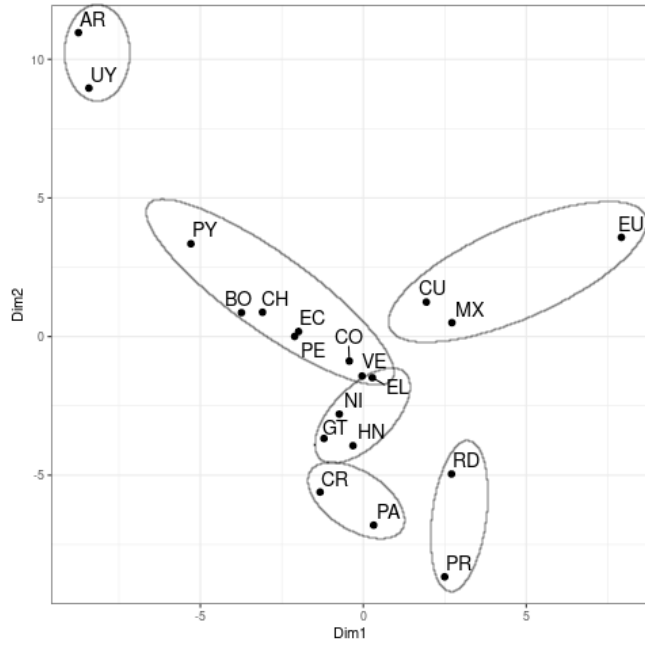
The 148 highest ranked features chosen by the support vector machine provides an analyses that more closely corresponds to extant divisions than does the outcome using all 592 features. The question now is exactly how these features are related to the dialect areas, and if it is possible to further eliminate some. To do this, the seven dialect regions proposed were correlated with the 148 features (Table 1).

**Table 1.** Features Correlated With Dialect Regions

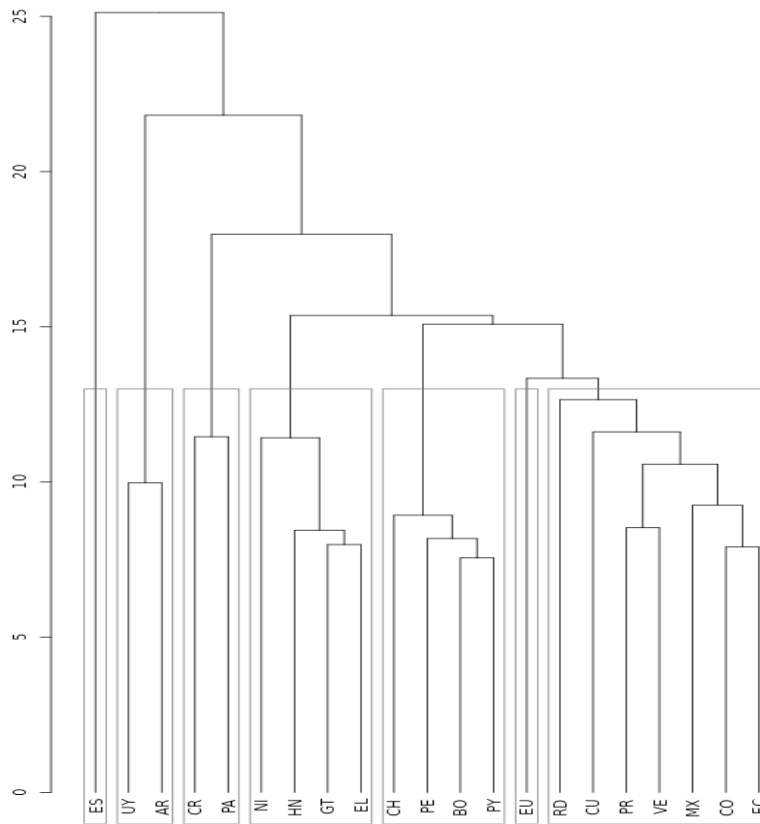
European	North Am.	Northern Central Am.	Southern Central Am.	Caribbean	Andean and Northern South Am.	Southern Cone
-ase past subj.	afiche (-)	auto (-)	acera	acera	arañazo (-)	alpargata
aero	ascensor (-)	autopista (-)	bar	anuncio	ascensor	atar
afiche (-)	bolígrafo	barca	bluyín	aro	bar (-)	auriculares
alubia	capaz	bocadillo	brassiere	autopista	barca (-)	auto
arañazo	colegio (-)	cantina	bus	brasiere	cartelón (-)	cabeza dura
armario	elevador	cartel (-)	cachetes	caldero	chancleta (-)	calesita
ático	habían with plural (-)	finca	cantina	cartelón	colegio	chance
auriculares	magnetófono	gasolinera	carro	cristal delantero	elevador (-)	colchón
balacera (-)	refrigerador	habían with plural	caterpillar	entretenimiento	escuela (-)	estación de servicio
barca	tú	magnetófono (-)	chancleta	escuela	falencia	frigorífico
bolígrafo		mico	escurridor	escurridor	gasolinera (-)	grabadora
cacahuete		quizás (-)	estola	estola	ocasión (-)	guardarropas
camarero		tablero	finca	goma de mascar	vitrina	habían with plural (-)
capaz			frigorífico	guineo		heladera
cazo			grasoso	ocasión		lavadora (-)
celular (-)			jeans	papel encerado		lavarropa
coger			lámina de queso	puerco		maleta (-)
constipado			lancha	salón		ocasión
escuela (-)			lavadora	sarten		parlante
estadía (-)			pasatiempo	ser consciente (-)		pasatiempo (-)
farola			poste de luz	sortija		perchero
gafas			queso americano	tiroteo		pizarrón
grasiento			rebanada de queso			present perfect (-)
guapa			sartén			reposera
lavadora			tiesto			ropero
lois			tirantes			ser consciente
mechero						sos
melocotón						sponsor
mofletes						valija
ordenador						veliz

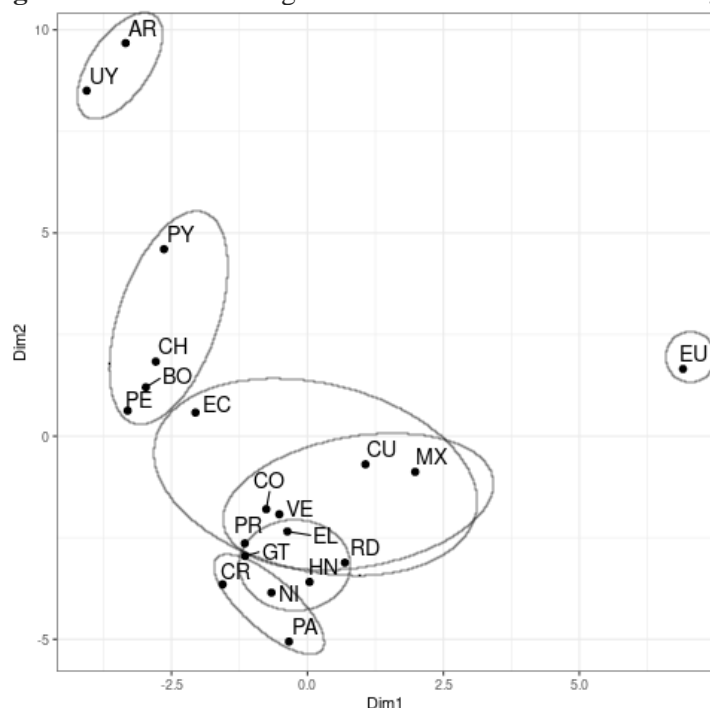


**Figure 7.** PCA Plot Using the 123 Highest Ranked SVM Features that Correlated with the 7 Proposed Dialect Areas



**Figure 8.** Hierarchical Clustering Dendrogram Using the Features from the One Rule Algorithm



**Figure 9.** PCA Plot Using the Features from the One Rule Algorithm

#### 4. Conclusions

The present study demonstrates how statistical and data mining methods can be applied to corpus data in order to group Spanish-speaking countries into dialect areas. In early attempts at delineating dialect areas researchers hand-picked features that formed the basis for their groupings. The variety of resulting clusterings may be due to the selection of features chosen. On the surface, using computational means to choose which features, and how many features to include in an analysis, seems to be a more objective method. In reality, it merely shifts the issue from cherry picking features to cherry picking the algorithm used to select the features.

As we saw in the final analysis, the hierarchical clustering dendrogram groups countries differently than the PCA plot does even when they both use the same data. An additional problem with PCA plots is that they allow varying interpretations. While some countries form clear clusters situated at a distance from other countries, when the countries are plotted closer together, clustering them into groups with ellipses by hand can become an extremely subjective task.

In this analyses described above, only two clustering algorithms were used to determine dialect boundaries. There are, however, myriads of different methods for clustering features and many algorithms for measuring the distance between entities such as countries, the combination of which can result in widely varying outcomes. As a result, the temptation to search through all of the combinatorial possibilities until one stumbles upon an analytical method that supports one's hypothesis becomes real.

The solution to this dilemma is to not rely on a single analysis. When many feature sets are evaluated using a number of different computational means, a consensus should begin to emerge, and it is that consensus, not the results of a single study that should be the focus of attention. With the limited number of analyses presented above, 15 of the countries consistently cluster into six dialect areas:

1. European (Spain)
2. Southern Cone (Uruguay, Argentina)
3. Southern Central America (Costa Rica, Panama)
4. Caribbean (Puerto Rico, Dominican Republic)
5. Northern Central America (Nicaragua, El Salvador, Guatemala, Honduras)
6. Andean South America (Bolivia, Paraguay, Chile, Peru)

What the present analyses do not firmly establish is which dialect regions Cuba, Ecuador, Mexico, Venezuela, Colombia, and the US belong to, nor whether there may be additional dialect zones not considered. Another limitation of the present study is that it relies solely on political boundaries. Dialects do not necessarily align with such boundaries. Studies that make use of geotagged tweets are in a good position to find dialect zones that transcend the limits of individual countries. In sum, using computational methods and large numbers of features to delineate dialect boundaries is an improvement over earlier methods, the use of these methods opens up another set of issues that must be dealt with.

### Acknowledgments

I appreciate the input by the reviewers as well as the suggestions provided by Earl Brown.

### References

- Alba, Orlando. 1992. Zonificación del español de America. In C. Hernández Alonso (ed.), *Historia y presente del español en America*, 63-84. Valladolid: Junta de Castilla y León.
- Aliaga Jiménez, José Luis. Dialectometría y léxico en las hablas de Teruel. 2003. *ELUA. Estudios de Lingüística* 17: 5-55.
- Brown, Earl. K. 2015. On the utility of combining production data and perceptual data to investigate regional linguistic variation: The case of Spanish experiential *gustar* ‘to like, to please’ on Twitter and in an online survey.” *Journal of Linguistic Geography* 3(2): 47-59. <https://doi.org/10.1017/jlg.2016.1>
- Burridge, J., Vaux, B., Gnacik, M., & Grudeva, Y. 2019. Statistical physics of language maps in the USA. *Physical Review E*, 99(3): 032305. <https://doi.org/10.1103/PhysRevE.99.032305>
- Armas y Céspedes, Juan Ignacio. 1882. *Orígenes del lenguaje criollo*. La Habana: Imprenta de la Viuda de Soler.
- Cahuzac, Philippe. 1980. La división del español de América en zonas dialectales. Situación etnolingüística o semántico-dialectal. *Lingüística Española Actual* 2: 385-461.
- Canfield, D. Lincoln. 1962. *La pronunciación del español en América*. Bogotá: Instituto Caro y Cuervo.
- Davies, Mark. 2017. Corpus del Español, Web/Dialects. <https://www.corpusdelespanol.org/web-dial/>

- Donoso, G., & Sánchez, D. 2017. Dialectometric analysis of language variation in Twitter. *Proceedings of the Fourth Workshop on NLP for Similar Languages, Varieties and Dialects*, 16-25 Valencia, Spain: Association for Computational Linguistics. 10.18653/v1/W17-1202
- Eddington, David Ellingson. 2021. A corpus analysis of some usage differences among Spanish-speaking countries” *Dialectologia* 27: 71-95.
- Frank, Eibe, Mark A. Hall, and Ian H. Witten. 2016. *Data Mining: Practical Machine Learning Tools and Techniques*, 4th Ed. San Francisco, CA: Morgan Kaufmann.
- Embleton, Sheila, Dorin Uritescu, & Eric S. Wheeler. 2013. Defining dialect regions with interpretations: Advancing the multidimensional scaling approach. *Literary and Linguistic Computing* 28: 13-22. <https://doi.org/10.1093/lc/fqs048>
- Henríquez-Ureña, Pedro. 1921. Observaciones sobre el español en América. *Revista de Filología Española* 8: 357-390.
- Holte, Robert C. 1993. Very Simple Classification Rules Perform Well on Most Commonly Used Datasets. *Machine Learning* 11.(1): 63-90.
- García Mouton, Pilar. 1991. *Dialectometría y léxico en Huesca*. Madrid: Consejo Superior de Investigaciones Científicas.
- Gonçalves, Bruno & David Sánchez. 2014. Crowdsourcing dialect characterization through Twitter. *PloS One*, 9(11): e112074. <https://doi.org/10.1371/journal.pone.0112074>; Data: <http://www.bgoncalves.com/languages/spanish.html>
- Gonçalves, Bruno and David Sánchez. 2016. Learning about Spanish dialects through Twitter. *Revista Internacional de Lingüística Iberoamericana* 14: 65-75.
- Grieve, Jack. 2011. A regional analysis of contraction rate in written Standard American English. *International Journal of Corpus Linguistics* 16(4): 514-546. <https://doi.org/10.1075/ijcl.16.4.04gri>
- 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(1): 39-72. <https://doi.org/10.1515/cllt-2012-0003>
- Grieve, Jack. 2014. A comparison of statistical methods for the aggregation of regional linguistic variation. In P. Auer, G. von Essen & W. Frick (eds.), *Aggregating dialectology, typology, and register analysis*, 53-88. Berlin: De Gruyter. <https://doi.org/10.1515/9783110317558.53>
- Guyon, Isabel, Jason Weston, Stephen Barnhill, & Vladimir Vapnik. 2002. Gene selection for cancer classification using support vector machines. *Machine Learning*, 46(1): 389-422.
- Henríquez-Ureña, P. H. 1921. Observaciones sobre el español en América. *Revista de Filología Española* 8: 357-390.
- Holte, Robert C. 1993. Very Simple Classification Rules Perform Well on Most Commonly Used Datasets. *Machine Learning* 11(1): 63-90.
- Huang, Yuan, Diansheng Guo, Alcie Kasakoff, & Jack Grieve. 2016. Understanding US regional linguistic variation with Twitter data analysis. *Computers, Environment*



and Urban Systems 59: 244-255.  
<https://doi.org/10.1016/j.compenvurbsys.2015.12.003>

The jamovi project. 2021. jamovi. (Version 1.6) [Computer Software.] Retrieved from <https://www.jamovi.org>

Leino, Antti, & Saara Hyvönen. 2008. Comparison of component models in analysing the distribution of dialectal features. *International Journal of Humanities and Arts Computing 2*: 73-187. DOI: 10.3366/edinburgh/9780748640300.001.0001

Manni, Franz, Wilbert Heeringa, Bruno Toupance, & John Nerbonne. 2008. Do surname differences mirror dialect variation. *Human Biology 80*: 41-64.

Moreno Fernández, Francisco. 1991. Morfología en el ALEANR: aproximación dialectométrica. In *I curso de geografía lingüística de Aragón*, 289-309. Zaragoza: Institución Fernando el Católico.

Moreno Fernández, Francisco, and Hiroto Ueda. 2018. Cohesion and particularity in the Spanish dialect continuum. *Open Linguistics 4*: 722-742.  
<https://doi.org/10.1515/opli-2018-0035>

Nagy, Naomi, Xiaoli Zhang, George Nagy, and Edgar W. Schneider. 2006. Clustering dialects automatically: A mutual information approach. *University of Pennsylvania Working Papers in Linguistics 12*: 12.

Nerbonne, John. 2009. Data-driven dialectology. *Language and Linguistics Compass 3*(1): 175-198. <https://doi.org/10.1111/j.1749-818X.2008.00114.x>

Quesada Pacheco, Miguel Ángel. 2014. División dialectal del español de América según sus hablantes Análisis dialectológico perceptual. *Boletín de Filología 49*(2): 257-309.

R Core Team 2020. R: A Language and environment for statistical computing. (Version 4.0) [Computer software]. Retrieved from <https://cran.r-project.org>. (R packages retrieved from MRAN snapshot 2020-08-24).

Resnick, Melvyn C. 1975. Phonological Variants and Dialect Identification in Latin American Spanish. Mouton: The Hague.

Rodríguez-Díaz, Carlos A., Sergio Jimenez, George Dueñas, Johnatan Estivan Bonilla, & Alexander F. Gelbukh. 2018. Dialectones: Finding statistically significant dialectal boundaries using twitter data. *Computación y Sistemas 22*(4): 1213-1222.

Rodríguez Vázquez, Paloma. 2019. La zonificación dialectal del español de América: propuestas clásicas y propuestas actuales. Document, Universidade de Dantiago de Compostela. <http://hdl.handle.net/10347/23567>

Rona, José Pedro. 1964. El problema de la división del español americano en zonas dialectales. In F. Moreno Fernández (ed.), *Presente y futuro de la lengua española*, vol. I, 215-226. Madrid: Ediciones Cultura Hispánica

Sato, Yo, and Kevin Heffernan. 2018. Creating Dialect Sub-corpora by Clustering: a case in Japanese for an adaptive method. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*, 3612-3616. Luxemburg: European Language Resources Association.

Sayce, David. n.d.. The Number of tweets per day in 2020. Accessed Feb. 2, 2022. <https://www.dsayce.com/social-media/tweets-day/>

Séguy Jean. 1971. La relation entre la distance spatiale et la distance lexicale. *Revue de Linguistique Romane* 35: 335–57.

Seol, Hyunsoo. 2020. *SnowCluster: Cluster Analysis*. [jamovi module]. <https://github.com/hyunsooseol/snowCluster>

Shackleton Jr, Robert G. 2005. English-American speech relationships: A quantitative approach. *Journal of English Linguistics* 33(2): 99-160. <https://doi.org/10.1177/0075424205279017>

Szmrecsanyi, Benedikt. 2011. Corpus-based dialectometry: a methodological sketch. *Corpora* 6(1): 45-76.

Tellez, Eric. S., Daniela Moctezuma, Sabino Miranda, & Mario Graff. 2021. A large scale lexical and semantic analysis of Spanish language variations in Twitter. *arXiv preprint arXiv:2110.06128*.

Tinoco, Antonio. R., & Hiroto Ueda. 2007. The VARILEX Project-Spanish Lexical Variation. *Linguistica Atlantica* 27: 117-121.

Ueda, Hiroto. 2009. Resultados y proyectos en las investigaciones sobre variación léxica del español. *Dialectologia* 2: 51-80.

Wagner, Max Leopold. 1920. Amerikanisch-Spanisch und Vulgärlatein. *Zeitschrift für Romanische Philologie* 40: 286-312; 385-404.

Wieling, Martijn & John Nerbonne. 2015. Advances in dialectometry. *Annual Review of Linguistics* 1(1): 243-264. <https://doi.org/10.1146/annurev-linguist-030514-124930>

Zamora, Juan y Jorge Guitart. 1988. *Dialectología hispanoamericana. Teoría, descripción, historia*. Salamanca: Almer.

## Appendix

Rankings of features by importance according to four algorithms

Random Forest	Information Gain	One Attribute	Support Vector Machine	Feature
79	293	175	1	<i>bolígrafo</i>
40	355	346	2	present perfect
370	342	234	3	<i>capaz</i>
191	254	468	4	<i>mofletes</i>
139	251	207	5	<i>arañazo</i>
453	390	375	6	<i>tú</i>
194	341	165	7	<i>auto</i>
153	441	192	8	<i>-ase subjunctive</i>
2	298	115	9	<i>estadía</i>
94	465	210	10	<i>armario</i>
304	188	23	11	<i>entretenimiento</i>
311	173	155	12	<i>ático</i>
4	397	389	13	<i>valija</i>

335	461	462	14	<i>maleta</i>
190	400	405	15	<i>vistazo</i>
359	323	224	16	<i>chance</i>
344	445	490	17	<i>lavadora</i>
248	324	69	18	<i>coger</i>
418	264	321	19	<i>sois</i>
565	49	503	20	<i>lois</i>
201	306	251	21	<i>celular</i>
178	185	191	22	<i>alpargata</i>
218	360	206	23	<i>amarrar</i>
234	231	545	24	<i>pómulo</i>
372	127	564	25	<i>pantalón vaquero</i>
252	137	361	26	<i>reposera</i>
81	294	173	27	<i>bocata</i>
369	241	301	28	<i>sandalia</i>
289	31	83	29	<i>grasiento</i>
518	75	127	30	<i>escurridor</i>
73	402	199	31	<i>acera</i>
69	388	308	32	<i>salón</i>
297	165	403	33	<i>vidriera</i>
161	422	423	34	<i>sujetador</i>
309	297	122	35	<i>estación de servicio</i>
182	327	225	36	<i>cartel</i>
155	417	198	37	<i>afiche</i>
3	52	347	38	<i>poste de luz</i>
206	157	535	39	<i>pasatiempo</i>
399	118	280	40	<i>cabeza dura</i>
371	108	349	41	<i>póster</i>
402	460	147	42	<i>finca</i>
430	418	370	43	<i>sos</i>
267	329	18	44	<i>el sarten</i>
7	414	302	45	<i>sartén</i>
54	463	107	46	<i>guapa</i>
320	381	221	47	<i>anuncio</i>
129	73	242	48	<i>cassette</i>
390	94	96	49	<i>guineo</i>
483	361	352	50	<i>quizá</i>
179	412	316	51	<i>se los las decir</i>
89	440	212	52	<i>ascensor</i>
530	379	577	53	<i>papel de plata</i>
245	93	303	54	<i>ropero</i>
293	362	353	55	<i>quizás</i>
86	450	118	56	<i>gafas</i>
510	432	102	57	<i>habían with singular</i>
253	32	478	58	<i>melocotón</i>
432	291	162	59	<i>auriculares</i>
138	98	477	60	<i>mejillas</i>
227	257	284	61	<i>caldero</i>
200	300	171	62	<i>bocadillo</i>
285	38	105	63	<i>guardarropas</i>
277	20	455	64	<i>mantecoso</i>
8	284	26	65	<i>eres</i>

316	398	391	66	<i>veliz</i>
405	344	235	67	<i>cantina</i>
323	33	432	68	<i>tio vivo</i>
12	312	65	69	<i>computadora</i>
364	376	568	70	<i>ocasión</i>
72	289	159	71	<i>autopista</i>
76	263	408	72	<i>tiroteo</i>
25	76	114	73	<i>estola</i>
53	332	6	74	<i>elevador</i>
70	141	145	75	<i>farola</i>
535	99	460	76	<i>magnetófono</i>
172	295	113	77	<i>escuela</i>
440	195	550	78	<i>pizarrón</i>
492	216	61	79	<i>constipado</i>
305	431	101	80	<i>heladera</i>
187	61	252	81	<i>caterpillar</i>
168	81	366	82	<i>rebanada de queso</i>
450	196	335	83	<i>sostén</i>
419	117	579	84	<i>papel encerado</i>
14	317	46	85	<i>closet</i>
215	276	268	86	<i>bus</i>
195	189	295	87	<i>cacahuete</i>
273	346	10	88	<i>distracciones</i>
522	410	330	89	<i>ser consciente</i>
125	154	431	90	<i>tiesto</i>
473	459	186	91	<i>atar</i>
90	279	180	92	<i>barca</i>
11	153	533	93	<i>parlante</i>
34	521	29	94	<i>engrasado</i>
156	12	465	95	<i>mechero</i>
472	348	229	96	<i>camarero</i>
62	190	71	97	<i>colchón</i>
407	229	572	98	<i>nube de polvo</i>
193	444	211	99	<i>aro</i>
88	240	160	100	<i>automóvil</i>
415	452	93	101	<i>grasoso</i>
96	442	484	102	<i>lámina de queso</i>
424	128	256	103	<i>cazo</i>
127	210	131	104	<i>falencia</i>
365	406	201	105	<i>aero</i>
549	486	176	106	<i>bluyín</i>
247	411	328	107	<i>seguidor</i>
114	419	417	108	<i>tablero</i>
23	149	285	109	<i>calesita</i>
395	303	21	110	<i>entrevistar</i>
528	495	90	111	<i>gancho de ropa</i>
517	504	336	112	<i>queso americano</i>
223	187	496	113	<i>jeans</i>
213	133	288	114	<i>cachetes</i>
160	53	205	115	<i>alubia</i>
348	282	120	116	<i>extrañar</i>
181	255	185	117	<i>audífonos</i>

126	115	78	118	<i>goma de mascar</i>
425	484	13	119	<i>cristal delantero</i>
343	177	266	120	<i>calzón</i>
224	455	475	121	<i>mico</i>
9	84	228	122	<i>cartelón</i>
254	15	511	123	<i>lavarropa</i>
319	449	137	124	<i>frigorífico</i>
445	163	434	125	<i>tirantes</i>
481	162	51	126	<i>chancleta</i>
174	123	253	127	<i>celofán</i>
19	244	491	128	<i>la calor</i>
422	443	487	129	<i>lancha</i>
130	328	226	130	<i>carro</i>
374	236	28	131	<i>ensalada de fruta</i>
346	1	563	132	<i>palomitas</i>
159	260	86	133	<i>gasolinera</i>
120	239	157	134	<i>balacera</i>
396	68	521	135	<i>perchero</i>
246	277	178	136	<i>bar</i>
366	259	77	137	<i>grabadora</i>
386	19	269	138	<i>brassiere</i>
124	378	560	139	<i>ordenador</i>
157	60	409	140	<i>sponsor</i>
100	405	406	141	<i>vitrina</i>
77	308	73	142	<i>colegio</i>
117	171	79	143	<i>grabador</i>
51	176	337	144	<i>puerco</i>
52	416	323	145	<i>sortija</i>
165	29	472	146	<i>miradita</i>
431	368	368	147	<i>refrigerador</i>
326	337	32	148	<i>encendedor</i>
459	67	435	149	<i>tetera</i>
177	357	314	150	<i>sala de estar</i>
290	314	72	151	<i>colectivo</i>
340	319	58	152	<i>cochino</i>
387	464	457	153	<i>maní</i>
189	122	262	154	<i>buhardilla</i>
489	491	112	155	<i>gafotas</i>
93	436	515	156	<i>lentes</i>
464	131	357	157	<i>resfrío</i>
103	424	570	158	<i>mozo</i>
87	371	548	159	<i>pizarra</i>
91	167	567	160	<i>ojeada</i>
39	104	522	161	<i>penthouse</i>
578	503	376	162	<i>tutifruti</i>
140	234	152	163	<i>azafata o</i>
147	426	469	164	<i>mono</i>
417	34	7	165	<i>edredón</i>
146	92	260	166	<i>calzoncillo</i>
361	338	2	167	<i>echar de menos</i>
400	206	273	168	<i>bombacho</i>
375	343	237	169	<i>camioneta</i>

151	4	177	170	<i>bidón</i>
204	492	541	171	<i>polvero</i>
347	64	276	172	<i>bonita</i>
154	334	166	173	<i>autobús</i>
436	235	195	174	<i>agarrar</i>
119	70	74	175	<i>escaparate</i>
501	304	290	176	<i>cacerola</i>
29	485	427	177	<i>tartera</i>
18	82	272	178	<i>botella grande</i>
345	385	143	179	<i>fanáticos</i>
57	435	99	180	<i>hamaca</i>
169	156	558	181	<i>paila</i>
203	44	373	182	<i>tumbona</i>
286	140	259	183	<i>camastro</i>
258	219	351	184	<i>preciosa</i>
202	261	573	185	<i>obstinado</i>
580	569	261	186	<i>bomba de nafta</i>
104	96	540	187	<i>polvareda</i>
149	215	364	188	<i>remezón</i>
216	50	377	189	<i>tiza</i>
166	377	571	190	<i>nevera</i>
406	415	324	191	<i>simio</i>
60	309	257	192	<i>cazuela</i>
410	347	11	193	<i>descongelar</i>
291	14	204	194	<i>altoparlante</i>
373	287	277	195	<i>bote</i>
452	250	82	196	<i>grapadora</i>
82	27	557	197	<i>movimiento sísmico</i>
48	126	363	198	<i>rasguño</i>
487	479	154	199	<i>backpack</i>
570	582	80	200	<i>gramola</i>
456	437	495	201	<i>hostess</i>
394	222	111	202	<i>gallinita ciega</i>
13	364	334	203	<i>ropa interior</i>
205	227	243	204	<i>catarro</i>
272	25	507	205	<i>macaco</i>
334	373	559	206	<i>oportunidad</i>
141	321	53	207	<i>chicle</i>
339	386	220	208	<i>anteojos</i>
217	3	305	209	<i>scotch</i>
435	23	384	210	<i>tocadiscos</i>
242	174	442	211	<i>terco</i>
268	273	506	212	<i>luneta</i>
563	474	322	213	<i>soquetes</i>
180	213	543	214	<i>polvorín</i>
520	470	320	215	<i>sobrecama</i>
15	396	436	216	<i>terremoto</i>
404	194	289	217	<i>cabezón</i>
312	446	494	218	<i>jugó</i>
460	59	291	219	<i>cabezota</i>
105	350	1	220	<i>zumó</i>
512	175	585	221	<i>pantufla</i>

17	325	149	222	<i>chancho</i>
266	253	33	223	<i>encerado</i>
420	524	44	224	<i>cosedora</i>
465	89	348	225	<i>poste eléctrico</i>
488	146	100	226	<i>headphones</i>
341	129	223	227	<i>bomba de gasolina</i>
508	518	587	228	<i>papel albal</i>
383	74	325	229	<i>silla plegable</i>
555	526	542	230	<i>polvoreda</i>
112	63	158	231	<i>autovía</i>
67	322	132	232	<i>fallo</i>
116	290	174	233	<i>boliche</i>
148	313	59	234	<i>computador</i>
337	243	367	235	<i>recibidor</i>
107	83	148	236	<i>chancla</i>
92	182	36	237	<i>endulzante</i>
444	203	241	238	<i>casquitos</i>
101	363	354	239	<i>rancho</i>
136	36	338	240	<i>pulpería</i>
162	370	552	241	<i>plátano</i>
46	214	215	242	<i>amplificador</i>
495	510	380	243	<i>tranchete</i>
232	403	392	244	<i>vosotros</i>
447	228	443	245	<i>terral</i>
376	191	526	246	<i>pileta</i>
240	62	388	247	<i>tozudo</i>
238	8	387	248	<i>topadora</i>
551	467	317	249	<i>seboso</i>
429	262	569	250	<i>movimiento telúrico</i>
398	88	345	251	present for past subjunctive
131	286	119	252	<i>expendio</i>
128	394	299	253	<i>sandwich</i>
360	209	109	254	<i>guardafango</i>
476	274	586	255	<i>papel de estaño</i>
546	106	67	256	<i>cóctel de fruta</i>
439	409	319	257	<i>sismo</i>
208	301	279	258	<i>caballitos</i>
457	404	399	259	<i>vos</i>
427	221	246	260	<i>cerilla</i>
324	383	528	261	<i>pc</i>
513	489	190	262	<i>albal</i>
158	457	139	263	<i>fósforo</i>
325	448	485	264	<i>la margen</i>
28	152	456	265	<i>mansarda</i>
433	382	218	266	<i>anillo</i>
377	428	446	267	<i>hinchas</i>
467	476	20	268	<i>zippo</i>
106	265	428	269	<i>tasca</i>
188	292	116	270	<i>estancia</i>
329	380	538	271	<i>piscina</i>
382	100	590	272	<i>papel de aluminio</i>
579	481	108	273	<i>guardabarro</i>

330	204	271	274	<i>bóxers</i>
338	202	396	275	<i>zancudo</i>
494	299	123	276	<i>estacionar</i>
251	256	505	277	<i>luminaria</i>
133	408	415	278	<i>tal vez</i>
503	224	146	279	<i>farolillo</i>
545	580	126	280	<i>escurridero</i>
211	30	476	281	<i>megáfono</i>
118	72	532	282	<i>parabrisa</i>
184	367	369	283	<i>refrigeradora</i>
263	186	523	284	<i>percha</i>
308	179	502	285	<i>lavaplatos</i>
239	296	172	286	<i>bocadito</i>
26	145	265	287	<i>buldócer</i>
438	352	141	288	<i>fanaticada</i>
287	453	209	289	<i>armador</i>
455	237	414	290	<i>tirita</i>
442	482	151	291	<i>azulón</i>
540	525	574	292	<i>pantalones tejanos</i>
31	531	419	293	<i>tabanco</i>
75	494	8	294	<i>edulcorante</i>
207	369	553	295	<i>platicar</i>
448	155	95	296	<i>gripe</i>
560	515	371	297	<i>soutien</i>
235	57	68	298	<i>colgador</i>
219	393	383	299	<i>tragamonedas</i>
496	184	240	300	<i>canoa</i>
143	208	24	301	<i>equivocación</i>
264	584	556	302	<i>papel sanitario</i>
468	501	555	303	<i>paquete postal</i>
363	21	97	304	<i>habichuela</i>
262	5	451	305	<i>machina</i>
303	425	531	306	<i>patrocinador</i>
261	164	471	307	<i>mosco</i>
45	192	181	308	<i>barman</i>
350	468	315	309	<i>scuela</i>
66	58	294	310	<i>cacahuate</i>
37	87	429	311	<i>testarudo</i>
1	69	514	312	<i>lavavajillas</i>
47	22	64	313	<i>corpiño</i>
408	339	3	314	<i>diversión</i>
336	389	518	315	<i>mosquito</i>
480	462	106	316	<i>guagua</i>
562	497	424	317	<i>surtidor de gasolina</i>
294	519	398	318	<i>wurlitzer</i>
58	399	402	319	<i>vereda</i>
328	318	43	320	<i>cintillo</i>
362	349	12	321	<i>damasco</i>
257	336	34	322	<i>encomienda</i>
186	178	286	323	<i>calzada</i>
5	375	566	324	<i>olla</i>
237	211	150	325	<i>chancha</i>



22	310	66	326	<i>comedor</i>
471	351	516	327	<i>letrero</i>
167	430	512	328	<i>linda</i>
588	591	449	329	<i>máquina de música</i>
582	533	103	330	<i>guardilla</i>
135	218	463	331	<i>matera</i>
292	2	372	332	<i>porfiado</i>
282	101	48	333	<i>chongo</i>
497	530	501	334	<i>lavalozza</i>
446	469	410	335	<i>tajador</i>
212	283	130	336	<i>falda</i>
16	80	329	337	<i>seísmo</i>
434	77	75	338	<i>escarpa</i>
505	144	170	339	<i>blue jean</i>
416	205	133	340	<i>fosforera</i>
461	577	270	341	<i>brasiel</i>
516	79	440	342	<i>tejanos</i>
544	201	264	343	<i>brik</i>
271	438	467	344	<i>mochila</i>
111	421	197	345	<i>afición</i>
411	148	247	346	<i>cerillo</i>
576	550	341	347	<i>queso en lonchas</i>
493	527	459	348	<i>magnetofón</i>
349	107	549	349	<i>pochoclo</i>
183	487	393	350	<i>yesquero</i>
6	447	488	351	<i>maceta</i>
276	65	482	352	<i>macetero</i>
332	439	486	353	<i>la puente</i>
20	161	168	354	<i>banano</i>
21	55	365	355	<i>rayón</i>
367	28	275	356	<i>bombona</i>
10	395	433	357	<i>tiradores</i>
121	281	278	358	<i>butaca</i>
44	483	591	359	<i>papel de baño</i>
55	11	466	360	<i>microcomputador</i>
99	387	310	361	<i>rótulo</i>
536	553	546	362	<i>polvazal</i>
525	549	342	363	<i>queso en rebanadas</i>
280	311	231	364	<i>casete</i>
389	132	200	365	<i>acolchado</i>
491	168	54	366	<i>chimpancé</i>
502	520	91	367	<i>gandula</i>
554	575	40	368	<i>cinta scotch</i>
171	466	194	369	<i>altillo</i>
307	230	182	370	<i>barrero</i>
548	517	430	371	<i>tierral</i>
358	170	52	372	<i>chango</i>
318	134	519	373	<i>movi</i>
225	330	283	374	<i>caldera</i>
123	413	304	375	<i>saya</i>
256	103	576	376	<i>pantaloncillo</i>
479	247	454	377	<i>máquina de lavar</i>

249	238	255	378	<i>cayuco</i>
553	490	309	379	<i>salita de estar</i>
317	454	138	380	<i>frijol</i>
176	407	202	381	<i>aeromoza o</i>
43	139	313	382	<i>sacarina</i>
56	180	450	383	<i>marrano</i>
63	112	169	384	<i>banqueta</i>
265	366	356	385	<i>resfriado</i>
506	124	161	386	<i>autocar</i>
113	10	312	387	<i>sacapuntas</i>
403	42	401	388	<i>vitrola</i>
152	433	98	389	<i>hacienda</i>
385	169	483	390	<i>mostrador</i>
515	587	447	391	<i>máquina excavadora</i>
228	340	233	392	<i>carretera</i>
577	252	30	393	<i>engrapadora</i>
220	233	179	394	<i>bicoca</i>
97	150	470	395	<i>morral</i>
504	120	63	396	<i>coriza</i>
260	456	474	397	<i>medias</i>
302	151	390	398	<i>vaqueros</i>
198	401	397	399	<i>zapatilla</i>
352	114	588	400	<i>papel aluminio</i>
331	480	385	401	<i>tolvanera</i>
484	475	500	402	<i>lavadora de platos</i>
284	458	479	403	<i>mesero</i>
210	54	216	404	<i>aparador</i>
32	199	282	405	<i>calcetín</i>
401	105	14	406	<i>cuarto de estar</i>
583	564	422	407	<i>stewardess</i>
122	326	230	408	<i>cascos</i>
380	508	124	409	<i>espónsor</i>
498	560	426	410	<i>tapabarros</i>
355	111	374	411	<i>trusa</i>
298	429	296	412	<i>hermosa</i>
466	496	539	413	<i>pitusa</i>
279	271	50	414	<i>chispero</i>
412	267	481	415	<i>macuto</i>
50	280	153	416	<i>azotea</i>
547	477	497	417	<i>jokey</i>
490	26	292	418	<i>cabezudo</i>
24	270	318	419	<i>silla reclinable</i>
475	555	300	420	<i>salveque</i>
241	423	421	421	<i>-stes for -ste preterite</i>
327	198	81	422	<i>gallina ciega</i>
214	166	208	423	<i>argolla</i>
80	46	381	424	<i>trancazo</i>
175	102	87	425	<i>garrafón</i>
30	493	509	426	<i>loncha de queso</i>
392	391	379	427	<i>ustedes</i>
270	45	110	428	<i>guardapolvo</i>
65	85	464	429	<i>matero</i>

441	269	222	430	<i>bolso de viaje</i>
574	472	41	431	<i>cinta adhesiva</i>
281	511	213	432	<i>arañón</i>
71	353	520	433	<i>poncho</i>
173	56	128	434	<i>esparadrapo</i>
185	159	70	435	<i>colcha</i>
313	41	340	436	<i>rollo de papel</i>
278	499	62	437	<i>contén</i>
59	24	395	438	<i>yola</i>
275	434	492	439	<i>la azúcar</i>
98	331	19	440	<i>emparedado</i>
384	471	412	441	<i>tajada de queso</i>
150	356	217	442	<i>andén</i>
391	285	167	443	<i>banana</i>
315	223	15	444	<i>cubrecama</i>
33	512	498	445	<i>jorongo</i>
573	590	489	446	<i>lasca de queso</i>
463	372	214	447	<i>aparcar</i>
109	245	187	448	<i>bolsón</i>
110	288	183	449	<i>bella</i>
229	121	343	450	<i>propiciador</i>
27	90	250	451	<i>cercha</i>
378	138	238	452	<i>canchita</i>
95	13	452	453	<i>mascar</i>
36	568	583	454	<i>papel de váter</i>
299	51	94	455	<i>gripa</i>
393	158	530	456	<i>patrocicante</i>
300	160	219	457	<i>añorar</i>
342	500	56	458	<i>cocaleca</i>
295	335	163	459	<i>auspiciador</i>
145	358	547	460	<i>pollera</i>
83	384	529	461	<i>pava</i>
209	40	524	462	<i>pianola</i>
164	359	554	463	<i>pluma</i>
108	278	27	464	<i>error</i>
428	71	534	465	<i>parquear</i>
192	135	188	466	<i>alberca</i>
314	392	378	467	<i>usted</i>
283	320	57	468	<i>coche</i>
486	592	592	469	<i>abrochadora</i>
409	345	17	470	<i>curita</i>
353	507	444	471	<i>pororó</i>
379	116	355	472	<i>reportear</i>
507	589	493	473	<i>juke box</i>
142	181	254	474	<i>catre</i>
356	200	164	475	<i>auspiciante</i>
49	183	227	476	<i>carrusel</i>
38	585	581	477	<i>papel higiénico</i>
310	220	129	478	<i>espejuelos</i>
397	17	156	479	<i>badén</i>
134	136	332	480	<i>silla de playa</i>
557	509	382	481	<i>traganíquel</i>

351	193	144	482	<i>farol</i>
414	217	439	483	<i>tecle</i>
566	566	411	484	<i>tablón de anuncio</i>
243	47	510	485	<i>livin</i>
274	225	416	486	<i>tallador</i>
322	48	575	487	<i>pantalón de mezclilla</i>
575	523	136	488	<i>friegaplatos</i>
244	6	5	489	<i>echar en falta</i>
584	581	89	490	<i>gallo ciego</i>
469	572	31	491	<i>engrapadora</i>
333	498	407	492	<i>tocacintas</i>
199	37	22	493	<i>entretención</i>
132	78	350	494	<i>pota</i>
137	43	49	495	<i>chola</i>
231	266	135	496	<i>fréjol</i>
539	571	16	497	<i>cubrelecho</i>
581	570	525	498	<i>papel toalet</i>
85	110	362	499	<i>repo</i>
585	586	582	500	<i>papel para cocinar</i>
538	532	400	501	<i>vellonera</i>
144	316	248	502	<i>cervecería</i>
388	86	386	503	<i>tombo</i>
454	522	458	504	<i>magnavoz</i>
569	232	360	505	<i>rodaja de queso</i>
41	354	344	506	<i>propaganda</i>
354	212	287	507	<i>calcetas</i>
269	514	307	508	<i>salpicadera</i>
42	109	189	509	<i>ajustador</i>
381	35	536	510	<i>patera</i>
163	113	551	511	<i>placard</i>
226	374	565	512	<i>ómnibus</i>
541	567	580	513	<i>papel estañado</i>
534	506	425	514	<i>suspensores</i>
306	535	88	515	<i>gallito ciego</i>
221	226	453	516	<i>masticar</i>
470	557	331	517	<i>silla de extensión</i>
521	528	37	518	<i>crispeto</i>
64	7	480	519	<i>mesonero</i>
357	207	267	520	<i>buldózer</i>
250	39	441	521	<i>temblor de tierra</i>
421	91	461	522	<i>mahón</i>
368	95	117	523	<i>excavadora</i>
170	16	232	524	<i>carrillos</i>
564	545	578	525	<i>papel de water</i>
426	451	84	526	<i>giz</i>
301	197	121	527	<i>estada</i>
550	268	404	528	<i>vidrio delantero</i>
532	534	9	529	<i>diurex</i>
519	540	281	530	<i>cacillo</i>
196	365	193	531	<i>altavoz</i>
462	576	47	532	<i>cínife</i>
561	573	125	533	<i>escurreplatos</i>

482	559	394	534	yins
529	544	274	535	bomba de bencina
74	305	258	536	celo
571	505	561	537	pala excavadora
511	427	298	538	hinchada
542	562	413	539	tajalápiz
485	249	236	540	cambur
572	548	333	541	rositas de maíz
233	142	60	542	concho
443	246	517	543	lighter
115	18	104	544	grúa
514	513	504	545	lonja de queso
556	272	55	546	chinela
451	125	142	547	foil
78	307	244	548	centro escolar
500	119	249	549	chalana
533	558	438	550	terregal
236	172	76	551	escondidas
288	9	306	552	salpicadero
478	478	263	553	breteles
259	302	4	554	durazno
437	541	45	555	clericó
413	539	184	556	bencinera
524	583	562	557	papel platina
84	66	92	558	garrafa
35	546	358	559	rocola
458	143	311	560	ruana
230	473	39	561	cotufa
586	561	437	562	taxibús
523	537	42	563	cinta pegante
559	488	35	564	endulzador
102	315	245	565	cerdo
61	333	25	566	enagua
477	547	359	567	roconola
296	258	508	568	macedonia
68	420	418	569	taberna
568	529	513	570	lavatrastos
449	516	85	571	gasolinería
587	588	499	572	judía verde
197	147	140	573	forofos
591	556	326	574	secaplatos
527	565	420	575	sutién
592	551	339	576	queso de sandwich
589	242	589	577	papel confor
567	502	448	578	máquina de lavar platos
526	552	527	579	pipocas
222	579	239	580	canguil
531	574	38	581	crispeta
552	578	293	582	cabrita de maíz
543	554	544	583	pomo plástico
423	275	203	584	aficionados
537	97	584	585	papel de inodoro

---

255	248	537	586	<i>piquera</i>
590	563	327	587	<i>silla de sol</i>
509	543	297	588	<i>hielera</i>
499	538	134	589	<i>freidero</i>
321	130	473	590	<i>microbús</i>
474	536	445	591	<i>poporopo</i>
558	542	196	592	<i>afilaminas</i>