

Data Analysis with R Philolab: stylometry applied to philology

Plan de Innovación y Mejora
Docente

ID2018/152
2018-2019

Claudia García-Minguillán
University of Salamanca
2018-2019



Stylometry applied to Philology

1. What is Stylometry?

2. Why? - Resource for
Humanists

3. How to use it?

4. Research Project sample

1. What is Stylometry?

WHAT?

- Analyze to authorial writing style

HOW?

- R: programming language for statistical computing
- Word Frequency - statistics
- Attractive interface: graphical visualization
- DATA possible exportation



...and... WHY STYLO?



UNIVERSITY OF
OXFORD

digital.humanities
@ oxford

2. Why?- Resource for Humanists

- Digital Humanities
- Traditional research topics
- Scientism to your humanist job!
- Spread your Creativity



European Summer University in Digital Humanities

16 de julio de 2017 ·

...

We're awaiting people from institutions of 27 countries and 4 continents
for #esudh2017. Have a safe trip to #Leipzig!



10



Claudia García-Minguillán
University of Salamanca 2018-2019

UNIVERSIDAD
DE SALAMANCA
CAMPUS DE EXCELENCIA INTERNACIONAL
1218 ~ 2018

UASal
800 AÑOS
1218 ~ 2018

Recent Advances in Computational Linguistics and their Application to Biblical Studies*

J. JOSÉ ALVIAR

Faculty of Theology, University of Navarra, 31080 Pamplona, Spain

Article

Hierarchical and Non-Hierarchical Linear and Non-Linear Clustering Methods to “Shakespeare Authorship Question”

Refat Aljumily

School of English Literature, Language and Linguistics, University of Newcastle, Newcastle upon Tyne, Tyne and Wear NE1 7RU, UK; E-Mail: refat.aljumily@ncl.ac.uk; Tel.: +44-191-208-6233

Academic Editor: Martin J. Bull

Received: 1 July 2015 / Accepted: 6 September 2015 / Published: 17 September 2015

[Did Shakespeare write his plays? - TED talk by Natalya St. Clair and Aaron Williams](#)

2. Resource for Humanists

REASSESSING THE APULEIAN CORPUS

669



FIGURE 17: Bootstrap Consensus Tree (3000-word samples)

Justin Stover & Mike Kestemont (2017)

REASSESSING THE APULEIAN CORPUS: A COMPUTATIONAL APPROACH TO AUTHENTICITY

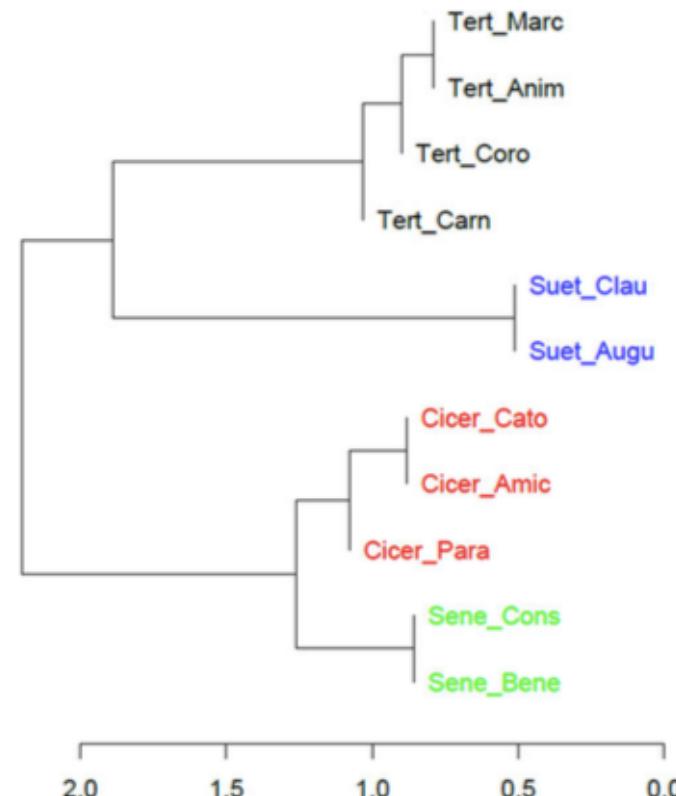


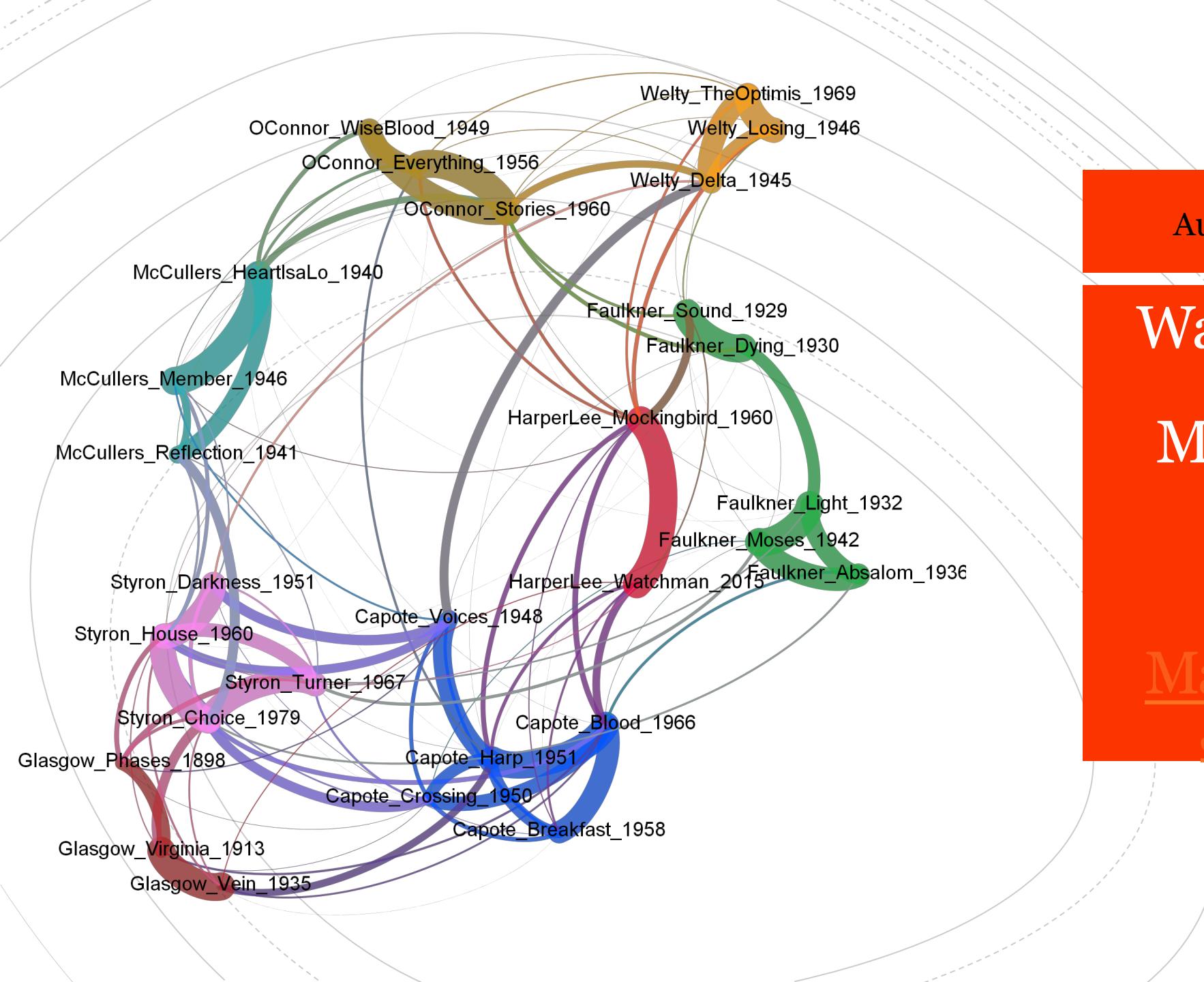
FIGURE 1: Cluster Analysis Dendrogram (no sampling), 100 MFW

Claudia García-Minguillán
University of Salamanca 2018-2019

PhiloLab: Laboratory of Philology

3. How to Stylo?

- Authorship verification
- Genres verification: **BIG** dendograms
- TOP Project!:
'Chemical Linguistics: between linguistics and organic chemistry'
- Functions:
 - `stylo()`
 - `rolling.classify()`
 - `rolling.delta()`
 - `stylo.network()`
 - `oppose()`
 - `rolling.delta()`



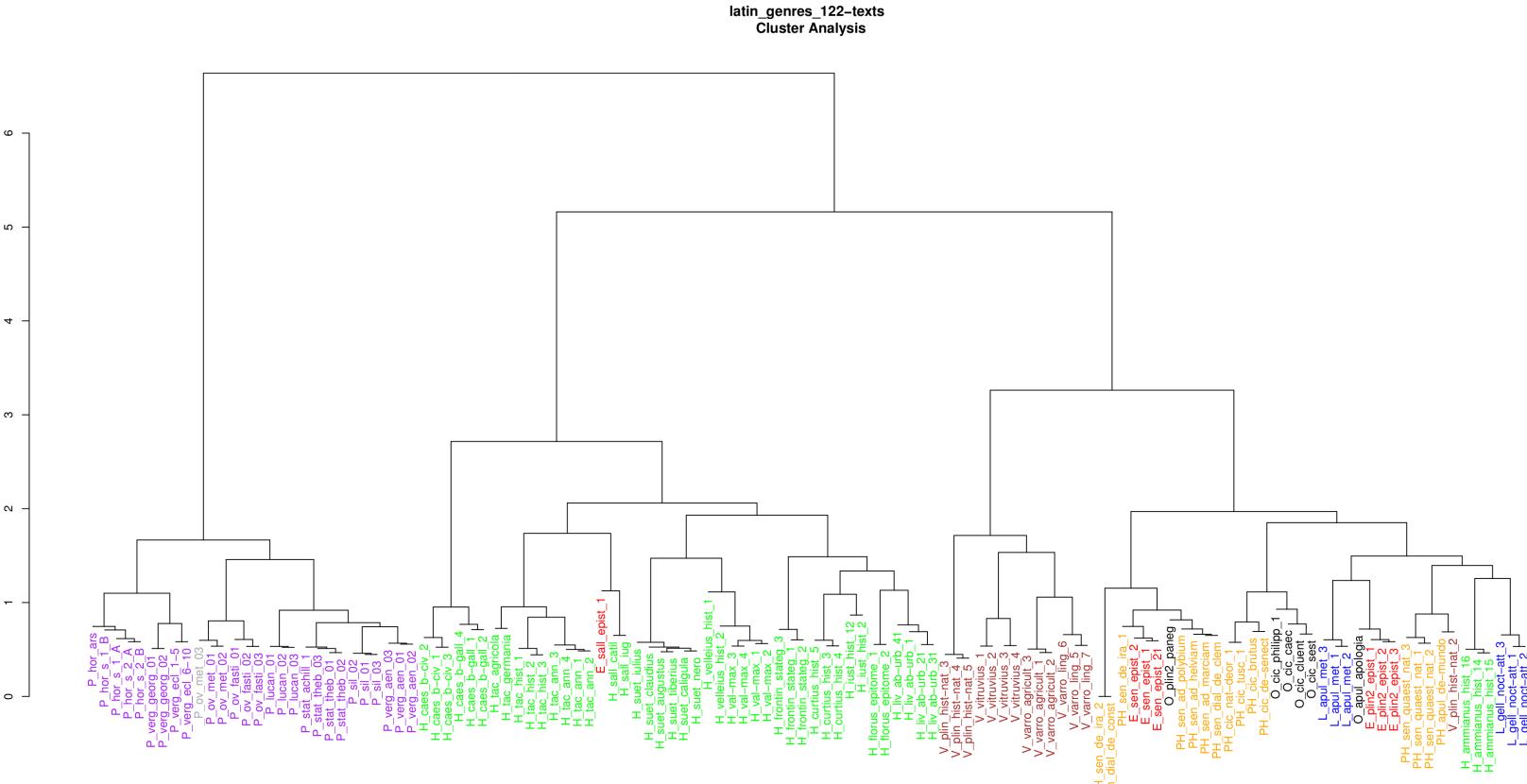
Authorship verification
‘Go Set A

Watchman while we Kill the Mockingbird in Cold Blood’

by
Maciej Eder and
Jan Rybicki

Claudia García-Minguillán
University of Salamanca
2018-2019



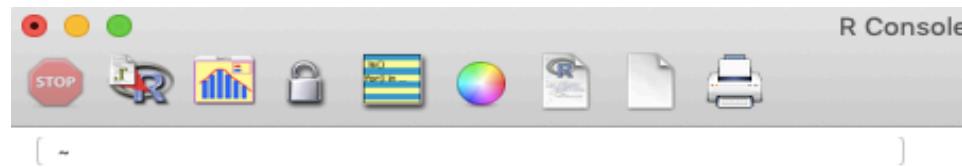


testing BIG dendograms

by Computational Stylistics Group

Claudia García-Minguillán
University of Salamanca





```
R version 3.5.3 Patched (2019-03-11 r76245) -- "Great Truth"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)
```

```
R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.
```

```
Natural language support but running in an English locale
```

```
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.
```

```
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.
```

```
[R.app GUI 1.70 (7543) x86_64-apple-darwin15.6.0]
```

```
[History restored from /Users/claudia/.Rapp.history]
```

```
> library.stylo()
Error in library.stylo() : could not find function "library.stylo"
> install.packages("stylo")
--- Please select a CRAN mirror for use in this session ---
Error in if (res > nrow(m)) { : argument is of length zero
> install.packages("stylo")|
```

stylo()

STATISTICS: Cluster Analysis MDS PCA (cov.) PCA (corr.) tSNE
Consensus Tree Consensus strength 0.5
DELTA DISTANCE: Classic Delta Manhattan Canberra Cosine Delta Eder's Delta Eder's Simple Euclidean Cosine Min-Max

INPUT: plain text xml xml (plays) xml (no titles) html
LANGUAGE: English English (contr.) English (ALL) Latin Latin (u/v > u)
Polish Hungarian French Italian Spanish
Dutch German CJK Other UTF-8

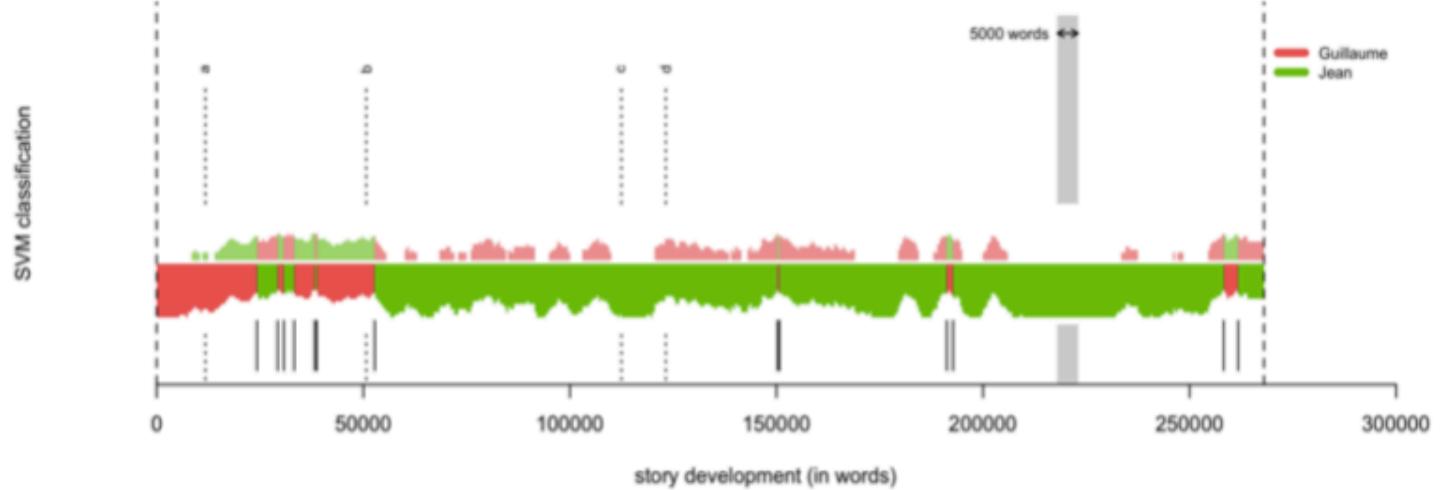
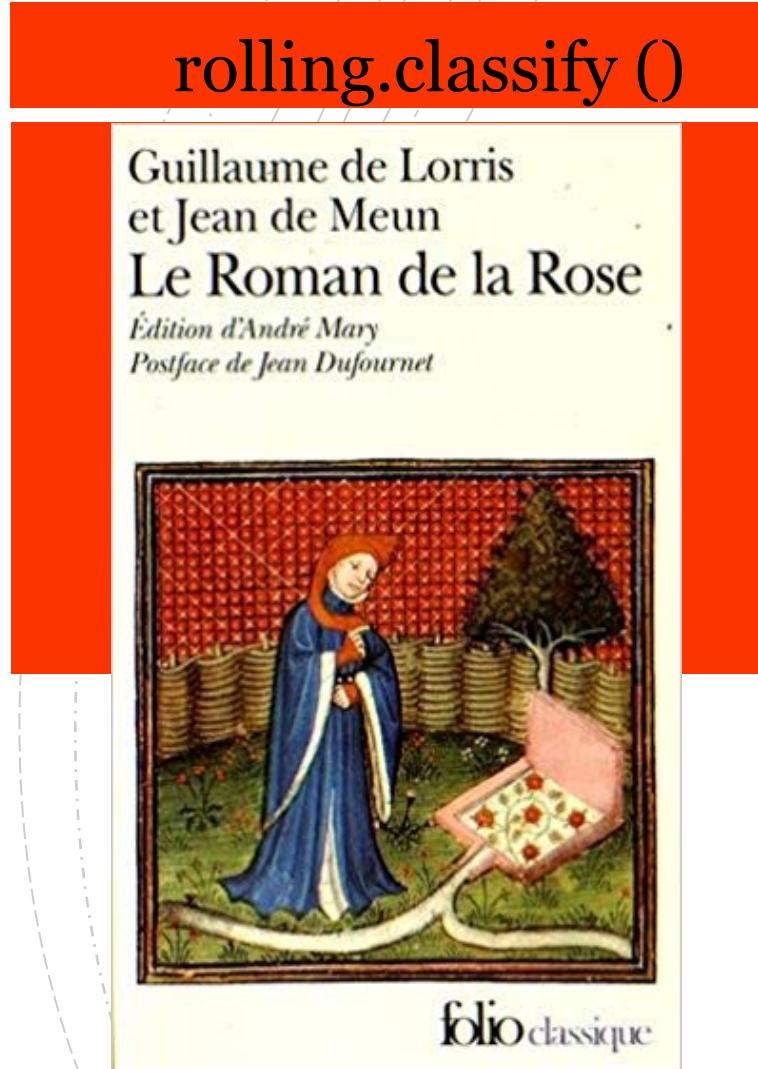


Figure 4: The Rolling Stylometry visualization. The medieval French allegoric story *Roman de la Rose* assessed using Rolling SVM and 100 MFWs; window size: 5,000 words, sample overlap: 4,500 words. Sections attributed to Guillaume de Lorris are marked red, those attributed to Jean de Meun are green. The level of certainty of the classification is indicated by the thickness of the bottom stripe. The commonly-accepted division into two authorial parts is marked with a vertical dashed line 'b'.

Source ↗
[Computational Stylistics Group](#)

Claudia García-Minguillán
University of Salamanca
2018-2019



```
/var/folders/ty/h9hk3p355c1gxfb0m4y4mm1c0000gn/T//RtmpD6MIHp downloaded_packages  
> install.packages("tsne")  
trying URL 'https://cran.rediris.es/bin/macosx/el-capitan/contrib/3.5/tsne_0.1-3.tgz'  
Content type 'application/x-gzip' length 21437 bytes (20 KB)  
  
'help.start()' for an HTML browser: downloaded 20 KB  
Type 'q()' to quit R.  
  
[R.app GUI 1.70 (7543) x86_64-app] The downloaded binary packages are in  
/var/folders/ty/h9hk3p355c1gxfb0m4y4mm1c0000gn/T//RtmpD6MIHp downloaded_packages  
[History restored from /Users/cla] > stylo()  
Error in stylo() : could not find function "stylo"  
> help()  
starting httpd help server ... do Error in stylo() : could not find function "stylo"  
> stylo()  
Error in stylo() : could not find library.stylo() : could not find function "library.stylo"  
> install.packages("tcltk2") > library(stylo)  
--- Please select a CRAN mirror for this session [1] (enter a number)  
trying URL 'https://cran.rediris.es/' ## stylo version: 0.6.9 ##  
Content type 'application/x-gzip'  
===== If you plan to cite this software (please do!), use the following reference:  
downloaded 946 KB  
Eder, M., Rybicki, J. and Kestemont, M. (2016). Stylometry with R:  
a package for computational text analysis. R Journal 8(1): 107-121.  
<https://journal.r-project.org/archive/2016/RJ-2016-007/index.html>  
  
The downloaded binary packages are in  
/var/folders/ty/h9hk3p355c1gx To get full BibTeX entry, type: citation("stylo")  
> install.packages("ape") Warning message:  
trying URL 'https://cran.rediris.es/bin/macosx/el-capitan/contrib/3.5/ape_2.8-1.tgz'  
Content type 'application/x-gzip' length 1030000 bytes (997 KB)  
running command '/usr/bin/otool' -L '/Library/Frameworks/R.framework/Resources/library/tcltk/libs/tcltk.so' had status 1  
===== sty  
downloaded 2.5 MB  
> stylo()  
using current directory...  
  
The downloaded binary packages are in  
/var/folders/ty/h9hk3p355c1gx !!!!!!!!  
> install.packages("class") Hey! The working directory should contain the subdirectory "corpus"  
trying URL 'https://cran.rediris.es/bin/macosx/el-capitan/contrib/3.5/class_7.3-15.tgz'  
Content type 'application/x-gzip'  
  
downloaded 92 KB  
Error in stylo() : Corpus prepared incorrectly  
>  
  
The downloaded binary packages are in  
/var/folders/ty/h9hk3p355c1gxfb0m4y4mm1c0000gn/T//RtmpAXumk downloaded_packages  
> install.packages("e1071")  
trying URL 'https://cran.rediris.es/bin/macosx/el-capitan/contrib/3.5/e1071_1.7-1.tgz'  
Content type 'application/x-gzip' length 896841 bytes (875 KB)  
  
downloaded 875 KB  
  
The downloaded binary packages are in  
/var/folders/ty/h9hk3p355c1gxfb0m4y4mm1c0000gn/T//RtmpAXumk downloaded_packages  
> install.packages("pamr")  
trying URL 'https://cran.rediris.es/bin/macosx/el-capitan/contrib/3.5/pamr_1.56.tgz'  
Content type 'application/x-gzip' length 817907 bytes (798 KB)  
  
downloaded 798 KB  
  
The downloaded binary packages are in  
/var/folders/ty/h9hk3p355c1gxfb0m4y4mm1c0000gn/T//RtmpAXumk downloaded_packages  
> install.packages("tsne")
```

help()



Claudia García-Minguillán
University of Salamanca
2018-2019



4. Research Project Sample

‘The many faces of an author’s style’

...Forthcoming publication

How a Computer Program Helped Show J.K. Rowling write A Cuckoo's Calling



Author of the *Harry Potter* books has a distinct linguistic signature



Hypothesis

Claudia García-Minguillán
University of Salamanca
2018-2019



Claudia García-Minguillán
Universidad de Salamanca 2018-2019

Fernando_Pessoa 0

- Heteronimia: Fernando Pessoa, Ricardo Reis, Bernardo Soares, Alvaro de Campos, Alberto Caeiro
 - Aims:
 - 1) Hypothesis
 - 2) Testing author's style
 - 3) Testing Stylo



Mind your corpus



1. Open Access Pessoa's work

Gutenberg Projects

OCR

Corpora—think about sharing it at GitHub!



2. Read and prepare it

Collatio – ‘Crítica Textual’

Txt.

Name your corpus



3. Remember the steps!

Parametres



4. Good practices

Be good and honest

 Caeiro_Guarda.txt

 Caeiro_Pastor.txt

 Caeiro_Poemas.txt

 Campos_Poemas.txt

 Campos_Ultim.txt

 Pessoa_Mensagem_I.txt

 Pessoa_Mensagem_II.txt

 Pessoa_Mensagem_III.txt

 Reis_poemas.txt

 Soares_Desassossego.txt

 Antologia Poética - Fernando Pessoa.txt

 Cancioneiro - Fernando Pessoa.txt

 Do Livro do Desassossego - Bernardo Soares.txt

 Ficções do interlúdio, para além do outro oc...o de Coelho Pacheco - Fernando Pessoa.txt

 Mensagem - Fernando Pessoa.txt

 O Banqueiro Anarquista - Fernando Pessoa.txt

 O Eu profundo e os outros Eus - Fernando Pessoa.txt

 O Guardador de Rebanhos - Alberto Caeiro.txt

 O Marinheiro - Fernando Pessoa.txt

 O pastor amoroso - Alberto Caeiro.txt

 Obra Completa - Fernando Pessoa.txt

 Poemas de Álvaro de Campos - Fernando Pessoa.txt

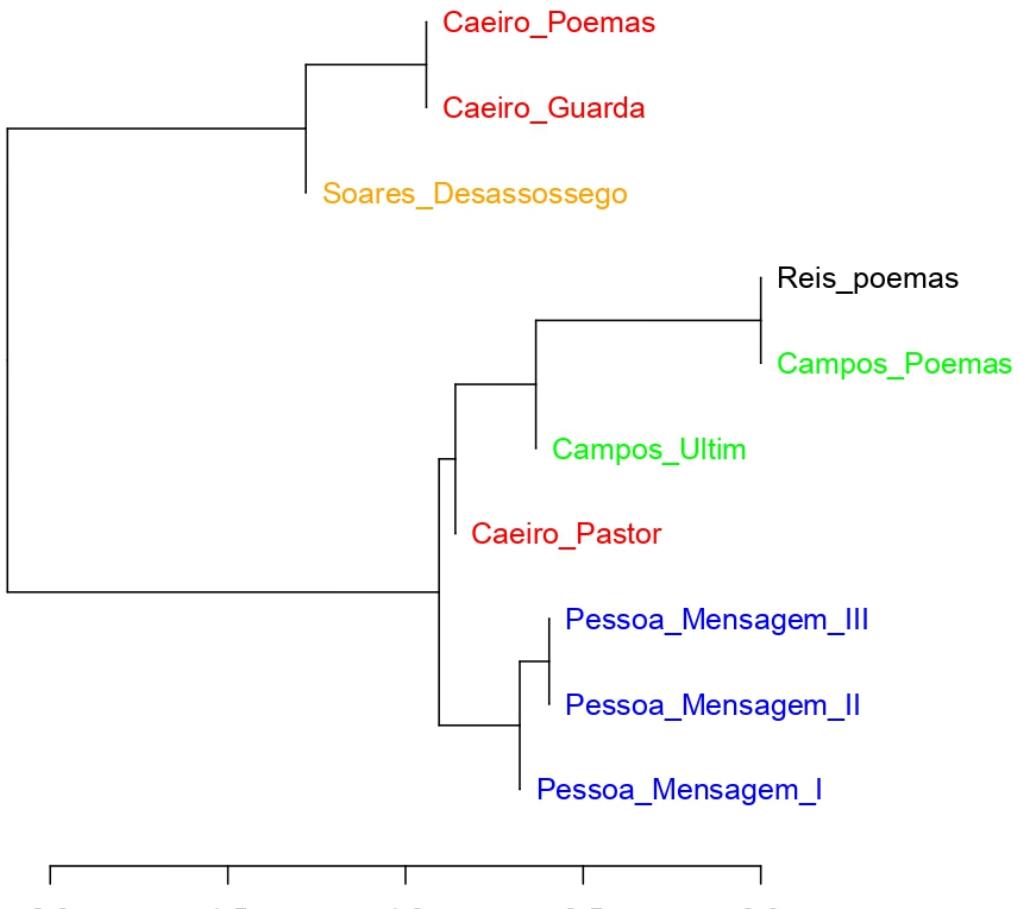
 Poemas de Fernando Pessoa.txt

 Poemas de Ricardo Reis - Fernando Pessoa.txt

 Poesias Inéditas - Fernando Pessoa.txt

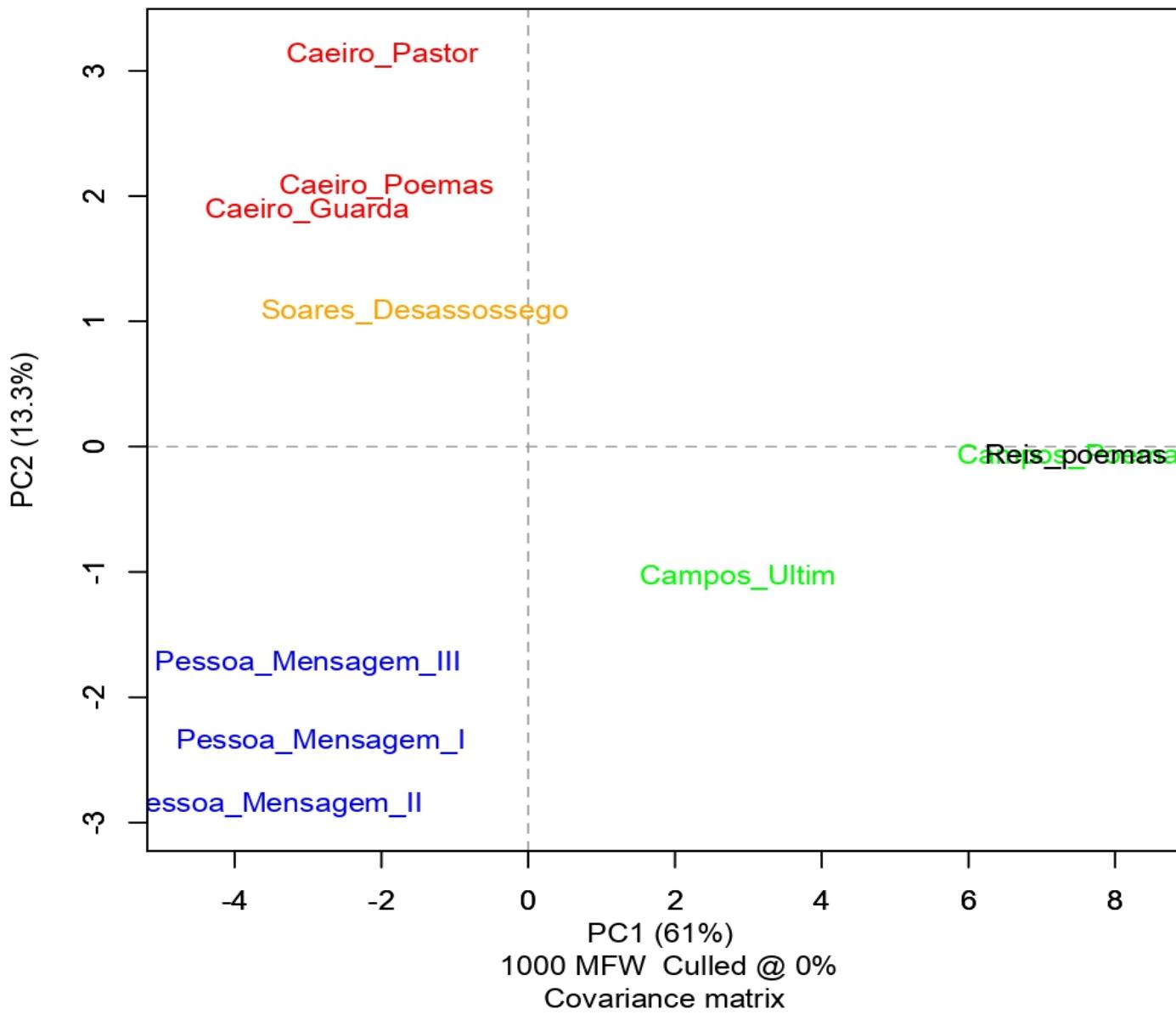
 Primeiro Fausto - Fernando Pessoa.txt

*Corpus Pessoa (I) Cluster Analysis

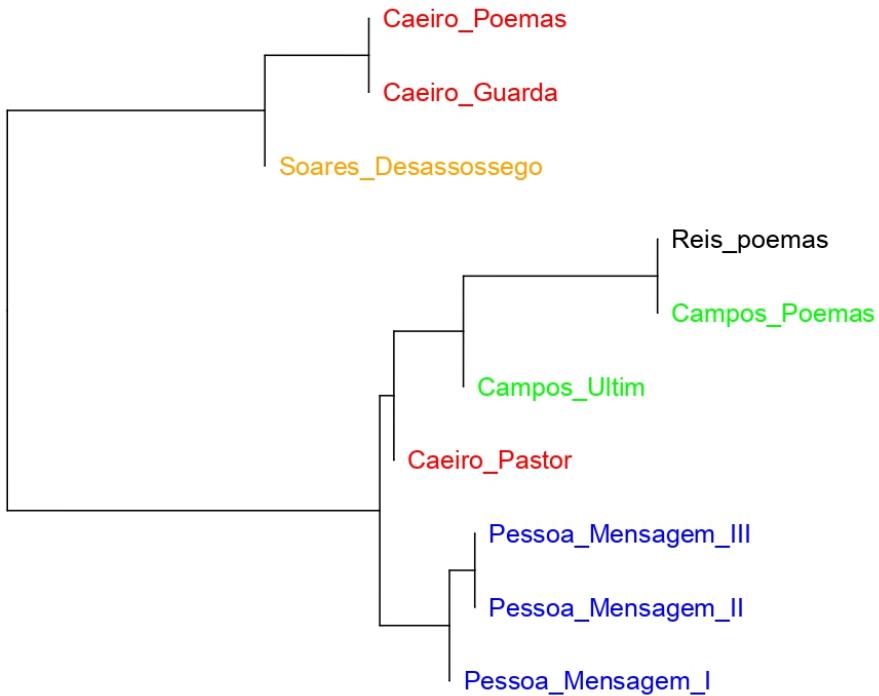


- Caeiro_Guarda.txt
- Caeiro_Pastor.txt
- Caeiro_Poemas.txt
- Campos_Poemas.txt
- Campos_Ultim.txt
- Pessoa_Mensagem_I.txt
- Pessoa_Mensagem_II.txt
- Pessoa_Mensagem_III.txt
- Reis_poemas.txt
- Soares_Desassossego.txt

*Corpus Pessoa (I) Principal Components Analysis

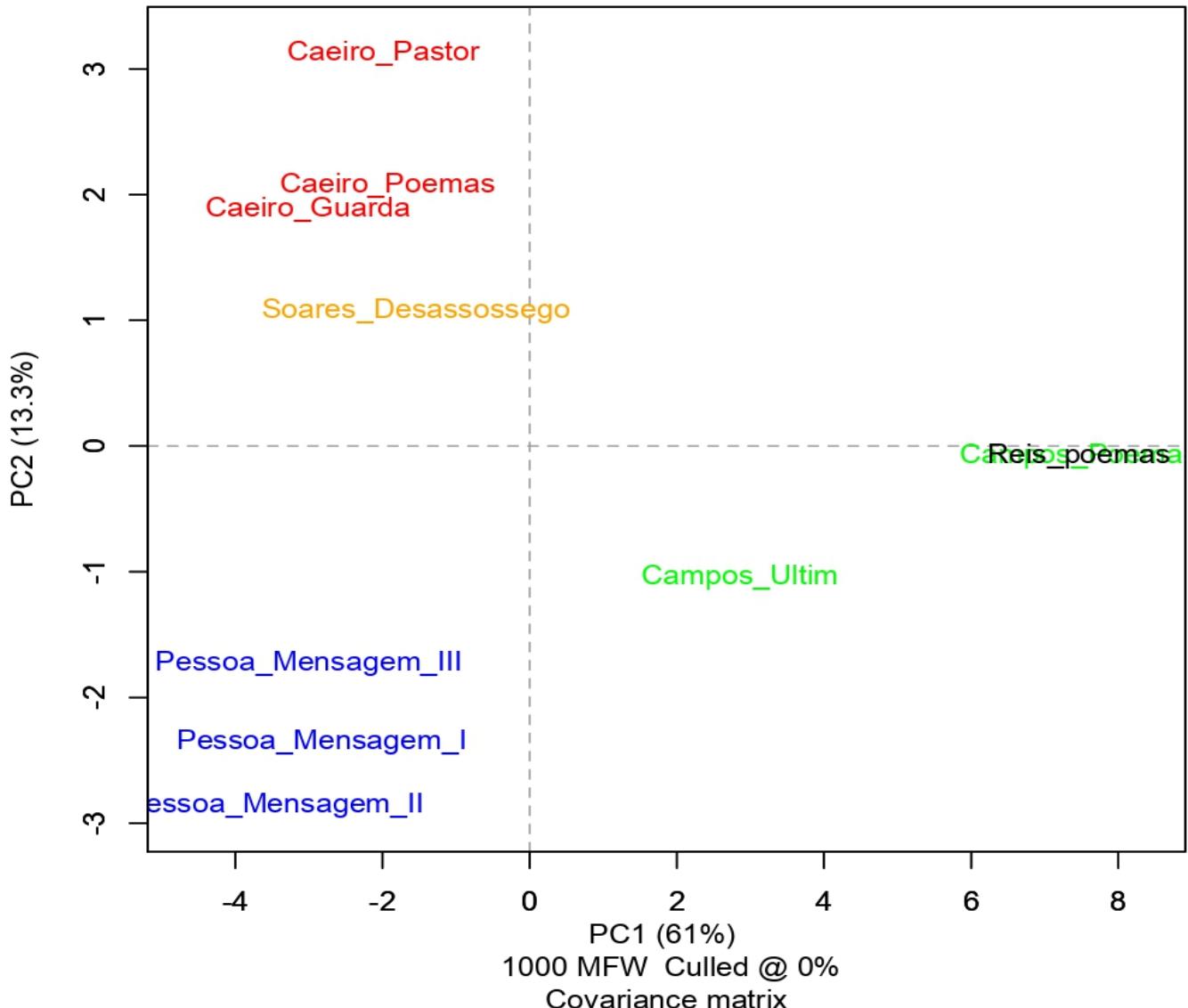


***Corpus Pessoa (I)**
Cluster Analysis

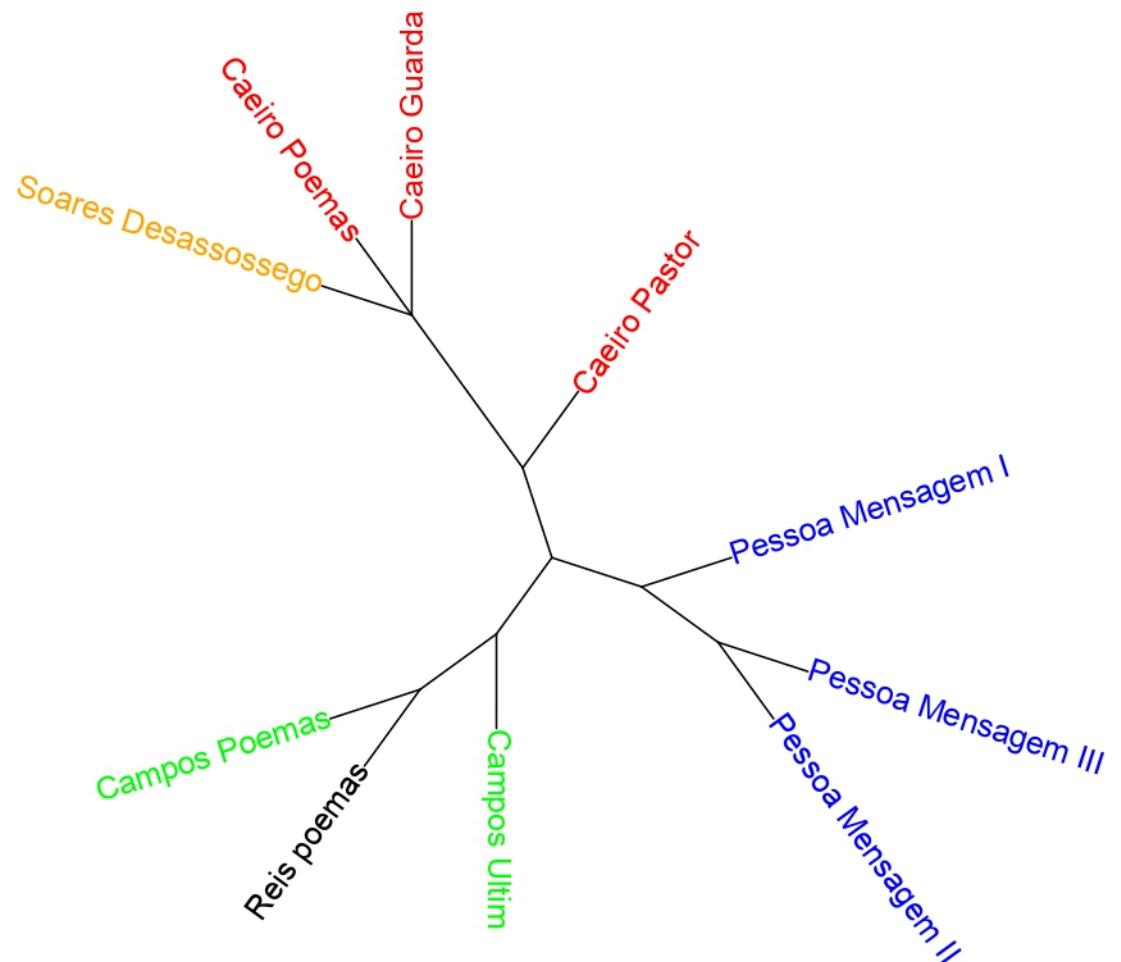


2.0 1.5 1.0 0.5 0.0
1000 MFW Culled @ 0%
Classic Delta distance

***Corpus Pessoa (I)**
Principal Components Analysis



***Corpus Pessoa (I)**
Bootstrap Consensus Tree



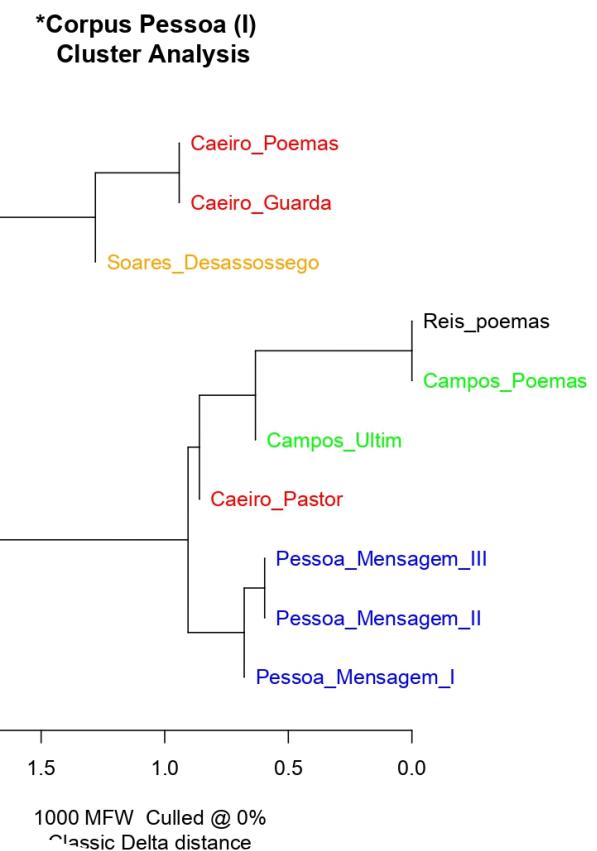
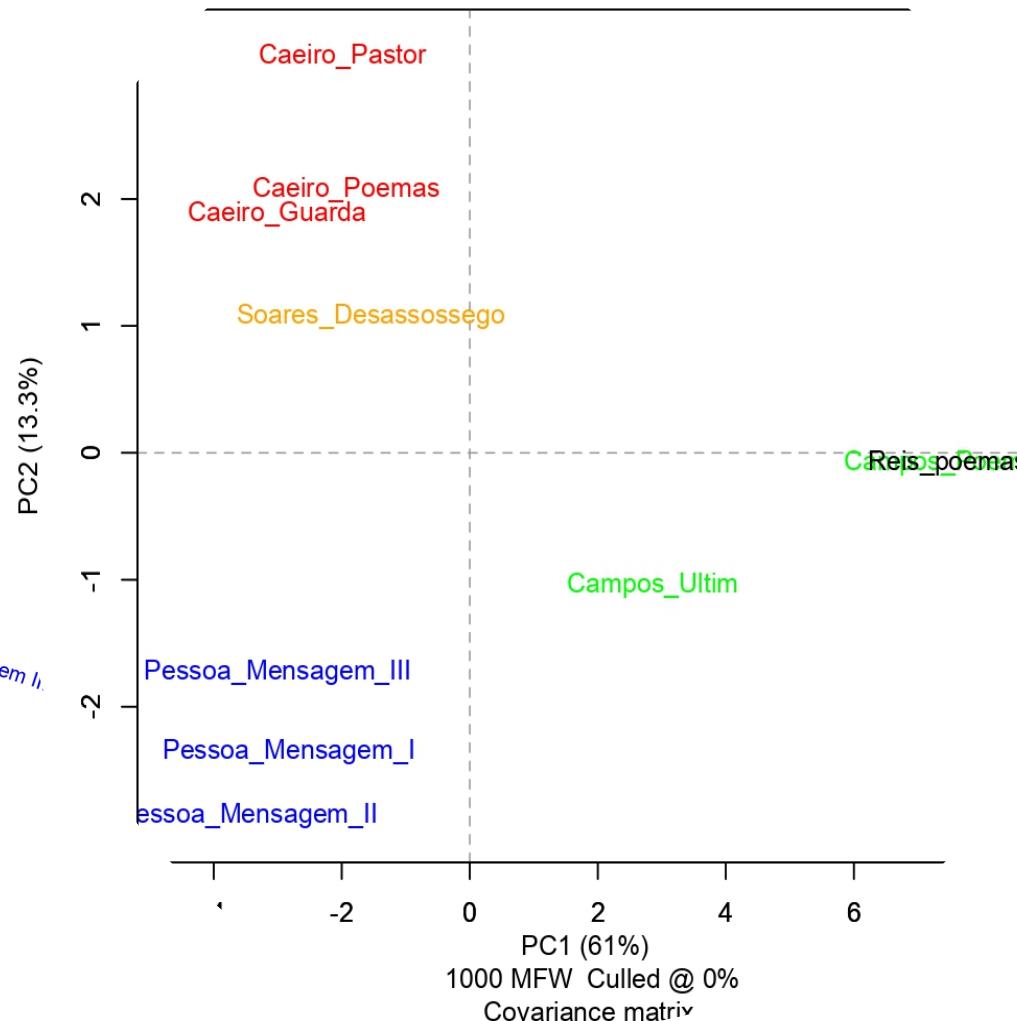
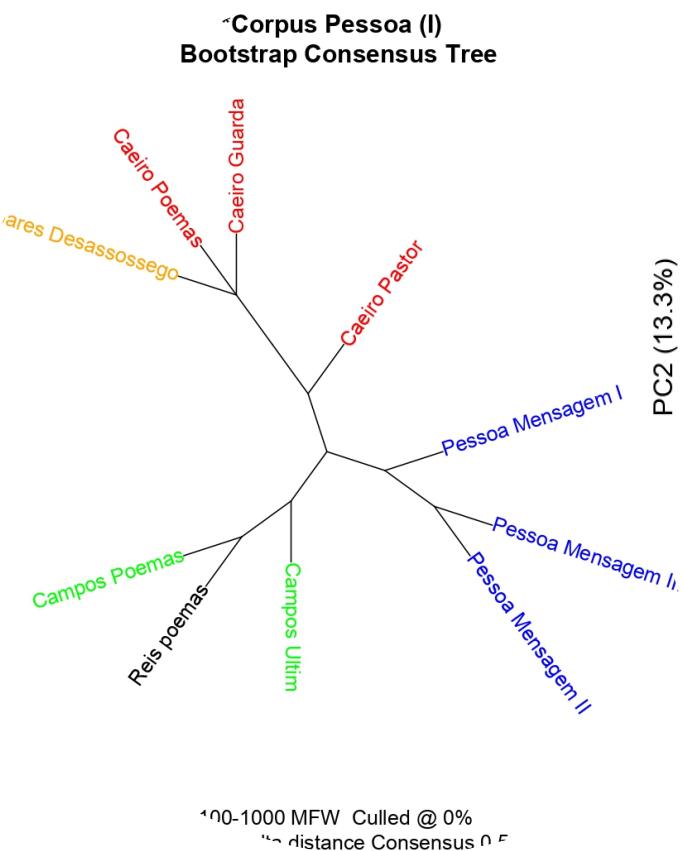
100-1000 MFW Culled @ 0%
Classic Delta distance Consensus 0.5



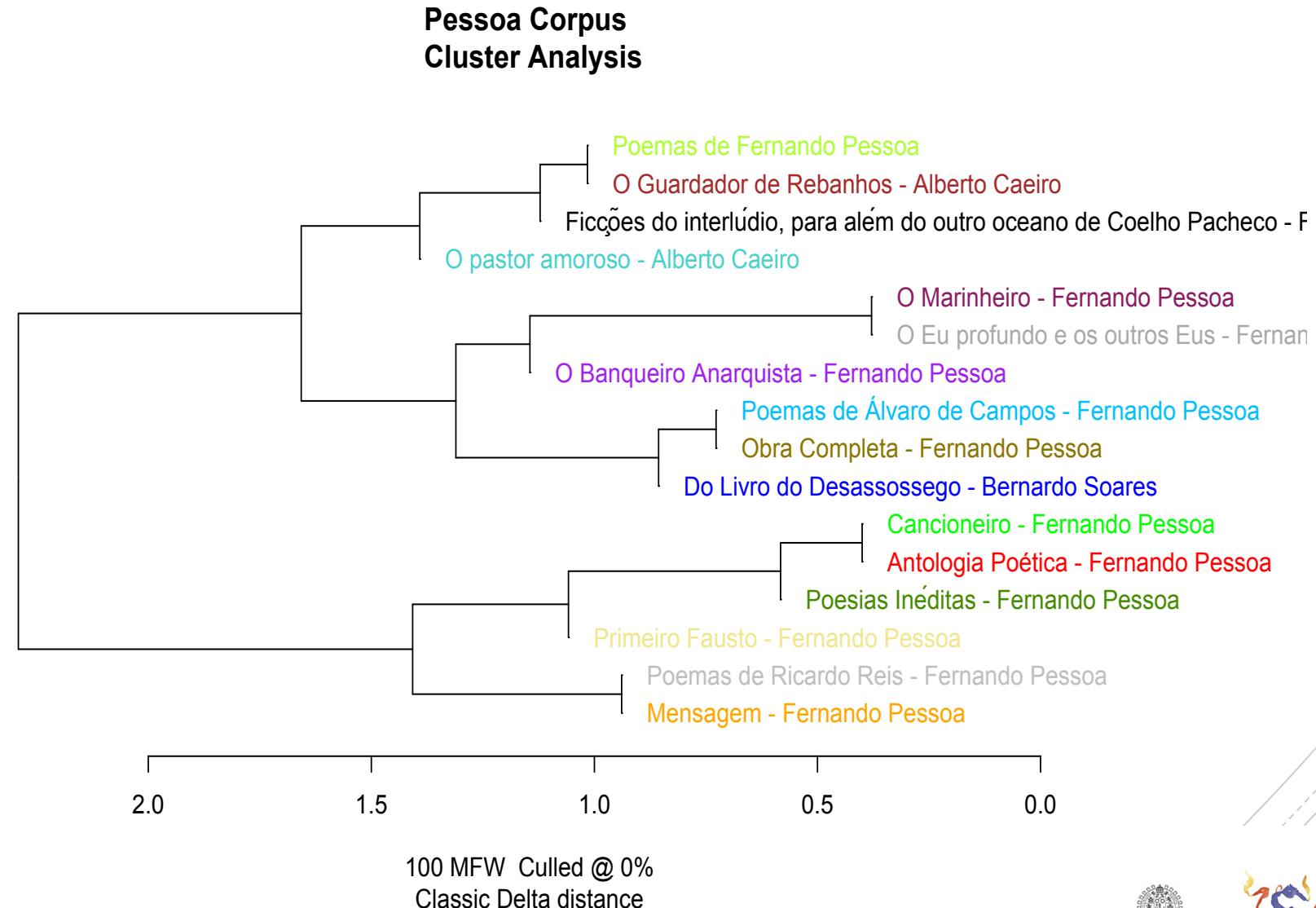
Claudia García-Minguillán
University of Salamanca
2018-2019



*Corpus Pessoa (I) Principal Components Analysis

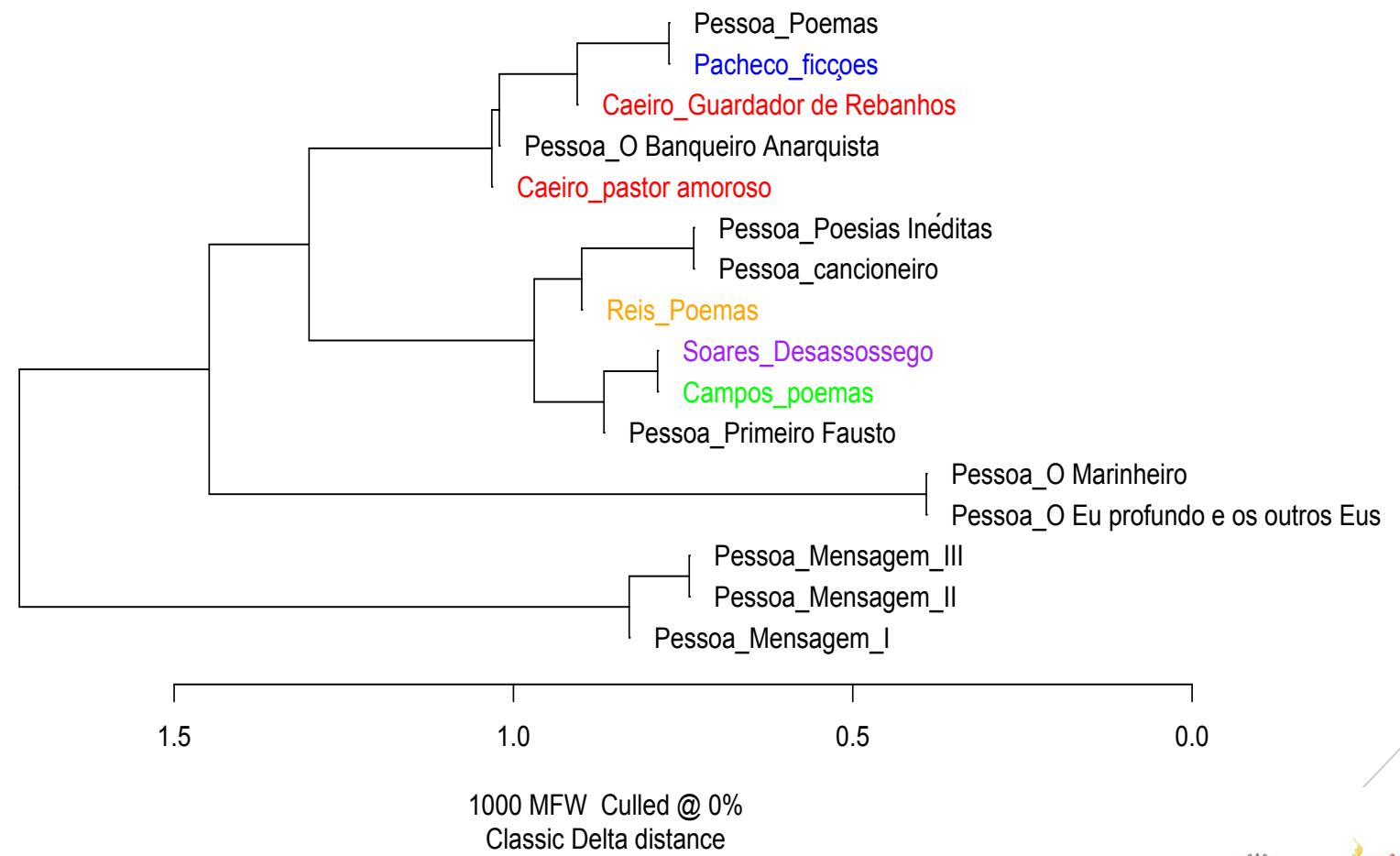


- Antologia Poética_Pessoa.txt
- Caeiro_Guarda.txt
- Caeiro_Pastor.txt
- Caeiro_Poemas.txt
- Campos_Poemas.txt
- Campos_Ultim.txt
- Cancioneiro_Pessoa.txt
- Ficções_Pessoa.txt
- O Banqueiro Anarquista_Pessoa.txt
- O Eu profundo e os outros Eus_Pessoa.txt
- O Marinheiro_Pessoa.txt
- Pessoa_Mensagem_I.txt
- Pessoa_Mensagem_II.txt
- Pessoa_Mensagem_III.txt
- Poesias Inéditas_Pessoa.txt
- Primeiro_Fausto_Pessoa.txt
- Reis_poemas.txt
- Soares_Desas.txt

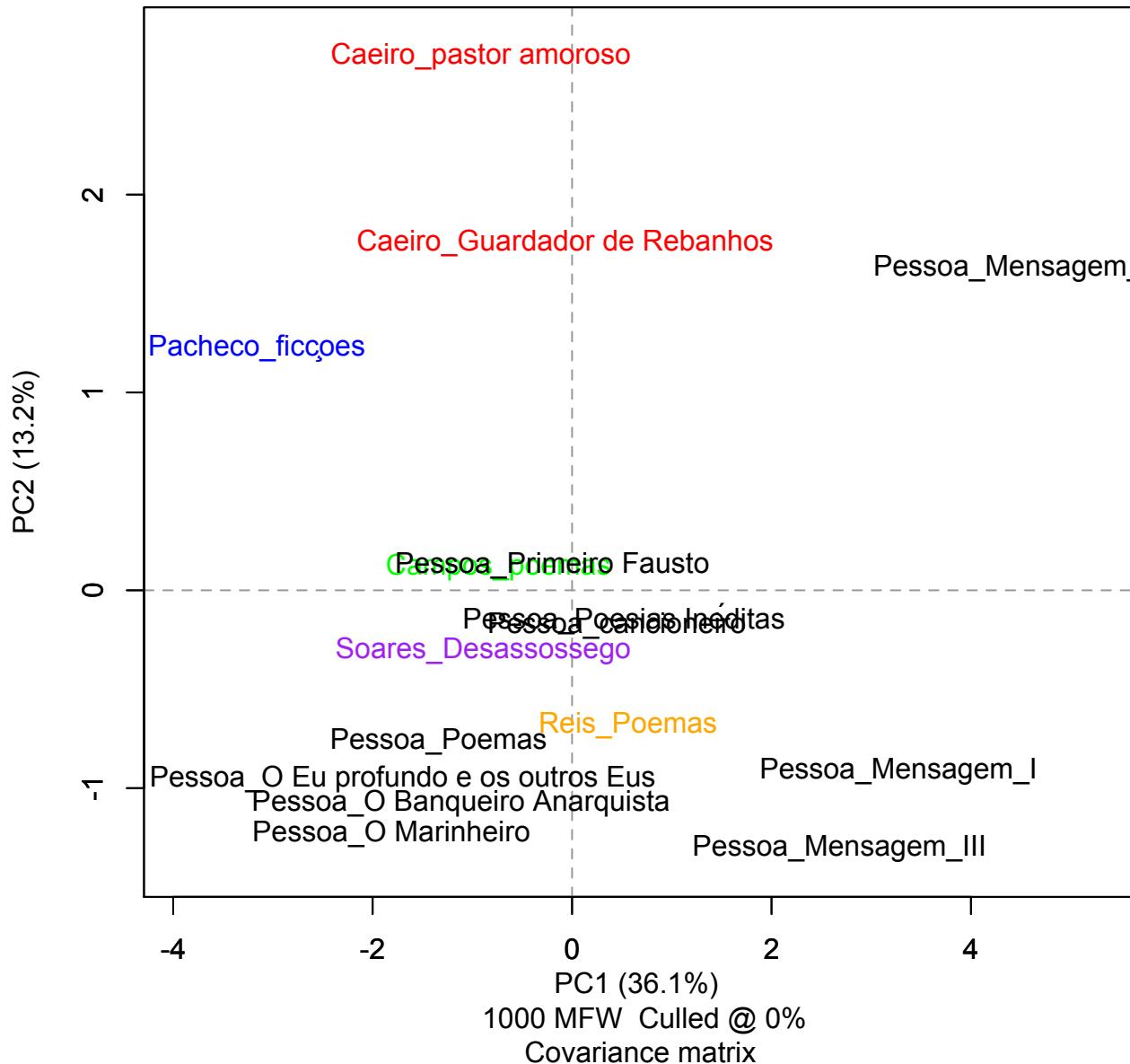


- Caeiro_Guardador de Rebanhos.txt
- Caeiro_pastor amoroso.txt
- Campos_poemas.txt
- Pacheco_ficções.txt
- Pessoa_cancioneiro.txt
- Pessoa_Mensagem_I.txt
- Pessoa_Mensagem_II.txt
- Pessoa_Mensagem_III.txt
- Pessoa_O Banqueiro Anarquista.txt
- Pessoa_O Eu profundo e os outros Eus.txt
- Pessoa_O Marinheiro.txt
- Pessoa_Poemas.txt
- Pessoa_Poesias Inéditas.txt
- Pessoa_Primeiro Fausto.txt
- Reis_Poemas.txt
- Soares_Desassossego.txt

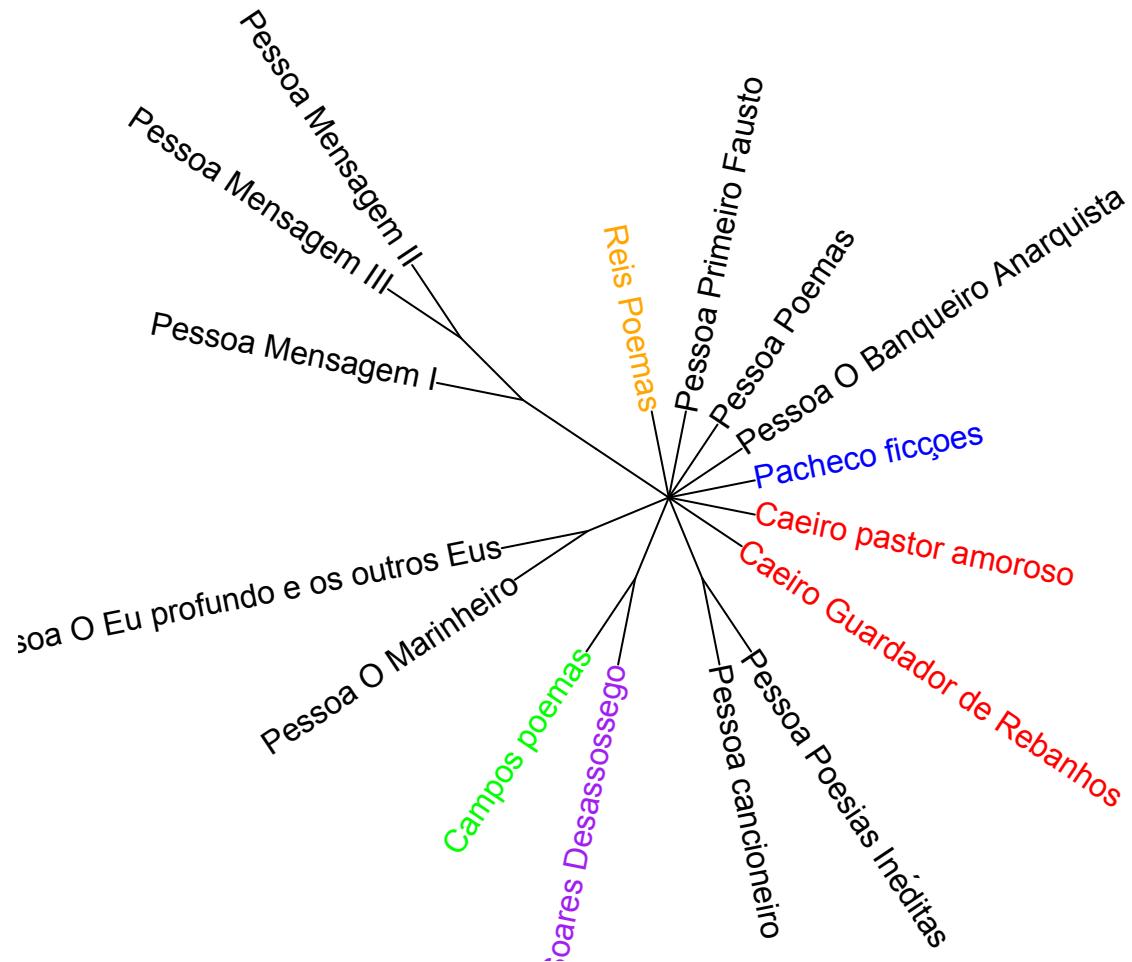
*Corpus_Pessoa(II)
Cluster Analysis



***Corpus_Pessoa(II)**
Principal Components Analysis



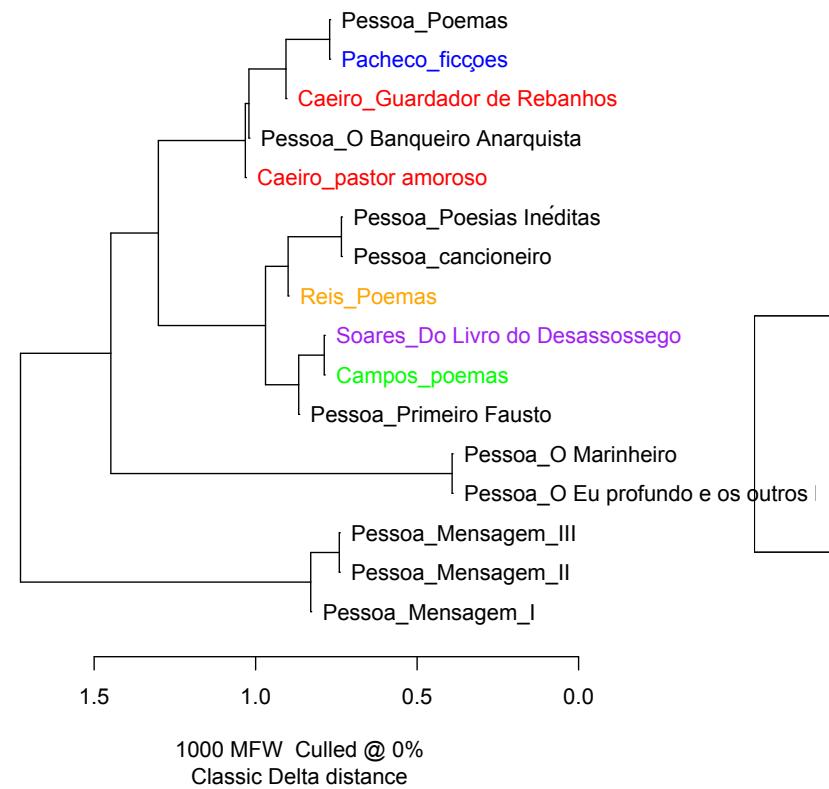
***Corpus_Pessoa(II)**
Bootstrap Consensus Tree



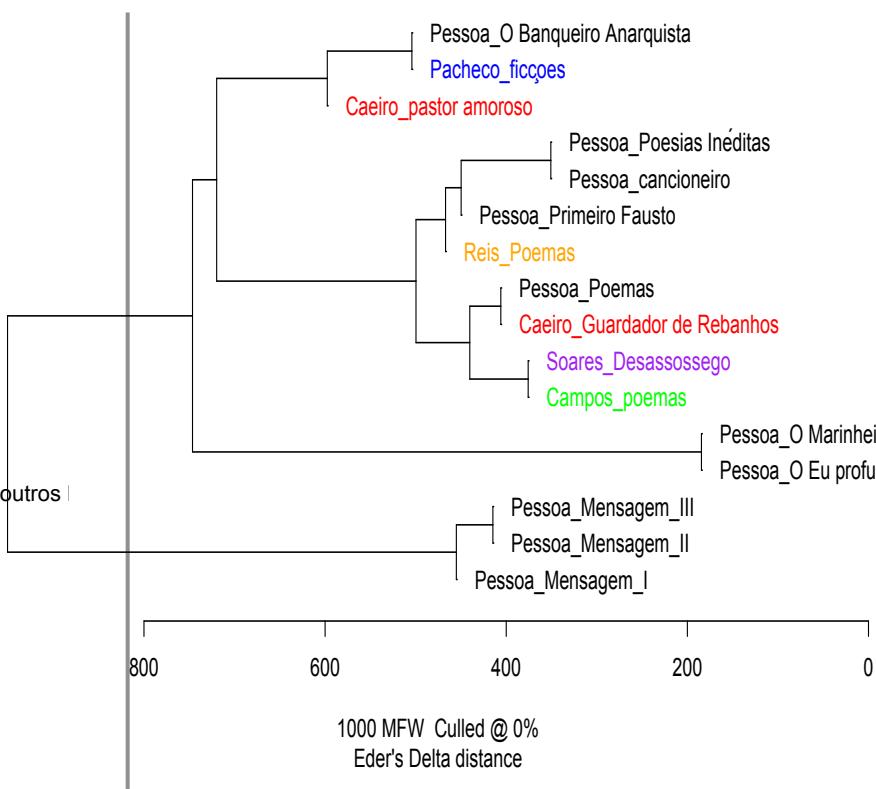
100-1000 MFW Culled @ 0%
Classic Delta distance Consensus 0.5

Claudia García-Minguillán
University of Salamanca
2018-2019

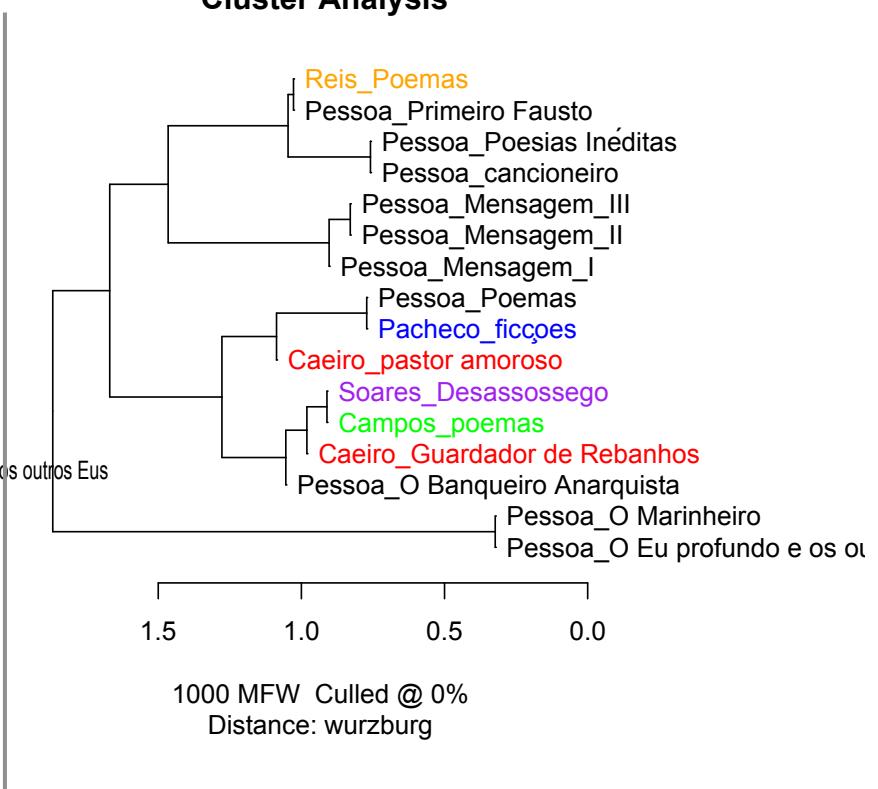
***Corpus_Pessoa(II)**
Cluster Analysis



***Corpus_Pessoa(II)**
Cluster Analysis



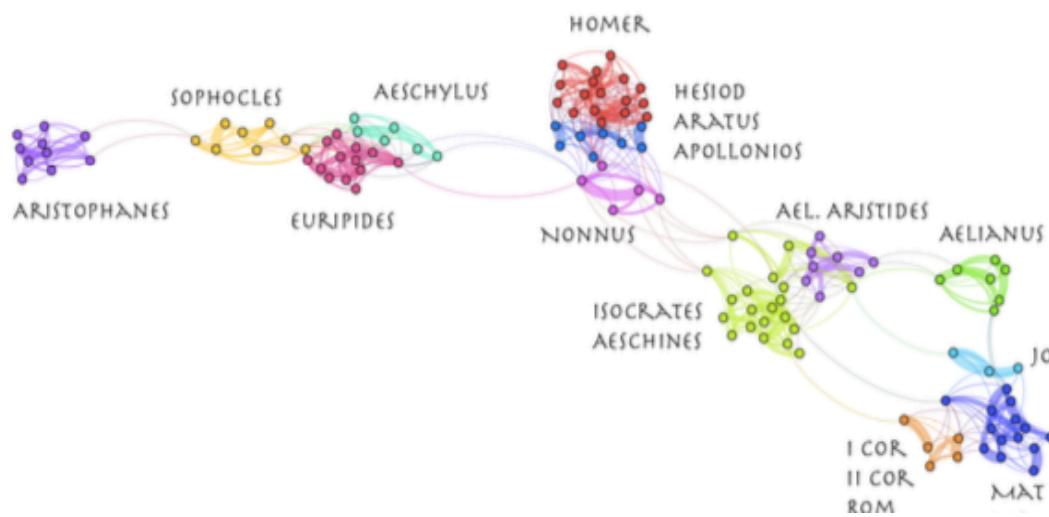
***Corpus_Pessoa(II)**
Cluster Analysis



3. Enjoy visualization

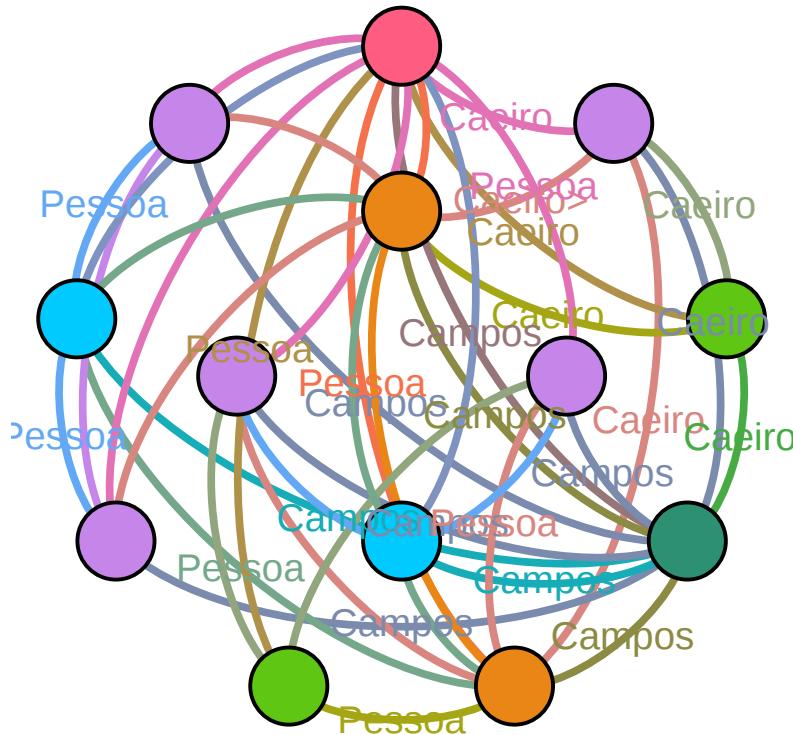
Gephi

visualization
software



Claudia García-Minguillán
University of Salamanca
2018-2019





References

- [Eder, M. ; Rybicki, J. & Kestemont, M. \(2016\): ‘Stylometry with R: A Package for Computational Text Analysis’, *The R Journal*, vol. 8](#)
- Best material for learning ➡ [Computational Stylistics Group](#)
- [‘Culture & Technology’ European Summer University in Digital Humanities, University of Leipzig](#)
- [Digital Humanities Oxford](#)
- [‘How a Computer Program Helped Show J.K. Rowling write A Cuckoo’s Calling’, by Patrick Juola, 08/20/2013](#)
- [Máster Patrimonio Textual y Humanidades Digitales](#)

Contact

cgmt@usal.es

[Claudia Garcia-Minguillan Academia.edu](http://Academia.edu)

[Claudia Garcia-Minguillan Researchgate](http://Researchgate)

