# STYLOMETRY WITH R DETAILED INSTRUCTIONS

JOANNA BYSZUK IJP PAN, JAN RYBICKI, UJ

2018

# STYLO
# MAIN FUNCTIONS

- stylo()
  - Calculates distances (differences) between series of most frequent words and draws graphs of those distances
    - CLUSTER ANALYSIS trees (for a single set of parameters)
    - BOOTSTRAP CONSENSUS trees (for multiple parameter settings)
    - MULTIDIMENSIONAL SCALING maps
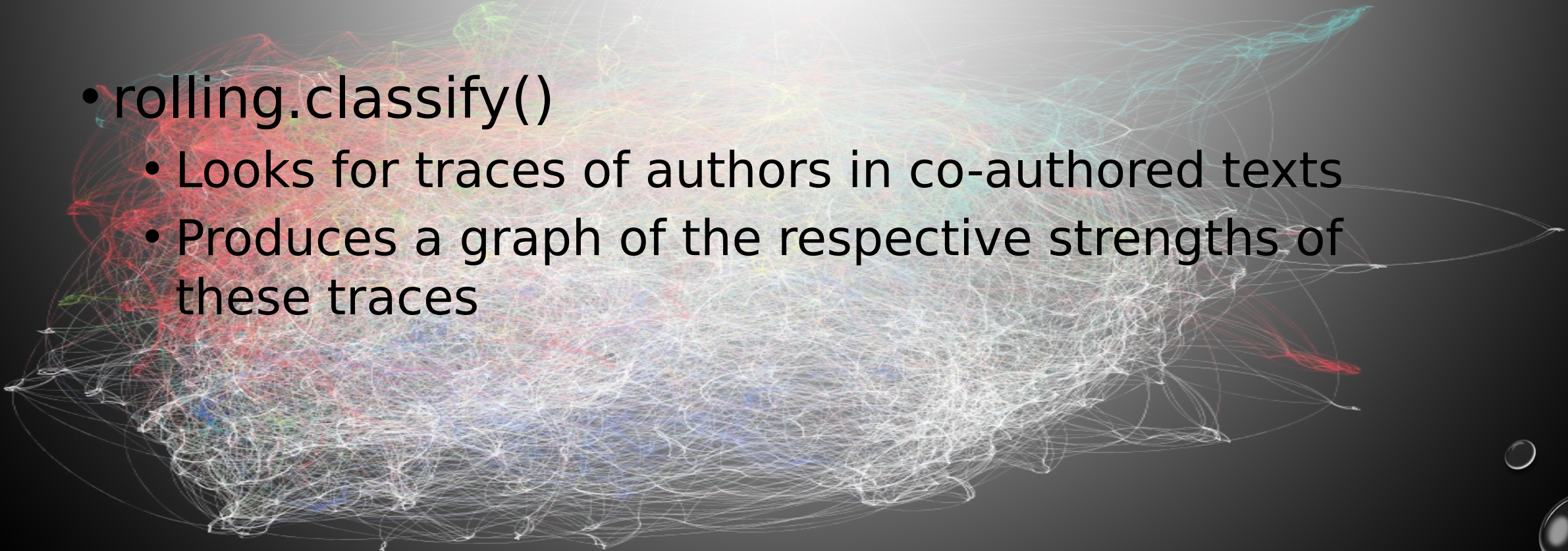    - PRINCIPAL COMPONENTS ANALYSIS maps

# STYLO
# MAIN FUNCTIONS

- oppose()
  - Cuts texts into equal-sized samples
  - Finds words characteristic for two (groups) of texts
    - These can be reused with stylo()
  - Produces a diagram of the use of each group's words

# STYLO
# MAIN FUNCTIONS

- rolling.classify()
  - Looks for traces of authors in co-authored texts
  - Produces a graph of the respective strengths of these traces

# RUNNING STYLO()

- Where are my texts?:
  - MENU:
  - FILE > CHANGE DIRECTORY >
  - E.G. English Benchmark etc

- (it contains the subfolder „cor

- but don't go there!)

- library(stylo) <ENTER>

- stylo() <ENTER>

# STYLO() PARAMETERS

- INPUT: STATE YOUR TEXTS' FO

- LANGUAGE

- DON'T PRESS „OK" YET!!!

# STYLO() PARAMETERS

- FEATURES: THINGS TO COUNT: WORDS OR CHARACTERS
  - ngram size: COUNT SINGLE FEATURES (1) OR THEIR CLUSTERS (>1)

- MFW SETTINGS: HOW MANY MOST FREQUENT WORDS TO USE

  UNLESS WE USE „bootstrap consensus tree" IN STATISTICS, Minimum=Maximum

# STYLO() PARAMETERS

- CULLING: MANIPULATING THE WORDLIST (0)
  - 0%: NO WORDS ARE REMOVED
  - 100%: ALL WORDS ARE REMOVED THAT DO NOT OCCUR IN ALL THE TEXTS

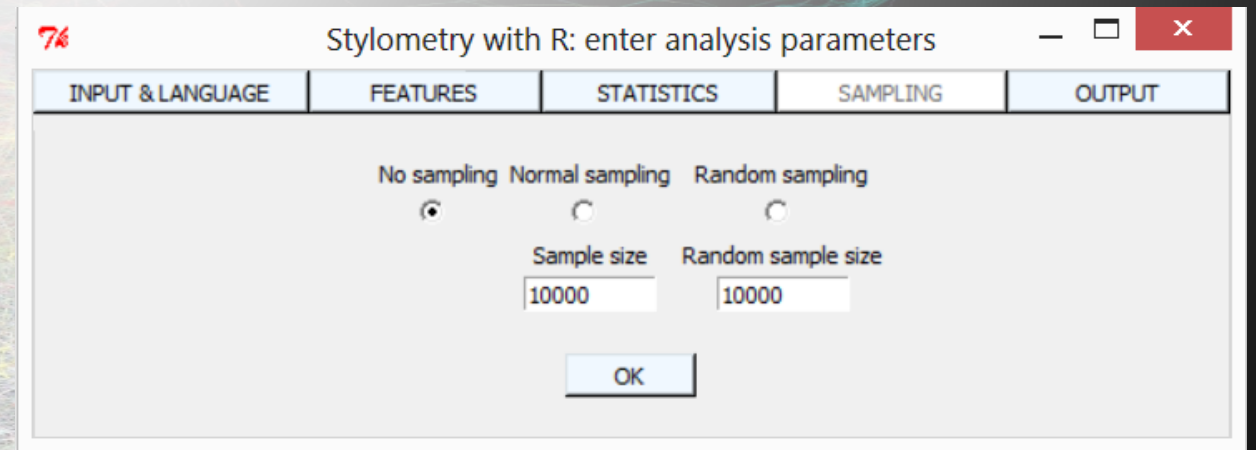- DELETE PRONOUNS?

- DON'T PRESS „OK" YET!!!

# STYLO() PARAMETERS

- STATISTICS: PICK STATISTICS METHOD (Cluster Analysis)

- DISTANCES: TYPE OF DISTANCE MEASURE (Classic Delta)

- DON'T PRESS „OK" YET!!!

# STYLO() PARAMETERS

- SAMPLING: (No sampling)
  - DO I WANT TO SAMPLE THE TEXTS
  - AND HOW

- DON'T PRESS „OK" YET!!!

# STYLO() PARAMETERS

- OUTPUT: (Onscreen)
  - GRAPH FORMAT ETC.
- PRESS „OK"!!!
- …WAIT FOR IT…

# RUNNING GEPHI

- FIRST WE NEED TO RUN STYLO:
  - It creates an OUTPUT FILE named e.g.

50 British Novels_100-1000_MFWs_Culled_100_Pronouns deleted_Classic Delta_C_0.5_EDGES.csv

# RUNNING GEPHI

- SELECT GEPHI>NEW PROJECT

- Data laboratory>Import Spreadsheet

- Import settings:
  - Separation: Comma
  - As table: Edges table
  - Charset: UTF-8? Windows-1252...

- Don't worry about this being somewhat illogical...

- Next

# RUNNING GEPHI



- Change „Weight" to „Double"

- Hit „Finish"

# RUNNING GEPHI

# RUNNING GEPHI



- Change „Edges merge strategy to „Average"
- Set „Append to existing workspace"

- Hit „OK"

# RUNNING GEPHI

- We need to get authors' names…

- Click

- And set:

- OK!

- Copy ID column to LABEL

# RUNNING GE[PHI]

- PREVIEW

- PARTITION
  - CLICK
  - SELECT: e.g. Author
  - Apply

- Show labels

- LAYOUT
  - ForceAtlas 2
    - Dissuade Hubs
    - Prevent Overlap
    - Edge Weight Influence 0.5
    - Scaling: 500
    - RUN!

- LAYOUT (CONT)
  - Expansion
    - RUN!

Gephi 0.9.1 - Network_Italiano_Consensus_100-100_MF

File  Workspace  Tools  Window  Help

Overview    Data Laboratory    Preview

Workspace 1  ✕

Appearance  ✕

Nodes  Edges

Unique  Partition  Ranking

---Choose an attribute

Graph  ✕

Dragging  (O

Arial Bold, 32

# RUNNING GEPHI

- OVERVIEW

- NODE LABELS
  - SHOW LABELS

- EDGES
  - SHOW EDGES
  - Thickness: np. 0.1, 0.01...

- REFRESH!

- Reset zoom

- SAVING
  - NETWORK:
  - File > Save
    - WITH .gephi EXTENSION
  - PICTURE:
  - File > Export > SVG/PDF/PNG
  - Options > Landscape
  - PHEW!

# RUNNING OPPOSE

- DIFFERENT SUBFOLDER STRUCTURE:
  - primary_set
  - secondary_set
  - test_set (OPTIONAL)

- library(stylo) <ENTER>

- oppose() <ENTER>

- What we get:
  - WORDS_PREFERRED characteristic for the primary_set texts
  - WORDS_AVOIDED 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)

# RUNNING ROLLING.CLASSIFY

- DIFFERENT SUBFOLDER STRUCTURE (AGAIN):
  - reference_set (individual writings)
  - test_set (collaborative text)

- library(stylo) <ENTER>

- rolling.classify(write.png.file = TRUE, classification.method = "delta",mfw=100, training.set.sampling = "normal.sampling", slice.size = 5000,slice.overlap = 4500)

- What we get:
  - Similarity graph