

Estilometridieval o cómo hacer números con textos medievales

José Manuel Fradejas Rueda
Universidad de Valladolid

IV Jornadas del español
Todo lo que un filólogo puede hacer con un ordenador
9 de noviembre de 2017

Twitter, Facebook, Wikipedia ... y NLP

J. M. Fradejas Rueda

III Jornadas de Lengua Española
Las vidas de las palabras

Grado en Español: Lengua y Literatura
Universidad de Valladolid
28-29 de marzo de 2017



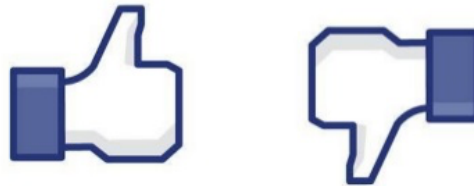
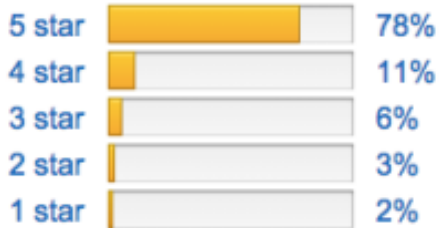
Likes & Dislikes



Customer Reviews

★★★★☆ 343

4.6 out of 5 stars



Likes and Dislikes



Análisis de sentimientos

¿Podemos medir / valorar el sentimiento que transmite una obra literaria?



Análisis de sentimientos

No me refiero a



DID NOT FINISH



A LONG SLOG



GOOD FOR WHAT IT WAS



AWESOME! READ IT!

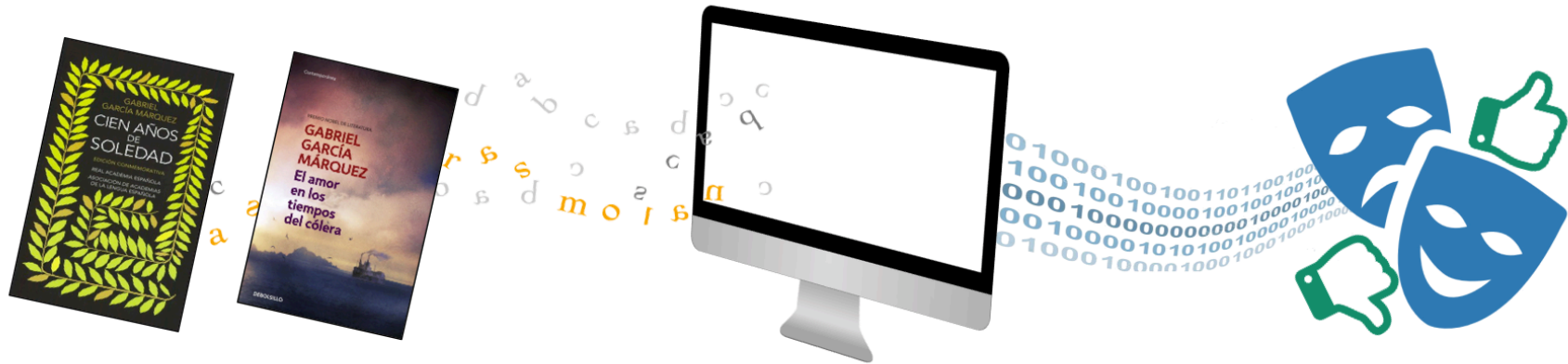


GOING ON MY FAVORITES SHELF!



Universidad de Valladolid

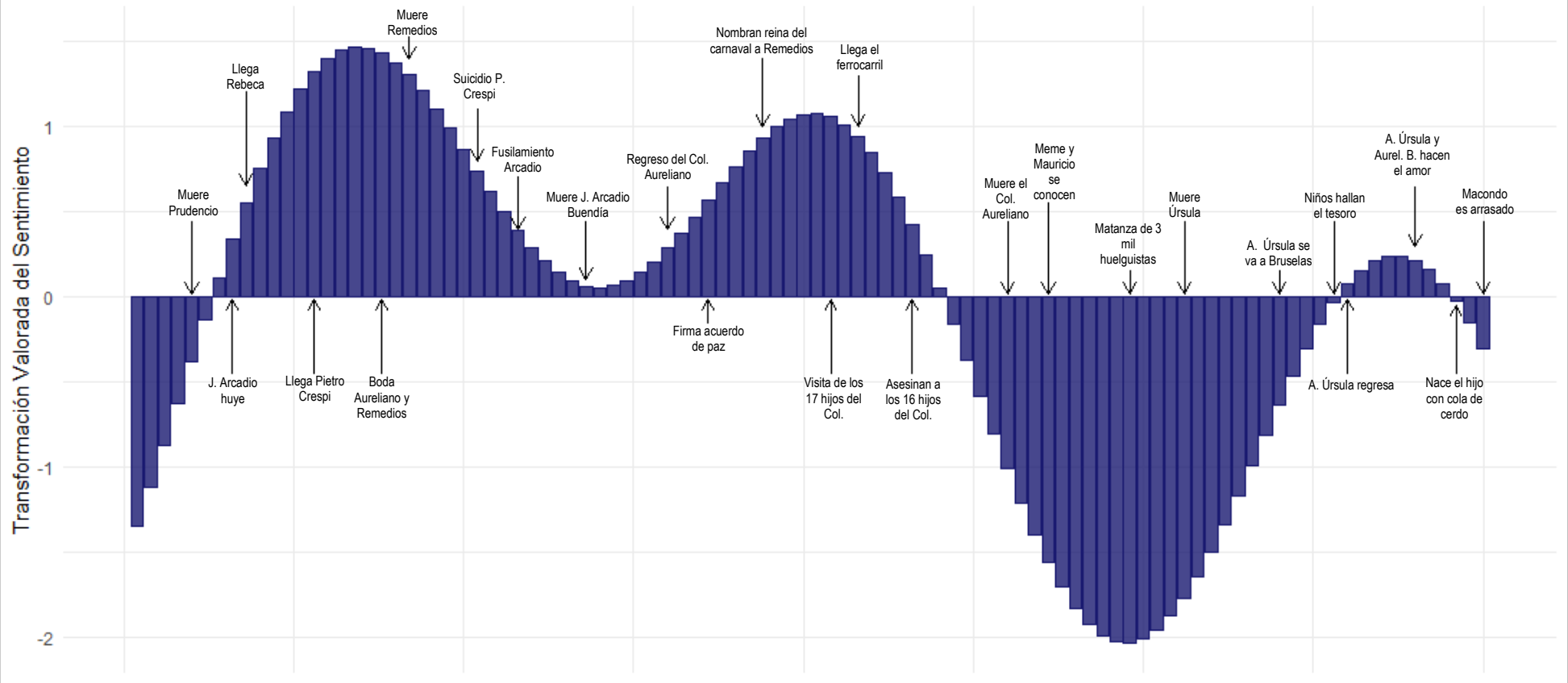
Análisis de Sentimientos de *Cien años de soledad* (1967) y *El amor en los tiempos del cólera* (1985) de Gabriel García Márquez



Autor: Danny Fernando Murillo Lanza
Tutor: Dr. José Manuel Fradejas Rueda

Valladolid, 24 de julio de 2017

Sentimientos en *Cien años de soledad*



¿Quién es Elena Ferrante?



Estilometría



Estilometría

Análisis estadístico del estilo literario

Estilometría

Trata de

- Identificar **semejanzas y diferencias** entre textos y
- **Agrupar** los textos de acuerdo con su características lingüísticas

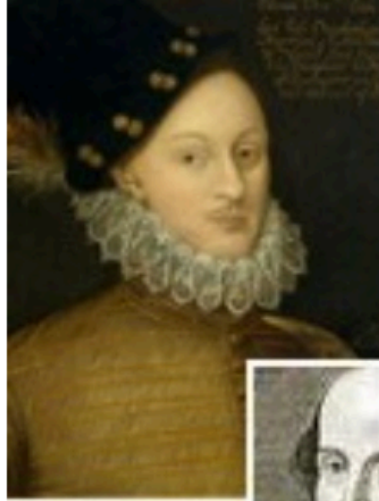
para

- **Detectar señales** en los textos (autoría, géneros, género (sexo), origen, estilo, etc.)



La cuestión chespiriana

Edward de Vere



Francis Bacon



Christopher Marlowe



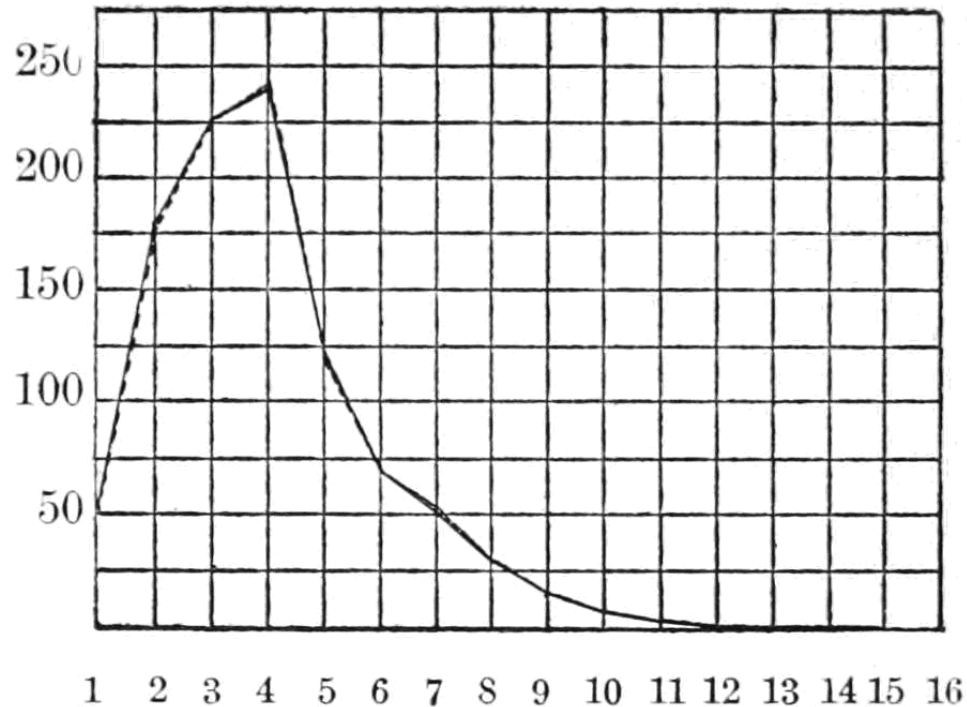
William Stanley



Estilometría

T. C. Mendenhall (1901)

“A Mechanical solution of a literary problem”



Estilometría

George Kingsley Zipf



American linguist,
philologist and statistics
professor for Harvard
University.

Moreno-Sánchez I, Font-Clos F, Corral Á (2016) Large-Scale Analysis of Zipf's Law in English Texts. PLoS ONE 11(1): e0147073. doi:10.1371/journal.pone.0147073 (publicado 22.01.16)

Estilometría

Federalist Papers

85 artículos

Publius

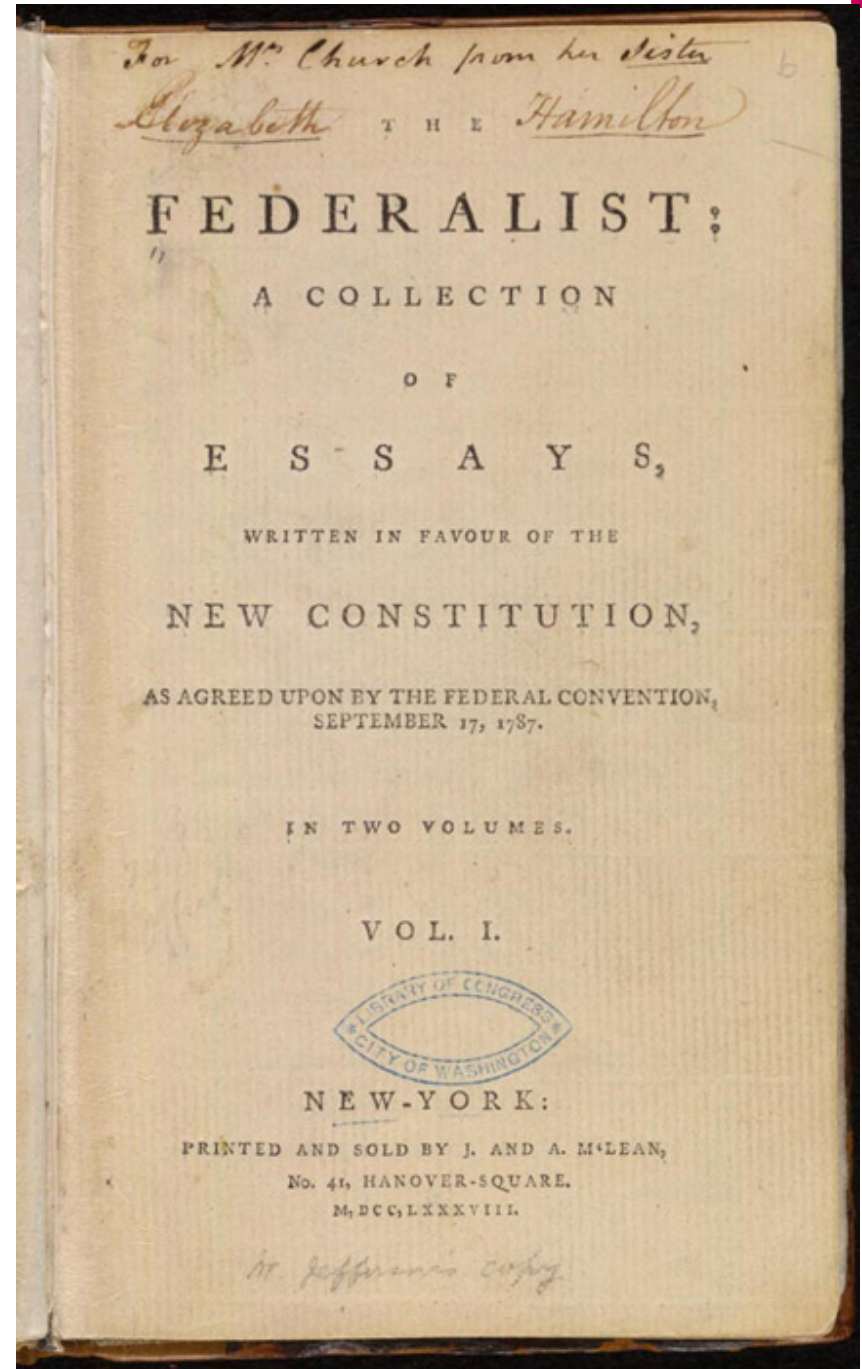
Alexander Hamilton

James Madison

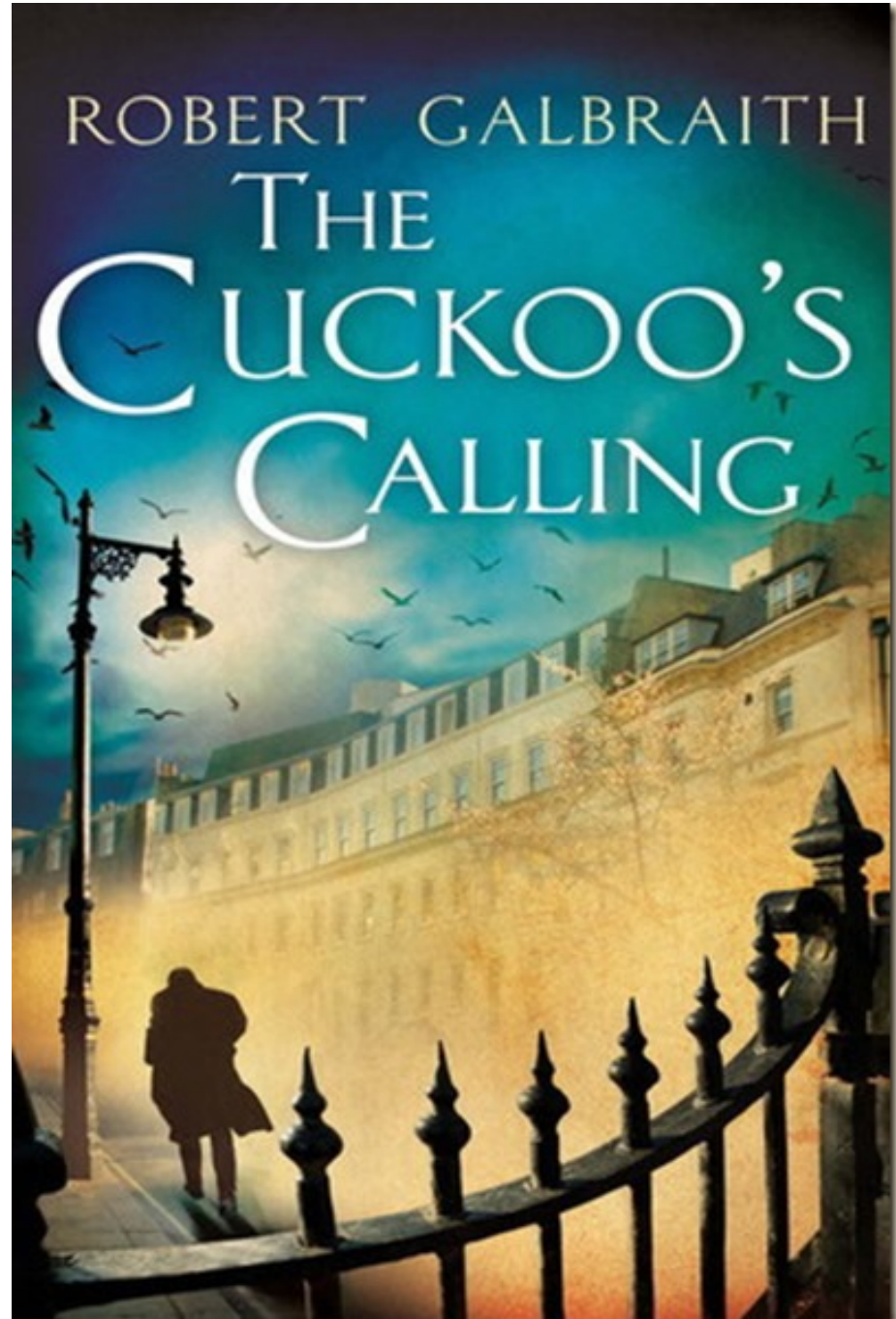
John Jay

Mosteller & Wallace 1964

http://avalon.law.yale.edu/subject_menus/fed.asp



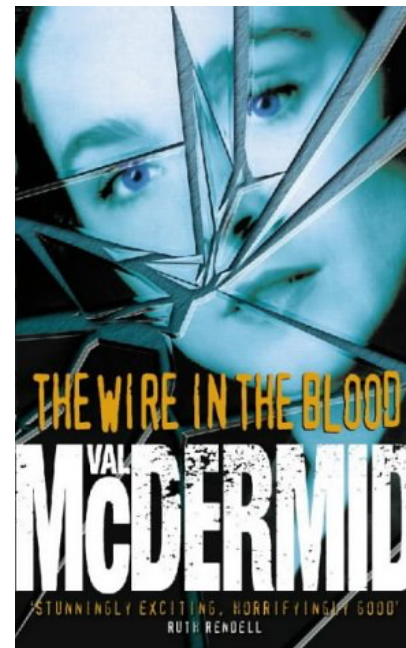
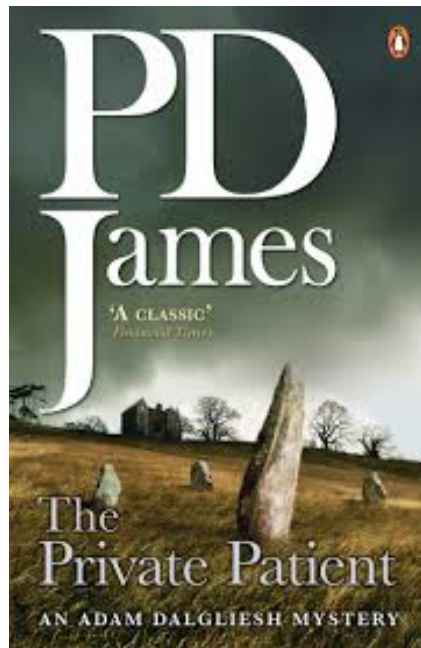
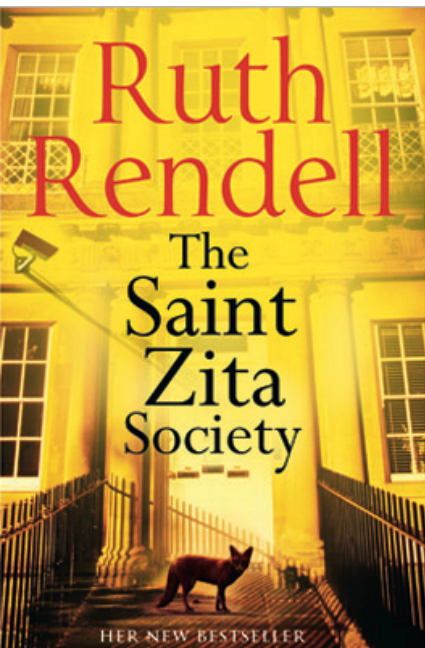
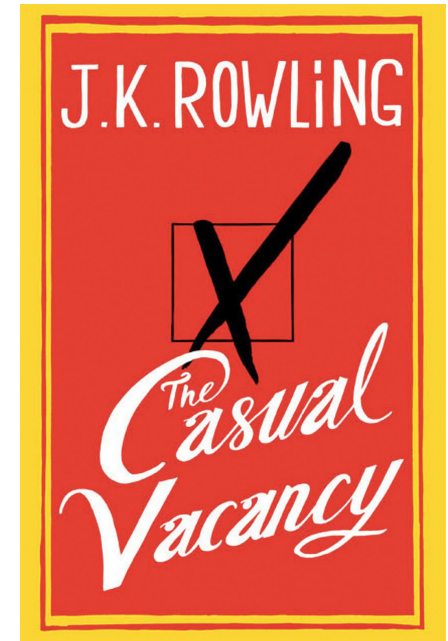
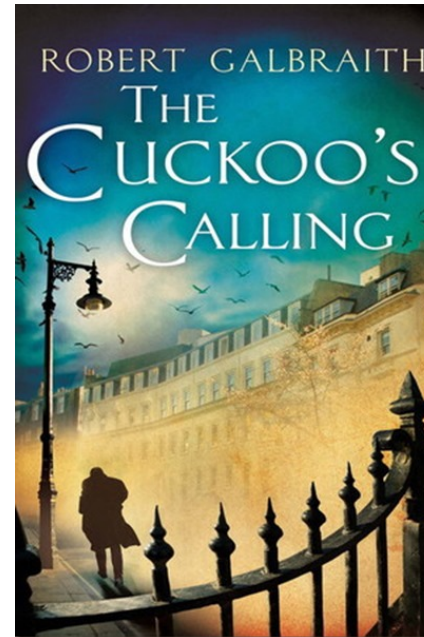
Estilometría



2013

Estilometría

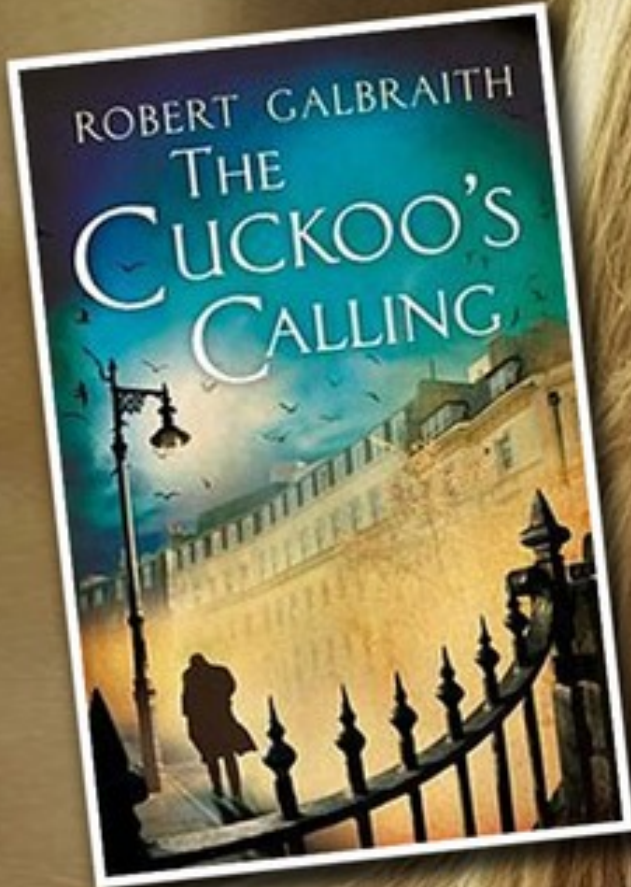
Patrick Juola



Estilometría

- Distribución de la longitud de las palabras
- Uso de las 100 palabras más comunes
- Distribución de 4-grams
- Distribución de bigramas

Estilometría



“The only person consistently suggested by every analysis was Rowlin, who showed up as the winner or the runner-up in each instance”

Minería de textos



Minería de textos



The image shows the RStudio interface. The left pane displays a script with the following code:

```
as.integer(a)
as.integer(as.roman(a))
a <- "xxi"
as.integer(as.roman(a))
library(stylo)
setwd("~/Desktop/LCA-LEXOMICA/LCA/palabras")
stylo()
```

The middle pane shows a dendrogram titled "palabras Cluster Analysis". The x-axis is labeled "800 MFW Culled @ 0% Classic Delta distance" and ranges from 2.0 to 0.0. The y-axis lists 48 samples, each labeled with "AIC_800_XX" where XX is a number from 1 to 48. The dendrogram shows hierarchical clustering of these samples.

The right pane is the console, showing the following output:

```
loading AIC_800_41.txt ...
loading AIC_800_42.txt ...
loading AIC_800_43.txt ...
loading AIC_800_44.txt ...
loading AIC_800_45.txt ...
loading AIC_800_46.txt ...
loading AIC_800_47.txt ...
loading AIC_800_48.txt ...
loading AIC_800_5.txt ...
loading AIC_800_6.txt ...
loading AIC_800_7.txt ...
loading AIC_800_8.txt ...
loading AIC_800_9.txt ...
slicing input text into tokens...
turning words into features, e.g. char n-grams (if applicable)...

Total nr. of samples in the corpus: 48
.....
The corpus consists of 38686 tokens

processing 48 text samples
....
combining frequencies into a table...

culling @ 0 available features (words) 4234
Calculating z-scores...

Calculating classic Delta distances...
MFW used: 800

Function call:
stylo()

Depending on your chosen options, some results should have been written
into a few files; you should be able to find them in your current
(working) directory. Usually, these include a list of words/features
used to build a table of frequencies, the table itself, a file containing
recent configuration, etc.

Advanced users: you can pipe the results to a variable, e.g.:
my.discovery = stylo()
this will create a class "my.discovery" containing some presumably
interesting stuff. The class created, you can type, e.g.:
summary(my.discovery)
to see which variables are stored there and how to use them.

for suggestions how to cite this software, type: citation("stylo")
```




mundo online | arte
cultura
lingüística
literatura
pensamiento

Revista de investigación para el análisis crítico de la cultura, el pensamiento y la sociedad digitales.

Presente en: ESCI (Thomson Reuters - WoK), ERIH Plus, MLAIB, Latindex, DIALNET, CIRC, ISOC, REDIB, DOAJ, ...



Caracteres

Caracteres. Estudios culturales y críticos de la esfera digital | ISSN: 2254-4496 | Salamanca

SÍGUENOS / FOLLOW US



¿QUÉ ES CARACTERES?

Es una publicación multilingüe, independiente e interdiscipli-

REVISTA/JOURNAL ▾ ACERCA DE / ABOUT ▾ NOVEDADES/UPDATES ▾ ADICIONALES/ADDITIONAL ▾ CONTACTAR/CONTACT



El análisis estilométrico aplicado a la literatura española: las novelas policíacas e históricas

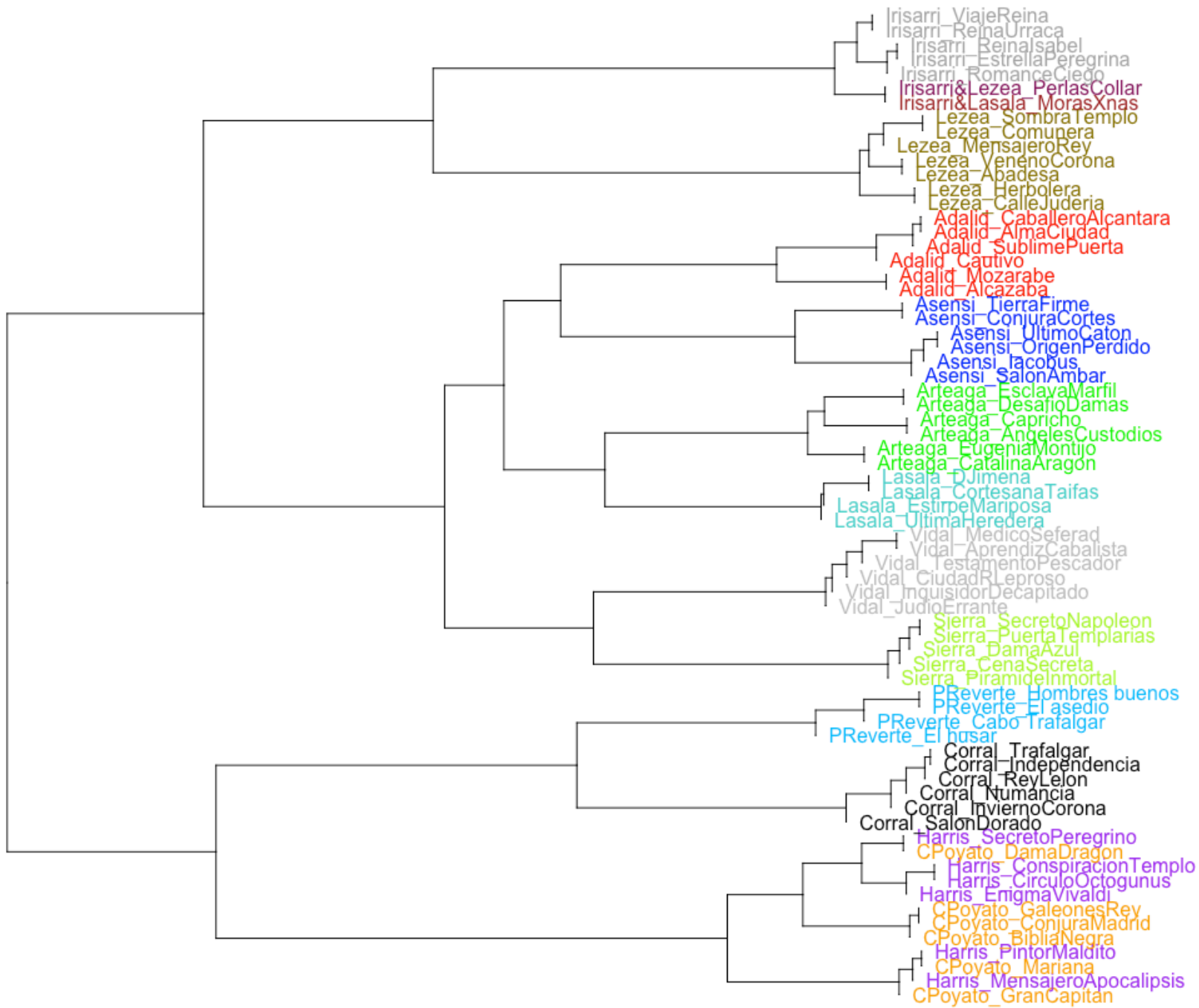
Stylometric Analysis Applied to Spanish Literature: Historical and Crime Fiction

José Manuel Fradejas Rueda (Universidad de Valladolid)

BOLETÍN / NEWSLETTER

[Suscribirse / Subscribe](#)

Histórica		Policíaca
69	textos	49
44 252 212	caracteres	22 732 809
7 659 249	tokens	4 010 728
122 038	tipos	91 789
195 805	párrafos	149 570
1.593342	Token-Tipo Ratio	2.288587
111 003.6	Media de tokens / novela	81 851.59
12 900.86	Media tipos / novela	10 645.39
2 837.754	Media párrafos / novela	3 052.449
641 336.4	Media caracteres / novela	463 934.9
39.11672	Media palabras/párrafo	26.81506
226.0014	Media caracteres / párrafo	151.9878
5.777618	Media caracteres / palabra	5.668001



```

56 texto.todo <- gsub("\n</p>", "</p>", texto.todo, perl = T) # Agrupa </p> aislados en una línea
57 texto.todo <- gsub("</p>\n<p>", "\n<lb break='no' rend='guion'/>", texto.todo, perl = T) # Marca fin de l
58 texto.todo <- gsub("<-\n</p>\n<p>", "\n<lb rend='guion'/>", texto.todo, perl = T) # Marca fin de línea cortad
59 texto.todo <- gsub("#</p>\n<p>", "\n<lb break='no'/>", texto.todo, perl = T) # Marca fin de línea cortado
60 texto.todo <- gsub("</p>\n<p>", "\n<lb/>", texto.todo, perl = T) # Marca fin de línea
61 texto.todo <- gsub("<-\n</p>\n<cb n='2'/>\n<p>", "\n<cb n='2'/>\n<lb break='no' rend='guion'/>", texto.todo,
62 texto.todo <- gsub("#</p>\n<cb n='2'/>\n<p>", "\n<cb n='2'/>\n<lb break='no'/>", texto.todo, perl = T) # M
63 texto.todo <- gsub("</p>\n<cb n='2'/>\n<p>", "\n<cb n='2'/>\n<lb/>", texto.todo, perl = T) # Marca fin de l
64 texto.todo <- gsub("<-\n</p>\n<cb n='(\n+)'>\n</fww type='encabezado'>(.*)</fww>\n</fww type='foliacion'>(\n+)</

```

```

1:1 CODIFICA TEI
R Script

```

```

~/Desktop/~STYLO/_policia/
loading Maluenda_Malatoscia.txt ...
loading Maluenda_RuidoCanerías.txt ...
loading Redondo_GuardianInvisible.txt ...
loading Redondo_LegadoHuesos.txt ...
loading Redondo_OfrendaTormenta.txt ...
loading Silva_AlquimistaImpaciente.txt ...
loading Silva_EstrategiaAgua.txt ...
loading Silva_LejanoPaisEstanques.txt ...
loading Silva_NieblaDoncella.txt ...
loading Silva_ReinaSinEspejo.txt ...
loading Villar_OjosAgua.txt ...
loading Villar_PlayaAhogados.txt ...
loading VMontalban_AsesinatoComiteCentral.txt ...
loading VMontalban_MaresSur.txt ...
loading VMontalban_SoledadManager.txt ...
loading VMontalban_Tatuaje.txt ...
loading VMontalban_YoMateKennedy.txt ...
loading Zanon_JohnnyThunders.txt ...
loading Zanon_NoLlamesACasa.txt ...
loading Zanon_TardeMaNunca.txt ...
slicing input text into tokens...
turning words into features, e.g. char n-grams (if applicable)...

Total nr. of samples in the corpus: 49
.....
The corpus consists of 4008394 tokens

processing 49 text samples
....
combining frequencies into a table...

culling @ 0 available features (words) 5000
Calculating z-scores...

Calculating classic Delta distances...
MFW used: 100

Function call:
stylo()

```

```

Environment History
To Console To Source
m <- addMarkers (m, lng=-1.726366, lat =38.919603, popup = "Chinchilla")
m
m <- addMarkers (m, lng=-0.860909, lat =38.631905, popup = "Villena")
m
m <- addMarkers (m, lng=-1.1300278, lat =37.986111, popup = "Murcia")
m
install.packages("tau")
install.packages(c("tidyverse", "rvest"))
library(stylo)
library(tidytext)
sentiments
library(syuzhet)
setwd("~/Desktop/~STYLO/_historica")
library(stylo)
stylo()
setwd("~/Desktop/~STYLO/_policia")
stylo()

```

Files Plots Packages Help Viewer

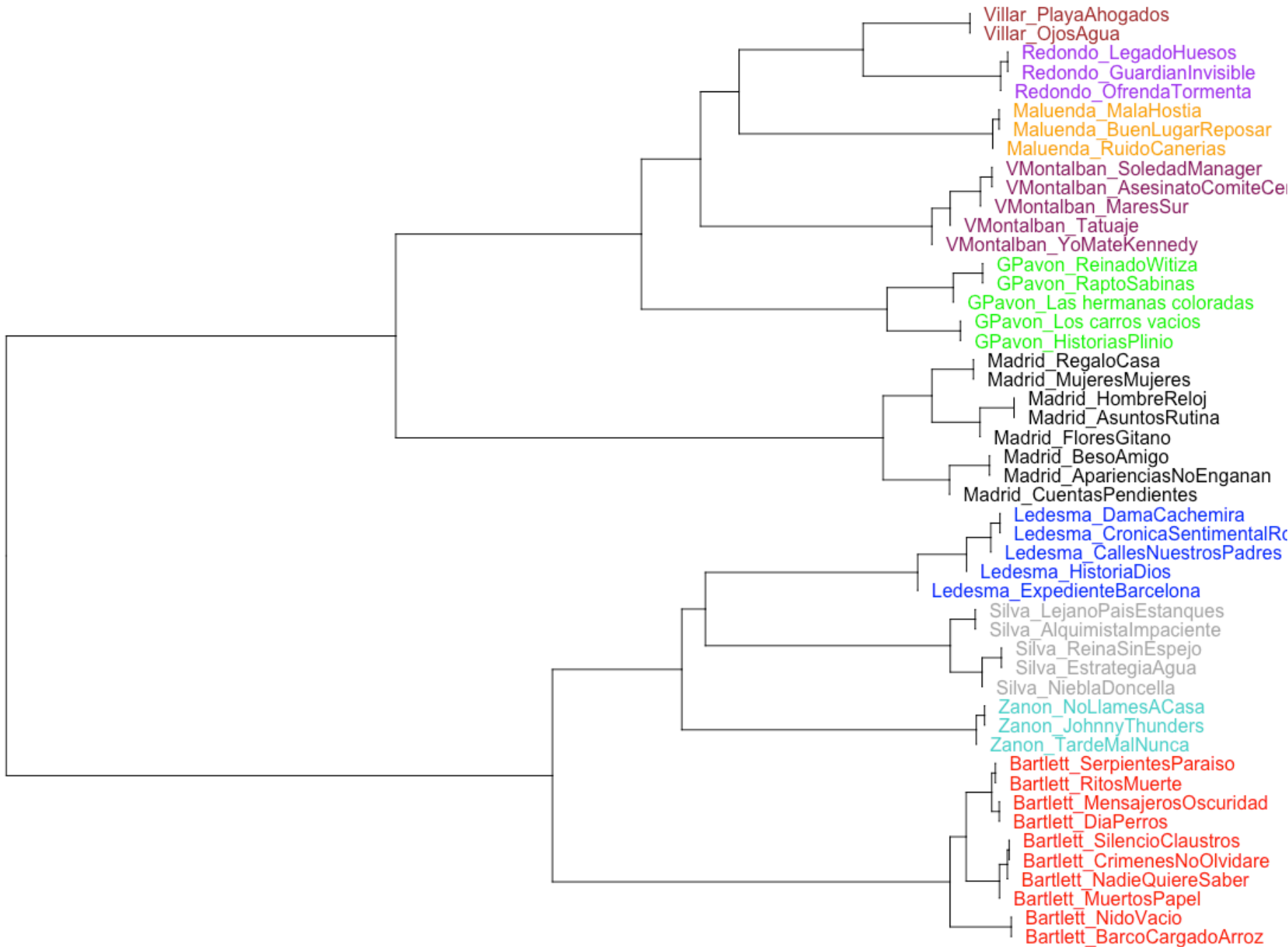
Zoom Export

Publish

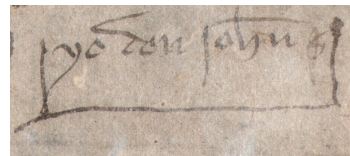
_policia Cluster Analysis

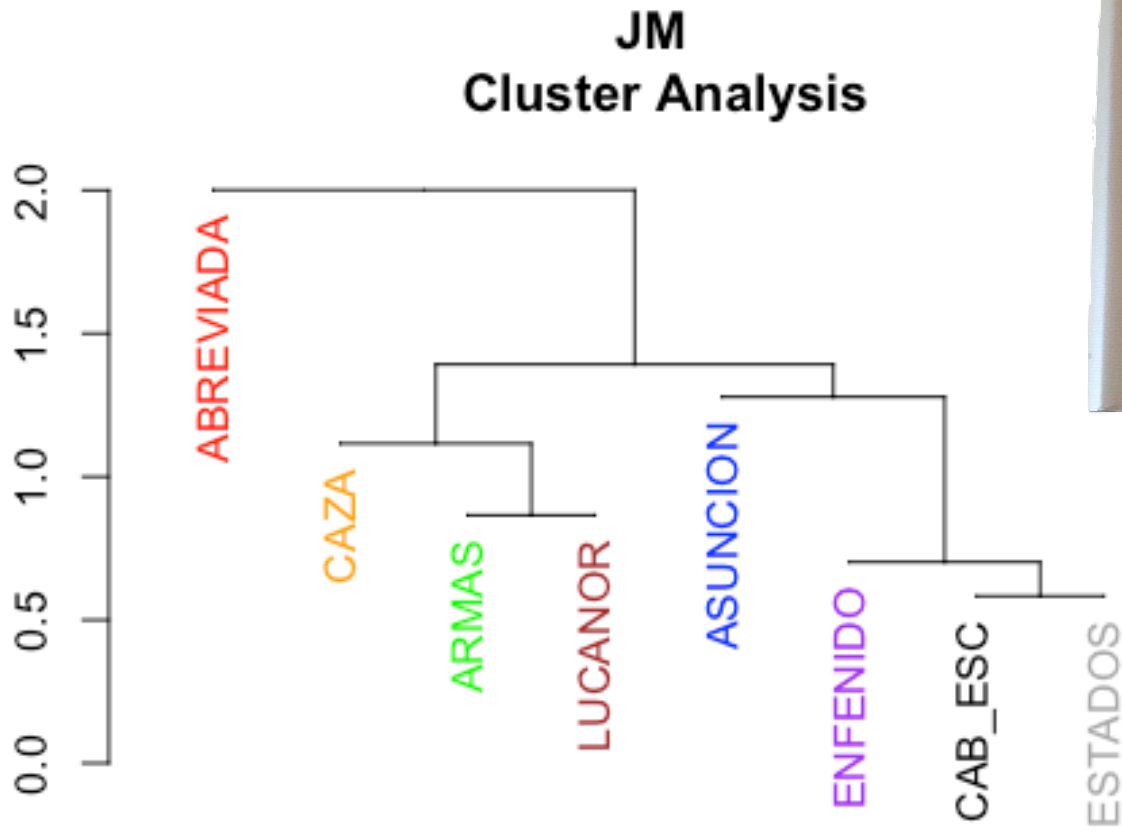
6 4 2 0

100 MFW Culled @ 0%
Classic Delta distance

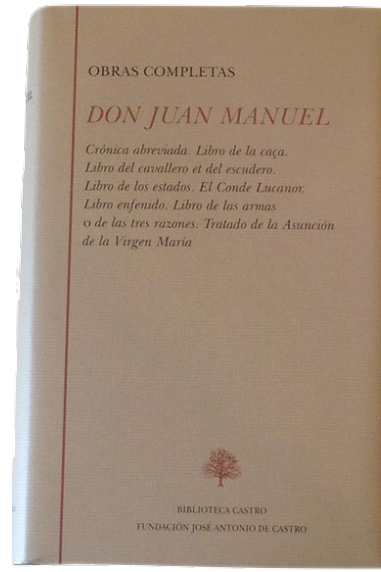


1. Crónica abreviada
2. Libro de la caza
3. Libro del caballero y del escudero
4. Libro de las armas
5. Libro de los estados
6. Libro enfenido
7. Conde Lucanor
8. Tratado de la Asunción





100 MFW Culled @ 0%
Classic Delta distance



Scriptorium alfonsí

21 textos

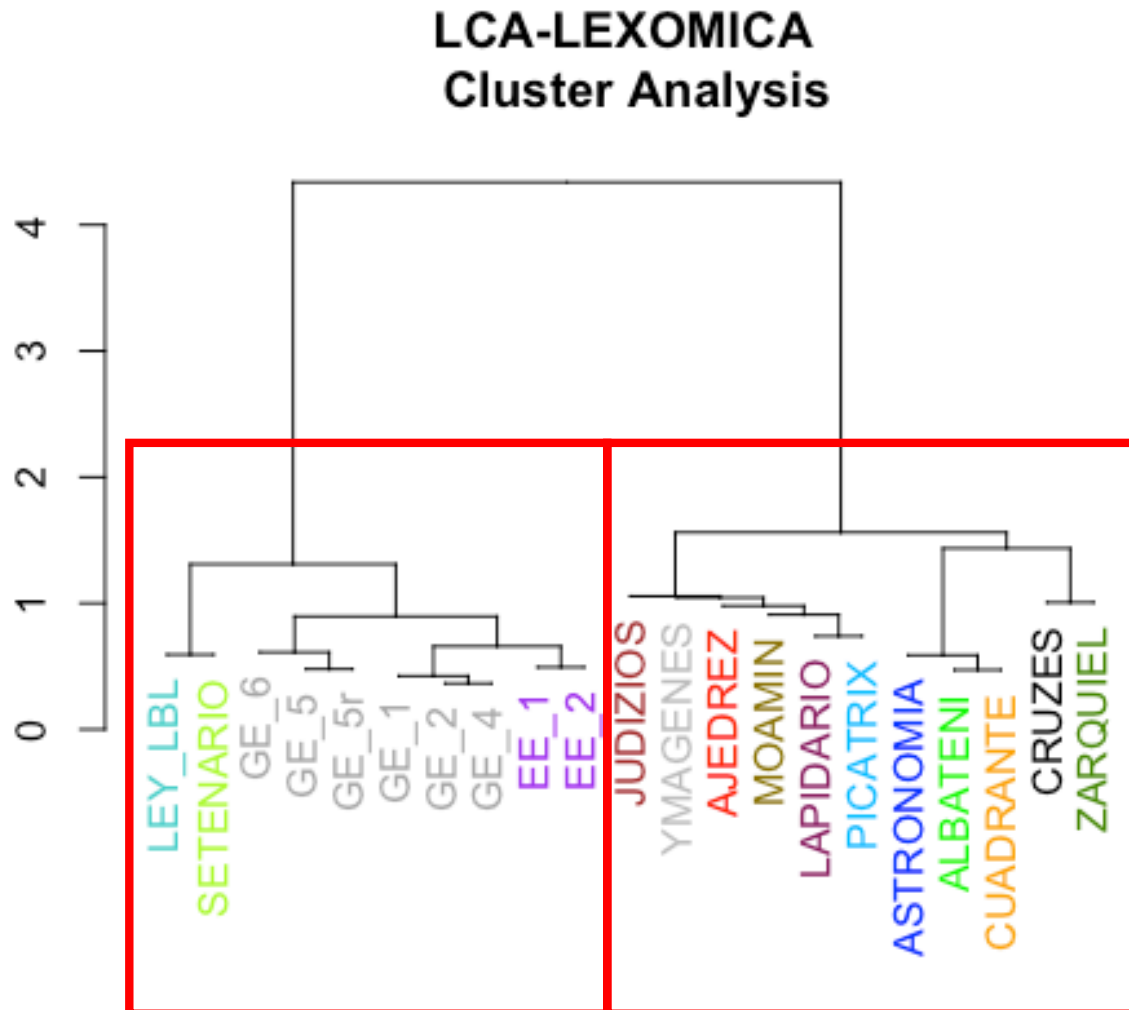
3.234.321 palabras



Scriptorium alfonsí

21 textos

3.234.321 palabras



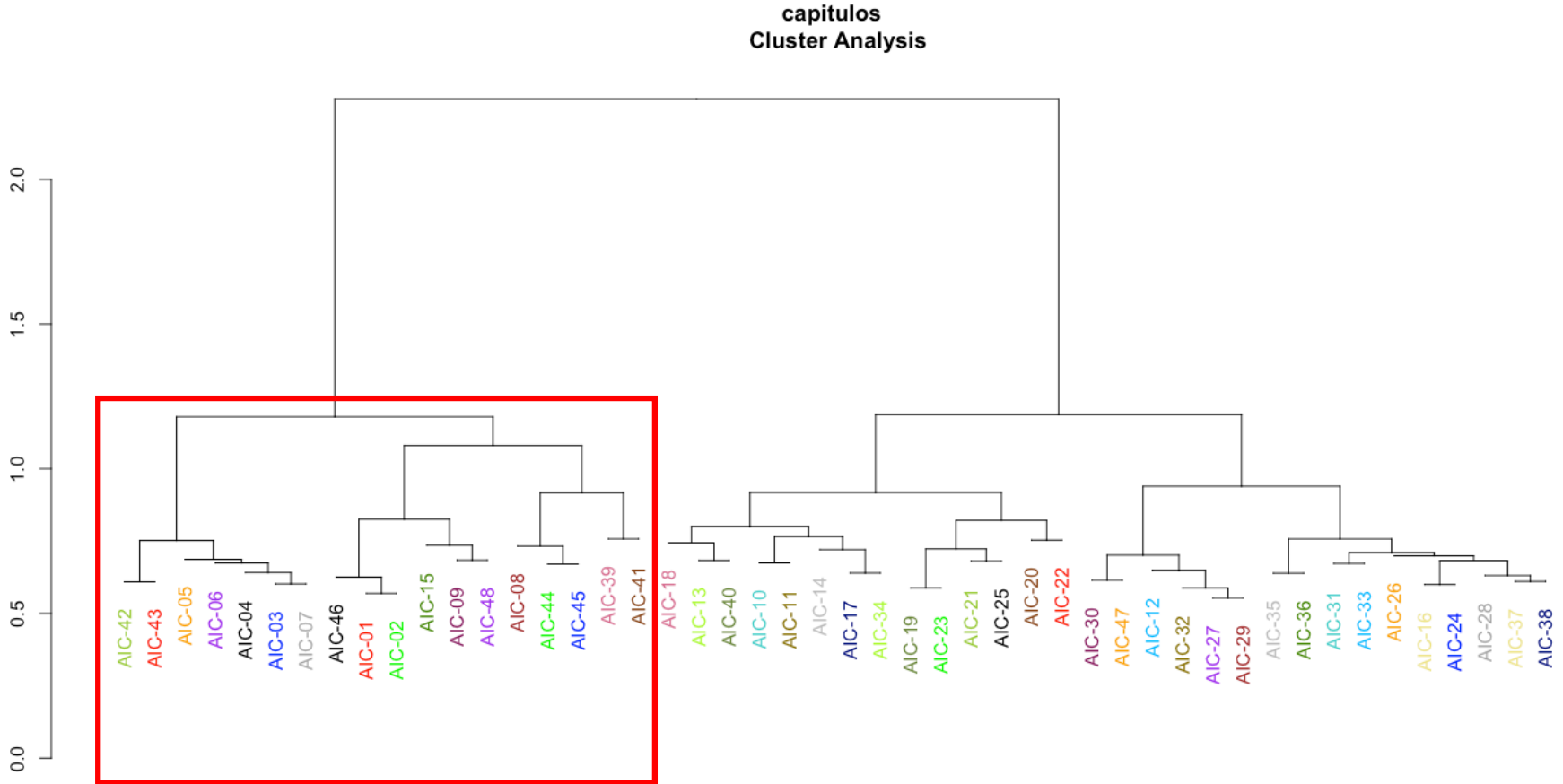


López de Ayala

el libro de la caza
de las aves

Libro de la caza de las aves

Pero López de Ayala



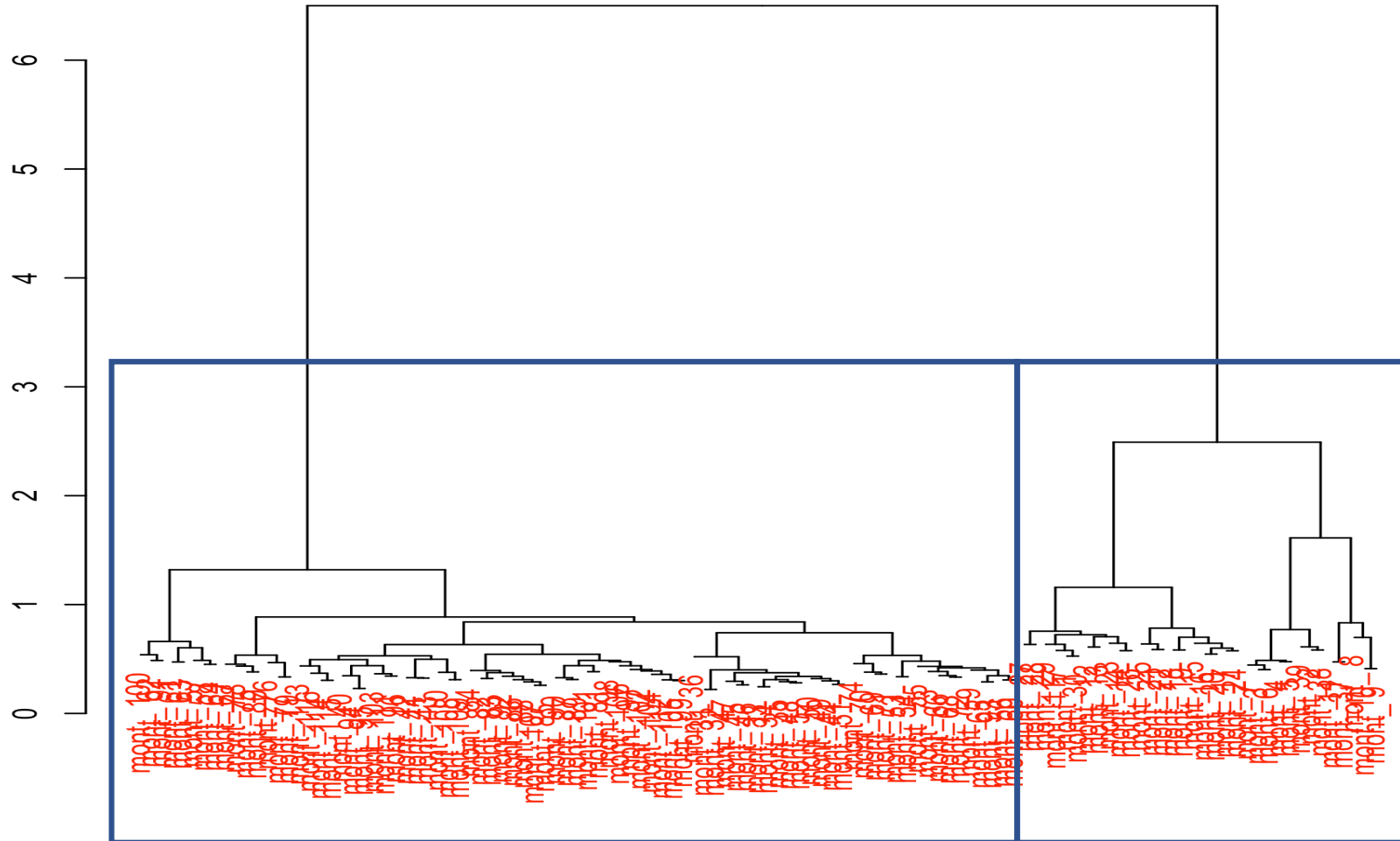
Libro de la montería



LIBRO
 MANDA
 MOSEA
 SIERA
 ROSEU
 ROBLE
 RUEVA
 DONALD
 RONSO

que fabla en to
 do lo q pertenece
 alas maneras
 dela Monteria
 E de parte de otros libros. El primero fabla del guisa y nexo que deve
 traer todo monterio. Anes sea de cavallo qe sea de pie. E en que manera
 deve pncipal e anes sus anes. Tambien de sabuelos. Como de almos.
 E de las feduras que deven anes y de las sermas lindos. Otrosi de las
 cosas que anes en de cada dia en el monte. Opueden anes. E que es
 lo que fagan en cada vna de ellas. E de la otra parte mientro del fuero de la li
 tera. e de los derechos que deven anes los monterios. Por que toman
 en que los ombres toman plazer. Conviene que sepan la Raye de ella. E
 el visio de ella para saberla mejor. Camas plazer avia ombres zuevos
 verro sefara en ella en cada un dia que non la en tendierro.

LCA-LEXOMICA Cluster Analysis



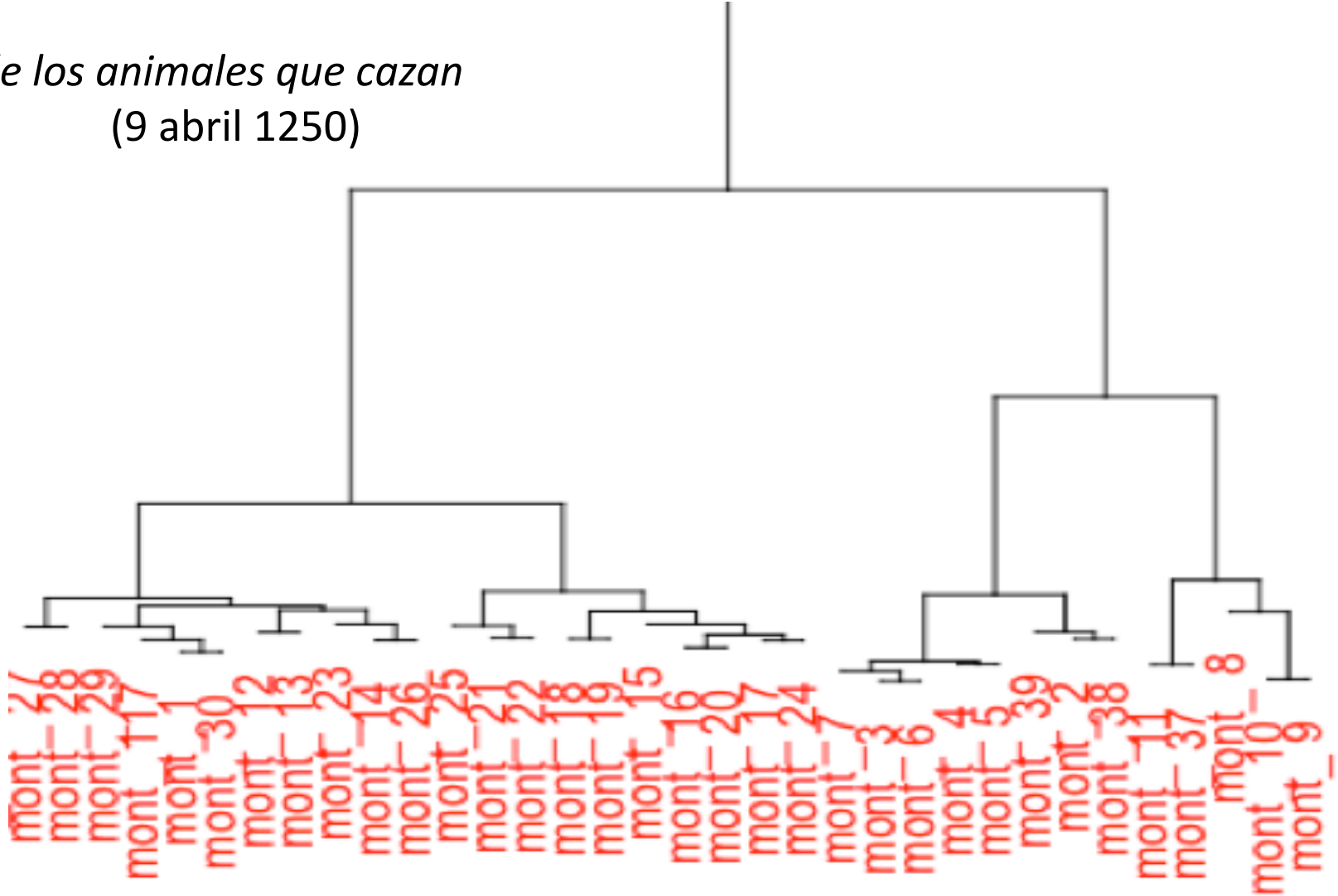
1000 MFW Culled @ 0%
Classic Delta distance

Libro de la montería (detalle)

1000 MFW Delta Classic

Lo procedente del

Libro de los animales que cazan
(9 abril 1250)



Argumento del primer auto desta comedia.

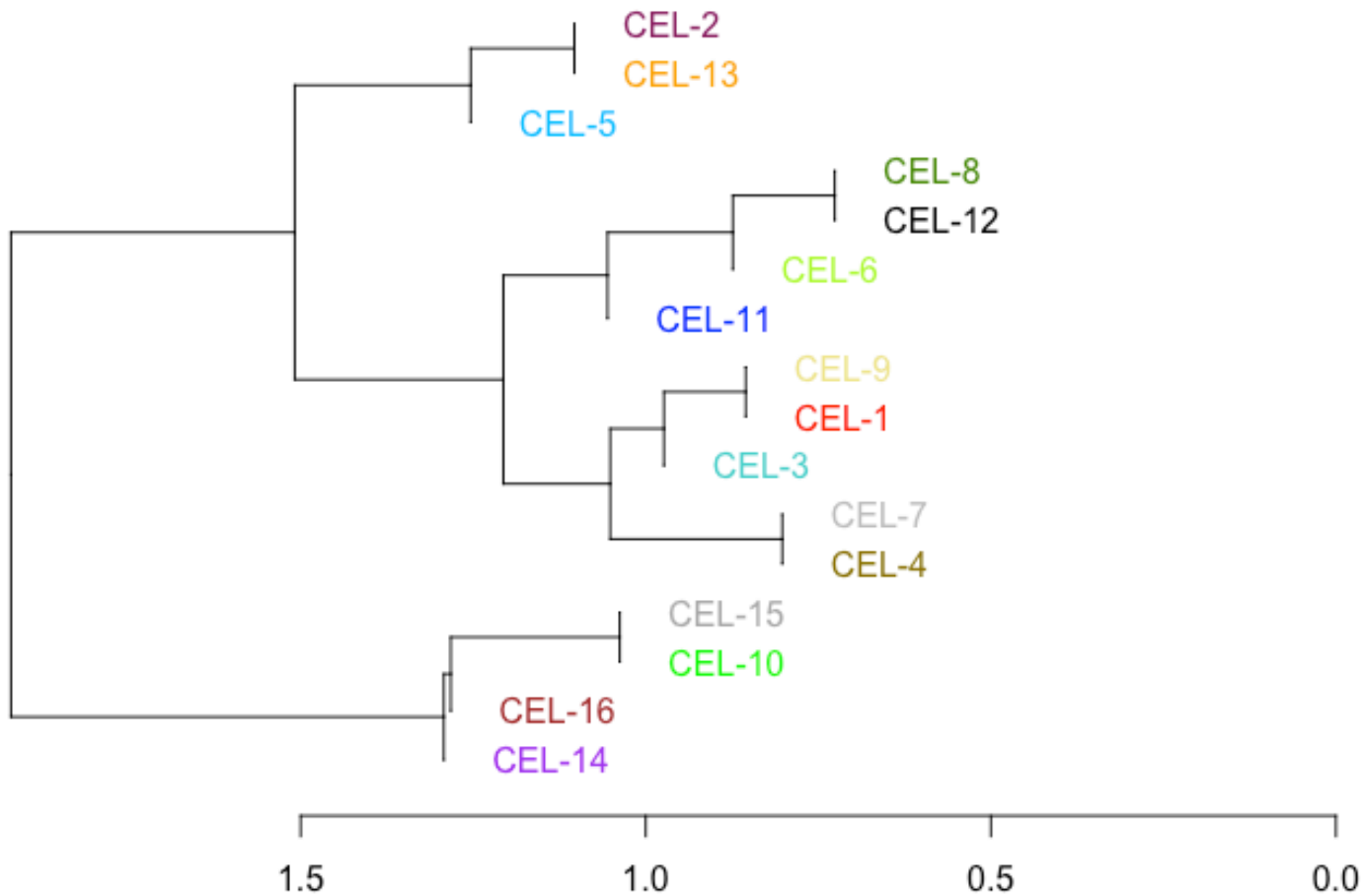
Adelibeas

Calisto



Enterado Calisto vna huerta empos d vn falcon fuyo fallo y a Adelibeas de cuyo amor preso comēcole de hablar: dela qual rigorosamēte despedito: fue para su casa muy sangustiado . hablo con vn criado fuyo llamado semponio. el qual despues de muchas razones le endereco a vna vieja llamada celestina: en cuya casa tenia el mesmo criado vna enamorada llamada elicia: la qual viniēdo sempnio a casa d celestina cōel negocio de su amo tenia a otro confesso llamado crito: al qual escondierō. Entre tanto q semponio esta negociado con celestina: calisto esta razonando cō otro criado fuyo por nōbre parmēno: el qual razonamiēto dura fasta q llega Semponio z celestina a casa de calisto. Parmēno fue conocido de celestina: la qual mucho le dize de los fechos z co

CELESTINA Cluster Analysis



100 MFW Culled @ 0%
Classic Delta distance

NEW APPROACHES TO TEXTUAL DISCUSSIONS IN HERRERA'S POETRY



Laura Hernández Lorenzo
FPU researcher and teaching fellowship
lhernandez1@us.es University of Seville



THE AUTHOR AND HIS WORKS



Fernando de Herrera, 'The Divine'

Algunas Obras, known as *H*



Fernando de Herrera (1534-1597)

- One of the most influential poets of Spanish Golden Age
- His poetic works published in life are known as *H*.
- Textual and authorship problems with the poetic works published after his death by the painter Pacheco and known as *P* (Cuevas 1985)



Versos de Fernando de Herrera, known as *P*

DIGITALISATION



Herrera, Fernando (1975). *Obra poética*. Volumes I and II. Annotated edition by José Manuel Blecuá. Madrid: Real Academia Española.

OCR

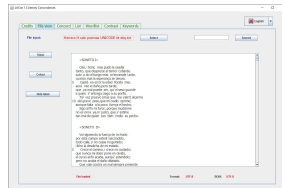
herrera2.txt = *H*
herrera3.txt = *P*
CORPUS

Digitalisation process:

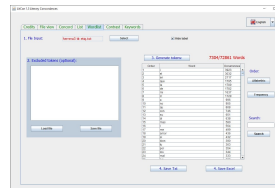
1. OCR through Blecuá's annotated edition
2. Generation of .docx and .txt files
3. Revision and correction of OCR and line break mistakes: à, â, ù
4. Resulting format: UTF-8

PROCESSING WITH LITCON

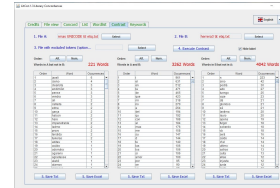
A Java software tool specifically developed to work with poetry



View of *H* corpus

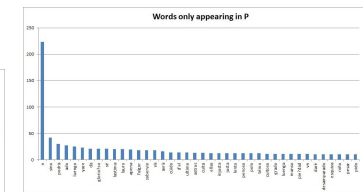
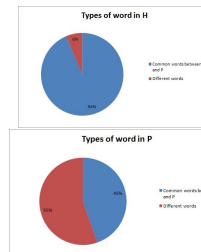
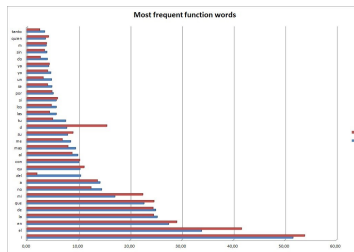


Generating a Wordlist of each of the corpus classified by frequency



Contrast function: *H* and *P* are processed and the software shows us common words between them and words which appear only in one of them, classified by frequency

SOME PRELIMINARY RESULTS



Anthony, I. (2013) 'A Critical Look at Software Tools in Corpus Linguistics', *Linguistic Research*, vol. 30, #2, pp. 141-161.
Blecuá, J. and Ruiz Urbón, C. (2009) 'Evaluación y cuantificación de algunas técnicas de "atribución de autoría" en textos españoles', *Casos: Estudios de literatura*, pp. 27-40.
Blecuá, J. M. (ed.) and Herrera, F. (1975) *Obra poética*. Volumes I and II. Madrid: Real Academia Española.
Cuevas, C. (ed.) and Herrera, F. (de) (1985). *Poesía castellanooriginal completa*. Madrid: Castalia.
Eder, M. (2010) 'Does size matter? Authorship attribution, small samples, big problems', *Proceedings of Digital Humanities*, pp. 132-135.

Holmes, D. I., Kardos, J. (2003). 'Who Was the Author? An Introduction to Stylistometry', *Chance*, vol. 16 #2, pp. 5-8.
Jockers, M. L. (2013) *Macroanalysis: digital methods and literary history*. Urbana: University of Illinois Press.
Kissel, A. D. (1966) *Wortstatistik der lateinischen Dichtung*. München: Beck.
Mauri, C. (1972) *Fernando de Herrera*. Madrid: Gredos.
McIntyre, T. and Hirste, A. (2002) *Corpus Linguistics: Method, Theory and Practice*. Cambridge: Cambridge University Press.
McIntyre, T. and Wilson, A. (2005) *Corpus Linguistics*. Edinburgh: Edinburgh University Press.
Moretti, F. (2007) *La letteratura vaide desde Ajeis*. Barcelona: Neobro Ediciones.

josemanuel.fradejas@uva.es

@JMFradeRue

7partidas.hypotheses.org

github.com/7partidasdigital



FFI2016-75014-P AEI-FEDER, EU