# Stylometry

Giovanni Pietro Vitali – University College Cork

giovannipietrovitali@gmail.com

https://github.com/digitalkoine

https://ucc-ie.academia.edu/GiovanniPietroVitali

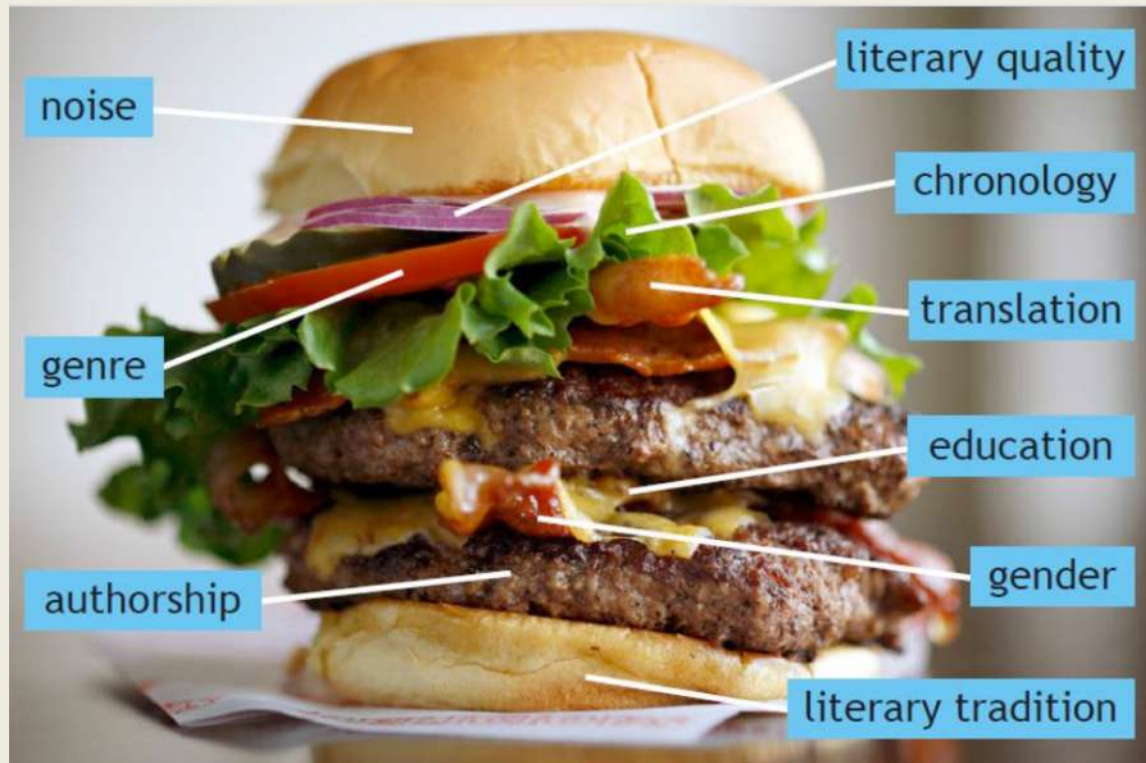# Lorenzo Valla and the *Constantini donatione*



Lorenzo Valla
(c. 1407–1457)

- The **Constantini donatione** was a forged decree where the emperor Constantine I transfers authority of the Roman Empire to the Pope.

- In **De falso credita et ementita Constantini donatione declamatio** (1517), Loreno Valla shows that the act was done in the eighth century by the same papal chancellery. Some grammatical forms could not have been used in the 4th century

- Act of immense value for the history of philology. The first instance of scholarly-based investigation of style

# STYLOmetry



- Author attribution identification of unknown authors
- Genre classification
- Historical study of language change
- Other applications
- Anonymity
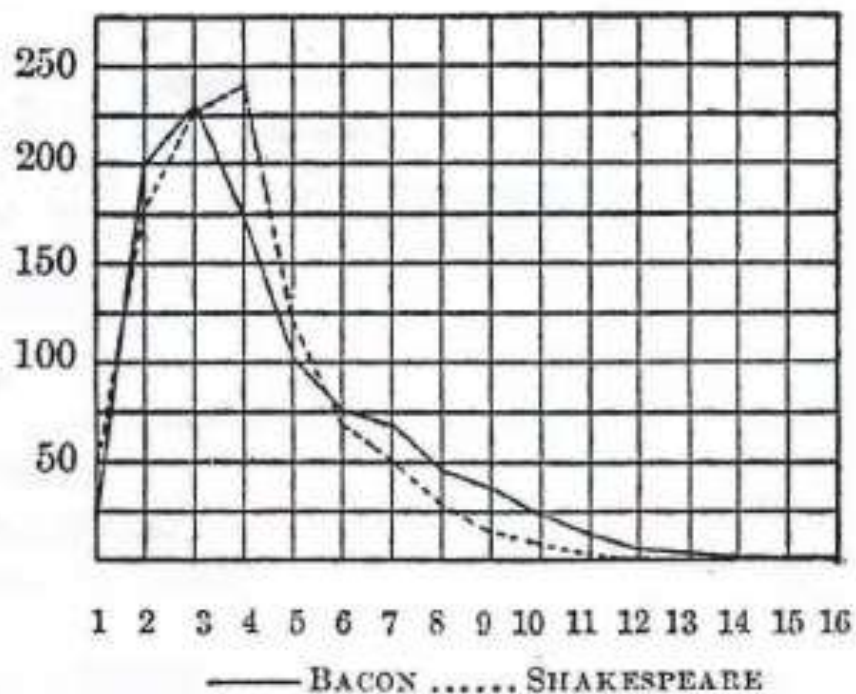- Plagiarism

# «Measuring» authorial style

Thomas Corwin Mendenhall (October 4, 1841 – March 23, 1924) was an American autodidact physicist and meteorologist. He was the first professor hired at The Ohio State University in 1873 and the superintendent of the U.S. Coast and Geodetic Survey from 1889 to 1894. Alongside his work, he was also an advocate for the adoption of the metric system by the United States.

He provided the first empirical evidence in favor of de Morgan's assumptions. In two subsequent studies, Mendenhall (1887, 1901) elaborated on de Morgan's ideas, suggesting that in addition to analy- ses "based simply on mean word-length" (1887: 239), one should attempt to graphically exhibit the peculiarities of style in composition: in order to arrive at such graphics, Mendenhall counted the frequency with which words of a given length occur in 1000-word samples from different authors, among them Francis Bacon, Charles Dickens, William M. Thackerey, and John Stuart Mill.
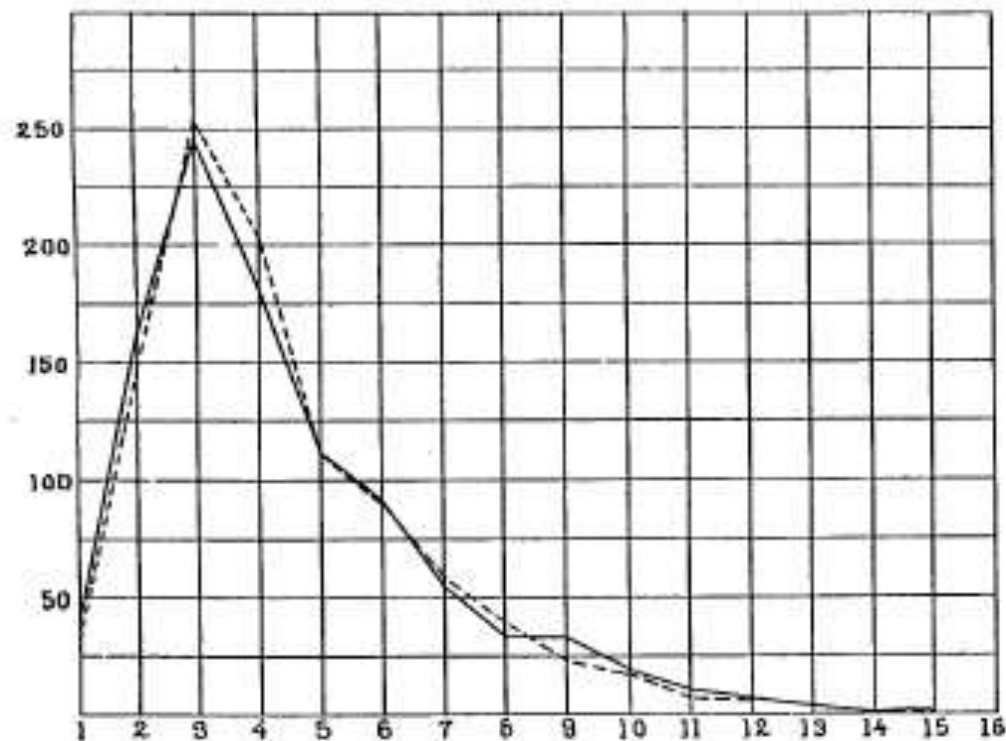
Robert E. Moritz, On The Significance Of Characteristic Curves Of Composition, Popular Science Monthly, volume 65, June 1904.

https://en.wikisource.org/wiki/Popular_Science_Monthly/Volume_65/June_1904/On_the_Significance_of_Characteristic_Curves_of_Composition

# Bacon, Shakespeare, Dickens (Mendenhall)



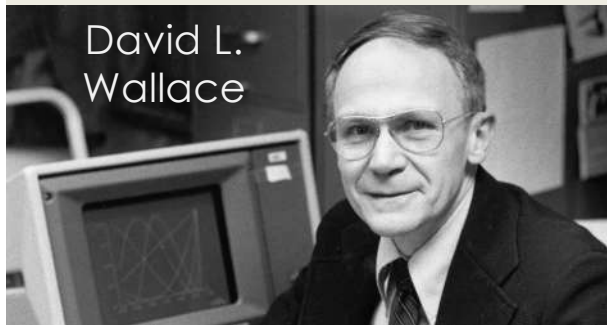Word Length Frequencies in Bacon's and Shakespeare's Texts (Mendenhall 1901)

Word Length Frequencies in Dickens' Oliver Twist (Mendenhall 1887)

# Mosteller and Wallace (1964)
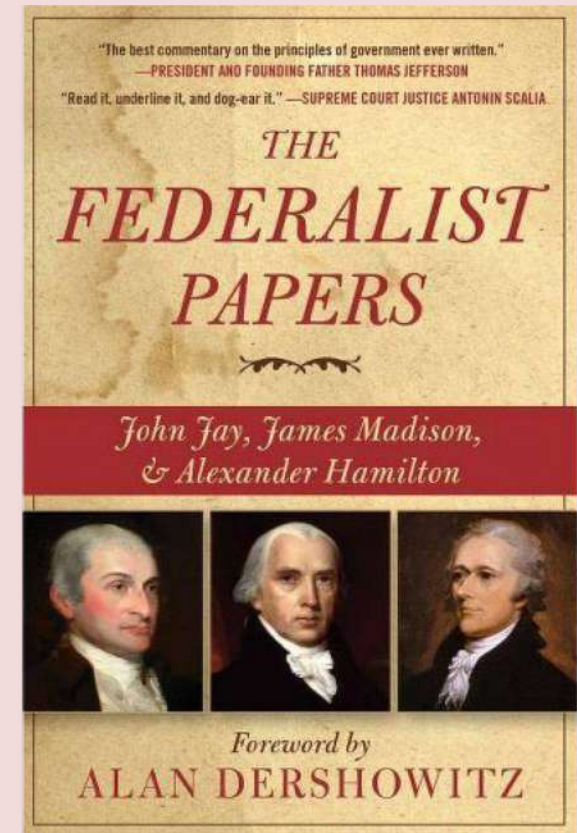
Frederick Mosteller

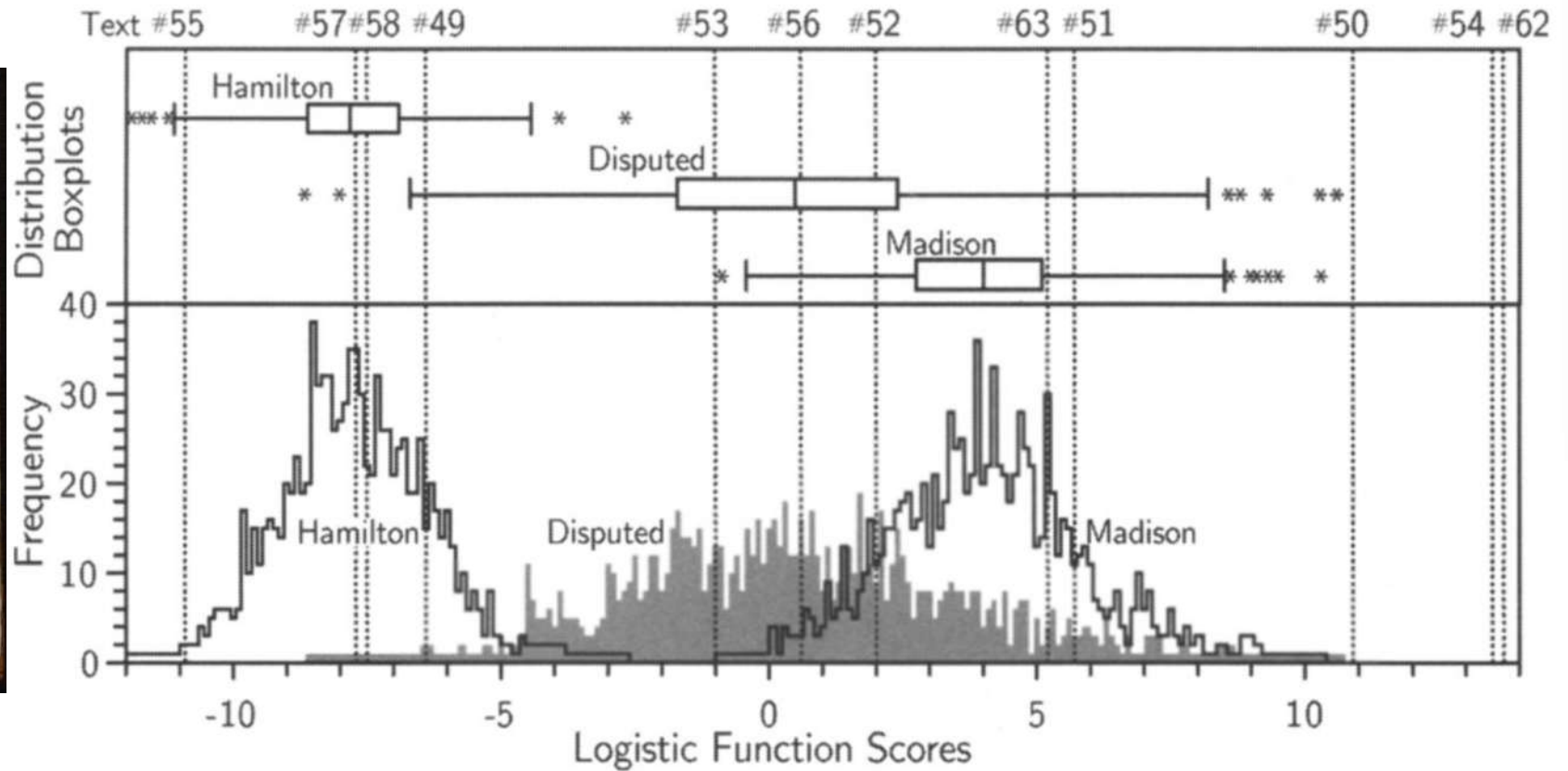David L. Wallace
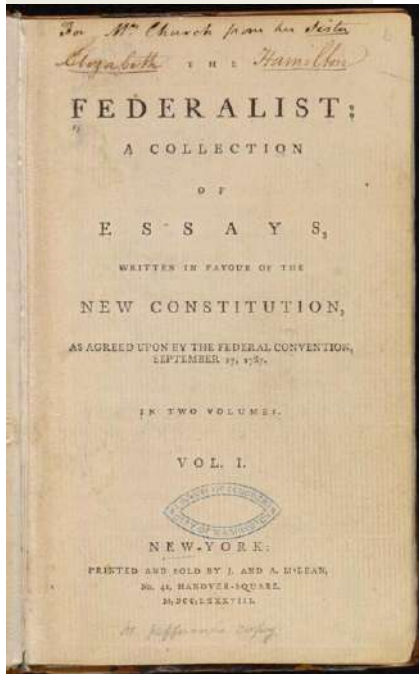
Case study: *The Federalist Papers* (1787-1788)

The Federalist Papers is a collection of 85 articles and essays written by Alexander Hamilton, James Madison, and John Jay under the pseudonym "Publius" to promote the ratification of the United States Constitution. The collection was commonly known as The Federalist until the name The Federalist Papers emerged in the 20th century.

The authors of The Federalist intended to influence the voters to ratify the Constitution.

# Statistical Approach

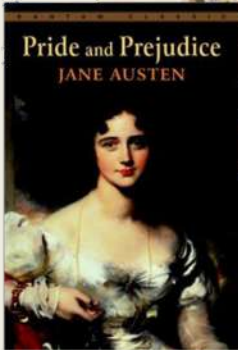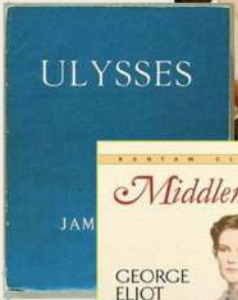|  | enough | while | whilst | upon |
|---|---|---|---|---|
| Hamilton | 0.59 | 0.26 | 0 | 2.93 |
| Madison | 0 | 0 | 0.47 | 0.16 |
| Disputed texts | 0 | 0 | 0.34 | 0.08 |
| Co-authored texts | 0.18 | 0 | 0.36 | 0.36 |

# The Methods For Stylometry And Authorship Attribution

- Character-level analysis
- Syntax-level analysis
- Multi-method analysis (e.g. JGAAP, PAN competition software...)
- ...and many others
- In this lesson, just two methods:
  - Zeta method (for the quantitative analysis of style)
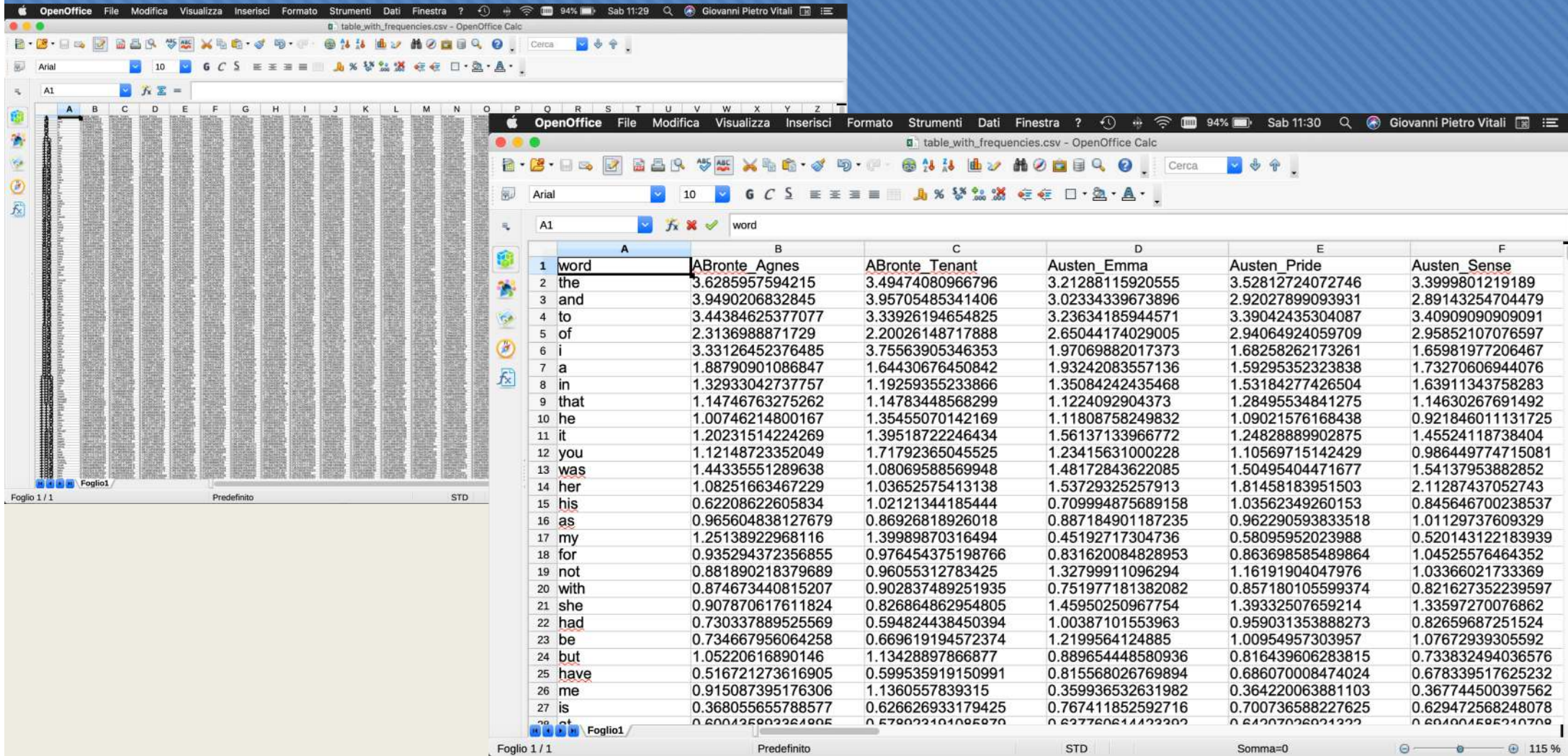  - Delta method (for authorship attribution)
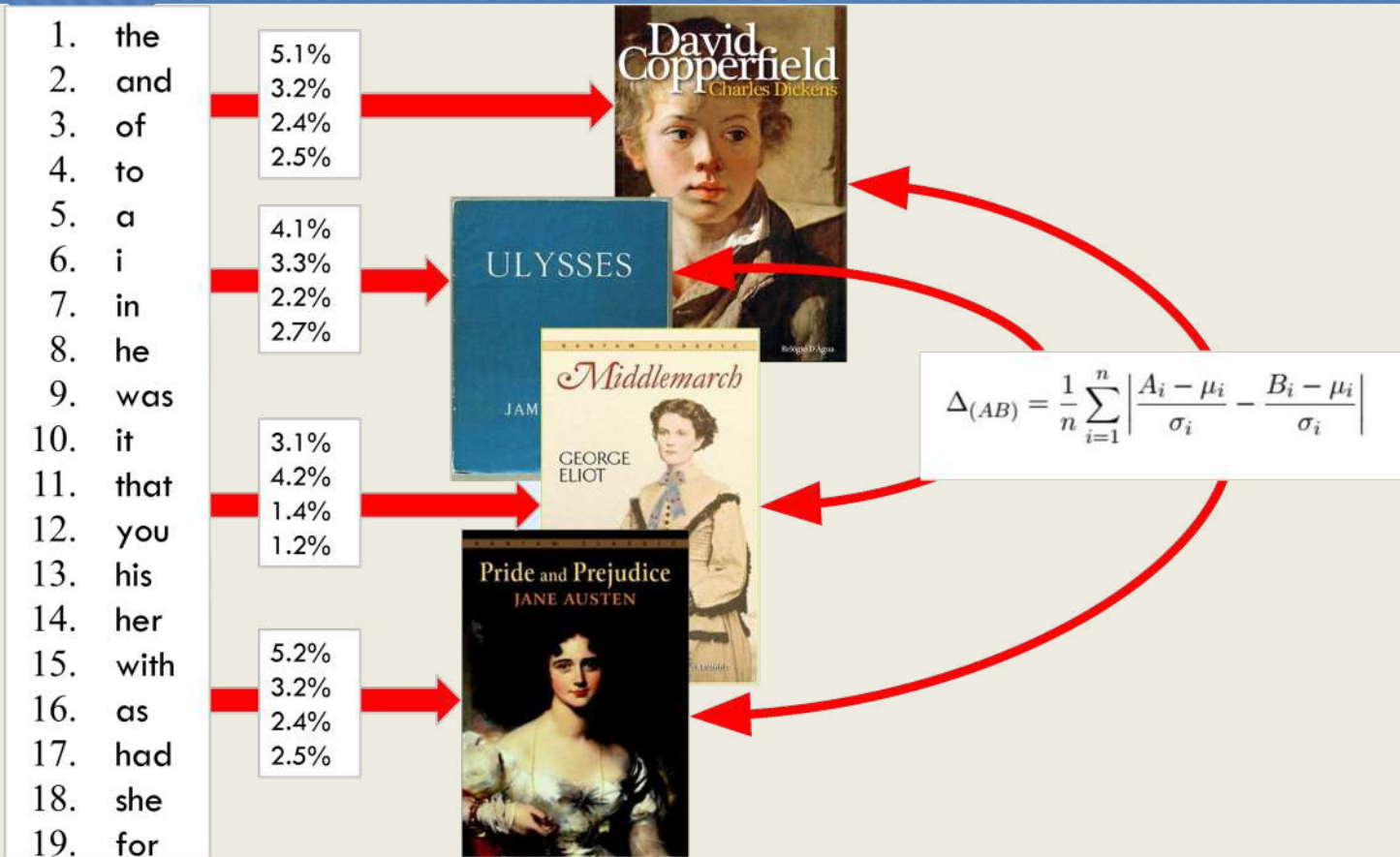
# Let's take an example: English Novels

# WORD-FREQUENCY BASED STYLOMETRY

# WORD-FREQUENCY BASED STYLOMETRY

1. the
2. and
3. of
4. to
5. a
6. i
7. in
8. he
9. was
10. it
11. that
12. you
13. his
14. her
15. with
16. as
17. had
18. she
19. for

5.1%
3.2%
2.4%
2.5%

4.1%
3.3%
2.2%
2.7%

3.1%
4.2%
1.4%
1.2%

5.2%
3.2%
2.4%
2.5%

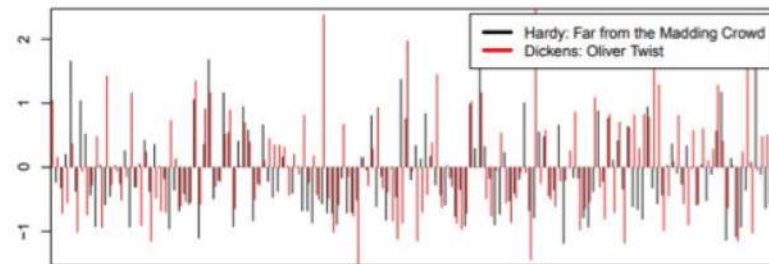$$\Delta_{(AB)} = \frac{1}{n} \sum_{i=1}^{n} \left| \frac{A_i - \mu_i}{\sigma_i} - \frac{B_i - \mu_i}{\sigma_i} \right|$$

# Burrows DELTA

Frequencies of 100 – 5,000 most frequent words (MFW) form a "fingerprint" of an author's style

Standardized to z-scores to give each word equal weight

$$\begin{array}{lcccccccccccc} & the & and & to & of & a & I & in & was & that & he & her \\ z(\text{Madding Crowd}) = ( & .53, & -.23, & -.32, & .20, & 1.66, & -.37, & 1.04, & .52, & -.44, & -.92, & .03, & \ldots) \\ z(\text{Tess of the d'U.}) = ( & .75, & -.48, & -.08, & .51, & -.24, & -.87, & .60, & .41, & -.14, & -.47, & 1.39, & \ldots) \\ z(\text{Oliver Twist}) = ( & 1.05, & .15, & -.71, & -.56, & .37, & -1.01, & -.06, & -.74, & -.28, & .48, & -.94, & \ldots) \end{array}$$

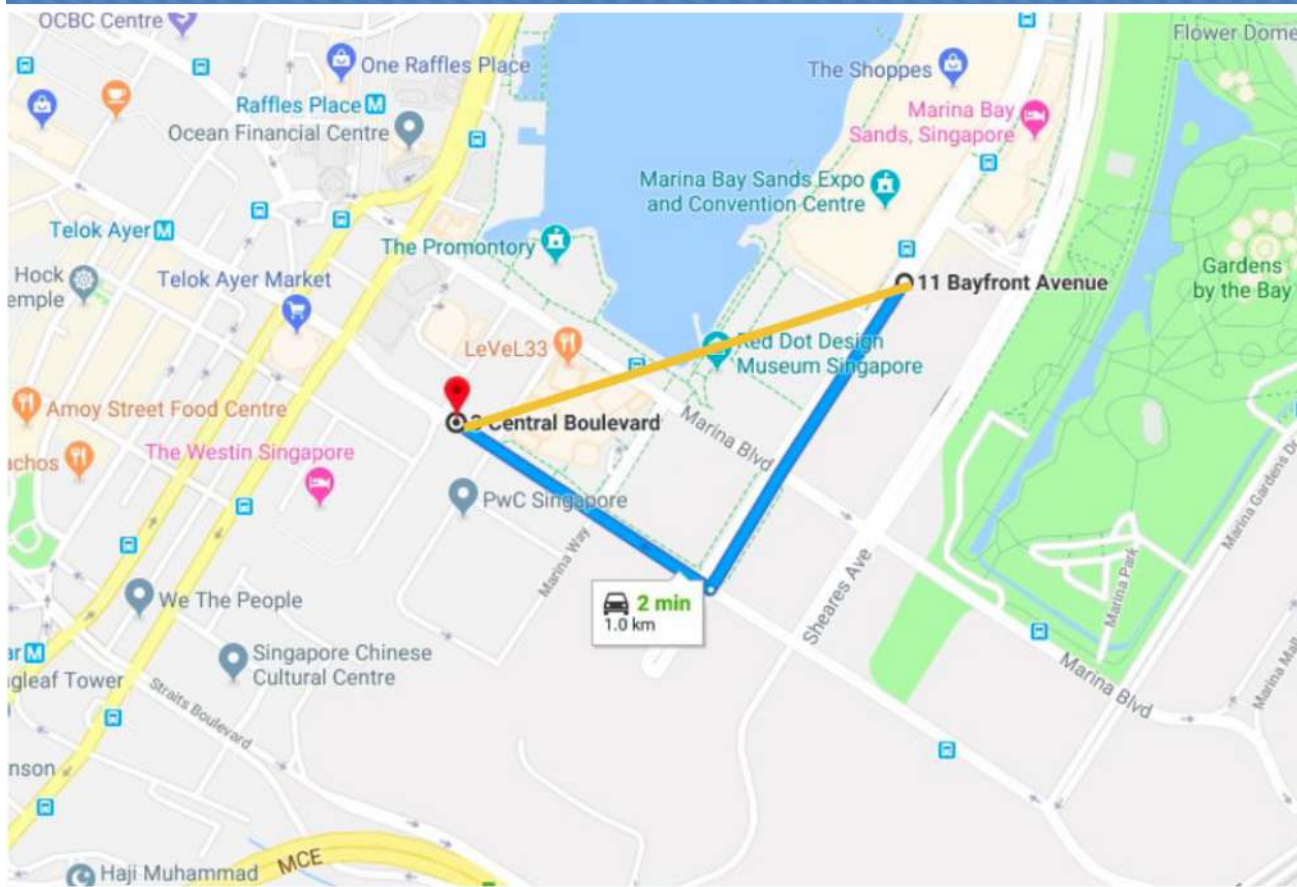— Hardy: Far from the Madding Crowd
— Dickens: Oliver Twist

# Zeta score

$$z = \frac{X - \bar{X}}{S}$$

- X - frequency of term
- Mean(X) - mean frequency of term
- S - standard deviation

# Distances



○ Time
○ Manhattan
○ Euclidean

# DELTA measures

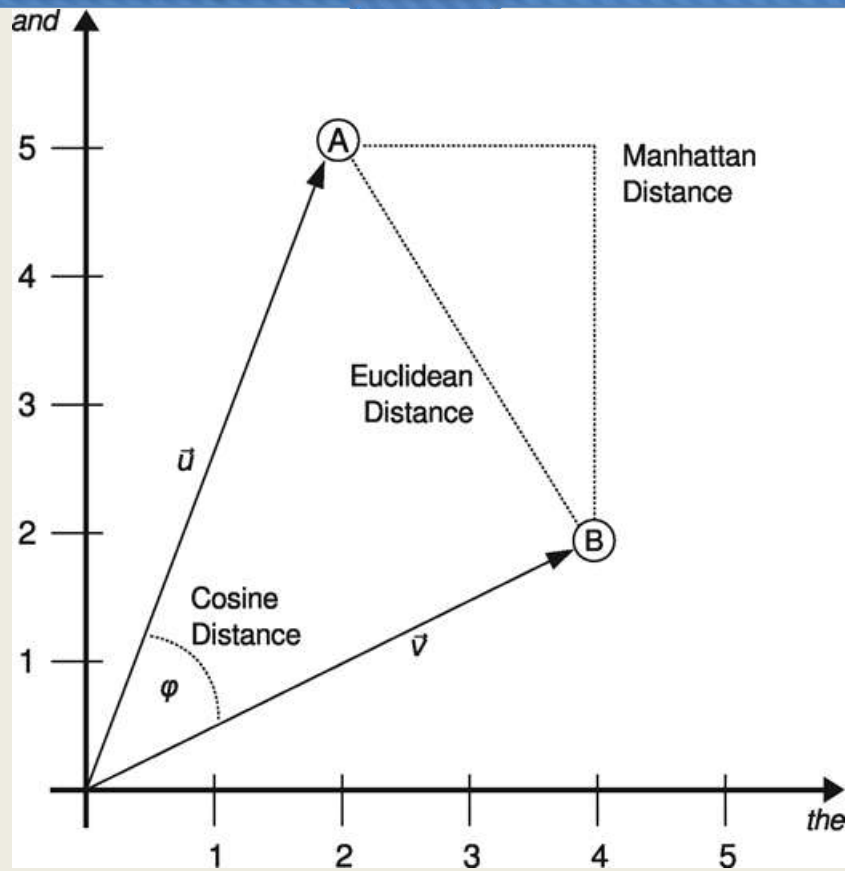Burrows's Delta = Manhattan distance (Burrows 2002)

$$\Delta_B(D, D') = \|\mathbf{z}(D) - \mathbf{z}(D')\|_1 = \sum_{i=1}^{n_w} |z_i(D) - z_i(D')|$$

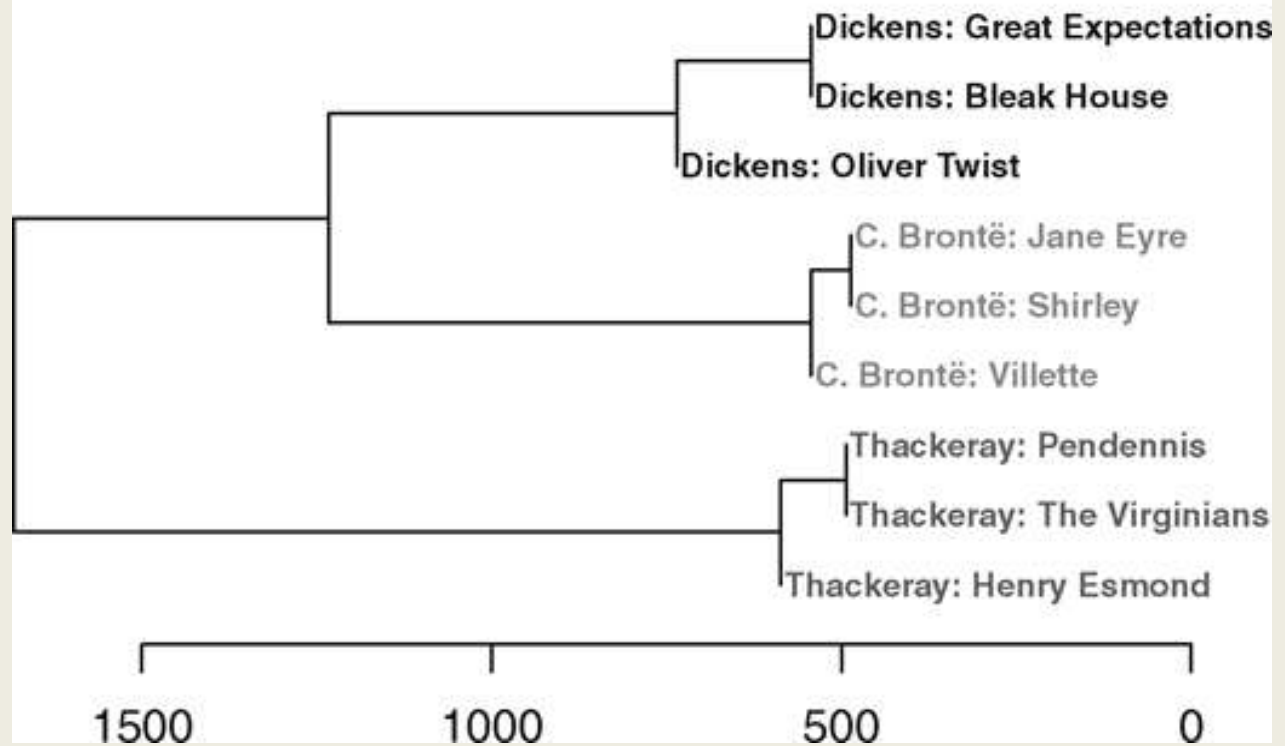Quadratic Delta = Euclidean distance (Argamon 2008)

$$\Delta_Q(D, D') = \|\mathbf{z}(D) - \mathbf{z}(D')\|_2^2 = \sum_{i=1}^{n_w} (z_i(D) - z_i(D'))^2$$

# Deltas & Distance

# Stylo package

https://cran.r-project.org/web/packages/stylo/index.html

https://sites.google.com/site/computationalstylistics/stylo

https://cran.r-project.org/web/packages/stylo/stylo.pdf

https://computationalstylistics.github.io/

10 Computational 01
01 Stylistics 0101000
11 Group 011010110

Stylo needs XQuartz to work on Mac

Command to Start

> install.packages("stylo")
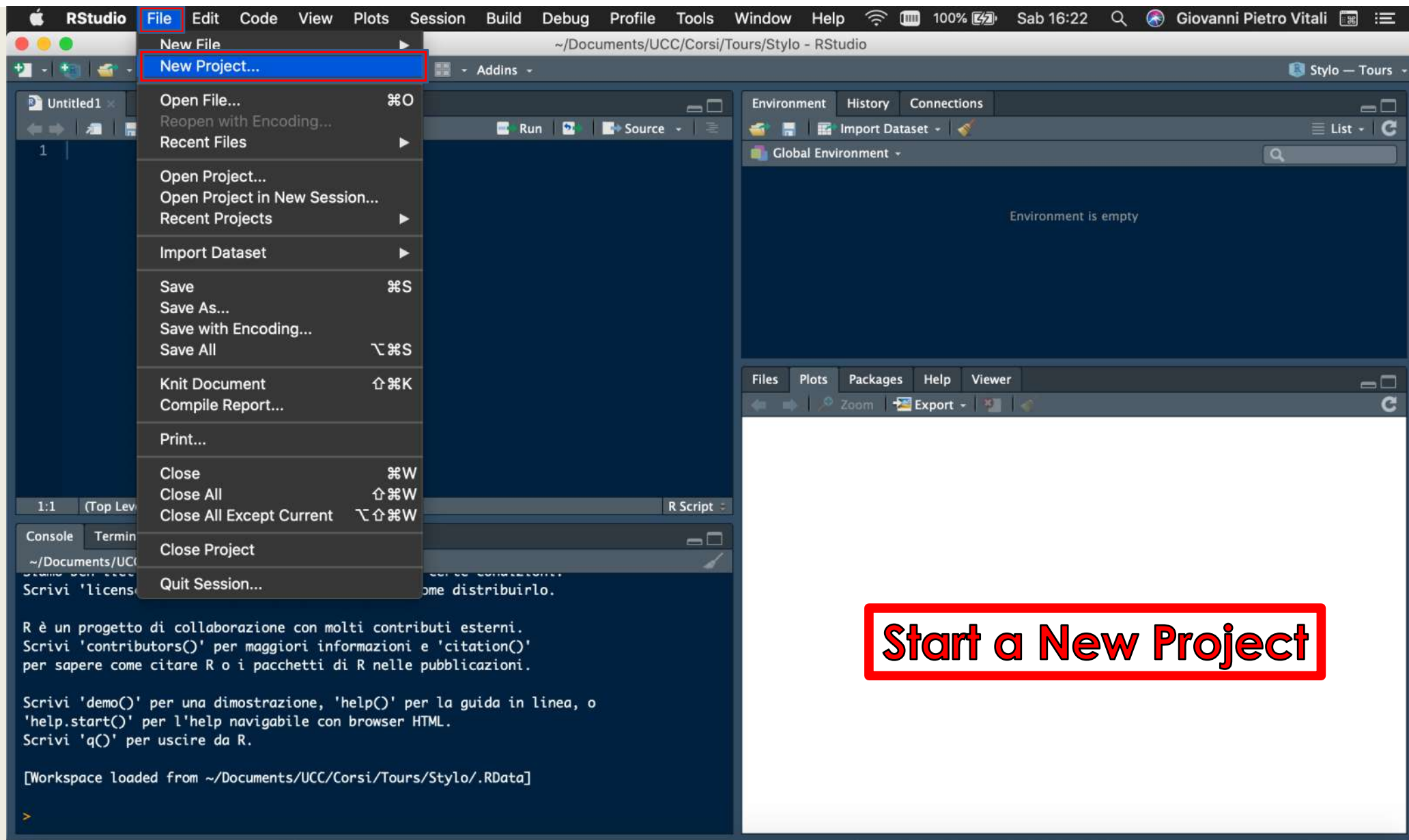
> library(stylo)

> stylo()

# Installing stylo

- run R
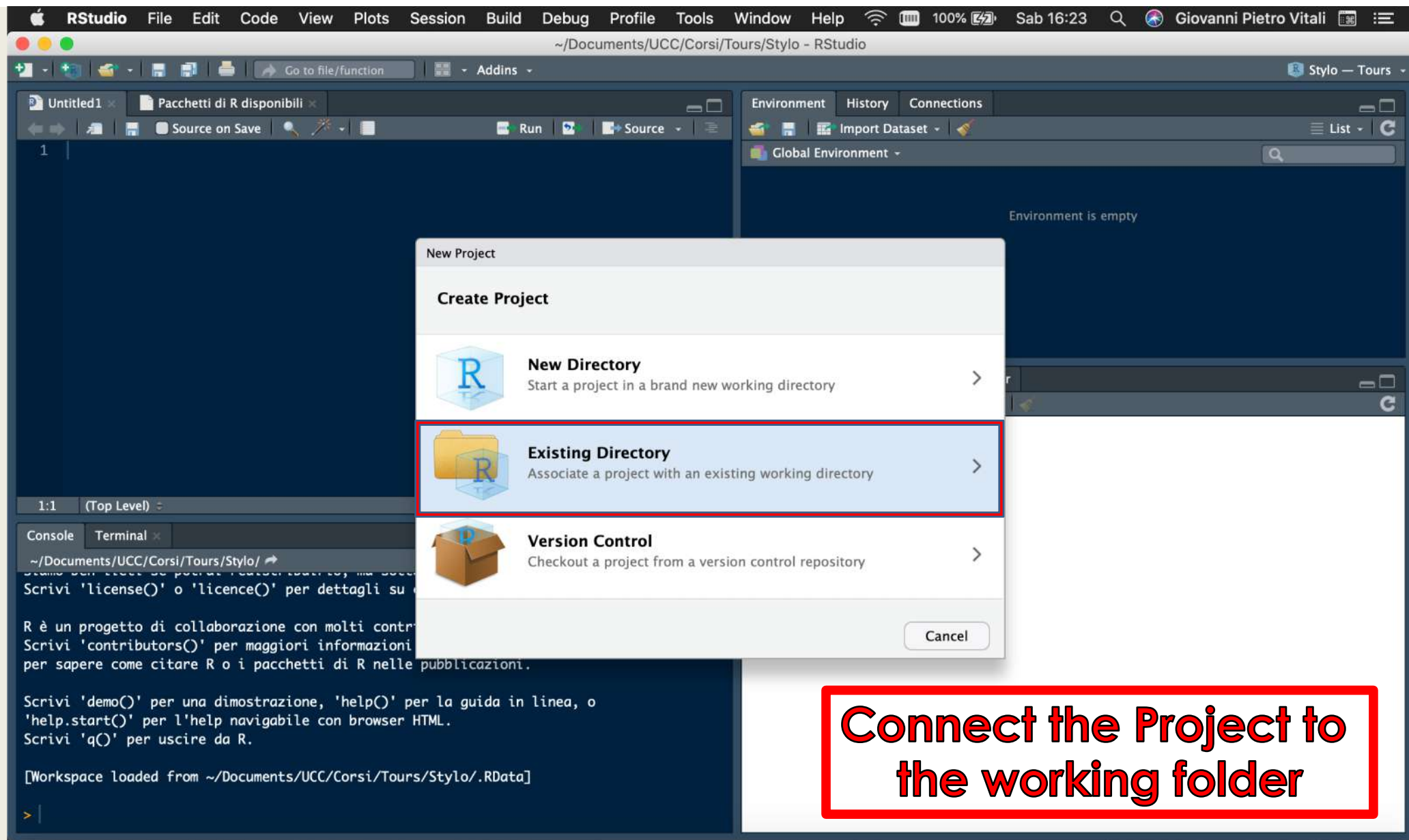- type install.packages("stylo")
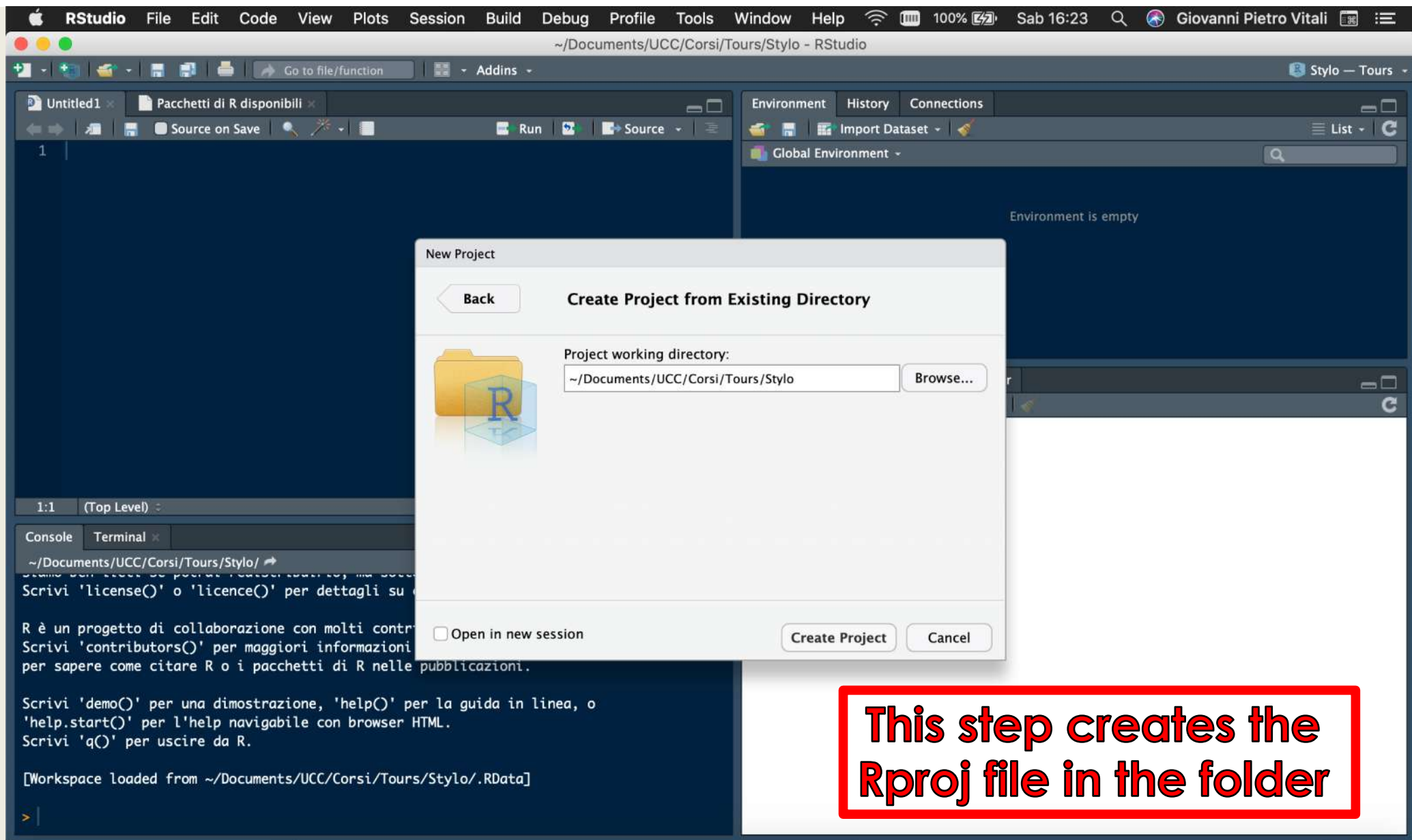- pick your R server
- click OK
- done!

# Some basic R functions

- to activate a package: library(stylo)
- to set working
  directory: setwd("path/to/my/stuff")
- to find your current location: getwd()
- to list files in your current location: list.files()
- to get help: help(<function>), e.g. help(stylo)
- to quit R: q()

# Main functions: stylo()

- It computes distances (differences) between texts, …
- … represented as rows of frequencies of most frequent words.
- Then it plots graphs of those distances:
    - **Cluster Analysis** plots (dendrograms)
    - **Multidimensional Scaling** scatterplots
    - **Principal Components Analysis** scatterplots
    - **Bootstrap Consensus Trees** plots (for multiple parameter settings)
    - **Bootstrap Consensus Networks** (other software will be needed to take over)
- The plots can be both displayed on screen and saved to a file (e.g. PNG).

Start a New Project

Connect the Project to the working folder

This step creates the Rproj file in the folder

# Main functions: stylo()

- It computes distances (differences) between texts, …
- … represented as rows of frequencies of most frequent words.
- Then it plots graphs of those distances:
  - **Cluster Analysis** plots (dendrograms)
  - **Multidimensional Scaling** scatterplots
  - **Principal Components Analysis** scatterplots
  - **Bootstrap Consensus Trees** plots (for multiple parameter settings)
  - **Bootstrap Consensus Networks** (other software will be needed to take over)
- The plots can be both displayed on screen and saved to a file (e.g. PNG).

Input & Language

# Features

○ words: words are used as the unit.

　○ characters: characters are used as the unit.

　○ *n*-gram size: this is where you can specify the

○ value of *n* for your *n*-grams

　○ preserve case: normally, all the words from

○ the input texts are turned into lowercase

| Book Number | The | Big-Data | Analytics | Tree | newbie | book | for | Girl | honest |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | Word Frequency | | | | |
| 1 | 120 | 80 | 60 | 20 | 1 | 5 | 120 | 0 | 0 |
| 2 | 110 | 0 | 0 | 100 | 10 | 20 | 100 | 40 | 10 |
| 3 | 130 | 0 | 0 | 10 | 11 | 30 | 110 | 20 | 10 |
| 4 | 100 | 0 | 0 | 2 | 20 | 40 | 100 | 10 | 100 |
| 5 | 90 | 0 | 0 | 10 | 30 | 20 | 100 | 100 | 40 |

# WHAT IS N-GRAM

# MFV (most-frequent-word) settings

- **Minimum**: this setting determines how many words (or features) from the top of the frequency list will be used

- **Maximum**: this setting determines how many words from the top of the word frequency list for the entire corpus will be used

- **Increment**: defines the value by which the value of Minimum will be increased at each subsequent run of your analysis

- **Start at freq. Rank**: how many words from the top overall frequency rank list to be skipped

# Culling

The culling values specify the degree to which words that do not appear in all the texts of your corpus will be removed. Thus, a culling value of 20 indicates that words that appear in at least 20% of the texts in the corpus will be considered in the analysis. A culling setting of 0 means that no words will be removed.

| Book Number | The | Big Data | Analytic | Tree | newbie | book | for | Girl | honest |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | Word Frequency | | | | |
| 1 | 120 | 80 | 60 | 20 | 1 | 5 | 120 | 0 | 0 |
| 2 | 110 | | 0 | 100 | 10 | 20 | 100 | 40 | 10 |
| 3 | 130 | 0 | 0 | 10 | 11 | 30 | 110 | 20 | 10 |
| 4 | 100 | 0 | 0 | 2 | 20 | 40 | 100 | 10 | 100 |
| 5 | 90 | 0 | 0 | 10 | 30 | 20 | 100 | 100 | 40 |

Sampling

Stylometry with R | stylo | set parameters

| INPUT & LANGUAGE | FEATURES | STATISTICS | SAMPLING | OUTPUT |

Stylo — Tours

GRAPHS:   Onscreen   PDF   JPG   SVG   PNG
☑        ☐     ☐     ☐     ☐

PLOT AREA:   Set default   Plot height   Plot width   Font size   Line width
☐        7        7        10        1

Colors   Grayscale   Black   Titles
◉        ○        ○        ☑

PCA/MDS:   Labels   Points   Both   Margins   Label offset
◉        ○        ○        2        0

PCA FLAVOUR:   Classic   Loadings   Technical   Symbols
◉        ○        ○        ○

VARIOUS:   Horizontal CA tree   Save distance table   Save features   Save frequencies   Dump samples
☑        ☐        ☐        ☐        ☐

OK

Environment   History   Connections
Import Dataset
Global Environment
Environment is empty

Files   Plots   Packages   Help   Viewer

broom        Convert Statistical Analysis Objects into Tidy
             Tibbles
callr        Call R from R

Output

Console   Terminal
~/Documents/UCC/Corsi/Tours/Stylo/

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


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


> stylo()
using current directory...

# Cluster analysis

Builds "tree" based on the most similar texts based on the MFV. It is not robust on the changes of the parameters.

To color the titles of novels of the same color is necessary to name the txt files before the title the name of the author and let it follow by underscore.

**Nome**
- .Rhistory
- .Rproj.user
- ABronte_Agnes.txt
- ABronte_Tenant.txt
- Austen_Emma.txt
- Austen_Pride.txt
- Austen_Sense.txt
- CBronte_Jane.txt
- CBronte_Professor.txt
- CBronte_Villette.txt
- Dickens_Bleak.txt
- Dickens_David.txt
- Dickens_Hard.txt
- EBronte_Wuthering.txt
- Eliot_Adam.txt
- Eliot_Middlemarch.txt
- Eliot_Mill.txt
- Fielding_Joseph.txt
- Fielding_Tom.txt
- Richardson_Clarissa.txt
- Richardson_Pamela.txt
- Sterne_Sentimental.txt
- Sterne_Tristram.txt
- Thackeray_Barry.txt
- Thackeray_Pendennis.txt
- Thackeray_Vanity.txt
- Trollope_Barchester.txt
- Trollope_Phineas.txt
- Trollope_Prime.txt



**Stylo**
**Cluster Analysis**

Eliot_Mill
Eliot_Adam
Eliot_Middlemarch
Dickens_David
Dickens_Bleak
Dickens_Hard
CBronte_Villette
CBronte_Professor
CBronte_Jane
EBronte_Wuthering
ABronte_Tenant
ABronte_Agnes
Sterne_Tristram
Sterne_Sentimental
Thackeray_Vanity
Thackeray_Pendennis
Thackeray_Barry
Austen_Sense
Austen_Pride
Austen_Emma
Fielding_Tom
Fielding_Joseph
Trollope_Prime
Trollope_Phineas
Trollope_Barchester
Richardson_Pamela
Richardson_Clarissa

3.5   3.0   2.5   2.0   1.5   1.0   0.5   0.0

100 MFW  Culled @ 0%
Classic Delta distance

Stylo
Bootstrap Consensus Tree

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

Uses many trees to discover unchanged patterns for different parameters.
It is more robust but harder to interpret

**Consensus Tree**

Stylo
Multidimensional Scaling

Multidimensional scaling is a statistical analysis technique often used to graphically show differences or similarities between elements of a set. It is a generalization of the concept of sorting: starting from a square matrix, containing the "similarity" of each row element with each column element, the multidimensional scaling algorithm assigns to each element a position in an N-dimensional space, with N established a priori..

If N is small enough, this space can be represented with a 3D graph or display. In practice this technique starts with a system with as many dimensions as the elements of the system, and reduces the dimensions up to a certain number N. In doing this then there is an inevitable loss of information (loss) and there are therefore different algorithms to do multidimensional scaling, which are better suited to different situations of use: in particular we distinguish between metric and non-metric algorithms.

# Stylometry

**Stylo**
**Principal Components Analysis**

# Functions: stylo.network()

**https://prismatic.phon.ox.ac.uk/index_network.html**

○ It is an extended version of the function stylo().

○ It performs Bootstrap Consensus Networks, or a network-like generalization of the Bootstrap Consensus Trees method.

○ It produces interactive visualizations in a web browser: to make it happen, you have to install an additional R package first.
Type: install.packages("networkD3")

Live Slides web content

To view

**Download the add-in.**
liveslides.com/download

**Start the presentation.**

# Main functions: classify()

- It trains a model for pre-defined groups of texts, e.g. authors.
- Then it computes distances (differences) between texts, …
- … represented as rows of frequencies of most frequent words.
- Finally, it compares the trained models with test texts, using:
  - **Delta** classifier (lazy learner introduced by Burrows)
  - **k-NN** classifier (lazy learner relying on >1 neighbors)
  - **Suppor Vector Machines**, a high-performance non-probabilistic classifier
  - **Naive Bayes**, a classical yet slightly outdated classifier
  - **Nearest Shrunken Centroids**, a classifier for high-dimensional datasets
- A final report of the classifier's performance is outputted.

# Main functions: oppose()

- Designed to compare two (groups of) texts
- It cuts input texts into equal-sized samples
- Finds words characteristic for two (groups) of texts
  - These can be reused with stylo() or classify()
- Produces a diagram of the use of each group's words

# Running oppose()

- Different subfolder structure:
  - primary_set
  - secondary_set
  - test_set (optional)
- Running the function:
  - library(stylo)
  - oppose()
- What we get:
  - words_preferred.txt characteristic for the primary_set texts
  - words_avoided.txt characteristic for the secondary_set texts
- word frequency graph

# oppose() parameters



- Slice length: size (in words) of the samples (5000)
- Slice overlap: (0)
- Method: (Craig's Zeta)
- Visualization: type of graph (Markers)

# oppose() parameters

Most of the parameters for this somewhat underdeveloped function are not on GUI. You can switch them on as command line parameters

- when your corpus contains non-Latin diacritics:
  - oppose(encoding = "UTF-8", corpus.lang = "Spanish")

# Men s Women - Words

# Men Vs Women - Markers

# Functions: rolling.classify()

- Looks for traces of authors in a co-authored text…

- … by sliding through this text sequentially in order to detect peculiarities.

- Produces a graph of the respective strengths of these traces.

# Le Roman de la Rose



Le Roman de la Rose (English: The Romance of the Rose) is a medieval French poem styled as an allegorical dream vision. It is a notable instance of courtly literature. The work's stated purpose is both to entertain and to teach others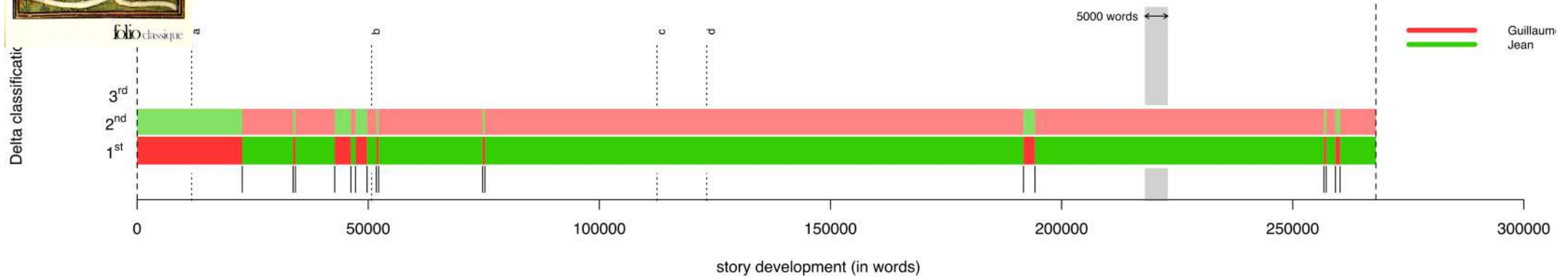 about the art of romantic love. Throughout the poem, Rose is used both as the name of the titular lady and as a symbol of female sexuality. The other characters' names also function both as regular names  as abstractions illustrating the various factors that are involved in a love affair.

The poem was written in two stages. The first 4,058 lines, written by Guillaume de Lorris circa 1230, describe the attempts of a courtier to woo his beloved. This part of the story is set in a walled garden (a locus amoenus), a traditional literary topos in epic and chivalric literature. Around 1275, Jean de Meun composed an additional 17,724 lines. In this enormous coda, allegorical personages (Reason, Genius, and so on) hold forth on love.

# Bibliography

BURROWS J. (2002), "'Delta': a Measure of Stylistic Difference and a Guide to Likely Authorship," in Literary and Linguistic Computing, XVII: 3, pp. 267-287.

EDER M. (2015), "Does size matter? Authorship attribution, small samples, big problem," in Digital Scholarship in the Humanities, XXX: 2, pp. 167-182.

EDER M. (2017), "Visualization in stylometry: Cluster analysis using networks," in Digital Scholarship in the Humanities, XXXII: 1, pp. 50-64.
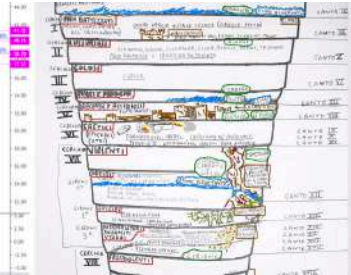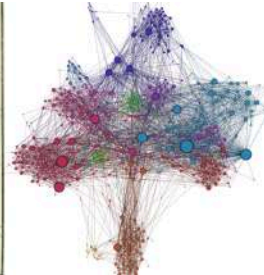
EVERT S., PROISL T., JANNIDIS F., REGER I., PIELSTRÖM S., SCHÖCH C., and VITT T. (2017), "Understanding and explaining Delta measures for authorship attribution" in Digital Scholarship in the Humanities, XXXII: suppl_2, pp. ii4-ii16.

MENDENHALL T. C. (1887), "The characteristic curves of composition," in Science, IX: 214, pp. 237-249.

MORTON A. Q. (1978), Literary Detection: How to Prove Authorship and Fraud in Literature and Documents, Scribner, New York.

MOSTELLER F. and WALLACE D. L. (1964), Inference and Disputed Authorship: The Federalist, Addison-Wesley, Reading Mass.

# Stylometry                                    _end

Giovanni Pietro Vitali – University College Cork

giovannipietrovitali@gmail.com

https://github.com/digitalkoine

https://ucc-ie.academia.edu/GiovanniPietroVitali