

Stylometry and Spanish Golden Age Theatre

An Evaluation of Authorship Attribution in a Control Group of One Hundred Undisputed Plays

Álvaro Cuéllar 

Abstract The aim of this study is to perform an evaluation of one hundred Spanish Golden Age theatre plays of undisputed authorship using the R package *stylo*, the stylometric analysis tool developed by Eder, Rybicki, and Kestemont (2016). In this paper we will determine which algorithms obtain best results on authorial classification (method, MFW, culling, and word n-grams). We will also evaluate the text length at which stylometry begins to be an effective diagnostic tool for authorship attribution in our corpus. This cross-validation evaluation can serve future analysis of similar corpora and will show the possibilities of applying stylometry to Spanish Golden Age theatre, which presents many cases of dubious authorship.

Keywords stylometry, Spanish Golden Age theatre, *stylo*, author identification, text length, most frequent words (MFW), culling, n-grams

1. Spanish Golden Age Theatre Authorship Issues

Menéndez Pidal (1949, xxviii–xxxii) explained that the frequency of authorship issues in Spanish literature can be attributed to three factors: its tendency to anonymity, its collectivism, and its collaborative writing. These three elements, though present throughout the history of Spanish literature, converge particularly in the Spanish Early Modern period, with canonical prosaic texts such as *La vida de Lazarillo de Tormes* or the apocryphal *Quijote* by Avellaneda serving as famous examples of this authorial

uncertainty.¹ Nevertheless, it is in the vast theatrical production of the Spanish Golden Age—over 10,000 texts are estimated to have been written for the stage during this period, although it is difficult to establish an exact number—where we find the highest quantity of authorship disputes that need to be resolved.

Issues of authorship in the field of Spanish Golden Age theatre are immeasurable and the academic literature surrounding them is substantial. For example, Tirso de Molina's famous prologue to the *Segunda parte de sus comedias* (1635) includes the following statement:

Destas doze comedias quatro, que son mias, en mi nombre, y en el de los dueños de las otras, ocho (que no se porque infortunio suyo, siendo hijas de tan ilustres padres, las echaron a mis puertas) las que restan.

“Of these twelve plays, only four are mine, in my name, and the remaining eight in the name of their owners (for I do not know why, being the daughters of such illustrious fathers, their misfortune put them at my doorstep).”²

This volume includes the play *El condenado por desconfiado*, crucial for its theological ideas, whose authorship has been in doubt precisely because of this statement by the playwright.

Calderón de la Barca is the greatest exponent of authorship confusion in Golden Age theatre, to the extreme that he is known to have even rejected comedies that he had actually written.³ We know that he wrote plays that he did not include in several lists of his authentic texts that he prepared toward the end of his life. Moreover, some of his plays appear in a list of apocryphal works that circulated in his name that Calderón included at the beginning of his *Parte cuarta* (1672), and which is inexact in several cases. This situation allows us to appreciate the complexities presented by this extreme case: we cannot even take for granted the information provided by authors regarding their own production. To this point, Germán Vega García-Luengos sums up the extent to which authorial uncertainty surrounds the attributed production of the two major playwrights of Spain's Golden Age theatre:

1 Some articles have recently been published addressing these issues with digital methods, especially using stylometry. See De la Rosa and Suárez (2016) for *La vida de Lazarillo de Tormes*, and Fradejas Rueda (2016), Rišler-Pipka (2016), and Rodríguez López-Vázquez (2018) for the apocryphal *Quijote*.

2 All translations are my own.

3 For an explanation about Calderón's textual issues and the problematics of his lists, see Coenen (2009a; 2009b; and 2019) or Vega García-Luengos (2008)

De las 506 comedias atribuidas en algún momento a Lope de Vega que estudian S. G. Morley y C. Bruerton en su *Cronología* (Morley y Bruerton 1968), 316 serían auténticas, 27 probablemente, 73 se califican como dudosas y 90 como ajenas. Calderón, el otro grande del teatro español, aún superaría estos porcentajes, tal como se apunta ya en los testimonios más tempranos: las 108 piezas de sus nueve partes se ven superadas por las 115 de las listas de ‘comedias supuestas’ que Vera Tassis dio en la *Verdadera quinta parte* (1682) y en la *Séptima* (1683).

“Of the 506 plays attributed at some point in time to Lope de Vega that S. G. Morley and C. Bruerton study in their *Cronología*, 316 would be authentic, 27 probably so, 73 dubious, and 90 spurious. Calderón, the other great writer of Spanish theatre, even exceeds these percentages, as already indicated in the earliest testimonies: the 108 texts of his nine volumes are exceeded by the 115 plays in the lists of ‘supposed comedies’ that Vera Tassis included in the *Verdadera quinta parte* (1682) and the *Séptima* (1683).” (2002, 16)

We may think that authorship problems affect only secondary works of Golden Age theatrical production, those that have exerted less sway of interest and influence in the literary tradition, effectively passing unnoticed until a scholar decides to recover them. However, some of the most celebrated plays fall into this category of dubious authorship, such as *Tan largo me lo fiáis* and *El burlador de Sevilla*, the two plays that started the Don Juan myth and whose authorship has been profusely debated.⁴

Among the possible explanations for the authorship problems that Spanish Golden Age theatre presents, the economic factor seems to be key. As a form of property, the text legally belonged to the theatrical director or the bookseller, not to the author of the play. The owners could change the real author’s name to attract audiences, or add or delete verses from the text. Although it may seem surprising today, Vega García-Luengos asserts that writers, “in general, did not make any great effort to resolve this state of affairs” (2009, 97; “en líneas generales, no se aprecia que dedicaran demasiados esfuerzos a solucionar tal estado de cosas.”) for two main reasons. First, it is worth noting that poets wrote their plays for the stage, and that publication was considered secondary. Second, the economic and ownership status of the theatrical text meant that textual stability was not legally protected or assured: playwrights did not keep the publishing rights when selling their original manuscripts with the final version of their plays to a theatrical company and, when they did recover those rights, there was no guarantee that the text would match the original.

⁴ Rodríguez López-Vázquez has been studying this issue for decades in numerous articles and books. For instance, see Rodríguez López-Vázquez (1983).

2. Stylometry

Stylometry is a technique that tries to automatically compare texts by their *style* and, among its many applications, it can be useful for authorship attributions of literary works. Within the category of *style*, which can be quite abstract, we are going to be working specifically with the concrete criterion of word usage frequency.⁵ Stylometry is based on this simple hypothesis: each writer uses some words more frequently than others, and these data can be used to establish relationships between different texts. For instance, in *La dama boba* by Lope de Vega, for each 100 words, 4.90 are *que*, 3.64 are *de*, and 2.90 are *y*, while in *Don Gil de las calzas verdes* by Tirso de Molina, these proportions are 4.86, 3.91 and 3.03. These tiny differences are not usually relevant to researchers; nevertheless, they are decisive for stylometry. Stylometry uses these proportions to establish relations of proximity between the different texts or classify them in a specific group through statistical models. Usually, stylometric analysis is performed on the basis of hundreds of these words and their proportions. Sometimes other parameters are more relevant, such as analyzing groups of words (n-grams) or accepting only words found in a certain percentage of texts, which is known as culling.

In the last years, stylometry has been extensively applied to authorship issues of Spanish Golden Age theatre. Some good examples are the attribution of *Siempre ayuda la verdad* to Lope de Vega (García-Reidy 2019); the study of *La conquista de Jerusalén* (Cerezo Soler and Calvo Tello 2019), where stylometric analysis seems to corroborate the consideration of Cervantes as the author of the text; the study of *La segunda Celestina*, a text attributed to Sor Juana Inés de la Cruz by Schmidhuber de la Mora some decades ago (Hernández Lorenzo and Byszuk 2019); the study of the *entremés* titled *Los mirones*, a short type of play that was performed in between acts of a play (Blasco 2019); and the stylometric study of Agustín Moreto y Cavana's production, which is one of the most complex repertoires of the period because of the author's tendency toward collaboration and rewriting (Ulla Lorenzo, Martínez Carro, and Calvo Tello 2020).

Professor Vega García-Luengos and I are currently developing the project *ETSO: Estilometría aplicada al Teatro del Siglo de Oro* (Stylometry applied to the Spanish Golden Age Theatre), <https://etso.es/> (Cuéllar and Vega 2017–2022), a collaborative project that currently includes dozens of scholars and seeks to build the biggest possible digital repository of Spanish Golden Age plays for the application of stylometric analysis. Our working philosophy is not to analyze a specific authorship issue within a controlled corpus. Instead, we use all the texts we have and let the corpus show unexpected connections. This is exactly what happened with *La monja alférez*, a play traditionally attributed to Juan Pérez de Montalbán, Lope de Vegas's disciple. We were not specifically

5 For an introduction to stylometry and the stylo package, see Fradejas Rueda (2019).

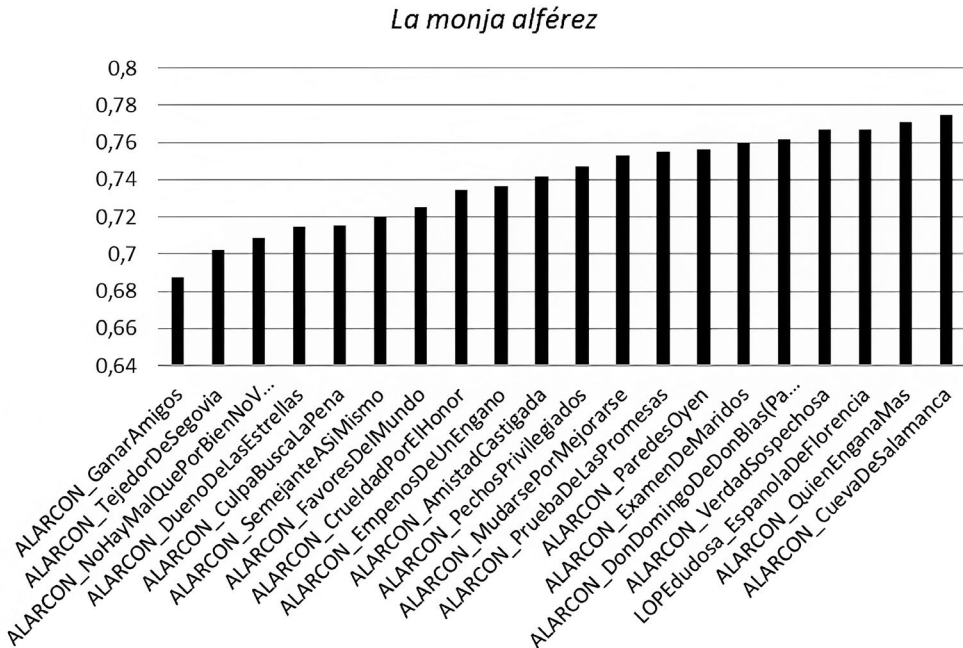


Fig. 1 20 most proximate plays to *La monja alférez* in a corpus of 1,028 texts (ETSO, September 1, 2020), Stylo (500 MFWs, Classic Delta, zero percent culled). (Vega García-Luengos 2021)

researching authorship issues surrounding this play, but when we ran an analysis of all the new texts we had collected, the following (Figure 1) were the plays that are most lexically proximate in a corpus of more than 1,000 texts.

As the results indicate, the closest plays were mostly written by Juan Ruiz de Alarcón, a quite prolific author born in New Spain. Thanks to this surprising result, Vega García-Luengos decided to isolate the historical and philological analysis of *La monja alférez* in a study that presents solid and variegated evidence confirming the attribution of the play to Juan Ruiz de Alarcón (2021). Vega García-Luengos and I are applying this same analytical procedure to other interesting cases that we hope to publish soon as a book-length anthology. We also use ETSo to assist scholars with specific authorship issues, such the collaborative writing in Moreto's *La adúltera penitente* (Moreto 2019) or *Empezar a ser amigos* (Demattè 2019), or the attribution dispute in *Las dos bandoleras* (Madrónal 2019).⁶

6 For a complete description of the collaboration of ETSo with different scholars, visit <http://etso.es/repercusion/>.

One of the most effective applications to perform stylometric analysis is *stylo*.⁷ This tool, which has to be installed as a package in R (R Core Team 2020), a programming language and free software environment for statistical computing, has been developed by the Computational Stylistics Group, a team formed by members of the universities of Krakow and Antwerp, headed by researchers Eder, Rybicki, and Kestemont (2016).

3. Corpus Selection and Homogenization

Before applying stylometry to texts that present attribution problems, it is necessary to ensure that our methods work with a corpus of controlled works. The aim of this article is to precisely verify that stylometry works with a control group of undisputed plays and determine which algorithms obtain better results.

Although it is impossible to ensure with total certainty that the plays chosen here were really written by the supposed authors—remember that they sometimes lie to us about their production—we have selected only works included in *Partes*, that is, volumes of usually twelve plays, whose publication was controlled by the authors or their relatives.⁸

To check the effectiveness of stylometry, we turn to a corpus of nine canonical authors of Golden Age theatre: Cervantes (1547–1616), Lope de Vega (1562–1635), Guillén de Castro (1569–1631), Tirso de Molina (1579–1648), Ruiz de Alarcón (1581?–1639), Calderón de la Barca (1600–1681), Pérez de Montalbán (1602–1638), Rojas Zorrilla (1607–1648), and Agustín Moreto (1618–1669). As we mentioned, we will only consider plays from *Partes de comedias* published by the authors or their relatives, which amounts to one hundred texts.

To lexically compare the texts, Vega García-Luengos and I had to homogenize the corpus. The original seventeenth-century orthography was quite unstable, and our texts come from different sources and editions, so we decided to modernize and homogenize the orthography to the current Spanish rules.⁹ In addition, we removed stage directions since we consider that they may negatively affect the authorship analysis, given our understanding that they are elements that were usually not written by playwrights but by theatrical producers and therefore are not related to the unconscious style of the authors' lines.¹⁰

7 Here we are using version 0.7.4.

8 For a complete description, see Vega García-Luengos (2010).

9 See Real Academia Española y Asociación de Academias de la Lengua Española (2010).

10 See Rodríguez-Gallego (2018) for a study about how Vera Tassis modified the stage directions in seven plays by Calderón de la Barca when he published them.

4. Corpus Evaluation

In this section, we will check to ensure that the corpus of undisputed works responds satisfactorily to a test with different settings. Only if the corpus shows reliable results when authorship is not in doubt can stylometry be considered a useful tool for works of dubious authorship.

We have a corpus of one hundred plays. The test we will apply is known as leave-one-out cross-validation, which consists of running an analysis of all the plays within the corpus as though each were of unknown authorship, forcing the algorithm in each instance to choose an author of the text in question, and making sure that the resulting automatic classification of each play in fact corresponds to the correct one of the nine authors included in the corpus.¹¹ In other words, we separate a text from the rest of the corpus, the work becomes our authorship *problem*, and it must be classified. The machine is asked to classify this work among the nine possible groups according to the frequency of its words and the parameters established. Since we have carefully selected our corpus to consist only of plays with undisputed authorship, this test will help us check if the tool is accurate and, if so, to what degree. The three parameters are:

1. Classification method. There are many methods to perform the analysis, such as SVM (Support Vector Machine), NSC (Nearest Shrunken Centroids), Naive Bayes, k-NN (k-nearest neighbors), Delta¹² with different distance measures (Classic, Eder, etc.), and other distances (Canberra, Cosine, etc.).¹³
2. Most frequent words (MFW). The application will take into account the most frequent words considering all the texts for the analysis. This means that function words, such as prepositions, conjunctions or auxiliary verbal forms are usually the most frequent ones, and words with meaning are secondary to them. Thus, an MFW of 100 means that the machine uses the 100 most frequent words in the texts as a whole. An MFW of 5,000 means that the machine uses the 5,000 most frequent words.
3. Culling. Once the machine has listed all the words by their frequency, we may want to dismiss those that occur in a tiny group of texts. We probably do not want to find stylometric closeness between two plays which share characters names, places, or plots. A culling of 30 means that the machine uses words that are at least in 30 percent of the texts. A culling of 100 implies that the machine uses words that appear in all the texts, which is usually a very small number.

11 For a step-by-step tutorial about how to program a cross-validation in *stylo*, see Cuéllar (2018).

12 The usage of distance measurements is not a classification method comparable to the others because it usually creates a matrix of distances between texts. Nevertheless, *stylo* uses this matrix to classify the texts based on the closest neighbors, so it can be eventually considered as a classification method.

13 For a deeper explanation, see Calvo Tello (2016) or Schöberlein (2017).

4.1 Results by Varying the Classification Method, MFW, and Culling

We can see in Figures 2, 3, 4 and 5 the classification results when varying the methods and the MFW. A result of 100 percent indicates a perfect classification: using leave-one-out cross-validation, all one hundred plays have been correctly attributed to their respective authors among the nine options given. A result of 94 percent means that the machine has correctly classified this percentage of texts.

By observing the figures, we can see that the global results are remarkably positive. Most of the methods offer classification results over 95 percent. Nevertheless, there are some methods that perform more poorly with our specific test set, such as Cosine, Euclidean and Naive Bayes. The second observation we can make is that the MFW variation between 100 MFW and 5,000 MFW does not seem to affect the classification results in a significant way.

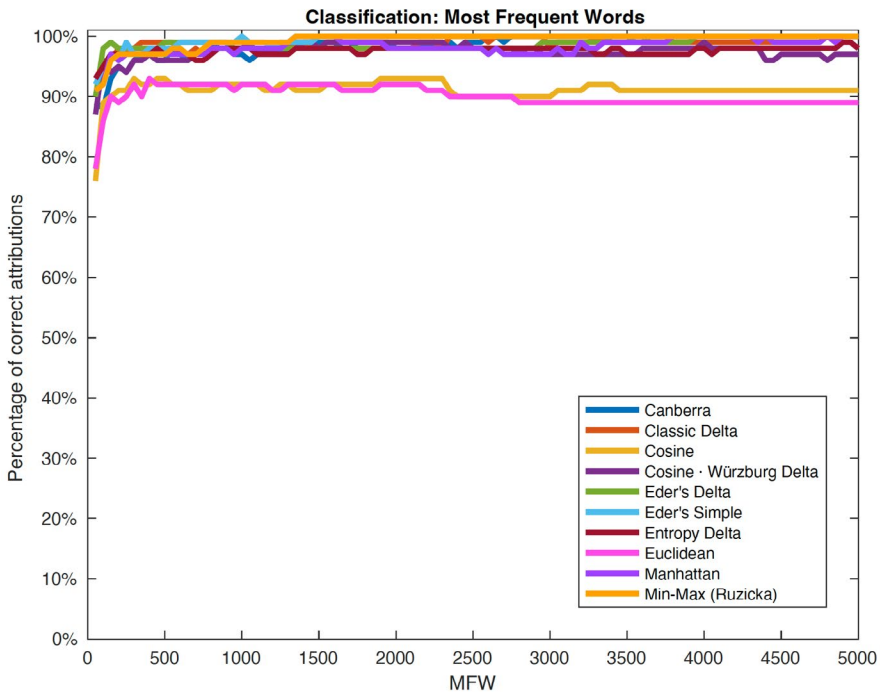


Fig. 2 Percentage of correct attributions with ten different distance methods by varying the MFW (Cuéllar, CC BY).

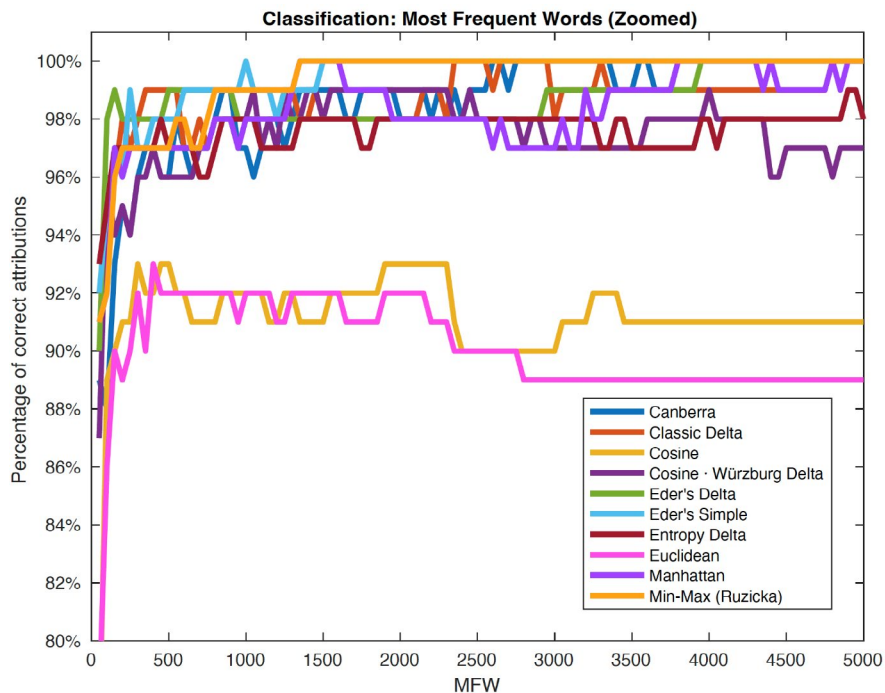


Fig. 3 Zoomed percentage of correct attributions with ten different distance methods by varying the MFW (Cuéllar, CC BY).

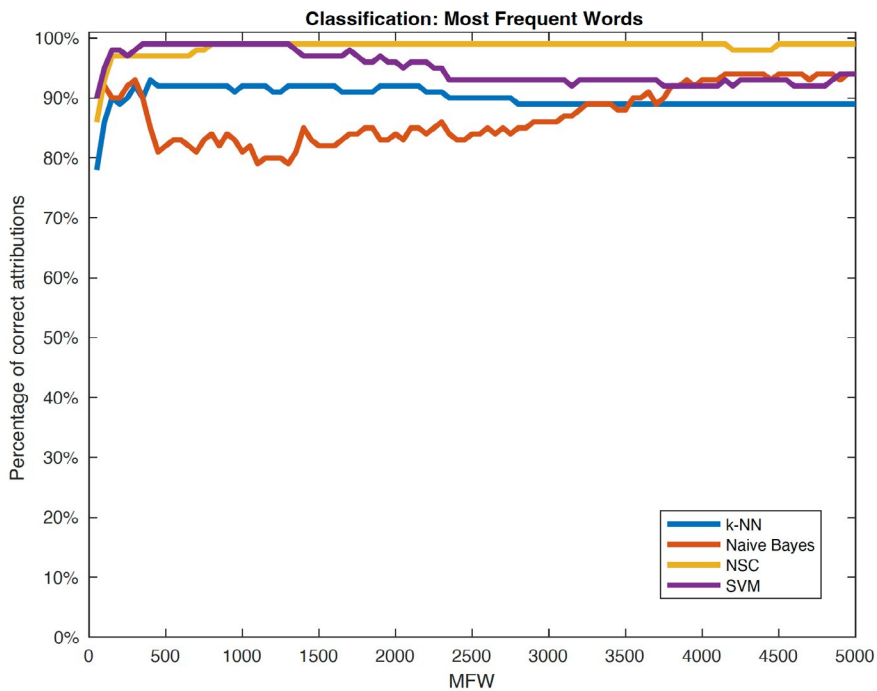


Fig. 4 Percentage of correct attributions with four different classification methods by varying the MFW (Cuéllar, CC BY).

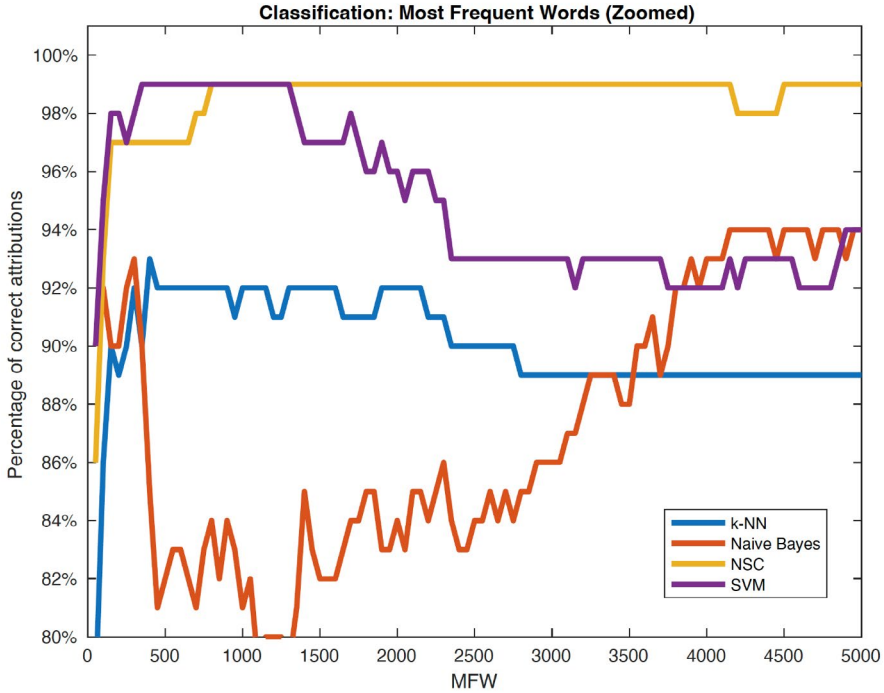


Fig. 5 Zoomed percentage of correct attributions with four different classification methods by varying the MFW (Cuéllar, CC BY).

We can now run an analysis choosing one method, such as Classic Delta, and varying the culling. We start with 5,000 MFW and a culling of 0 percent. We are going to increase the culling by steps of 0.5 percent until we reach a 100 percent of culling. This culling increment is going to affect the number of MFW available because there are not enough common words, so this number is going to reduce consequently. For example, with a 20 percent culling we perform the analysis with 3,597 words and with a 80 percent culling we perform the analysis with 603 words. It is complex, therefore, to understand if the results we obtain are due the culling or to the reduction of parameters. Nonetheless, we can see the results in Figure 6, they remain almost constant when varying the culling.¹⁴

The results do not seem to change significantly when we vary the MFW or the culling. It does not seem to matter if we use the 100, 500 or 5,000 most frequent

¹⁴ Results for the other methods are quite similar: culling does not seem to significantly affect the results of the automatic classification.

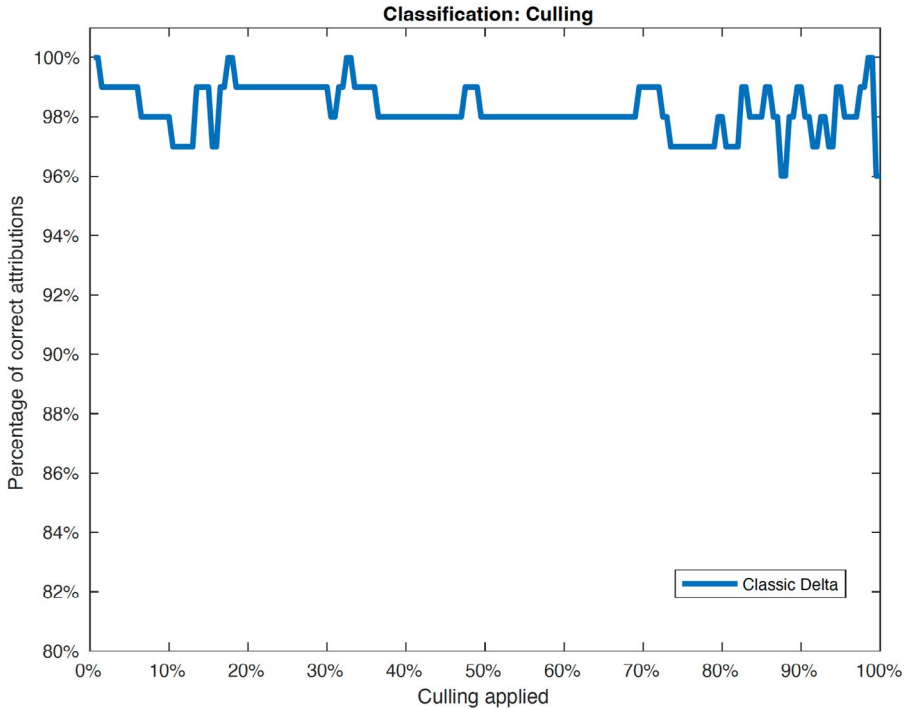


Fig. 6 Percentage of correct attributions by varying the culling. Classic Delta (Cuéllar, CC BY).

words, or if we use words that have to appear in a specific percentage of texts: the results are still extraordinary.

We can get a glimpse of the potential that these results suppose: if we now introduce a play with dubious authorship and if it is classified with one of these nine authors, we can be on the trail of an attribution. Of course, it is possible that it does not belong to any of the authors or that the system has failed, so we must be cautious and take into account additional evidence.

4.2 Results by Varying the Word N-grams

For this test, we will compare the results of using just one word (e.g. *que*), and groups of two (e.g. *lo que*), three (e.g. *ha de ser*), four (e.g. *qué es lo que*), and five words (e.g. *qué es lo que dices*). In Figure 7, we are using Classic Delta and zero percent culling. Again, the results present the percentage of success: 100 percent means a perfect classification;

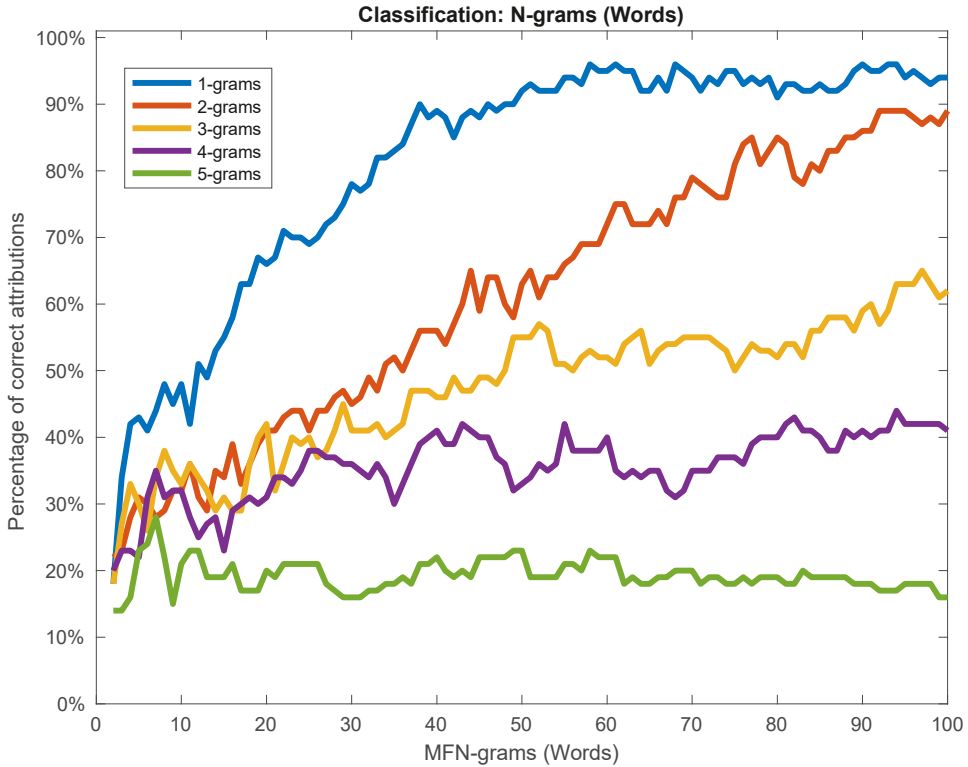


Fig. 7 Percentage of correct attributions when using different N-grams (words). Classic Delta (Cuéllar, CC BY).

77 percent means that this proportion of texts has been correctly classified among the nine groups when performing the leave-one-out cross-validation classification.

The results in Figure 7 are straightforward: they decline when the quantity of words in the n-grams increases. It seems that for authorship verification in corpora similar to ours, using words instead or bigrams, trigrams, etc., may be the best option.

4.3 Results by Varying the Text Length

In the field of Spanish Golden Age theatre, we are not limited to working with entire plays. Collaborative writing was frequent; thus, it was not uncommon for multiple writers to collaborate on the same play, each of them writing one act and then, with minimal corrections, combining them. Then there are short theatrical texts such as

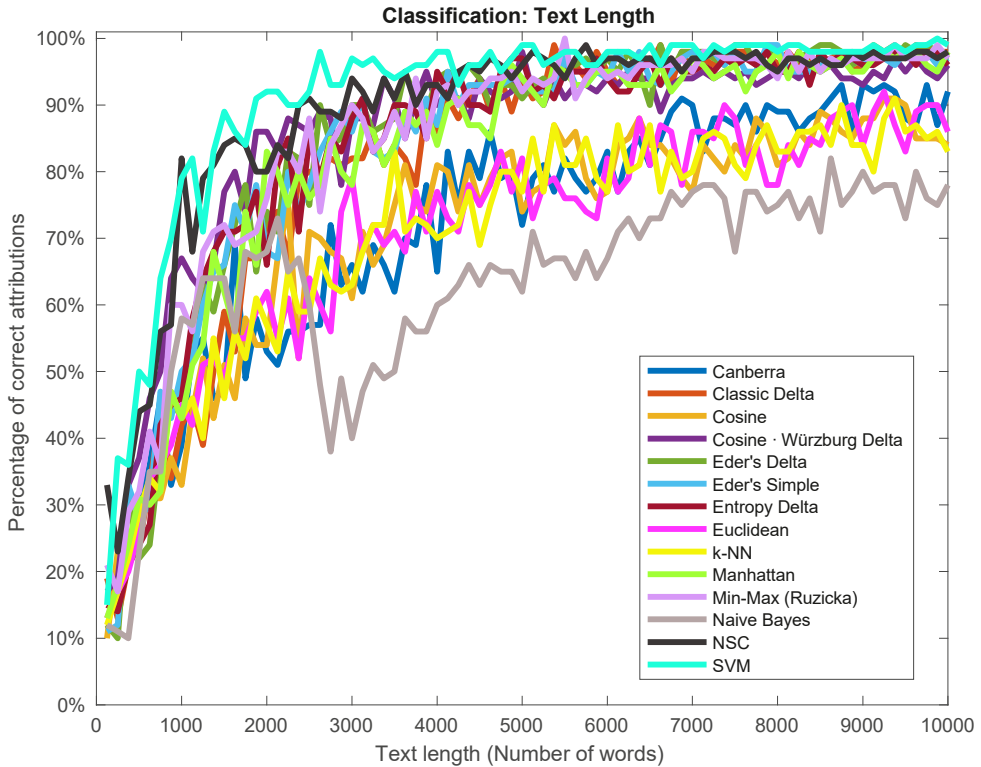


Fig. 8 Percentage of correct attributions with 14 different methods by varying the text length (Cuéllar, CC BY).

entremeses, loas, bailes, jácaras, mojigangas, etc. that also present authorship issues. Finally, some scholars are interested in exploring only passages of plays, such as a specific sonnet whose authorship is in doubt. The texts' length is therefore an element that must be taken into consideration.

In Figure 8, we are going to try to answer the length of text—how many words—where stylometry starts to work with acceptable success on my specific corpus.¹⁵ We will compare the 14 methods, randomly sampling our corpus from 50 to 10,000 words, and we will apply 500 MFW and 0 percent culling.

The results show that trying to use stylometry on texts with less than 1,000 words is quite problematic with our methods. Around the 2,500-word mark, most of methods, such as Classic Delta, Cosine, Würzburg Delta, Eder's Delta, Eder's Simple,

15 For a study on text length and authorial classification, see Eder (2010).

Entropy Delta, Manhattan, Min-Max (Ruzicka), NSC and very specially SVM, start working with quite a high success rate, above 80 percent. The other methods need a higher overall word count to improve their results. From 7,500 words onward, most of the methods start working as they did with the entire plays (which usually have 15,000–20,000 words). It seems clear that the larger the texts are, the better the results that we obtain.¹⁶

5. Conclusion

In this paper we have tested whether stylometry is an effective tool for authorship identification in Spanish Golden Age theatre with a corpus of one hundred plays of undisputed authorship. We have done so using *stylo* with 14 different methods, and parameters such MFW, culling and n-grams. The global results have been very satisfactory, reaching more than 95 percent of correct attributions. Surprisingly, results do not vary significantly when changing the method, the number of MFW or the percentage of culling, but they tend to be worse when using a number of n-grams larger than single words. We also have tested from which text length stylometry starts working effectively, obtaining results of above 80 percent of correct attributions for texts from 2,500 words with most methods. Thus, applying stylometry to shorter text can be problematic and we have to be extremely careful in those cases.

In conclusion, stylometry appears to be a promising technique for shedding light on the abundant authorship challenges that the Spanish Golden Age theatre presents. Its hypothesis is quite simple: each author uses unconsciously the words in different frequencies and, if we are able to measure these differences, we can establish relations of proximity between the texts, or classify them satisfactorily. The hypothesis seems to be supported by studies like this one, but we need to continue exploring this technique and always use it in combination with other approaches when dealing with a specific attribution.

16 These results are close to those offered by Hernández Lorenzo (2019) on Early Modern Spanish poetry. She concluded that with texts of 2,000 words or more, stylometry starts working satisfactorily.

Acknowledgements

This article is part of the project *Sound and Meaning in Spanish Golden Age Literature* (FWF Austrian Science Fund 32563). Special thanks to Germán Vega García-Luengos (Universidad de Valladolid) for his help in the elaboration of this paper. The corpus of frequencies that allows the replication of these experiments is placed on <https://github.com/AlvaroCuellar> or by request to the author.

ORCID®

Álvaro Cuéllar  <https://orcid.org/0000-0002-9934-6321>

References

- Blasco, Javier. 2019. "Atribuciones cervantinas desde la estilometría: el entremés de *Los mirones*." In *Cartografía literaria en homenaje al profesor José Romera Castillo*, edited by Guillermo Laín Corona and Rocío Santiago Nogales, 151–68. Madrid: Visor Libros. <http://uvadoc.uva.es/handle/10324/37760>.
- Calvo Tello, José. 2016. "Entendiendo Delta desde las Humanidades." *Caracteres: Estudios culturales y críticos de la esfera digital* 5 (1): 140–76.
- Cerezo Soler, Juan, and José Calvo Tello. 2019. "Autoría y estilo. Una atribución cervantina desde las Humanidades Digitales. El caso de *La conquista de Jerusalén*." *Anales Cervantinos* 51: 231–50. <https://doi.org/10.3989/anacervantinos.2019.011>.
- Coenen, Erik. 2009a. "En los entresijos de una lista de comedias de Calderón." *Revista de filología española* 89 (1): 29–56. <https://doi.org/10.3989/rfe.2009.v89.i1.63>.
- Coenen, Erik. 2009b. "Las atribuciones de Vera Tassis." *Castilla. Estudios De Literatura*: 111–33. <https://doi.org/10.24197/cel.o.2009.111-133>.
- Coenen, Erik. 2019. "Everett Hesse, Vera Tassis y el texto de las comedias de Calderón." *Bulletin of the Comediantes* 71 (1/2): 87–102. <https://doi.org/10.1353/boc.2019.0006>.
- Cuéllar, Álvaro. 2018. "La necesidad de la validación cruzada en Stylo y cómo programarla en R." *Caracteres. Estudios culturales y críticos de la esfera digital* 7 (2): 301–20. <https://dialnet.unirioja.es/servlet/articulo?codigo=7104985>.
- Cuéllar, Álvaro, and Germán Vega García-Luengos. 2017–2022. *ETSO: Estilometría aplicada al teatro del Siglo de Oro*. <https://etso.es/>.
- De la Rosa, Javier, and Juan-Luis Suárez. 2016. "The Life of Lazarillo de Tormes and of His Machine Learning Adversities." *Lemir: Revista de Literatura Española Medieval y del Renacimiento* 20: 373–438. <http://arxiv.org/abs/1611.05360>.

- Demattè, Claudia. 2019. "Una nueva comedia en colaboración entre ¿Calderón?, Rojas Zorrilla y Montalbán: *Empezar a ser amigos* a la luz del análisis estilométrico." *Rilce* 35 (3): 852–74. <https://doi.org/10.15581/008.35.3.852-74>.
- Eder, Maciej. 2010. "Does Size Matter? Authorship Attribution, Small Samples, Big Problem." *Digital Scholarship in the Humanities* 30 (2): 167–82. <https://doi.org/10.1093/llc/fqt066>.
- Eder, Maciej, Jan Rybicki, and Mike Kestemont. 2016. "Stylometry with R: A Package for Computational Text Analysis." *R Journal* 8 (1): 107–21. <https://journal.r-project.org/archive/2016/RJ-2016-007/index.html>.
- Fradejas Rueda, José Manuel. 2016. "El análisis estilométrico aplicado a la literatura española: las novelas policíacas e históricas." *Caracteres. Estudios culturales y críticos de la esfera digital* 5 (2): 196–45.
- Fradejas Rueda, José Manuel. 2019. *Cuentapalabras. Estilometría y Análisis de Textos Con R Para Filólogos*. <http://www.aic.uva.es/cuentapalabras>.
- García-Reidy, Alejandro. 2019. "Deconstructing the Authorship of *Siempre ayuda la verdad*: A Play by Lope de Vega?" *Neophilologus* 103 (4): 493–510. <https://doi.org/10.1007/s11061-019-09607-8>.
- Hernández Lorenzo, Laura. 2019. "Poesía áurea, estilometría y fiabilidad: Métodos supervisados de atribución de autoría atendiendo al tamaño de las muestras." *Caracteres. Estudios culturales y críticos de la esfera digital* 8 (1): 189–28. <https://dialnet.unirioja.es/servlet/articulo?codigo=7105381>.
- Hernández-Lorenzo, Laura, and Joanna Byszuk. 2019. "Challenging Stylometry: The Authorship of the Baroque Play *La Segunda Celestina*." *Digital Humanities Conference*. <https://github.com/JoannaBy/La-Segunda-Celestina>.
- Madroñal, Abraham. 2019. "Entre la historia y la leyenda. A propósito de *Las dos bandoleras*, comedia atribuida a Lope de Vega." *Anuario Lope de Vega* 25: 281–310. <https://doi.org/10.5565/rev/anuariolopedevga.298>.
- Menéndez Pidal, Ramón. 1949. "Caracteres primordiales de la literatura española. Con referencias a las otras literaturas hispánicas, latina, portuguesa y catalana." In *Historia general de las literaturas hispánicas*, edited by Guillermo Díaz-Plaja, xiv–lix. Barcelona: Editorial Barna.
- Moreto, Agustín. 2019. *La adúltera penitente*, edited by Fernando Rodríguez-Gallego. Alicante: Biblioteca Virtual Miguel de Cervantes. <http://www.cervantesvirtual.com/nd/ark:/59851/bmc0942634>.
- Morley, Sylvanus Griswold, and Courtney Bruerton. 1968. *Cronología de las comedias de Lope de Vega*. Barcelona: Gredos.
- R Core Team. 2020. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Real Academia Española y Asociación de Academias de la Lengua Española. 2010. *Ortografía de la lengua española*. Madrid: Espasa.
- Rißler-Pipka, Nanette. 2016. "Avellaneda y los problemas de la identificación del autor. Propuestas para una investigación con nuevas herramientas digitales." In *El otro Don Quijote. La continuación de Fernández de Avellaneda y sus efectos*, edited by Hanno Ehrlicher, 27–51. Augsburg: Universität Augsburg, Institut für Spanien-, Portugal- und Lateinamerika-Studien (ISLA).

- Rodríguez López-Vázquez, Alfredo. 1983. "La autoría de 'El Burlador de Sevilla': Andrés de Claramonte." *Castilla: Estudios de Literatura* 32 (5): 87–108. <http://uvadoc.uva.es/handle/10324/16144>.
- Rodríguez López-Vázquez, Alfredo. 2018. "Cervantes y helecho de procusto: notas críticas al uso de la Estilometría en obras de atribución dudosa y en obras apócrifas." *EHumanista* 41: 193–201. <https://dialnet.unirioja.es/servlet/articulo?codigo=6946620>.
- Schöberlein, Stefan. 2016. "Poe or not Poe? A Stylometric Analysis of Edgar Allan Poe's Disputed Writings." *Digital Scholarship in the Humanities* 32 (3): 643–59. <https://doi.org/10.1093/lc/fqw019>.
- Ulla Lorenzo, Alejandra, Elena Martínez Carro, and José Calvo Tello. 2021. "Las comedias de dudosa atribución de Agustín Moreto: nuevas perspectivas estilométricas." *Neophilologus* 105: 57–73. <https://doi.org/10.1007/s11061-020-09649-3>.
- Vega García-Luengos, Germán. 2002. "Atribución." In *Diccionario de la comedia del Siglo de Oro*, edited by Frank P. Casa, Luciano García Lorenzo, and Germán Vega García-Luengos, 16–19. Madrid: Editorial Castalia.
- Vega García-Luengos, Germán. 2008. "Consideraciones sobre la configuración del legado de comedias de Calderón." *Criticón* 103–104: 249–71. <https://doi.org/10.4000/criticon.12179>.
- Vega García-Luengos, Germán. 2009. "Los problemas de autoría en el teatro español de los Siglos de Oro." In *Unidad y multiplicidad: tramas del hispanismo actual. VIII Congreso Argentino de Hispanistas*, edited by Mariana Genoud de Fourcade and Gladys Granata de Egües, 95–135. Mendoza: Zeta Editores.
- Vega García-Luengos, Germán. 2010. "Sobre la identidad de las partes de comedias." *Criticón* 108: 57–78. <https://doi.org/10.4000/criticon.14266>.
- Vega García-Luengos, Germán. 2021. "Juan Ruiz de Alarcón recupera 'La monja alférez'." In *Sor Juana Inés de la Cruz y el teatro novohispano: XLII Jornadas de teatro clásico*, edited by Rafal González Cañal and Almudena García González, 89–159. Castilla-La Mancha: Ediciones de la Universidad de Castilla-La Mancha. <https://ruidera.uclm.es/xmlui/handle/10578/28570>.

