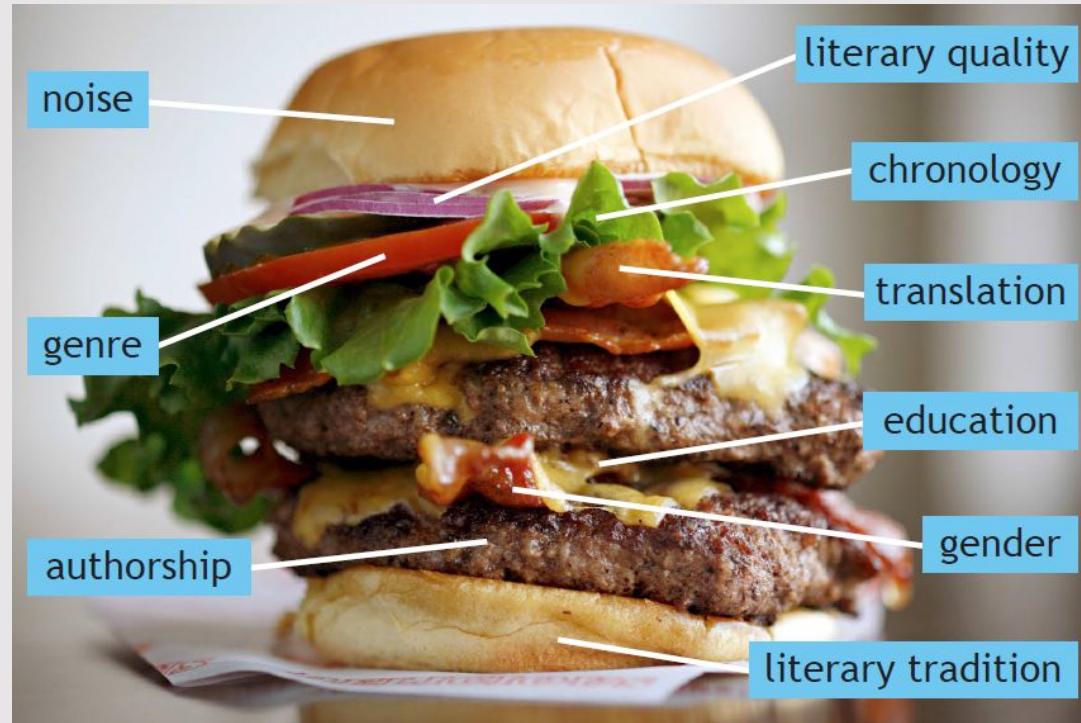


STYLOmetry



The art of measuring style

Agenda

- Software installation
- Why to analyze text?
- History of stylometry?
- What is STYLO?
- STYLO()
- CLASSIFY()
- ROLLING.CLASSIFY()
- Way Forward

About me

- Data science with math background
- Hiking
- Psychology
- People
- yoga



LinkedIn: <https://www.linkedin.com/in/janidziak/>

Count number of F's

FINISHED FILES ARE THE RESULT OF
YEARS OF SCIENTIFIC STUDY COMBINED
WITH THE EXPERIENCE OF YEARS.

Count number of F's

FINISHED FILES ARE THE RESULT OF
YEARS OF SCIENTIFIC STUDY COMBINED
WITH THE EXPERIENCE OF YEARS.

Environment preparation



Install R

<https://cran.r-project.org/bin/windows/base>

R-3.5.1 for Windows (3)

[Download R 3.5.1 for Windows](#) (62 megabytes, 32/64 bit)

[Installation and other instructions](#)
[New features in this version](#)

If you want to double-check that the package you have downloaded matches the package distributed by CRAN, you can compare the [md5sum](#) of the .exe to are available.

Frequently asked questions

- [Does R run under my version of Windows?](#)
- [How do I update packages in my previous version of R?](#)
- [Should I run 32-bit or 64-bit R?](#)

Please see the [R FAQ](#) for general information about R and the [R Windows FAQ](#) for Windows-specific information.

Other builds

- Patches to this release are incorporated in the [r-patched snapshot build](#).
- A build of the development version (which will eventually become the next major release of R) is available in the [r-devel snapshot build](#).
- [Previous releases](#)

Note to webmasters: A stable link which will redirect to the current Windows binary release is
<<CRAN MIRROR>/bin/windows/base/release.htm>.

Last change: 2018-07-02

Install RStudio

www.rstudio.com/products/rstudio/download



Products

Resources

Pricing

About Us

Blogs



Choose Your Version of RStudio

RStudio is a set of integrated tools designed to help you be more productive with R. It includes a console, syntax-highlighting editor that supports direct code execution, and a variety of robust tools for plotting, viewing history, debugging and managing your workspace. [Learn More about RStudio features.](#)



RStudio Desktop
Open Source License

FREE

RStudio Desktop
Commercial License

\$995 per year

RStudio Server
Open Source License

FREE

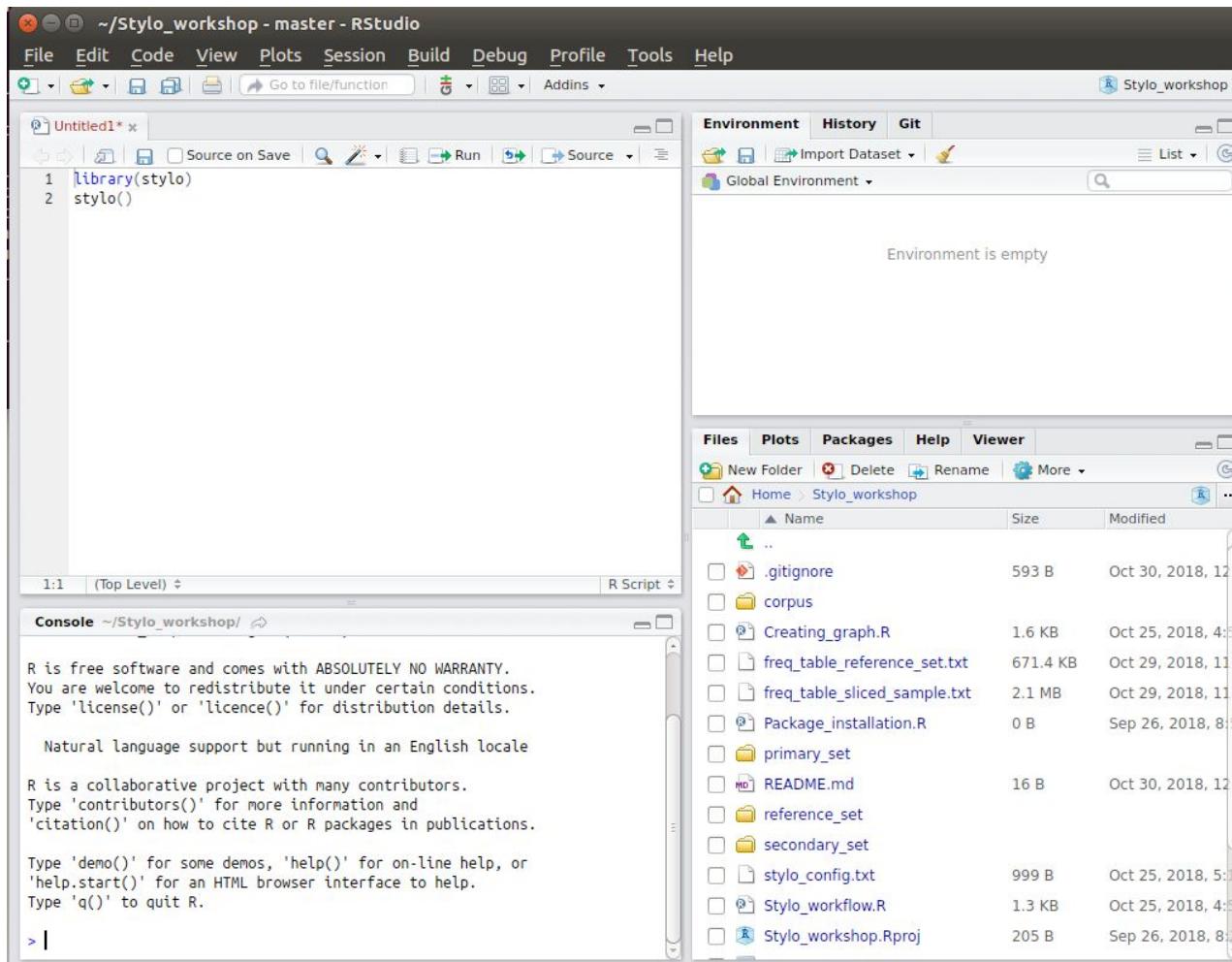
RStudio Server Pro
Commercial License

\$9,995 per year

RStudio Server Pro +
RStudio Connect
Commercial License

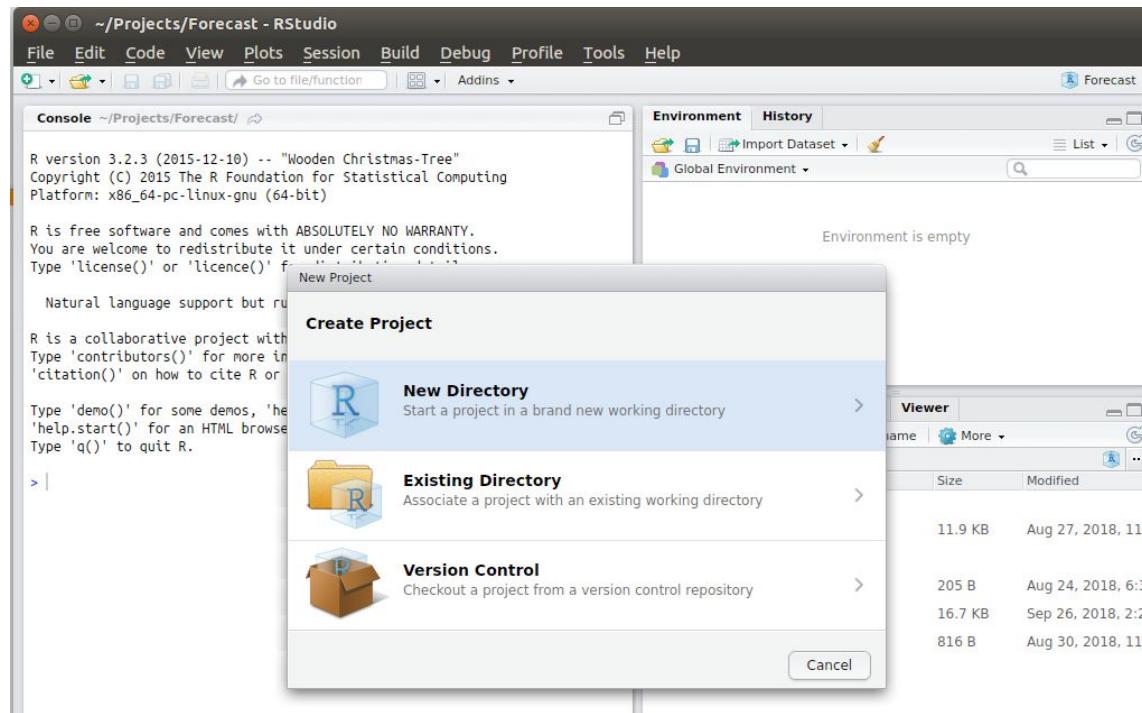
\$29,995 per
year

RStudio



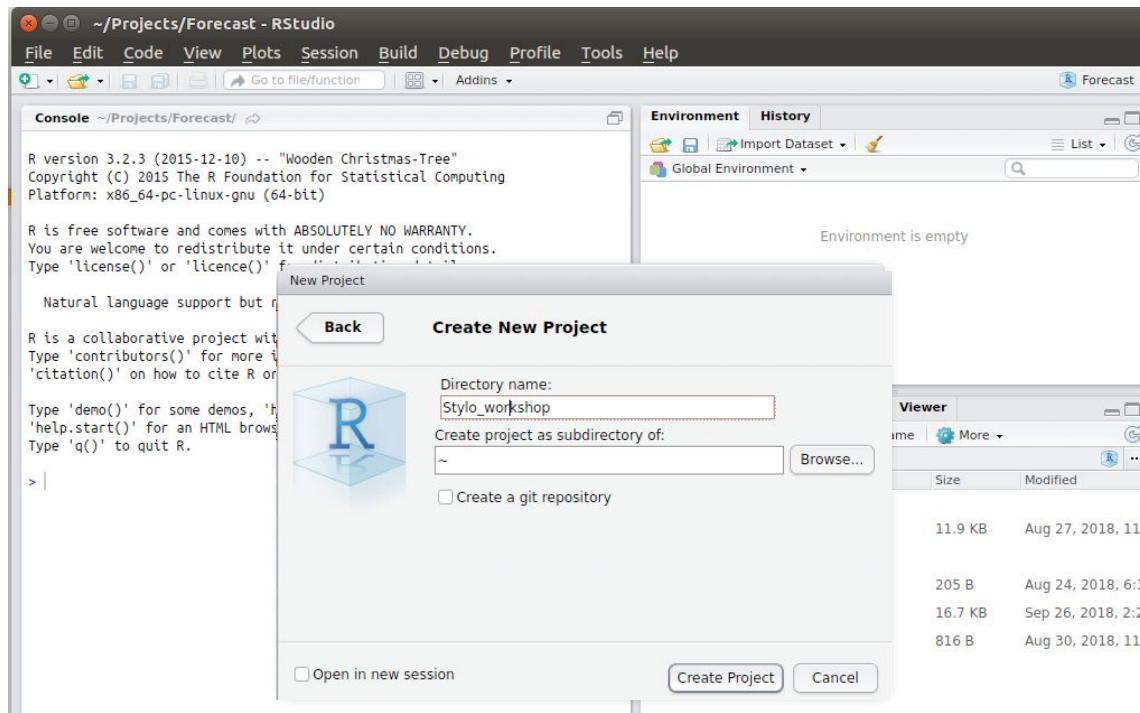
Create project directory

File >> New Project... >> New Directory >> Empty Project



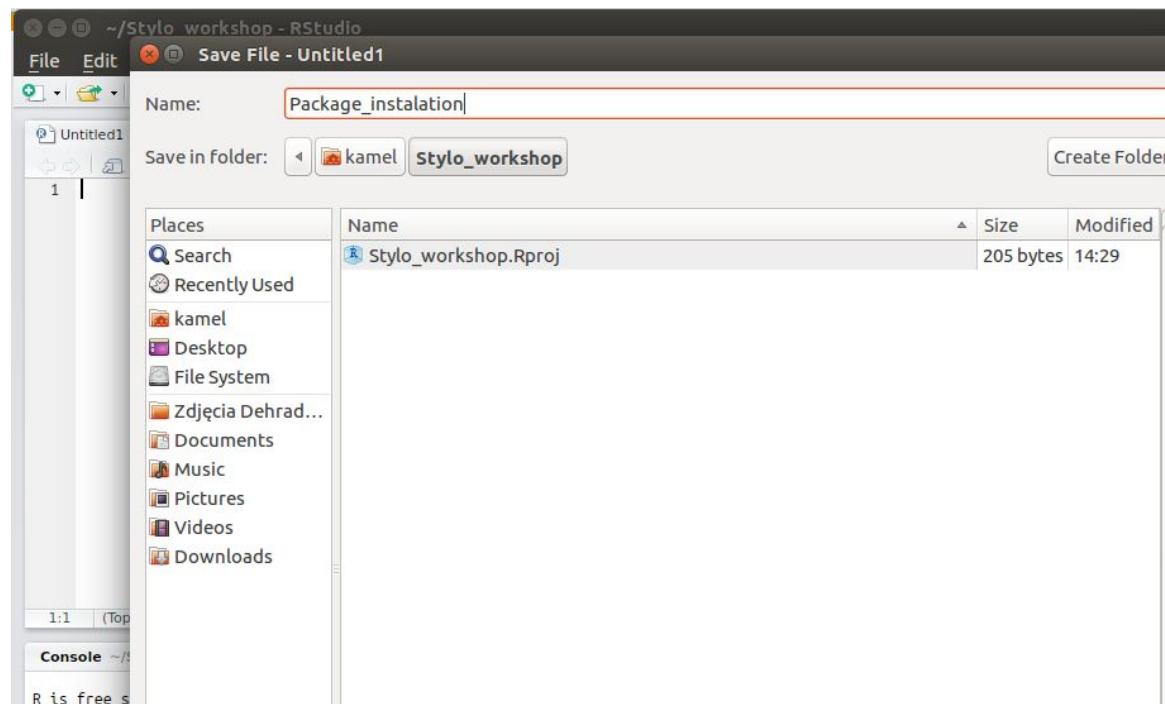
Project Name

Stylo_workshop



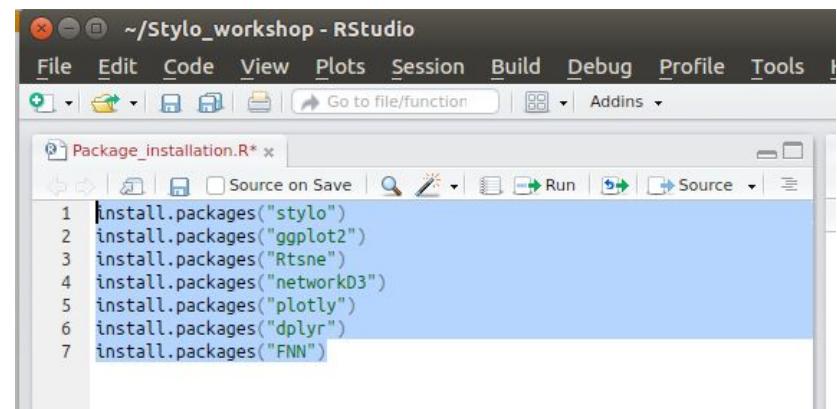
Create script

File >> New File >> R Script
File >> Save >> Package_installation



Install packages

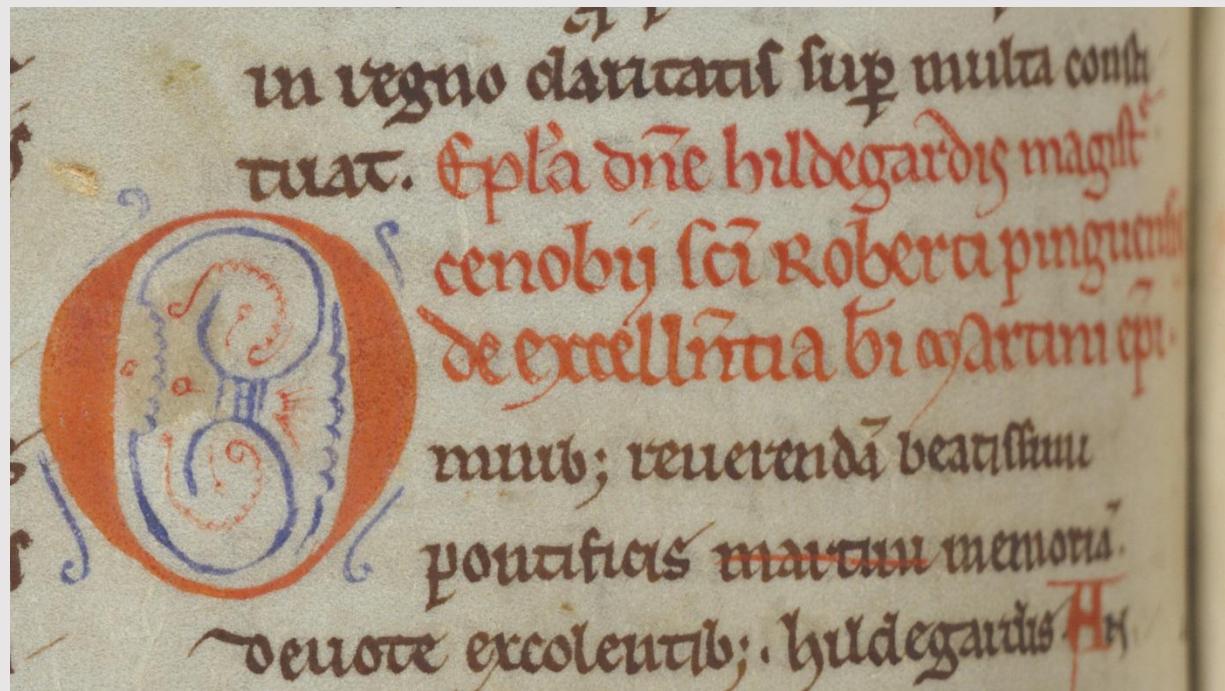
```
install.packages("stylo")
install.packages("ggplot2")
install.packages("Rtsne")
install.packages("networkD3")
install.packages("plotly")
install.packages("dplyr")
install.packages("FNN")
```



The screenshot shows the RStudio interface with the title bar "~/Stylo_workshop - RStudio". The code editor window displays a file named "Package_installation.R" containing the following R code:

```
1 install.packages("stylo")
2 install.packages("ggplot2")
3 install.packages("Rtsne")
4 install.packages("networkD3")
5 install.packages("plotly")
6 install.packages("dplyr")
7 install.packages("FNN")
```

SHORT HISTORY OF STYLOMERTY



Lorenzo Valla (c. 1407–1457)



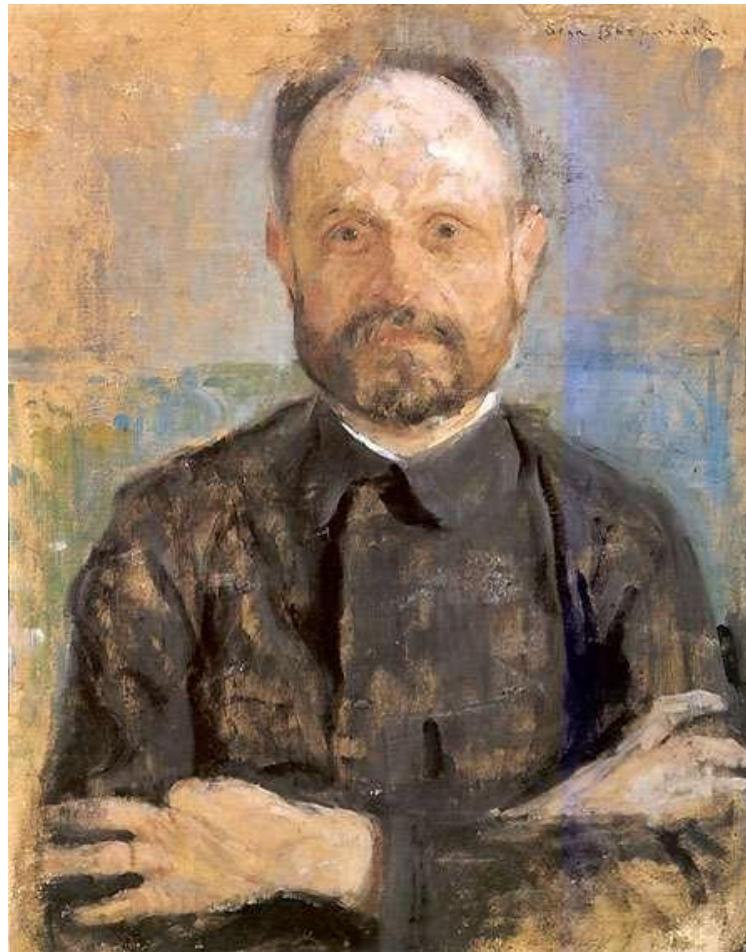
Constantini Donatione declamatio



Constantini Donatione declamatio

- A forged decree where the emperor Constantine I transfers authority of the Roman Empire to the Pope
- Written probably in the 8th century, claimed to be authentic for many centuries
- The first instance of scholarly-based investigation of style:
- Some grammatical forms could not have been used in the 4th century

Wincenty Lutosławski (1863–1954)



Wincenty Lutosławski

- 1890: chronology of Plato's dialogues
- The basics of stylometry were set out by Polish philosopher Wincenty Lutosławski in *Principes de stylométrie* (1890)
- Pioneer of yoga in Poland

John Burrows

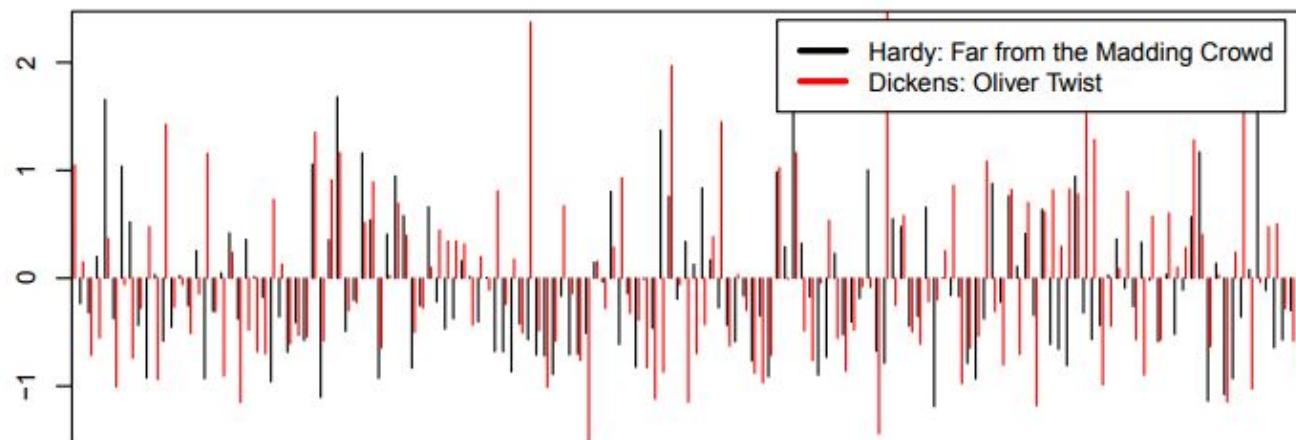


Burrows DELTA

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

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



Z Score

- X - frequency of term
- $\text{mean}(X)$ - mean frequency of term
- S - standard deviation

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

DELTA measures

Burrows's Delta = Manhattan distance (Burrows 2002)

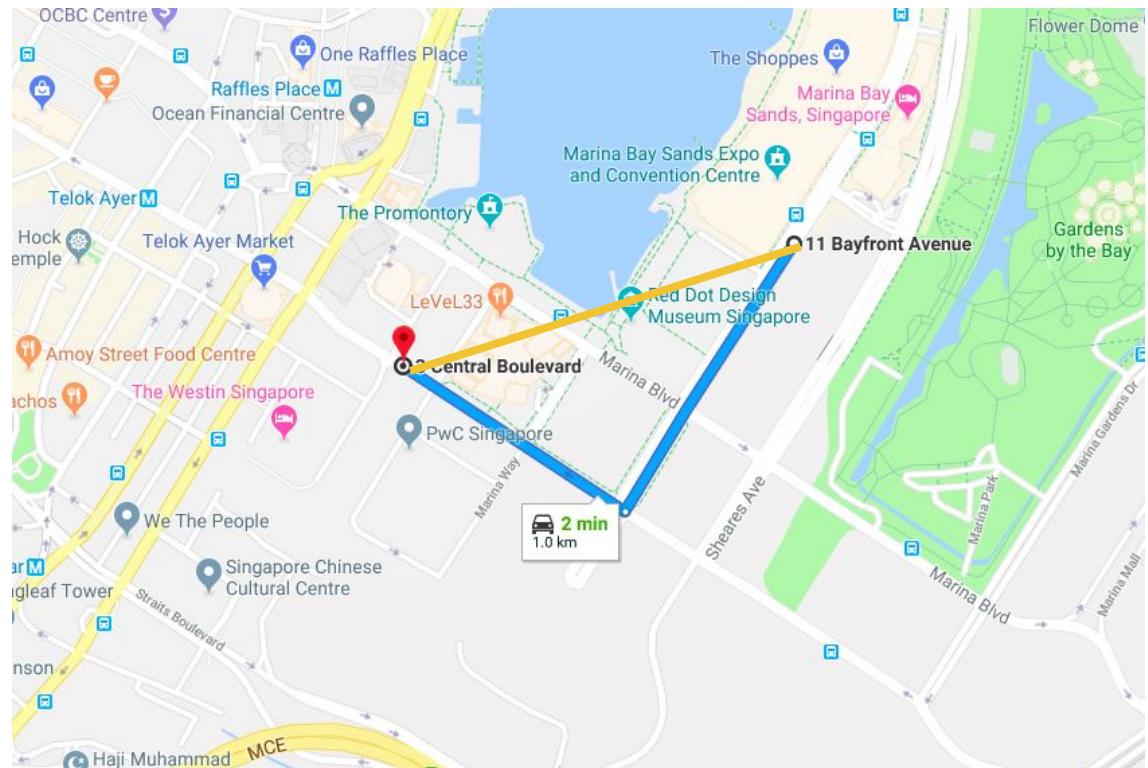
$$\Delta_B(D, D') = \|\mathbf{z}(D) - \mathbf{z}(D')\|_1 = \sum_{i=1}^{n_w} |z_i(D) - z_i(D')|$$

Quadratic Delta = Euclidean distance (Argamon 2008)

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

Distances

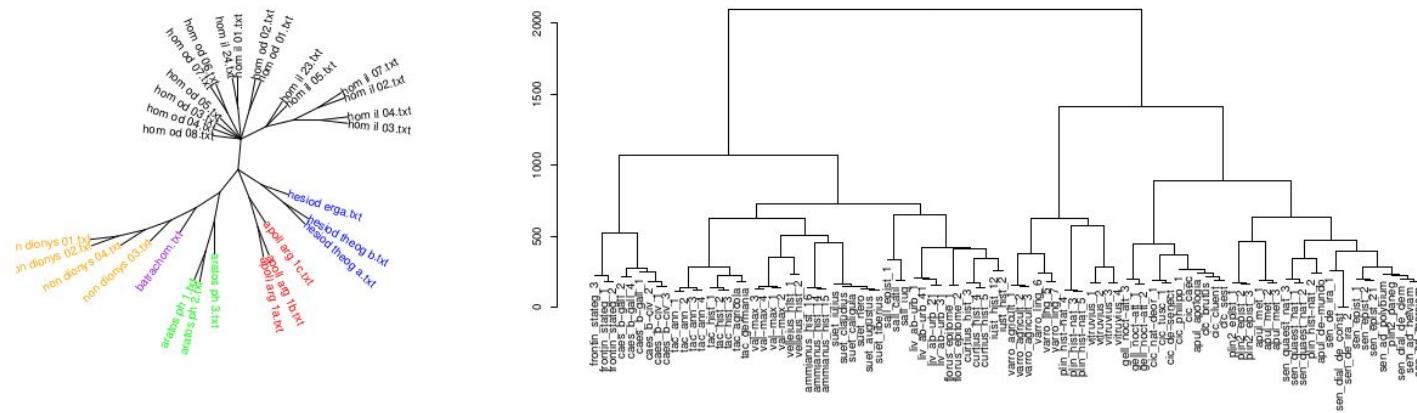
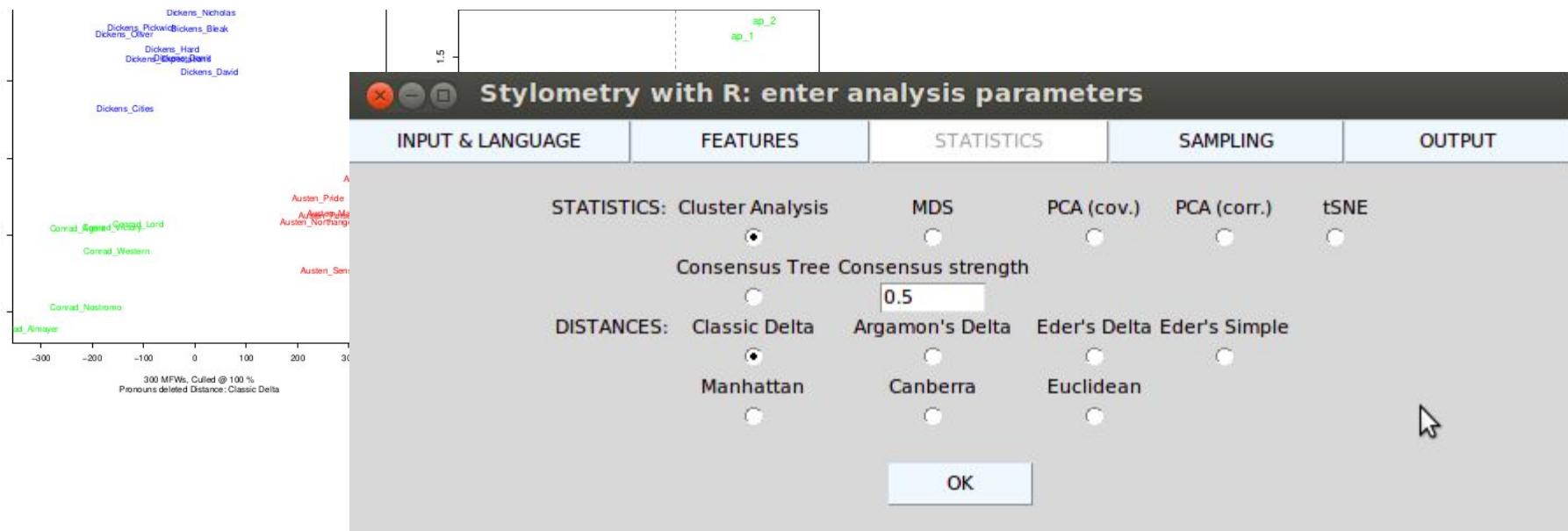
- Time
- Manhattan
- Euclidean



Maciej Eder



Stylo R software



COMPUTATIONAL STYLISTICS GROUP

<https://computationalstylistics.github.io/>

About Blog People Projects Publications Resources

Computational Stylistics Group

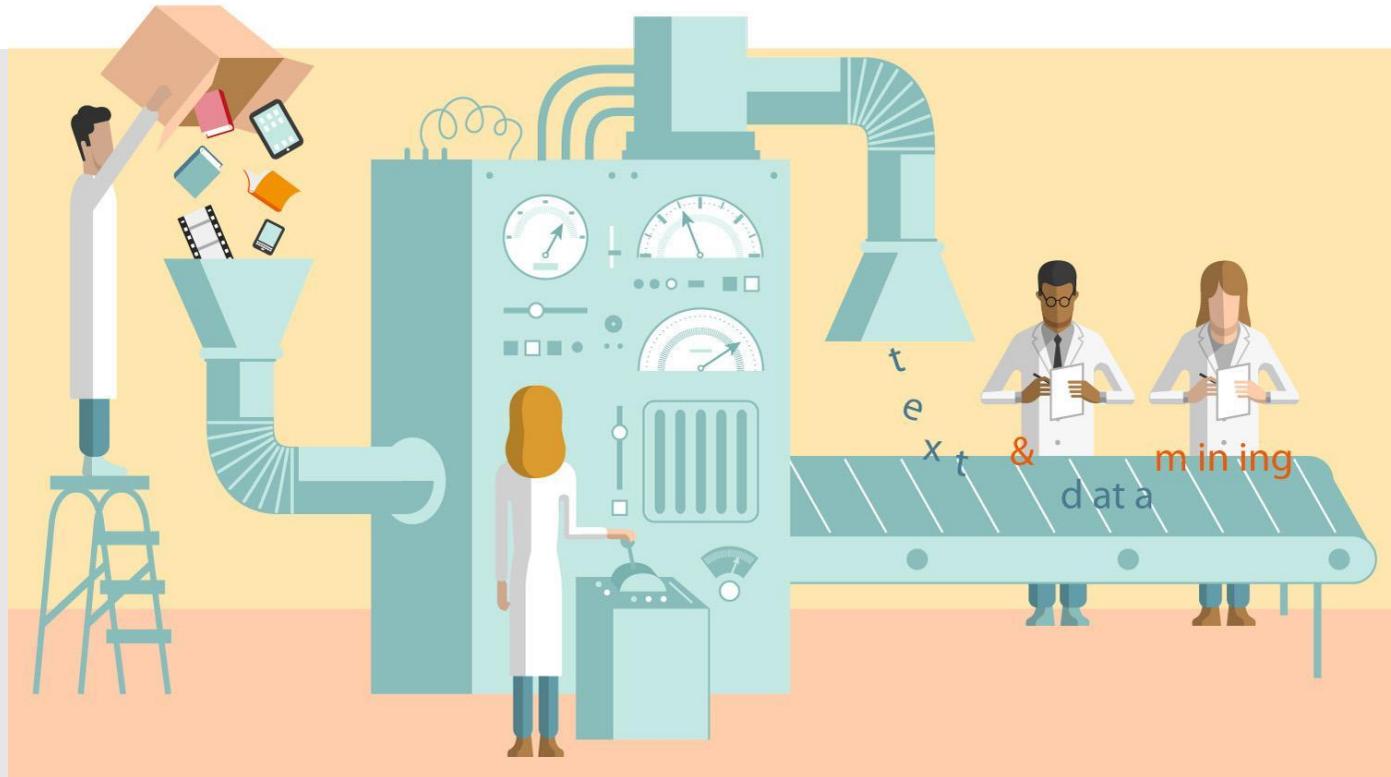
A cross-institutional research team focused on computer-assisted text analysis.

Computational Stylistics Group is a cross-institutional research team focused on computer-assisted text analysis, stylometry, authorship attribution, sentiment analysis, and the like stuff. The research projects conducted by the team members could be described as an intersection of linguistics, literary criticism, and computer sciences – however the best name here would be “Digital Humanities”. The group is based mostly in Kraków, at the Institute of Polish Language (Polish Academy of Sciences), but also at the Jagiellonian University and the University of Antwerp.



Even if the Group has been involved in several research projects (some of them are listed on this website, on the [Projects](#) subpage), it is probably known – at the first place – for the R package [stylo](#), which is a comprehensive collection of functions written in the programming language R, for performing a variety of experiments in computational stylistics. More information about the package can be found [here](#). Also, please check the [discussion list](#) dedicated to various issues in stylometry and beyond.

Why to analyze text?



Applications of stylometry

Digital Humanities

- Author attribution identification of unknown authors
- Genre classification
- Historical study of language change

Other applications

- Anonymity
- Plagiarism
- Criminal civil security- (detection if someone's suicide note is real)

Agatha Christie

5 September 1890 – 12 January 1976

66 detective novels and 14 short story collections

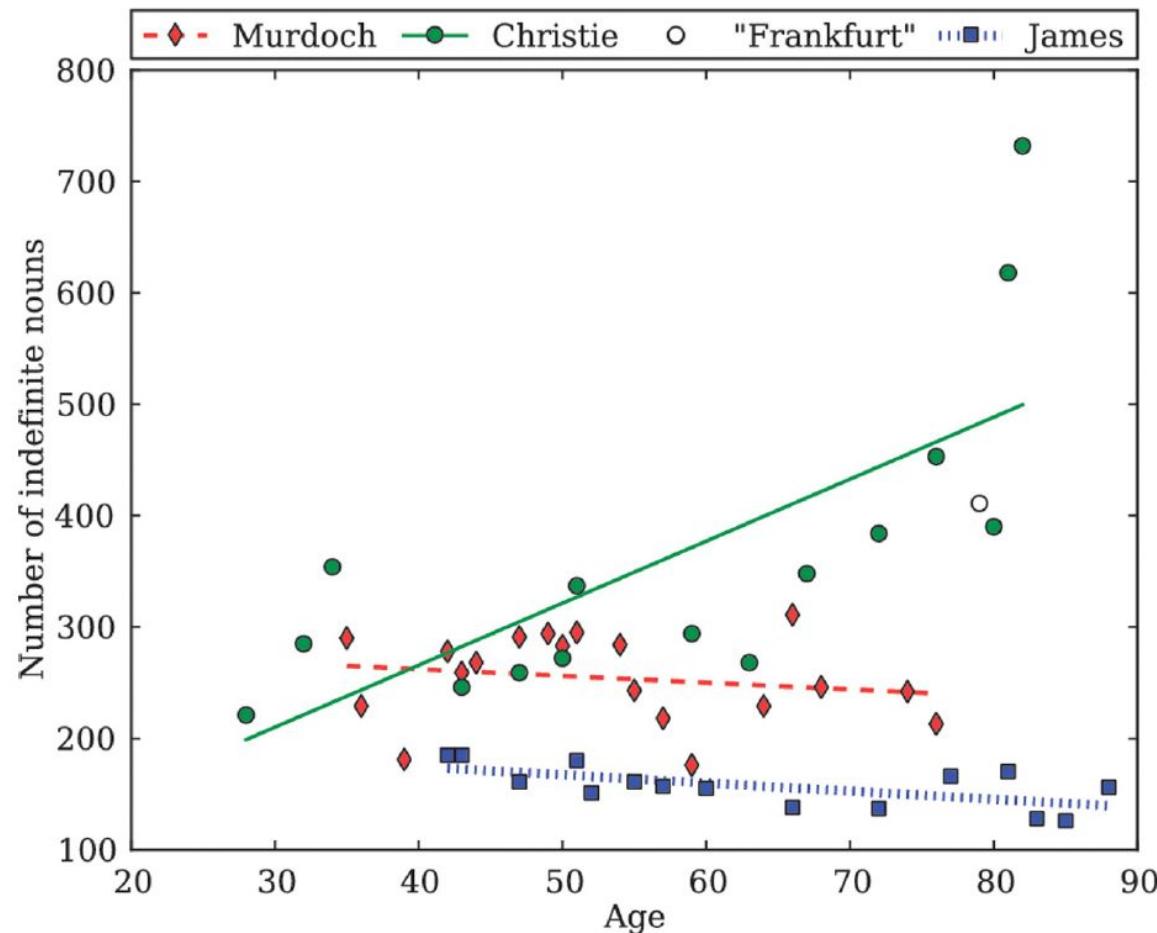
from 1971 to 1974, Christie's health began to fail...

...although she continued to write.

Four indefinite nouns

- thing(s),
- something,
- anything,
- nothing

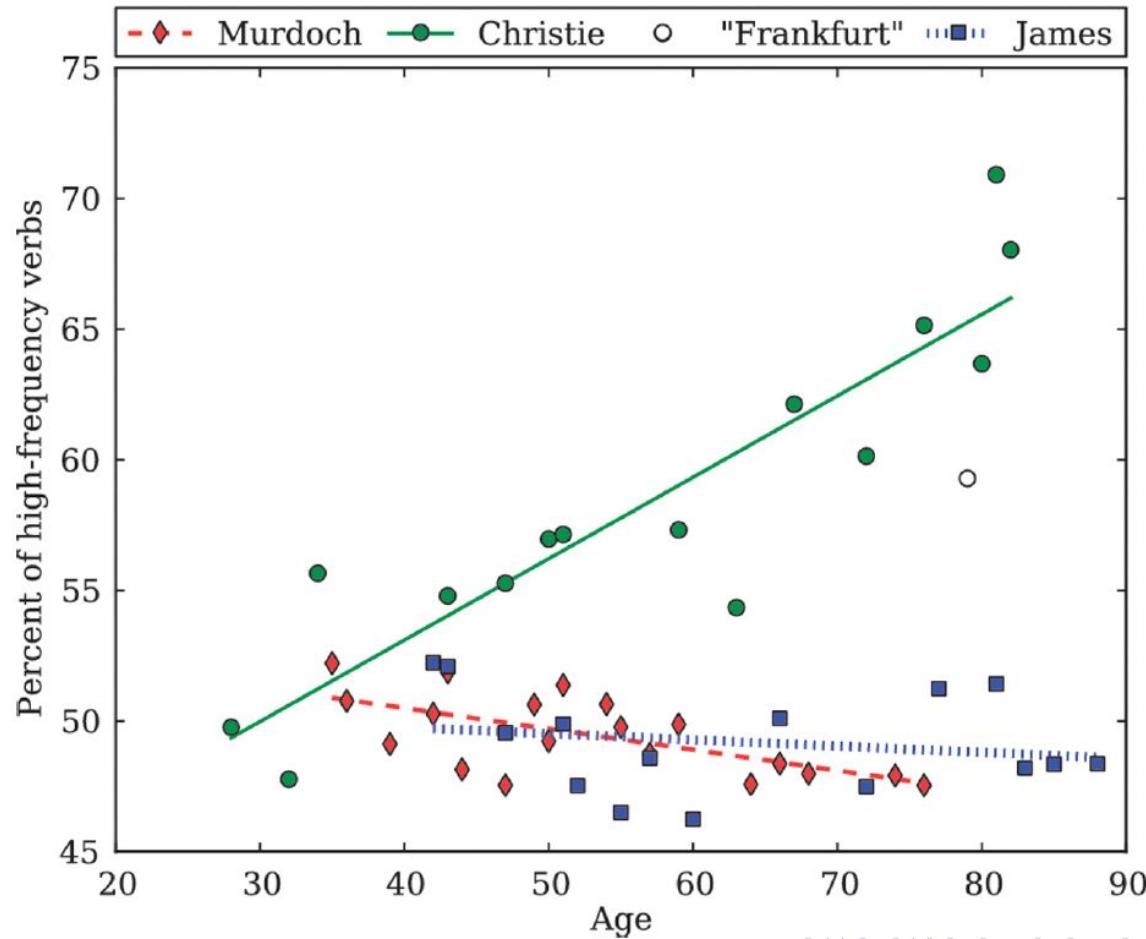
Number of indefinite noun occurrences



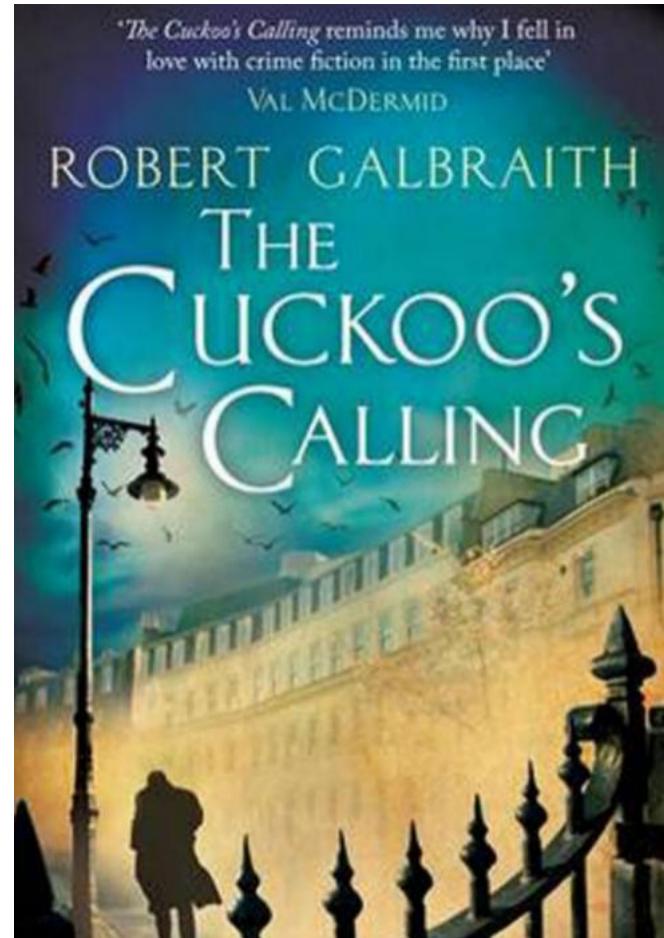
High-frequency verbs of low specificity

be, come, do, get, give, go, have, know, look,
make, see, tell, think, want, ask, feel, find,
forget, happen, hear, like, live, mean, meet, put,
remember, run, say, seem, speak, suppose, take,
use, walk, wonder

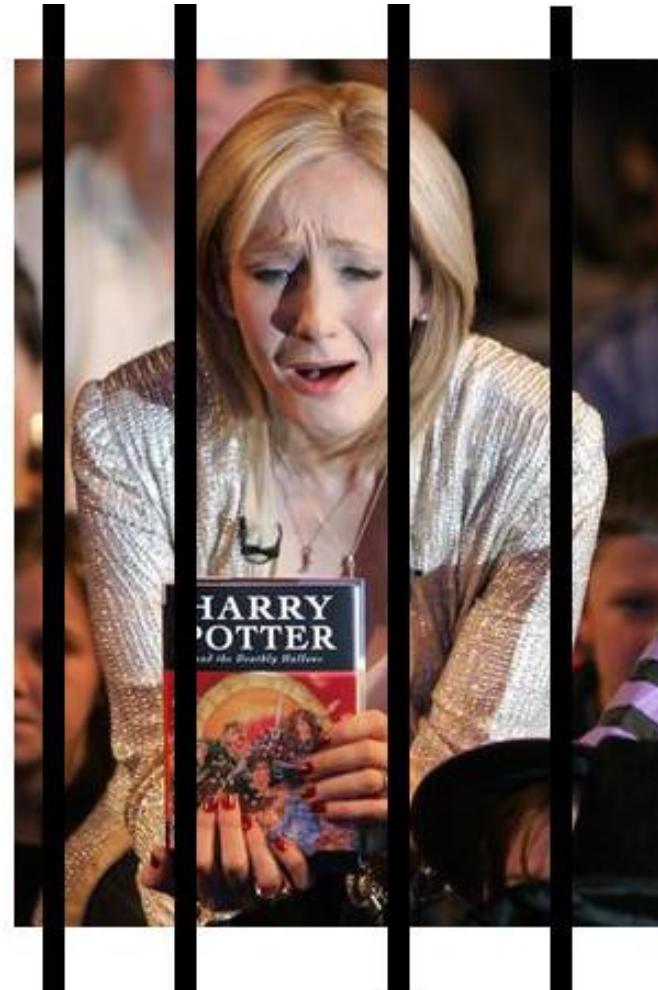
Proportion of 35 high-frequency verbs



Who is Robert Galbraith?



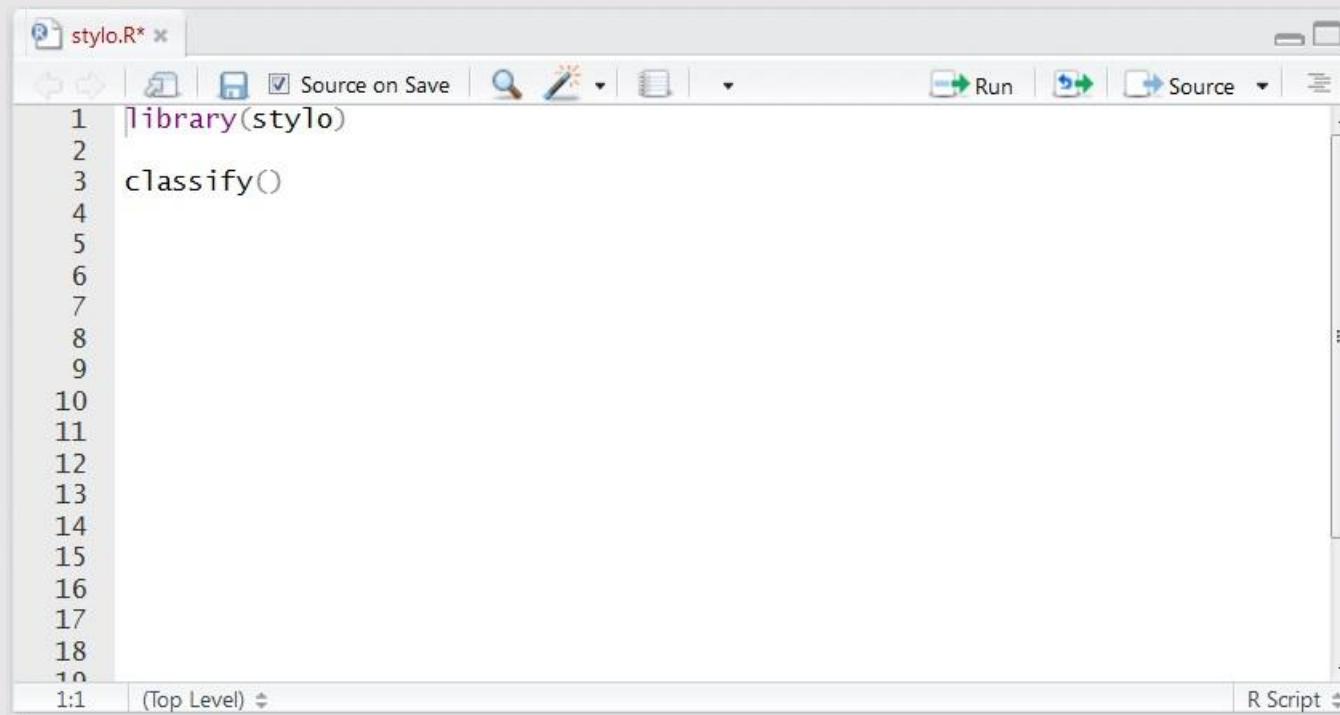
J. K. Rowling . . .



SIRI



How to use stylo?



The screenshot shows an RStudio interface with a file named "stylo.R" open. The code in the editor is:

```
1 library(stylo)
2
3 classify()
```

The RStudio toolbar at the top includes icons for file operations, search, and run. The status bar at the bottom shows "1:1 (Top Level)" and "R Script".

FUNCTIONS PROVIDED

- `stylo()` - an all-in-one tool for a variety of experiments in computational stylistics.
- `classify()` -Function that performs a number of machine-learning methods for classification used in computational stylistics
- `rolling.classify()` - Function that splits a text into equal-sized consecutive blocks (slices) and performs a supervised classification of these blocks against a training

Analysis Workflow

1. Data preparation
 - a. Naming Convention
 - b. Directory
2. Language and Input Formatting
3. Text Characteristics (Features)
4. Summary Type (Statistics)
5. Output

COMPUTATIONAL STYLISTICS GROUP

<https://computationalstylistics.github.io/>

About Blog People Projects Publications Resources



Computational Stylistics Group

A cross-institutional research team focused on computer-assisted text analysis.

Computational Stylistics Group is a cross-institutional research team focused on computer-assisted text analysis, stylometry, authorship attribution, sentiment analysis, and the like stuff. The research projects conducted by the team members could be described as an intersection of linguistics, literary criticism, and computer sciences – however the best name here would be “Digital Humanities”. The group is based mostly in Kraków, at the Institute of Polish Language (Polish Academy of Sciences), but also at the Jagiellonian University and the University of Antwerp.



Even if the Group has been involved in several research projects (some of them are listed on this website, on the [Projects](#) subpage), it is probably known – at the first place – for the R package [stylo](#), which is a comprehensive collection of functions written in the programming language R, for performing a variety of experiments in computational stylistics. More information about the package can be found [here](#). Also, please check the [discussion list](#) dedicated to various issues in stylometry and beyond.

Data preparation

10 Computational 01
01 Stylistics 0101000
11 Group 011010110

About Blog People Projects Publications Resources

Resources

Materials prepared by the Group. (More to be added soon...).

Corpora

The following selection of links is but a tip of an iceberg when it comes to the corpora (text collections) suitable for text analysis. The corpora listed below, however, are compiled by the members of CSG, and checked for compatibility with commonly known stylometric software.

- [A Small Collection of British Fiction](#)
- [100 Polish Novels](#)
- [100 English Novels](#)
- [68 German Novels](#)



Data preparation

computationalstylistics / A_Small_Collection_of_British_Fiction

Watch 1 Star 1 Fork 1

Code Issues 0 Pull requests 0 Projects 0 Insights

Join GitHub today

GitHub is home to over 28 million developers working together to host and review code, manage projects, and build software together.

Dismiss Sign up

A selection of 28 classic British novels from the 19th century (including a few late 18th-century items). Full text versions, in plain text format, harvested from trustworthy public domain sites.

5 commits 1 branch 0 releases 2 contributors

Branch: master New pull request Find file Clone or download

computationalstylistics Update README.md

corpus 28 novels added to the repository 6 months ago

README.md Update README.md 6 months ago

overview Create overview 6 months ago

README.md



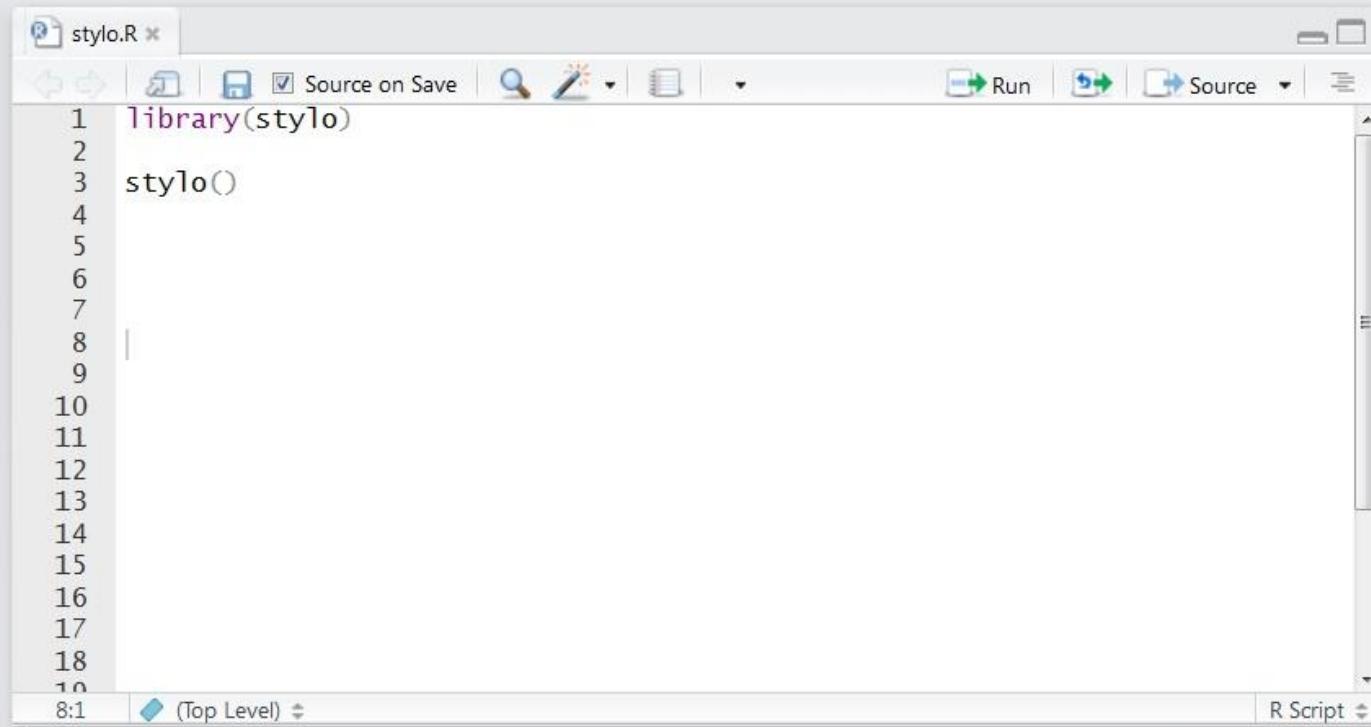
Data preparation

- https://github.com/computationalstylistics/100_english_novels
- https://github.com/computationalstylistics/A_Small_Collection_of_British_Fiction

Exercise - project preparation

- Create new project for the Exercises
- Name the project Stylometry_Excercise
- Create script with named exercise.R

STYLO



The screenshot shows the RStudio interface with a script file named "stylo.R" open. The code in the file is:

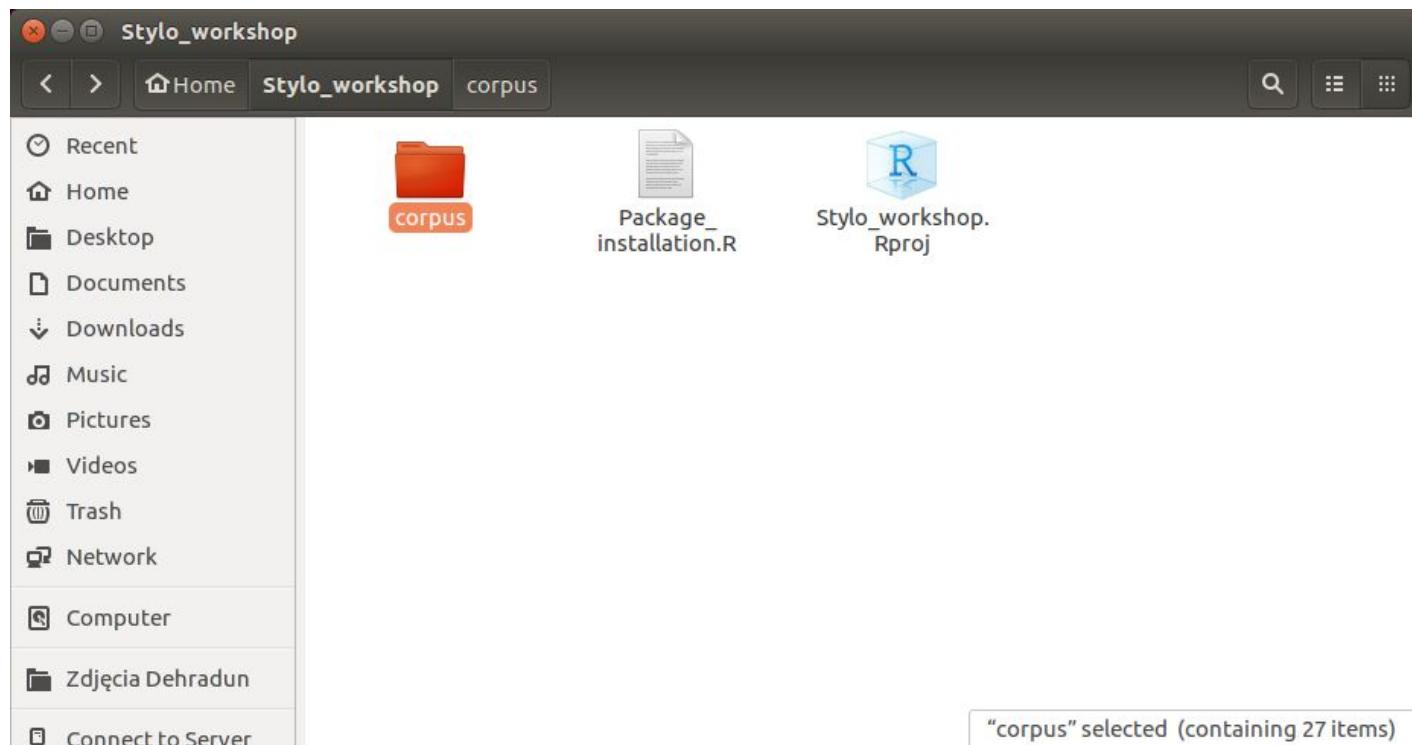
```
1 library(stylo)
2
3 stylo()
```

The "stylo" library is loaded at line 1, and the function "stylo" is called at line 3. The RStudio toolbar at the top includes icons for file operations, search, and run. The status bar at the bottom shows "8:1" and "(Top Level)".

What is Stylo?

Stylo: Comprehensive collection of functions written in the programming language R, for performing variety of experiments in computational stylistics.

Data preparation for Stylo()



Corpus preparation - labels

Colors on graphs are assigned according to filenames: the sequence of letters before “ - ” (underscore) is assumed to be the label of the author (genre, etc.).

- ▶ ABronte_Agnes.txt
- ▶ Austen_Emma.txt
- ▶ Austen_Pride.txt
- ▶ Conrad_Lord.txt
- ▶ Dickens_Pickwick.txt
- ▶ ...

LIVE DEMO

The screenshot shows the RStudio interface with a script editor window titled "stylo.R" and a "Stylometry with R | stylo | set parameters" dialog box.

stylo.R content:

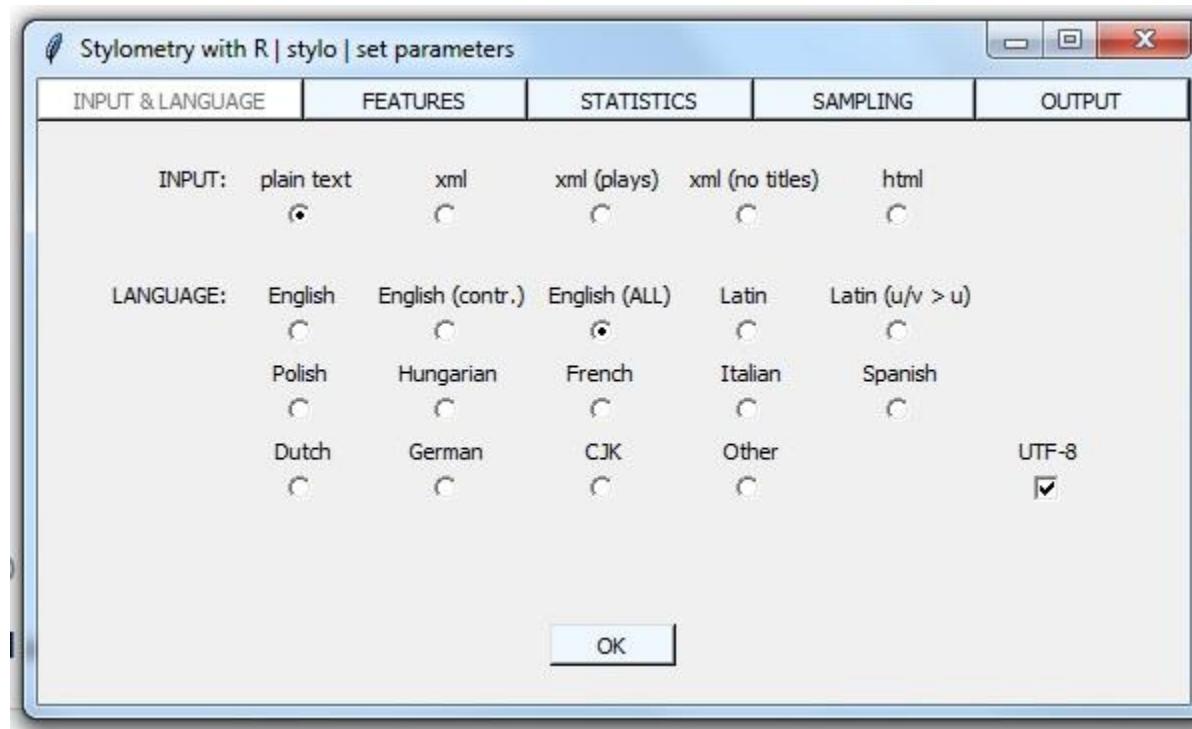
```
1 Library(stylo)
2
3 stylo()
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
```

Stylometry with R | stylo | set parameters dialog box:

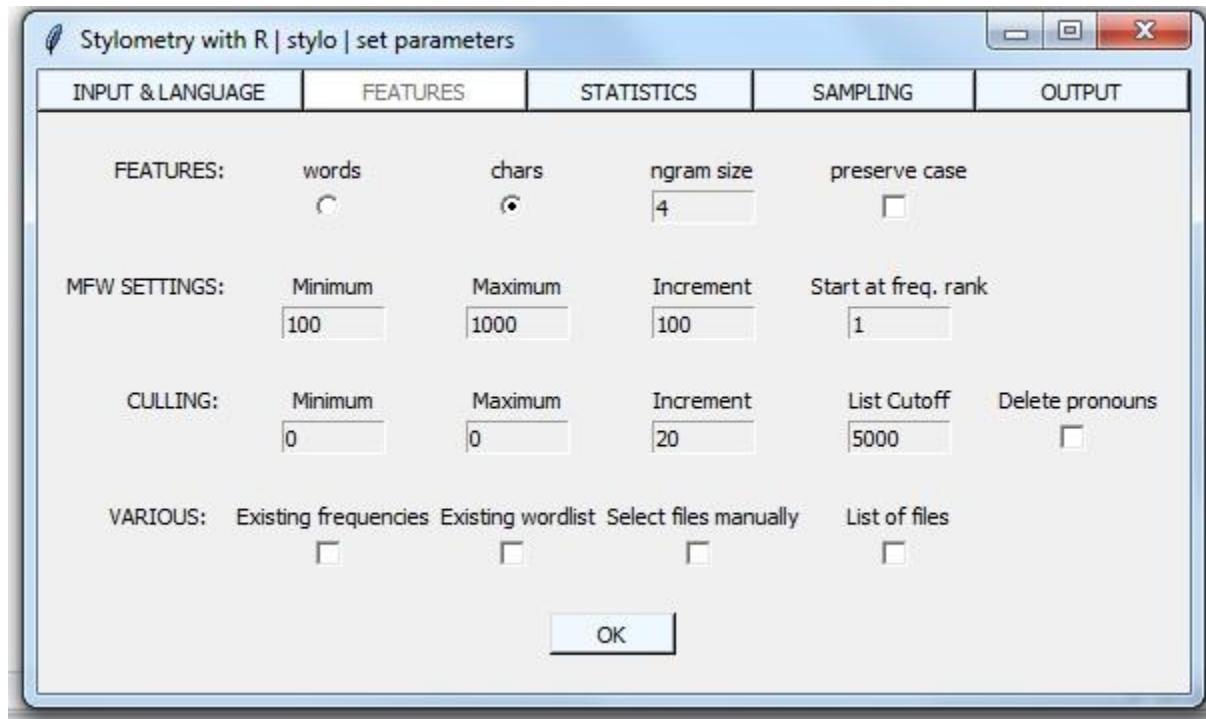
	INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT
INPUT:	plain text <input checked="" type="radio"/>	xml <input type="radio"/>	xml (plays) <input type="radio"/>	xml (no titles) <input type="radio"/>	html <input type="radio"/>
LANGUAGE:	English <input type="radio"/>	English (contr.) <input type="radio"/>	English (ALL) <input checked="" type="radio"/>	Latin <input type="radio"/>	Latin (u/v > u) <input type="radio"/>
	Polish <input type="radio"/>	Hungarian <input type="radio"/>	French <input type="radio"/>	Italian <input type="radio"/>	Spanish <input type="radio"/>
	Dutch <input type="radio"/>	German <input type="radio"/>	CJK <input type="radio"/>	Other <input type="radio"/>	UTF-8 <input checked="" type="checkbox"/>

OK

INPUT & LANGUAGE



FEATURES



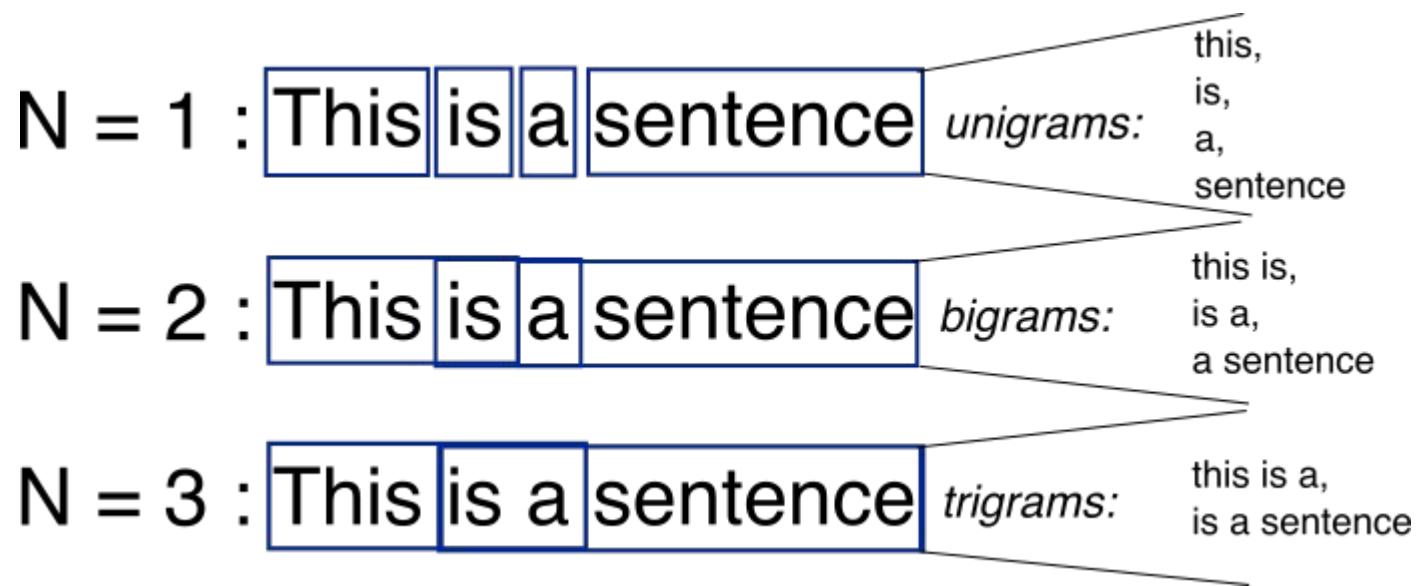
FEATURES

- 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.

FEATURES

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

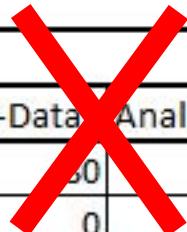
settings

- 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

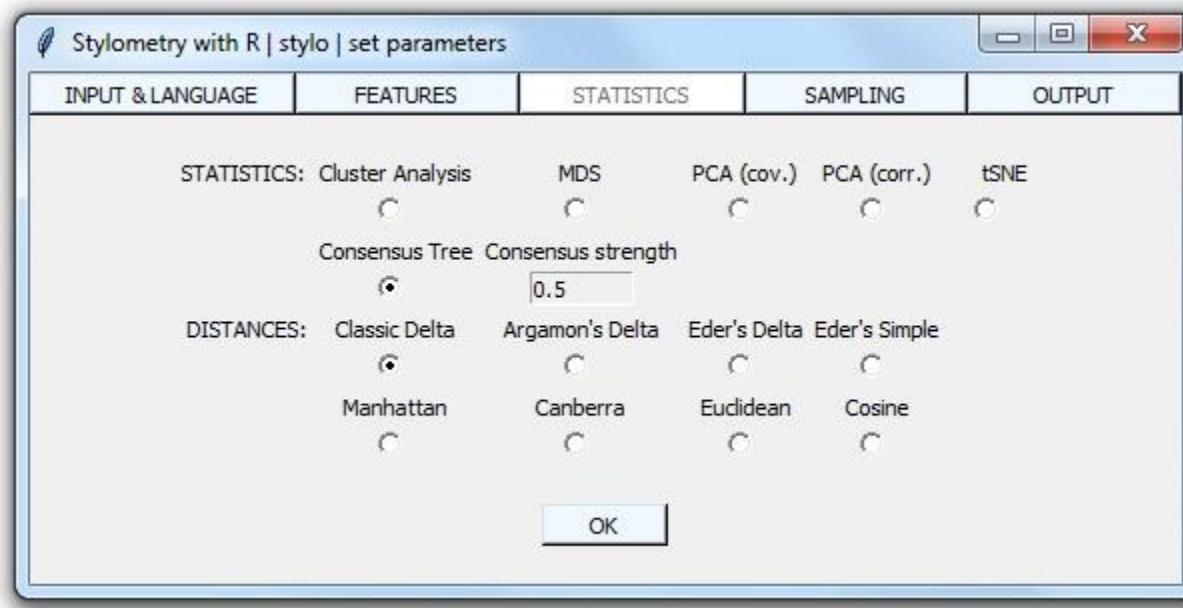
- The culling values specify the degree to which words that do not appear in all the texts of your corpus will be removed. Thus, a culling value of 20 indicates that words that appear in at least 20% of the texts in the corpus will be considered in the analysis. A culling setting of 0 means that no words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

Culling

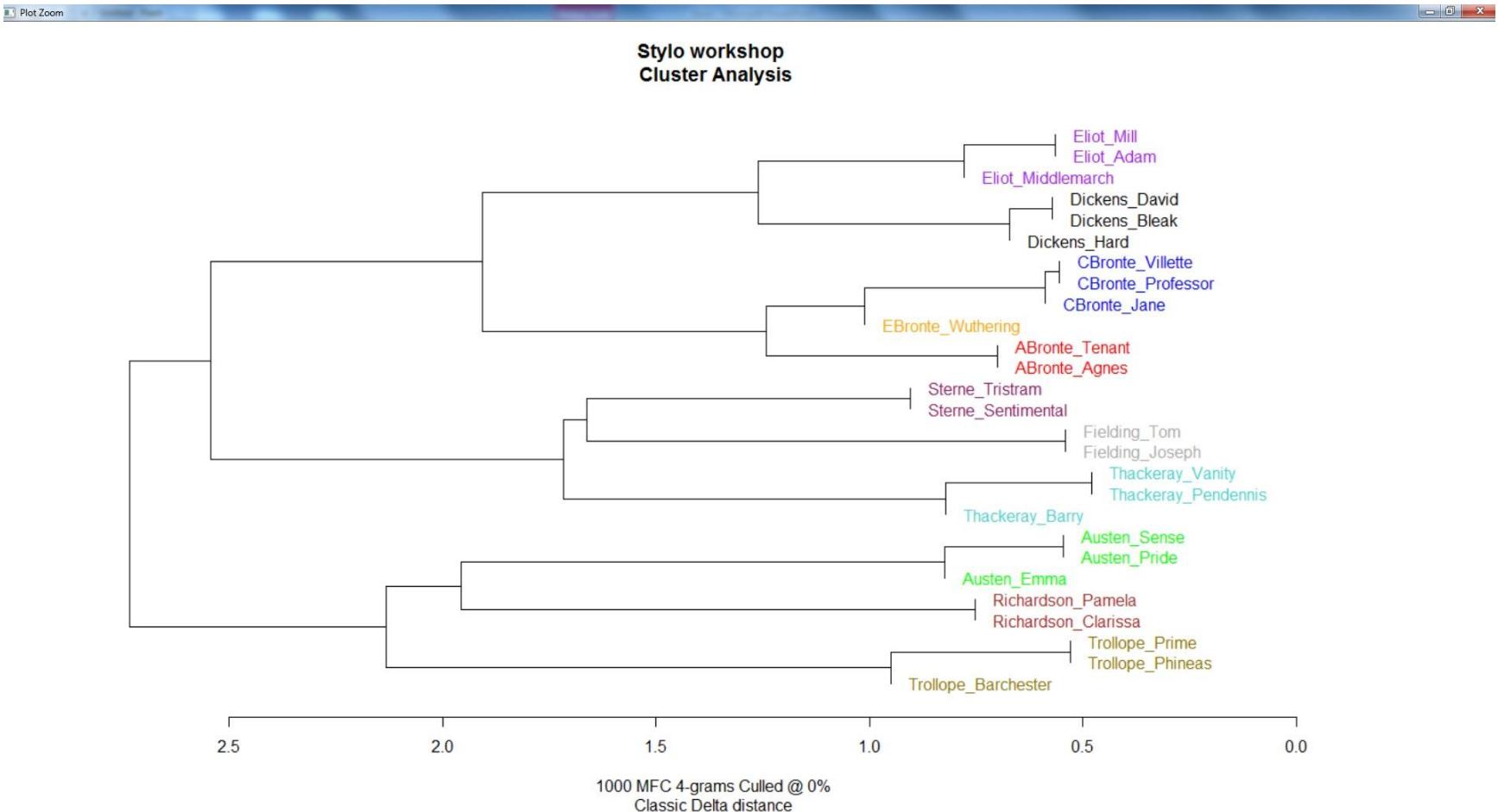


Book Number	Word Frequency									
	The	Big-Data	Analytics	Tree	newbie	book	for	Girl	honest	
1	120	30	60	20	1	5	120	0	0	0
2	110	0	0	100	10	20	100	40	10	
3	130	0	0	10	11	30	110	20	10	
4	100	0	0	2	20	40	100	10	100	
5	90	0	0	10	30	20	100	100	40	

STATISTICS



CLUSTER ANALYSIS

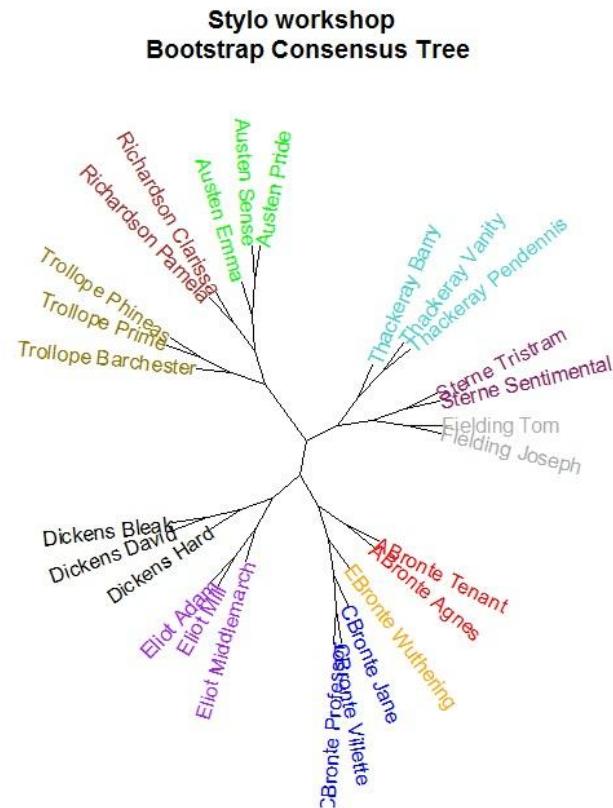


Cluster Analysis

Builds “tree” based on the most similar texts based on the MFV.

It is not robust on the changes of the parameters.

Consensus tree



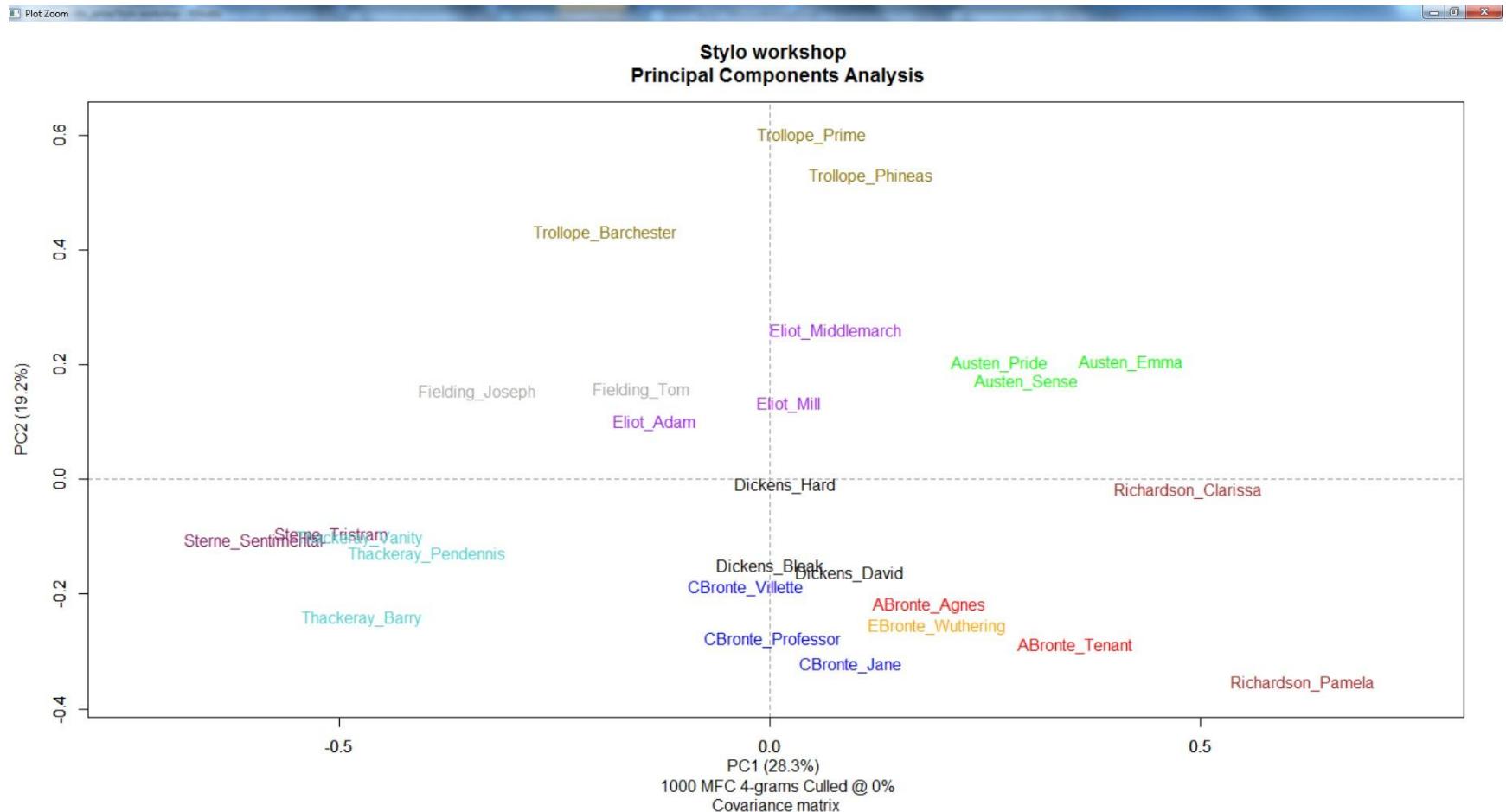
100-1000 MFC 4-grams Culled @ 0%
Classic Delta distance Consensus 0.5

Consensus Tree

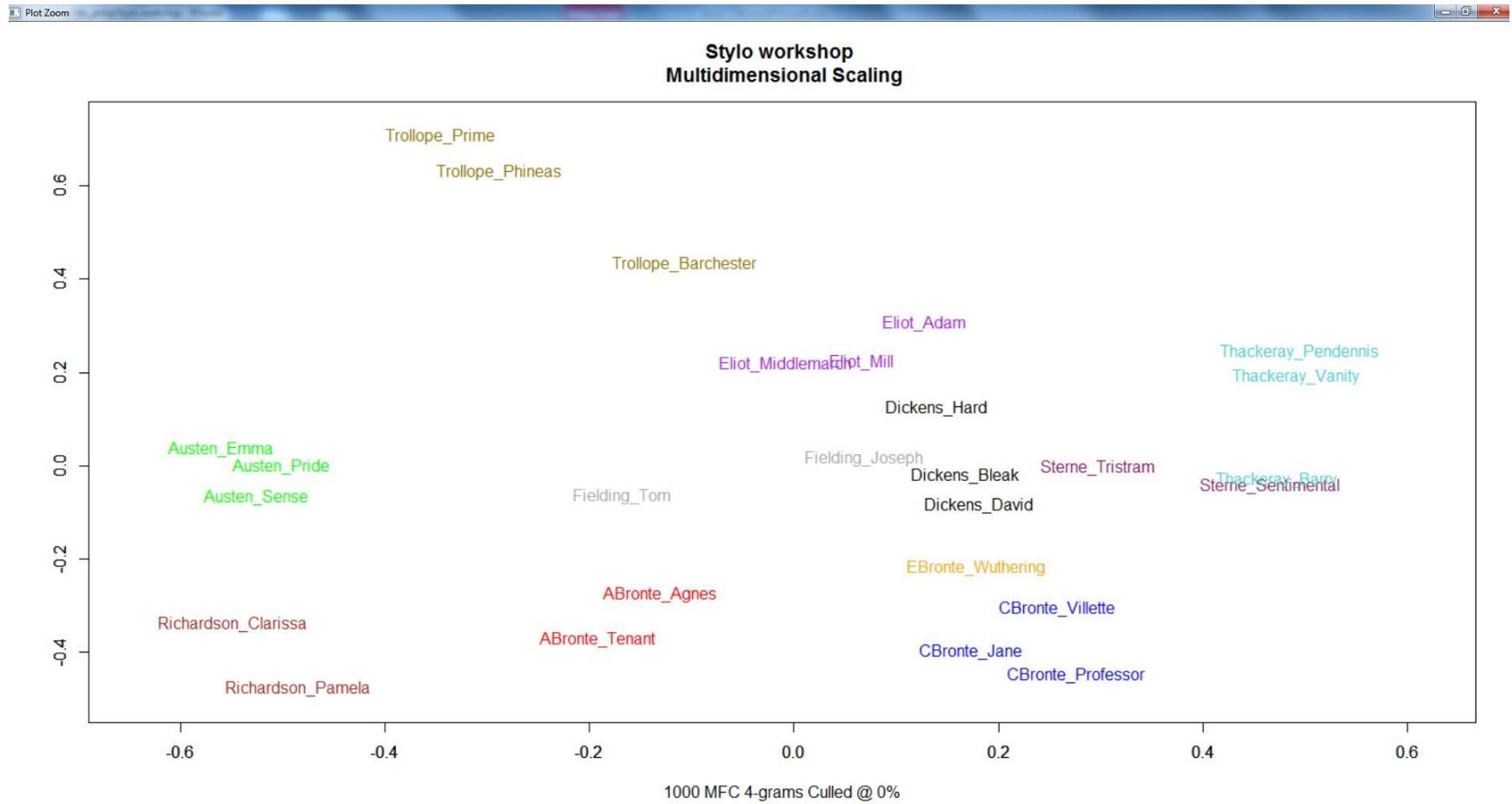
Uses many trees to discover unchanged patterns
for different parameters

It is more robust but harder to interpret.

PCA (principal component analysis)



MDS (multidimensional scaling)

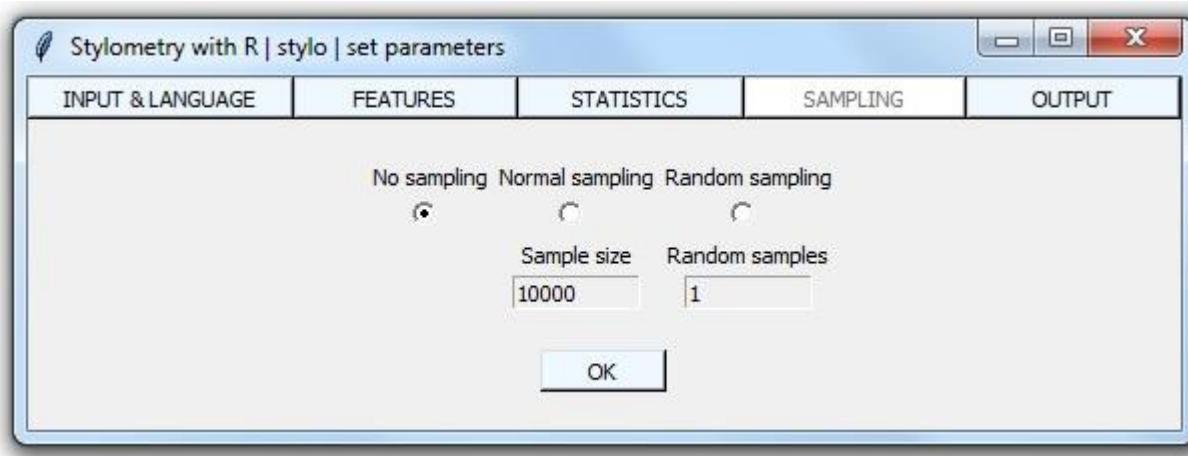


Dimension reduction

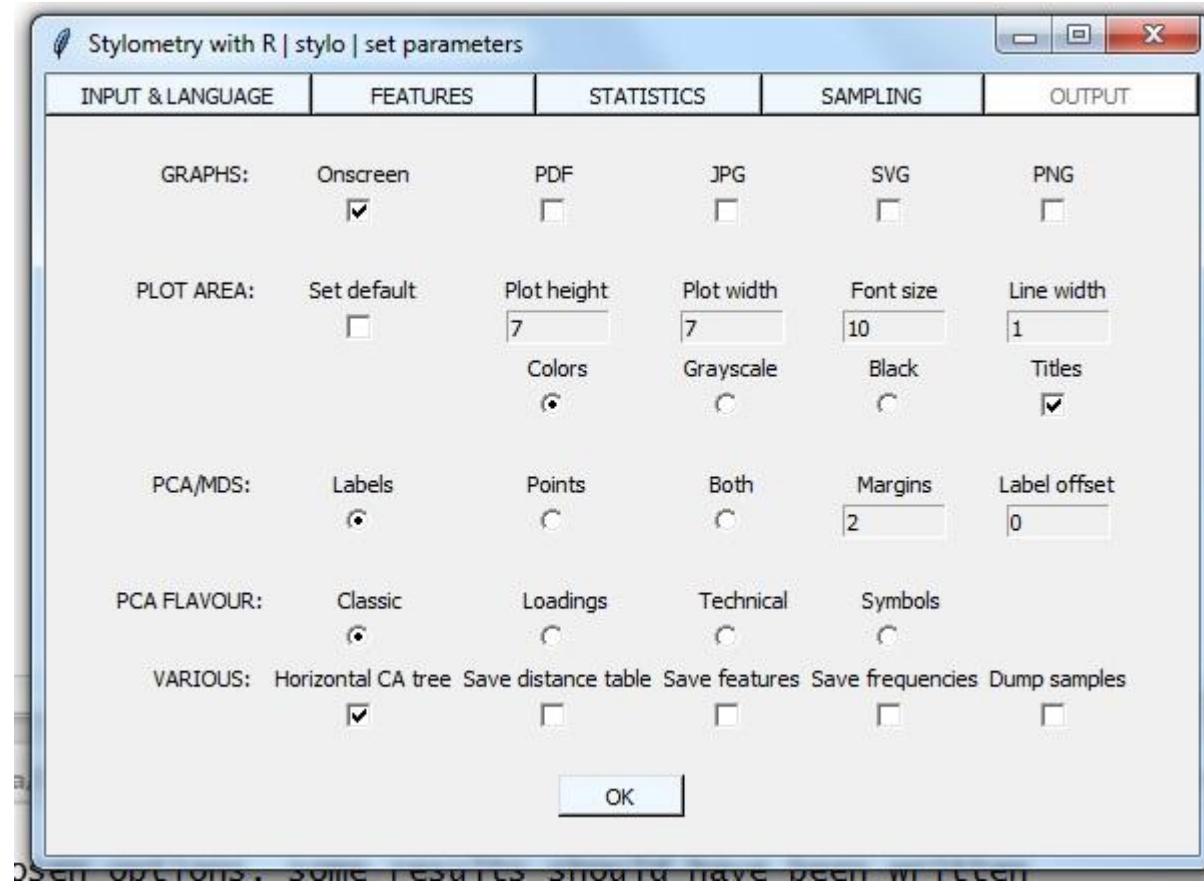
Methods to reduce high dimensions into easy to interpret 2D or 3D space.

It is easy to interpret and visualization gives good intuitions on the similarity between texts.

SAMPLING



OUTPUT



Graphs

- Onscreen: the graph on the screen
- PDF: a PDF file with your graph
- JPG: graph in JPEG format
- SVG: a SVG vector file
- PNG: graph in PNG format

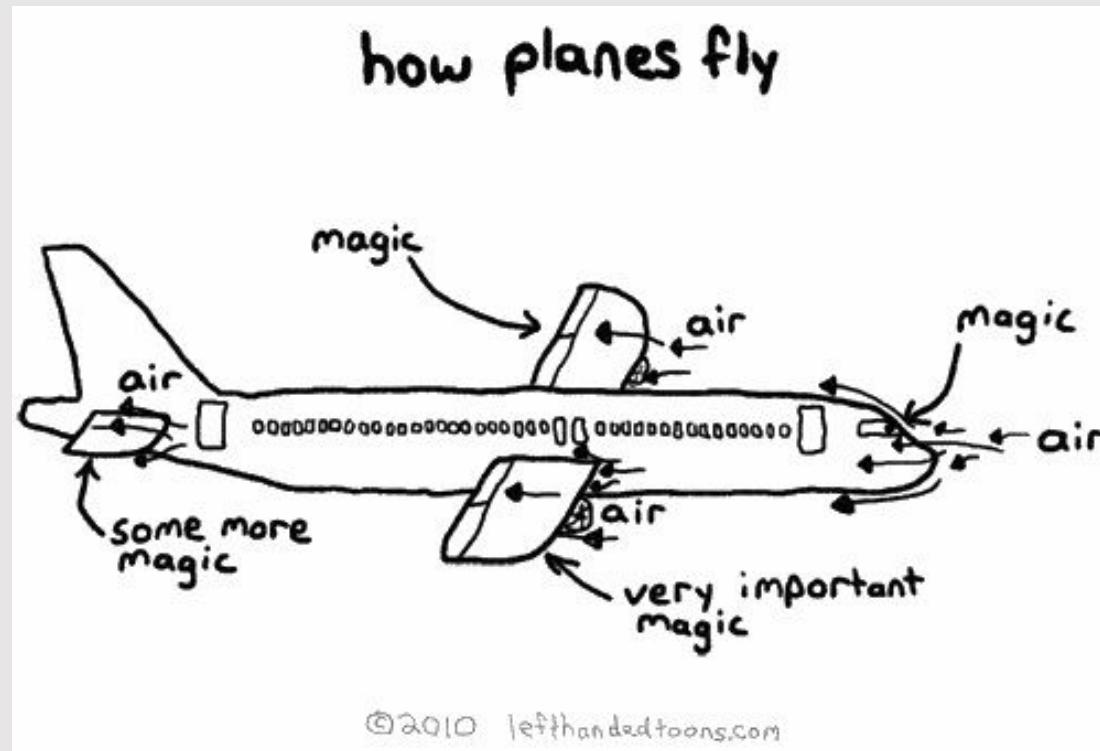
Lets get it to work

```
> library(stylo)
> stylo()
using current directory...
Performing no sampling (using entire text as sample)
loading ABronte_Agnes.txt      ...
loading ABronte_Tenant.txt     ...
loading Austen_Emma.txt ... 
loading Austen_Pride.txt       ...
loading Austen_Sense.txt       ...
loading CBronte_Jane.txt       ...
loading CBronte_Professor.txt  ...
loading CBronte_Villette.txt   ...
loading Dickens_Bleak.txt      ...
loading Dickens_David.txt     ...
loading Dickens_Hard.txt       ...
loading EBronte_Wuthering.txt  ...
loading Eliot_Adam.txt        ...
loading Eliot_Middlemarch.txt  ...
loading Eliot_Mill.txt         ...
loading Fielding_Joseph.txt    ...
loading Fielding_Tom.txt       ...
loading Richardson_Clarissa.txt ...
```

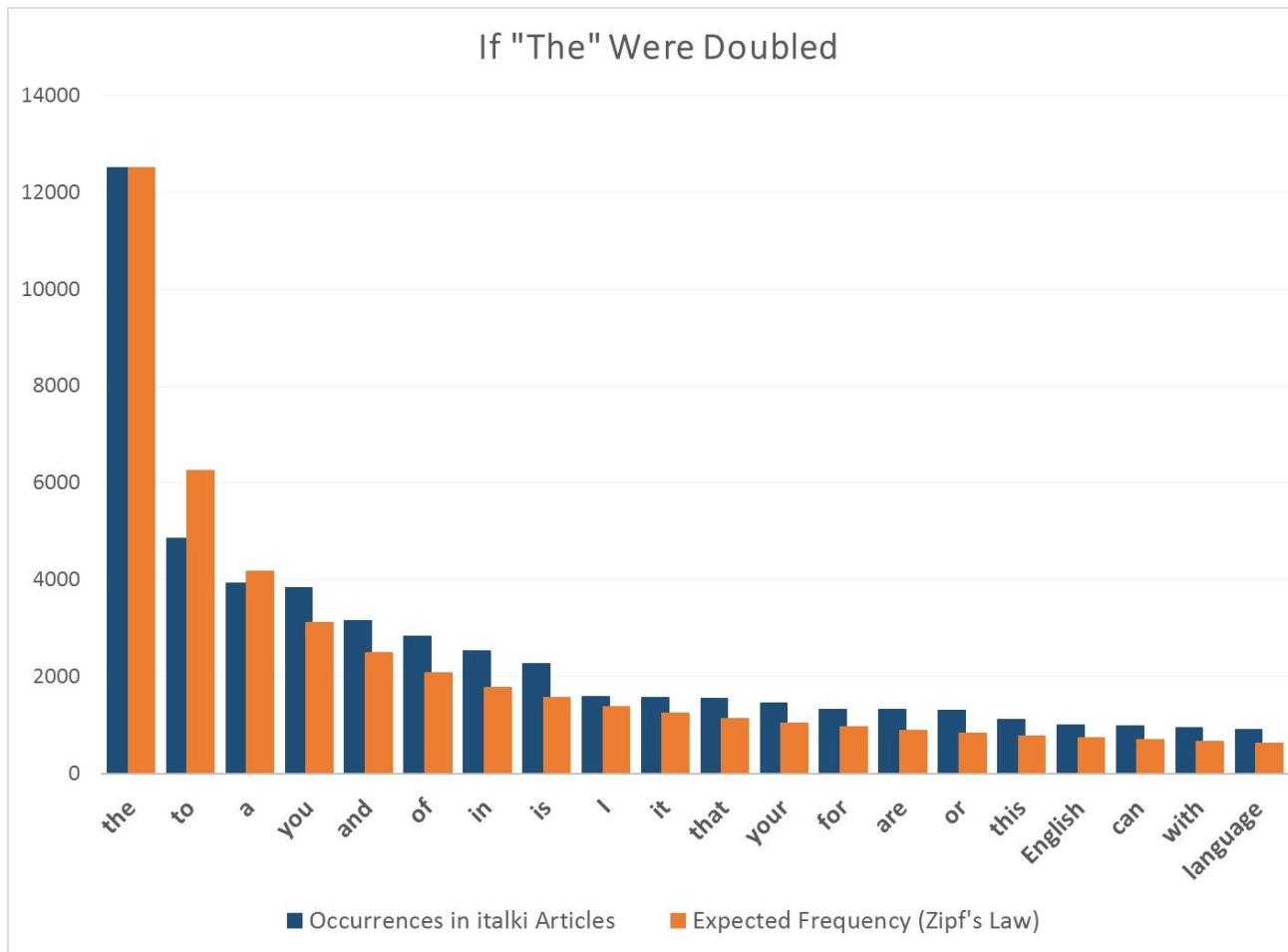
Stylo - Exercise

- Use 100 English Novels for stylo() analysis
- Use Eder's and Manhattan Distance
- Cluster analysis, PCA, Consensus tree
- Use different values for culling and MFW

Intuitions how it is done



Zpif Law I



Zpif Law II

1. i	15. po	3600. przystani
2. się	16. za	3601. racji
3. w	17. tak	3602. twardo
4. nie	18. już
5. na	19. jej	80725. dzieciobójstwo
6. z	20. od	80726. dzieciiskami
7. do	80727. dzieciisków
8. to	3593. ludzki	80728. dzieciński
9. a	3594. mgnienie	80729. dziedzicami
10. że	3595. pacierz	80730. dziedzicowe
11. ale	3596. pióra	80731. dziedzicowym
12. jak	3597. podała	80732. dziedzictwem
13. o	3598. postawie	80733. dziedziczko
14. co	3599. pragnęła	80734. dziedzicznej

Zpif Law III

The more frequent word the more meanings it has

Term document frequency - novel into numbers

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

Z Score

- X - frequency of term
- $\text{mean}(X)$ - mean frequency of term
- S - standard deviation

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

Z Score WHY SO MUCH MATH?

- X - frequency of term

Too much math

- $\text{mean}(X)$ -

- S - standard deviation



$$Z = \frac{\text{_____}}{S}$$

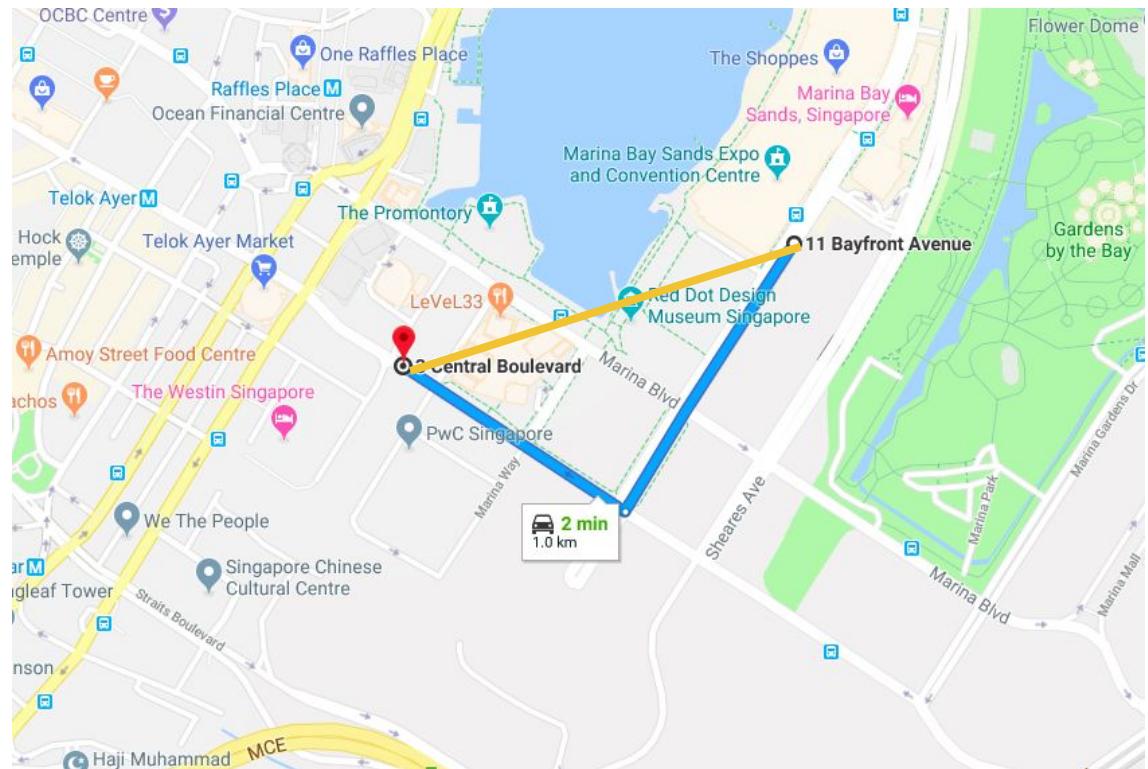
S

Term document frequency - novel into numbers

```
"EBronte_Wuthering" "Eliot_Mill_rep"
"the" 3.9190127861008 4.19596902207898
"and" 4.0682268396635 2.97705517340901
"to" 2.98771127938188 2.92943383520131
"of" 1.90548061503632 2.43013131944778
"i" 3.07775424273868 1.32810621001491
"a" 1.98866316213736 2.31709076915677
"in" 1.25631372683538 1.49405935831449
"that" 1.03249264649133 1.41805762662947
"he" 1.67479911843651 1.10106306219635
"it" 1.12425071391207 1.15493770744144
"you" 1.48956787953109 1.16792534513445
"was" 0.966461140029671 1.46471691760065
"her" 1.33006320158476 1.38390494973303
"his" 1.22115409352462 0.948097551589783
"as" 0.802668701923489 0.968781567174948
"my" 0.948452547358311 0.372312280532974
"for" 0.722058810918352 0.856222040502189
"with" 0.691186937767449 1.0063014093992
"not" 0.803526253955459 0.671027947472221
"be" 0.618295015050038 0.683053537928712
"had" 0.590853350027013 1.0014911732166
"she" 1.10023925701692 0.946654480735004
"have" 0.542830436236719 0.538265428832556
"me" 0.91243536201559 0.417047477031122
"but" 0.589138245963074 0.771080860070229
"is" 0.560839028908079 0.427148973014575
"at" 0.672320793064119 0.682091490692193
"him" 0.794093181603794 0.513733224301313
"so" 0.307861179477065 0.347780076001732
"this" 0.253835401462984 0.381932752898167
```

Distances

- Time
- Manhattan
- Euclidean



Distance

- dictionary (the, dog, cat, to be, sun)
- text 1 (12, 4, 4, 5, 1)
- text 2 (21, 1, 7, 5, 2)
- text 3 (17, 3, 6, 3, 2)

$$\begin{aligned} d(\text{text1}, \text{text2}) &= |12-21| + |4-1| + |4-7| + |5-5| + |1-2| \\ &= 9 + 3 + 3 + 0 + 1 = \\ &= 16 \end{aligned}$$

Distance - Exercise

- dictionary (the, dog, cat, to be, sun)
- text 1 (12, 4, 4, 5, 1)
- text 2 (21, 1, 7, 5, 2)
- text 3 (17, 3, 6, 3, 2)

Exercise:

Calculate $d(\text{text2}, \text{text3})$.

Distances

- Non-standard words weights adjusted
- Non-standard words weights adjusted
- More importance to MFV
- Modification of Manhattan distance
- All words equal weights

- ▶ Classic Delta:

$$\delta_{(AB)} = \frac{1}{n} \sum_{i=1}^n \left| \frac{f_i(A) - f_i(B)}{\sigma_i} \right|$$

- ▶ Argamon's Delta:

$$\delta_{(AB)} = \frac{1}{n} \sum_{i=1}^n \left| \frac{\sqrt{f_i(A)^2 - f_i(B)^2}}{\sigma_i} \right|$$

- ▶ Eder's Delta:

$$\delta_{(AB)} = \frac{1}{n} \sum_{i=1}^n \left(\left| \frac{f_i(A) - f_i(B)}{\sigma_i} \right| \times \frac{n - n_i + 1}{n} \right)$$

- ▶ Eder's Simple:

$$\delta_{(AB)} = \sum_{i=1}^n \left| \sqrt{f_i(A)} - \sqrt{f_i(B)} \right|$$

- ▶ Manhattan:

$$\delta_{(AB)} = \sum_{i=1}^n |f_i(A) - f_i(B)|$$

CLASSIFY

The screenshot shows an RStudio interface with a single open script file named "stylo.R". The code in the file is as follows:

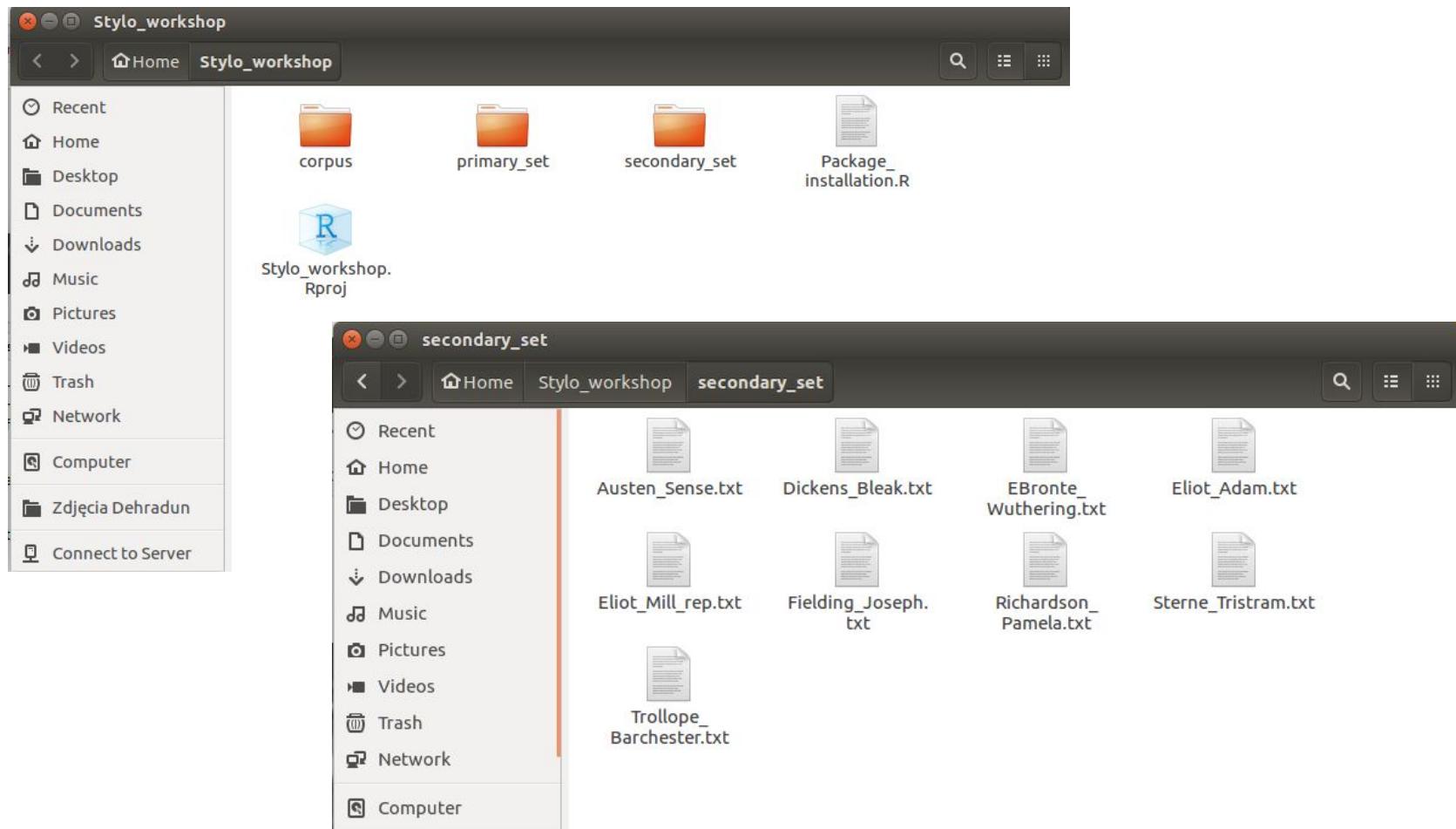
```
1 library(stylo)
2
3 classify()
```

The RStudio interface includes a toolbar with various icons for file operations, search, and help. Below the toolbar is a menu bar with "Run" and "Source" options. The status bar at the bottom shows "1:1 (Top Level)" and "R Script".

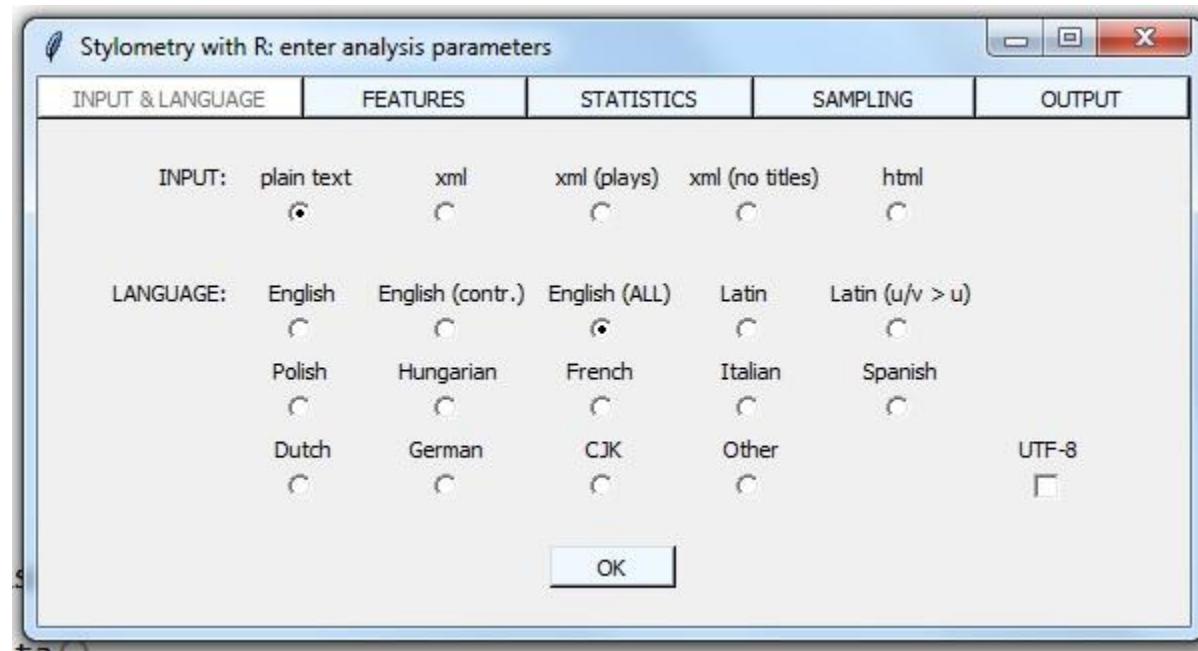
What is classification?



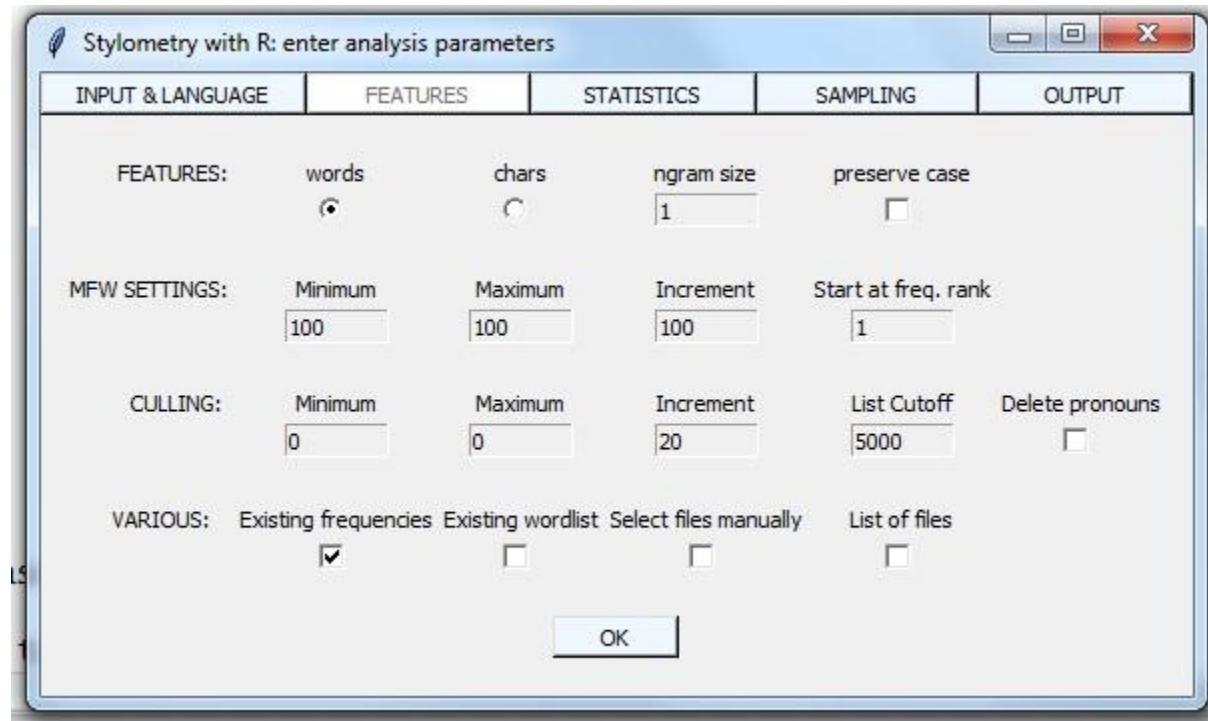
CORPUS PREPARATION FOR CLASSIFY



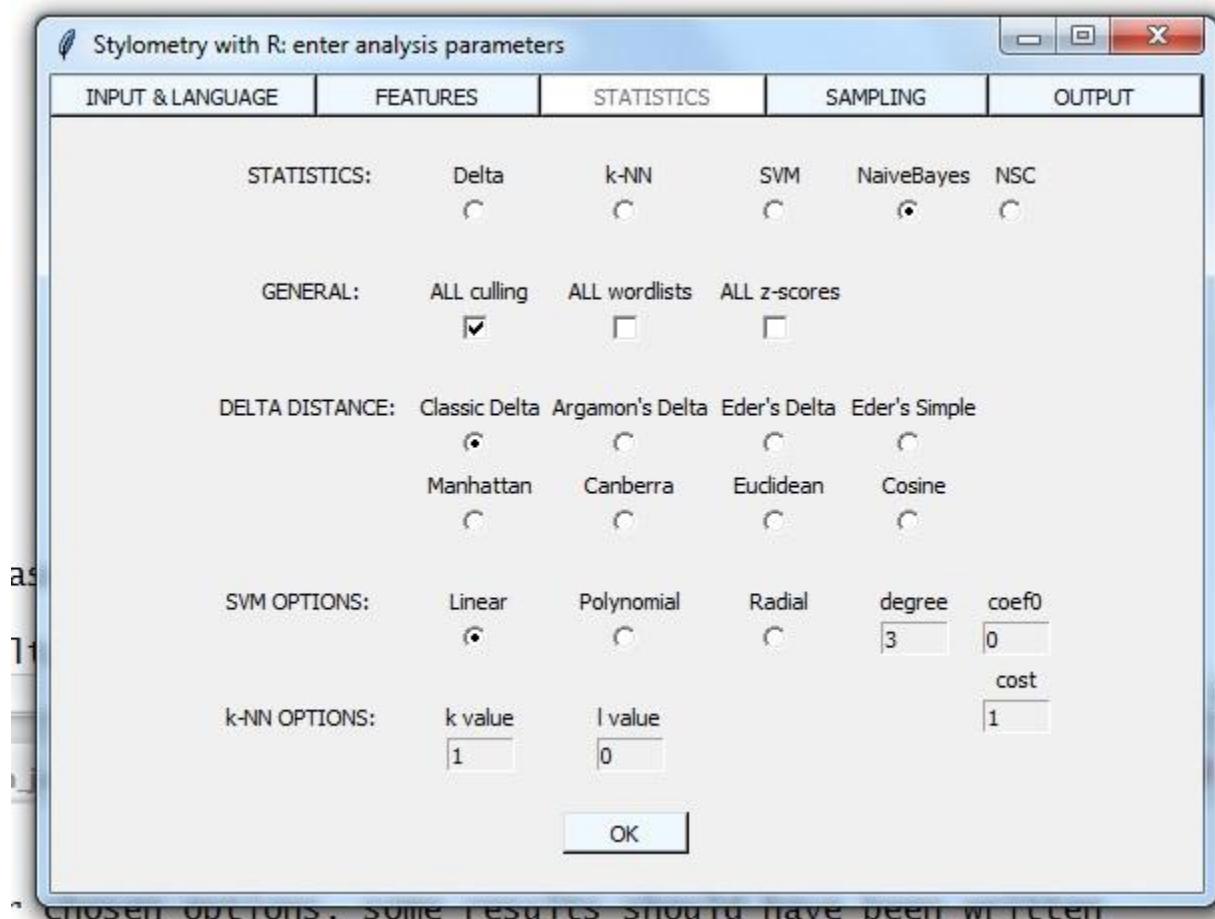
INPUT & LANGUAGE



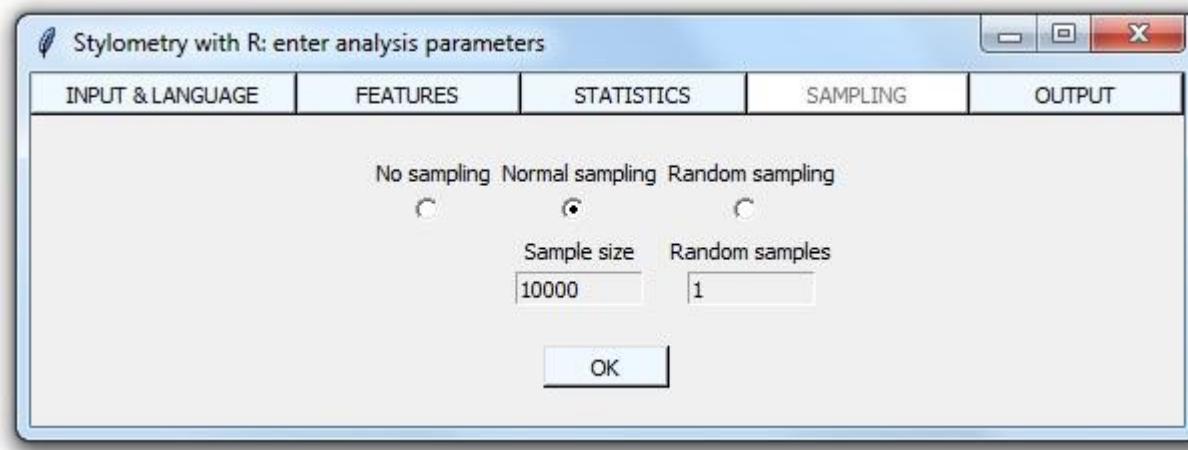
FEATURES



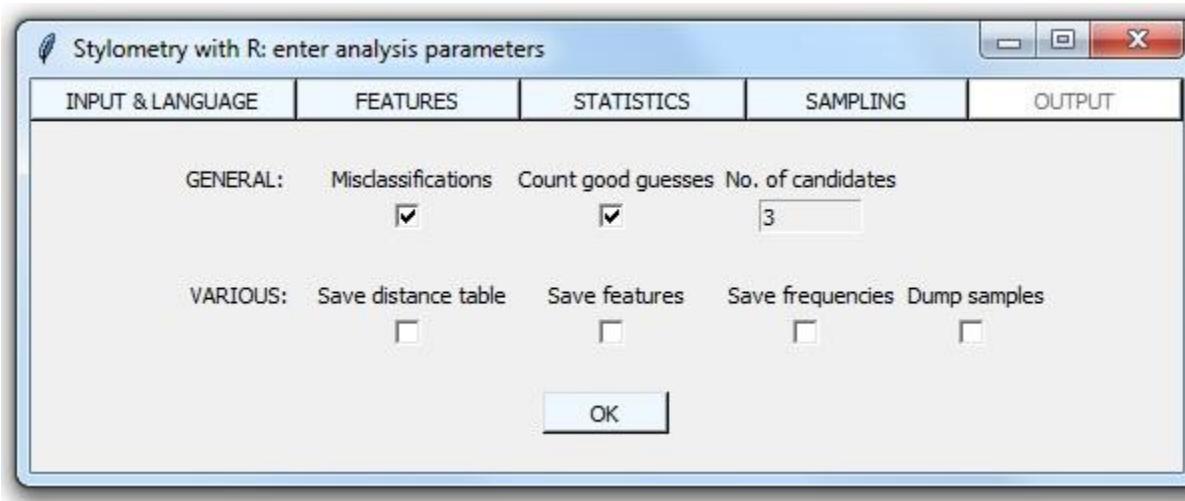
STATISTICS



SAMPLING



OUTPUT



RESULTS FOR CLASSIFY

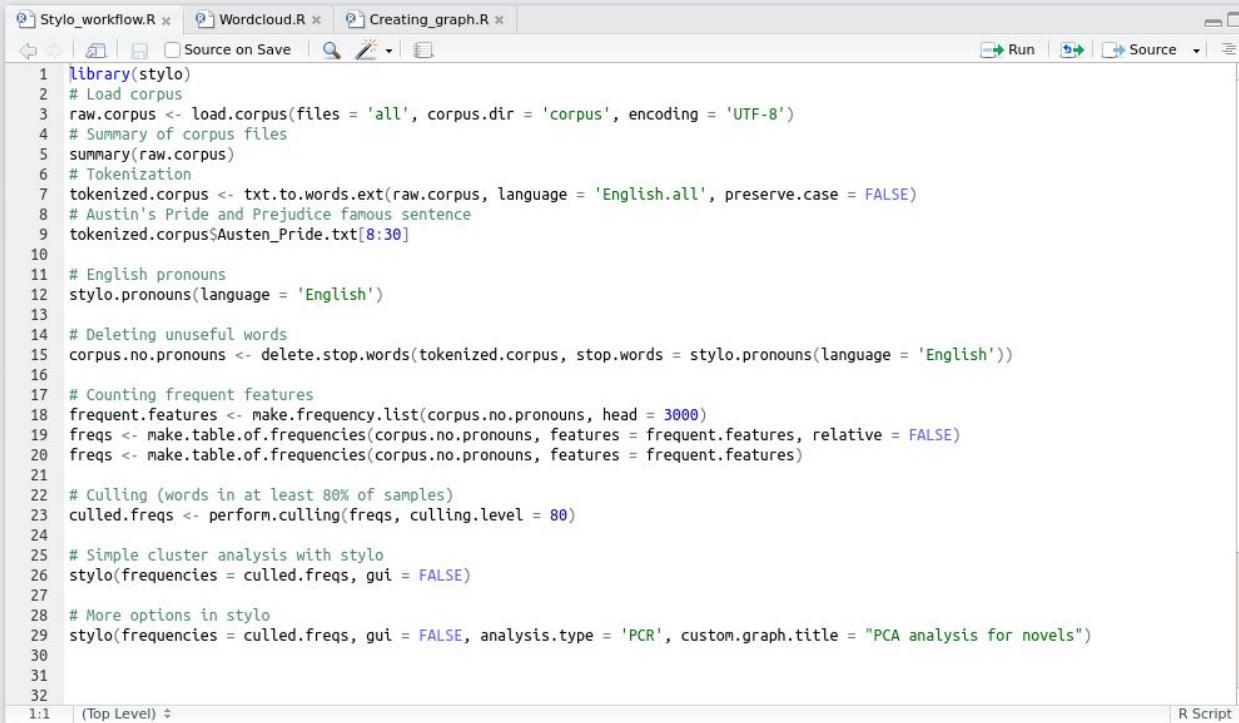
```
final_results - Notepad
File Edit Format View Help

Dickens_Bleak    -->    Eliot
Dickens_David    -->    Eliot
Dickens_Hard     -->    Eliot
Fielding_Joseph   -->    Eliot
Fielding_Tom      -->    Austen
Richardson_Clarissa --> ABronte
Richardson_Pamela  --> ABronte
Sterne_Sentimental --> CBronte
Sterne_Tristram   --> CBronte
Thackeray_Barry   --> CBronte
Thackeray_Pendennis --> CBronte
Thackeray_Vanity   --> CBronte
Trollope_Barchester --> Eliot
Trollope_Phineas   --> Eliot
Trollope_Prime     -->    Eliot

100 MFW , culled @ 0%,  0 of 0  (NaN%)
General attributive success:  0 of 0 (NaN%)
MFWs from 100 to 100 @ increment 100
```

Only misclassified

WAY FORWARD

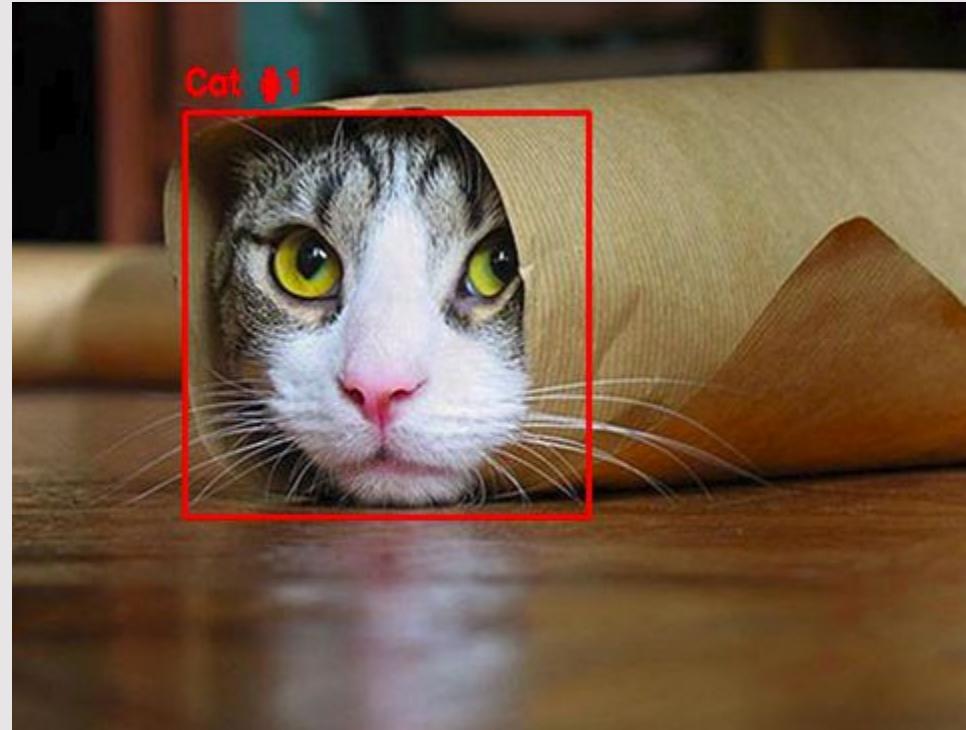


The screenshot shows the RStudio interface with the 'Stylo_workflow.R' script open. The code performs the following steps:

- Imports the 'stylo' library.
- Loads a corpus from files.
- Creates a summary of the corpus files.
- Performs tokenization on the raw corpus.
- Selects a specific sentence from Austin's Pride and Prejudice.
- Identifies English pronouns in the selected sentence.
- Deletes stop words from the tokenized corpus.
- Counts frequent features in the corpus.
- Creates frequency lists and tables.
- Performs culling on the frequencies.
- Performs cluster analysis using the 'stylo' function.
- Creates a PCA analysis graph titled "PCA analysis for novels".

```
1 library(stylo)
2 # Load corpus
3 raw.corpus <- load.corpus(files = 'all', corpus.dir = 'corpus', encoding = 'UTF-8')
4 # Summary of corpus files
5 summary(raw.corpus)
6 # Tokenization
7 tokenized.corpus <- txt.to.words(ext(raw.corpus, language = 'English.all', preserve.case = FALSE)
8 # Austin's Pride and Prejudice famous sentence
9 tokenized.corpus$Austen_Pride.txt[8:30]
10
11 # English pronouns
12 stylo.pronouns(language = 'English')
13
14 # Deleting unuseful words
15 corpus.no.pronouns <- delete.stop.words(tokenized.corpus, stop.words = stylo.pronouns(language = 'English'))
16
17 # Counting frequent features
18 frequent.features <- make.frequency.list(corpus.no.pronouns, head = 3000)
19 freqs <- make.table.of.frequencies(corpus.no.pronouns, features = frequent.features, relative = FALSE)
20 freqs <- make.table.of.frequencies(corpus.no.pronouns, features = frequent.features)
21
22 # Culling (words in at least 80% of samples)
23 culled.freqs <- perform.culling(freqs, culling.level = 80)
24
25 # Simple cluster analysis with stylo
26 stylo(frequencies = culled.freqs, gui = FALSE)
27
28 # More options in stylo
29 stylo(frequencies = culled.freqs, gui = FALSE, analysis.type = 'PCR', custom.graph.title = "PCA analysis for novels")
30
31
32
```

More intuitions on classification



Naive Bayes

What is the probability that this is ABronte given there are such words used?

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

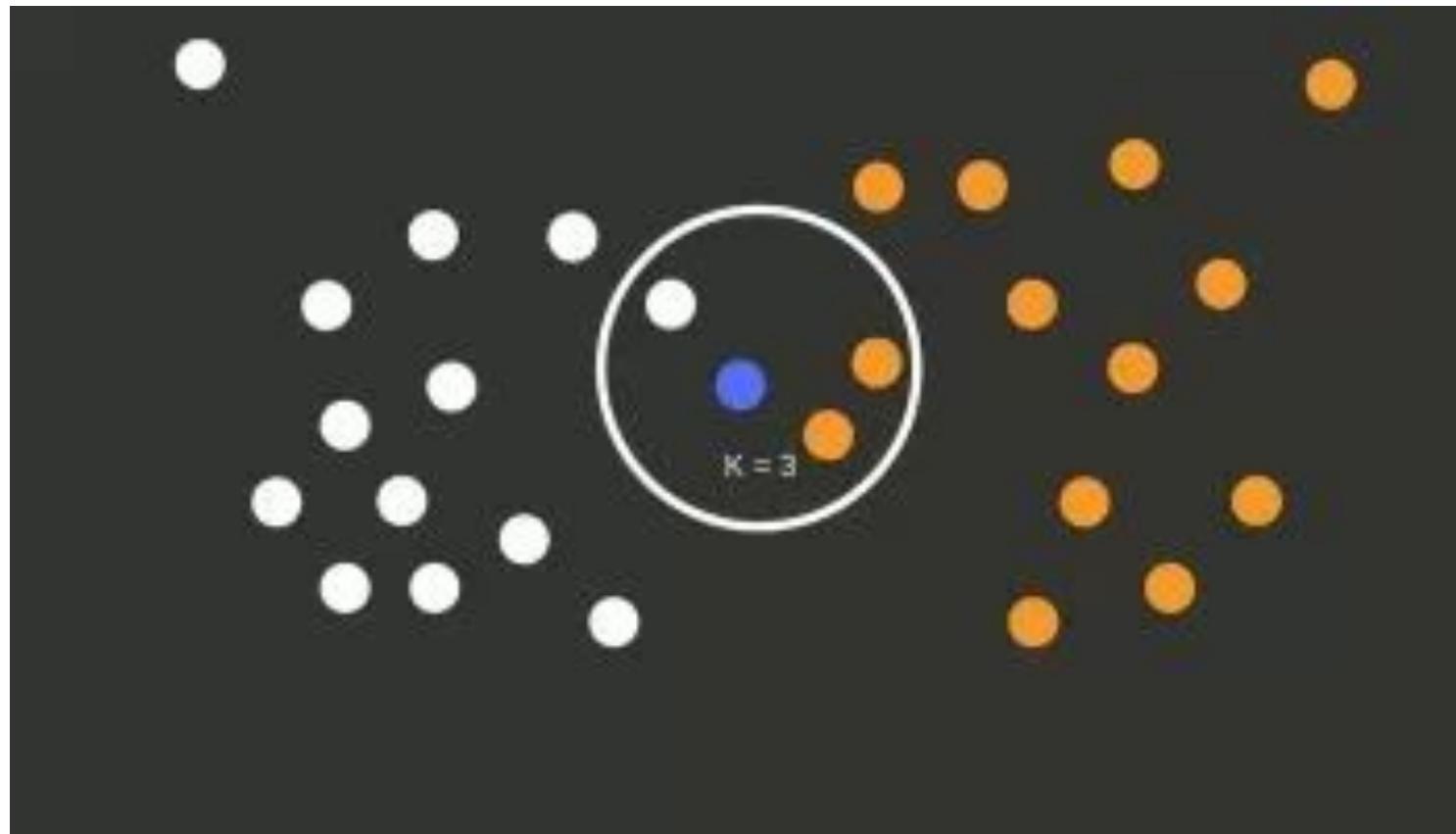
Diagram illustrating the components of the Naive Bayes formula:

- Likelihood: $P(x|c)$
- Class Prior Probability: $P(c)$
- Posterior Probability: $P(c|x)$
- Predictor Prior Probability: $P(x)$

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \cdots \times P(x_n|c) \times P(c)$$

K-nn

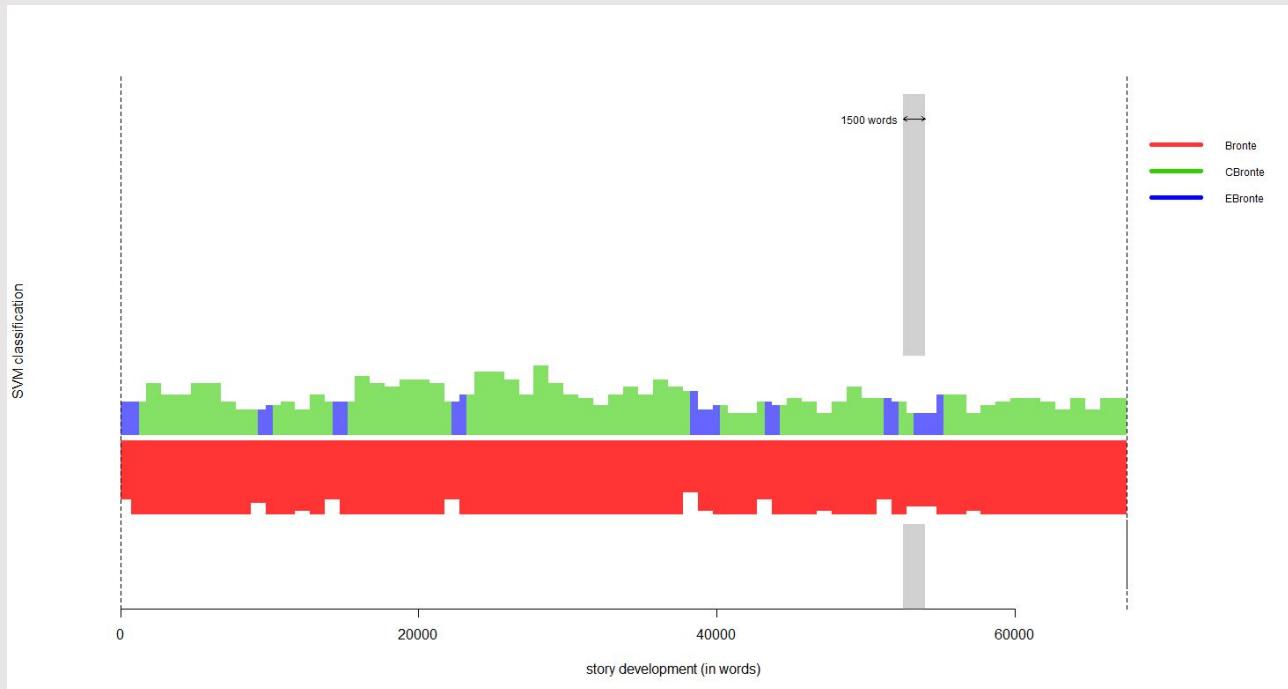
What are the closest observations to ABronte?



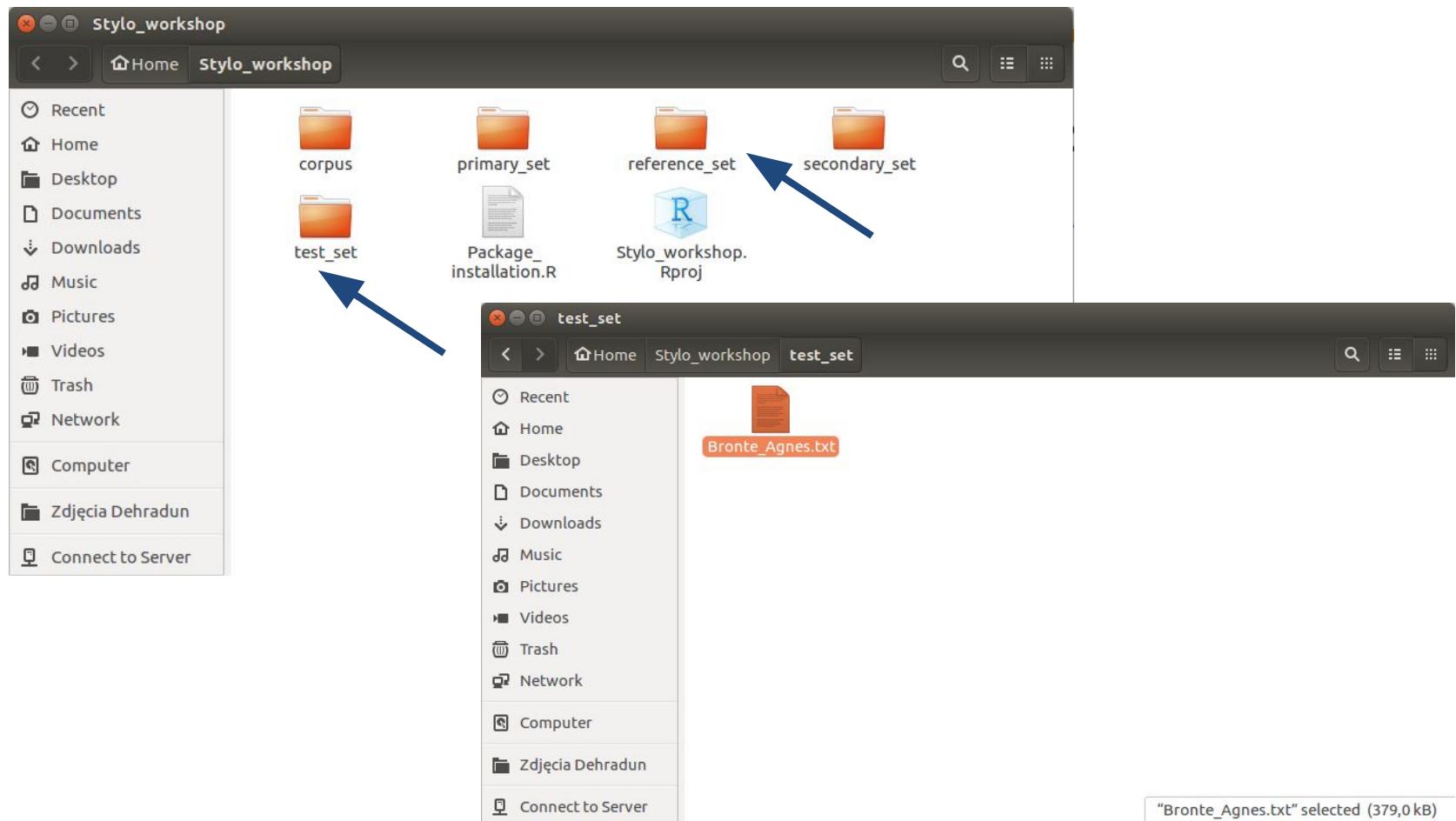
Classify - Exercise

- Remember about correct data set up
- Use 100 English Novels for classify() analysis
- KNN (for different k values), Naive Bayes,
- Use different values for culling and MFW

ROLLING CLASSIFY



Corpus prep for Rolling.classify



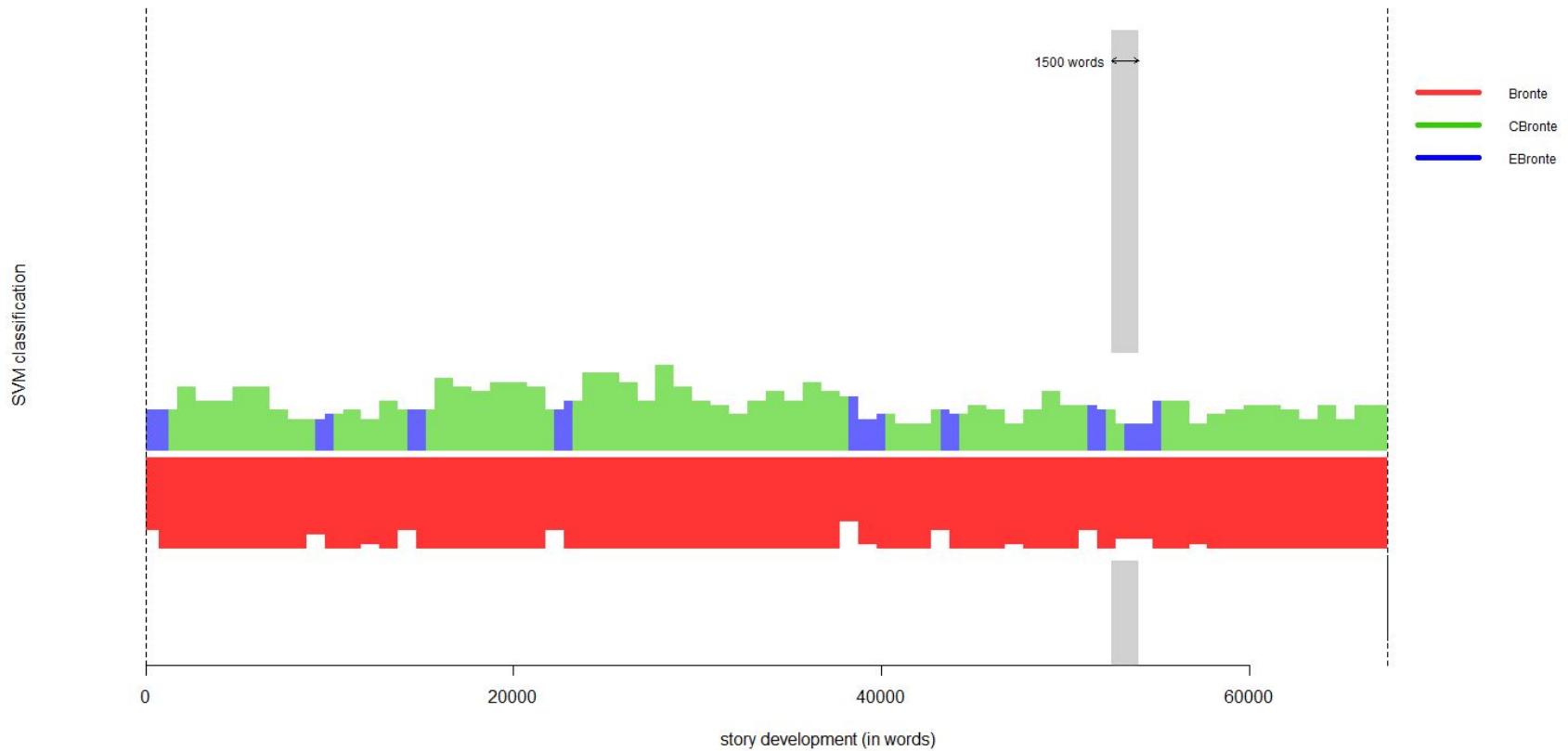
How to use it

```
rolling.classify(slice.size = 1500, mfw = 100,  
slice.overlap = 500, plot.legend = T,  
classification.method = "svm")
```

Rolling.classify

```
> rolling.classify(slice.size = 500, mfw = 100, slice.overlap = 100, plot.legend = T)
using current directory...
The subcorpora will be loaded from text files...
loading Bronte_Wuthering.txt ...
slicing input text into tokens...
turning words into features, e.g. char n-grams (if applicable)...
loading Bronte_Wuthering.txt ...
slicing input text into tokens...
Bronte_Wuthering
  - text length (in words): 116612
  - nr. of samples: 291
  - nr. of words dropped at the end of the text: 212
turning words into features, e.g. char n-grams (if applicable)...
loading Bronte_Agnes.txt ...
loading Bronte_Jane.txt ...
loading Bronte_Professor.txt ...
loading Bronte_Tenant.txt ...
loading Bronte_Villette.txt ...
loading CBronte_Jane.txt ...
loading CBronte_Professor.txt ...
loading CBronte_Villette.txt ...
loading EBronte_wuthering.txt ...
slicing input text into tokens...
turning words into features, e.g. char n-grams (if applicable)...
```

Results of rolling classify



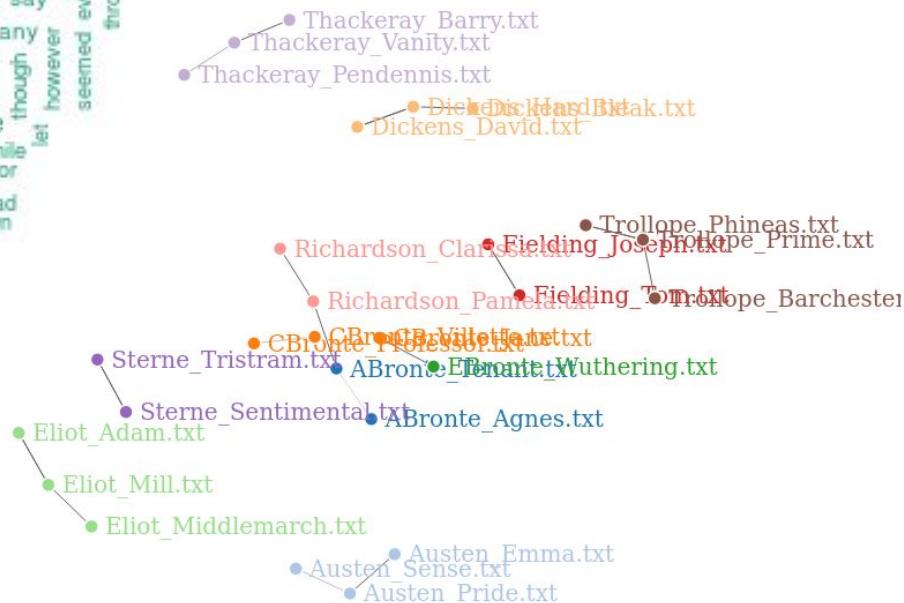
Rolling.classify - Exercise

- Remember about correct data set up
- Use 100 English Novels for rolling.classify()
- Choose 2 authors for reference set and only one novel
- Change values slice.size, mfw, and slice.overlap

Way Forward

- Workflow using code
- Word clouds
- Graph data visualizations
- ...

Visualizations



Download R files

https://github.com/jandziak/Stylo_workshop

Workflow using code

- # Load corpus
- # Tokenization
- # Deleting useless words
- # Creating features
- # Culling
- # Simple cluster analysis with stylo
- # More options in stylo

Load corpus

```
# Load corpus  
raw.corpus <- load.corpus(files = 'all', corpus.dir  
= 'corpus', encoding = 'UTF-8')
```

```
# Summary of corpus files  
summary(raw.corpus)
```

Tokenization

```
# Tokenization  
tokenized.corpus <-  
txt.to.words.ext(raw.corpus, language =  
'English.all', preserve.case = FALSE)
```

```
# Austin's Pride and Prejudice famous sentence  
tokenized.corpus$Austen_Pride.txt[8:30]
```

Deleting useless words

English pronouns

```
stylo.pronouns(language = 'English')
```

Deleting useless words

```
corpus.no.pronouns <-  
  delete.stop.words(tokenized.corpus,  
  stop.words = stylo.pronouns(language =  
  'English'))
```

Creating features

```
# Counting frequent features
frequent.features <-
make.frequency.list(corpus.no.pronouns, head
= 3000)
freqs <-
make.table.of.frequencies(corpus.no.pronouns,
features = frequent.features, relative = FALSE)
freqs <-
make.table.of.frequencies(corpus.no.pronouns,
features = frequent.features)
```

Culling

```
# Culling (words in at least 80% of samples)
culled.freqs <- perform.culling(freqs,
culling.level = 80)
```

Simple cluster analysis with stylo

```
# Simple cluster analysis with stylo  
stylo(frequencies = culled.freqs, gui = FALSE)
```

More options in stylo

```
# More options in stylo  
stylo(frequencies = culled.freqs, gui = FALSE,  
analysis.type = 'PCR', custom.graph.title = "PCA  
analysis for novels")
```

Word clouds

A word cloud centered around the word "the". The word "the" is the largest and most prominent word in the center. Other large words include "and", "of", "be", "with", "it", "for", "had", "that", "to", "a", "was", "as", and "at". The words are colored in various shades of gray, black, red, green, blue, purple, orange, yellow, and pink. The background is white with a subtle grid pattern.

Workflow using code

- # Prepare environment
- # Plot Word cloud

Prepare Environment

```
# Install missing packages  
install.packages("wordcloud") # word-cloud  
generator  
install.packages("RColorBrewer") # color  
palettes  
library("wordcloud")  
library("RColorBrewer")
```

Plot Word cloud

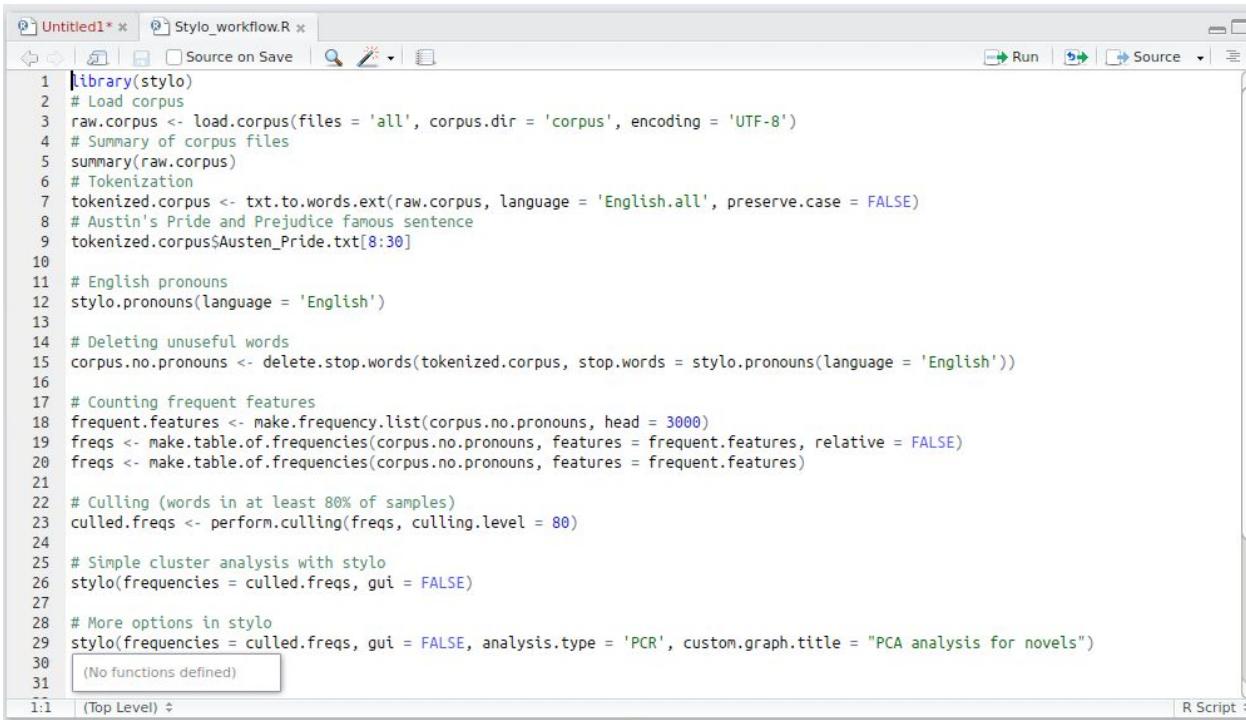
```
# Ensure reproducibility  
set.seed(1234)
```

```
wordcloud(words = names(df), freq =  
as.numeric(df[1,1:2543]), min.freq = 0.01,  
max.words=200, random.order=FALSE,  
rot.per=0.35,  
colors=brewer.pal(8, "Dark2"))
```

Graph data visualizations



Just go to the code



The screenshot shows an RStudio interface with the following details:

- Title Bar:** Untitled1* | Stylo_workflow.R
- Toolbar:** Includes icons for file operations, search, and run.
- Code Editor:** Displays the following R script:

```
1 library(stylo)
2 # Load corpus
3 raw.corpus <- load.corpus(files = 'all', corpus.dir = 'corpus', encoding = 'UTF-8')
4 # Summary of corpus files
5 summary(raw.corpus)
6 # Tokenization
7 tokenized.corpus <- txt.to.words(ext(raw.corpus, language = 'English.all', preserve.case = FALSE)
8 # Austin's Pride and Prejudice famous sentence
9 tokenized.corpus$Austen_Pride.txt[8:30]
10
11 # English pronouns
12 stylo.pronouns(language = 'English')
13
14 # Deleting unuseful words
15 corpus.no.pronouns <- delete.stop.words(tokenized.corpus, stop.words = stylo.pronouns(language = 'English'))
16
17 # Counting frequent features
18 frequent.features <- make.frequency.list(corpus.no.pronouns, head = 3000)
19 freqs <- make.table.of.frequencies(corpus.no.pronouns, features = frequent.features, relative = FALSE)
20 freqs <- make.table.of.frequencies(corpus.no.pronouns, features = frequent.features)
21
22 # Culling (words in at least 80% of samples)
23 culled.freqs <- perform.culling(freqs, culling.level = 80)
24
25 # Simple cluster analysis with stylo
26 stylo(frequencies = culled.freqs, gui = FALSE)
27
28 # More options in stylo
29 stylo(frequencies = culled.freqs, gui = FALSE, analysis.type = 'PCR', custom.graph.title = "PCA analysis for novels")
30
31 (No functions defined)
```

The code performs the following steps:

- Loads the `stylo` package.
- Loads a corpus from files named 'all' in a directory named 'corpus' using UTF-8 encoding.
- Sums up the files in the corpus.
- Performs tokenization on the corpus.
- Selects a specific sentence from the tokenized corpus.
- Identifies English pronouns in the corpus.
- Deletes words that are considered unuseful (stop words).
- Counts the most frequent features (words) in the corpus.
- Culls words that appear in at least 80% of the samples.
- Performs a simple cluster analysis using the `stylo` function.
- Provides more options for the cluster analysis, including PCA analysis for novels.

Way forward - Exercise

- Apply stylo code analysis on new data
- Use 100 English Novels
- Produce new word cloud using different novel
- Create graph for novels