# Stylometry with R

mike.kestemont@gmail.com
www.mike-kestemont.org
University of Antwerp

NEH Summer Institute, Pittsburgh
July 2017

# Installing R

- **R** (http://www.r-project.org/)

- Open-source statistical software

- Runs on all major platforms

- Install instructions: http://cran.freestatistics.org/

# Sublime 3

- For viewing files today

- If you don't have a good text editor

- (esp. if you are on Windows)

- Install Sublime 3

- Free download

- Install: http://www.sublimetext.com/2

# Stylo

- "Stylometry with R"

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

- Free package for easy stylometric analysis in R

- Graphical user interface (no coding!)

# There's no I in team



Maciej

Jan

# Install Stylo

- Install from within R

- Launch R: double-click icon (e.g. in `Applications`)

- To download and install, type in the console:

    - `install.packages("stylo")`

- Every time you restart R, import Stylo:

    - `library(stylo)`

# Download course material

- Download course materials from:

  - tinyurl.com/y73tc2es

- Unzip the folder (pitt17)

- Place it e.g. on your Desktop

# Medieval French Genres

- Jean Bodel (French poet, late 12th C.)

- Famous quote *Chanson de Saisnes*:

   *Ne sont que 3 matières à nul homme atandant,*
   *De France et de Bretaigne, et de Rome la grant.*

- Distinguishes 3 *matières* or "genres":
   1. *Matière de France* (chansons de geste; Charlemagne)
   2. *Matière de Bretaigne* (romans arturiens; King Arthur)
   3. *Matière de Rome* (romans antiques; e.g. Troie)

- Question: can we distinguish these using stylometry?

# Clustering in Stylo

- Let's do a clustering experiment on our genres

- Create a folder `corpus` under `pitt17/data/genres/`

- Copy all `bre_*` and `fra_*` texts to this folder

# Run stylo

- Stylo needs to know where our data is. Type in R:

  - `setwd(`"`~/Desktop/pitt17/data/genres/`"`)`

  - (You can use tab to navigate!)

  - It has to see `corpus` (and not be inside it!)

- Make sure stylo is loaded:

  - `library(stylo)`

- Run command:

  - `stylo()`

- The GUI should load…

# Stylo GUI

# Adjust parameters and hit OK

# And you should get a tree...

# OK… What happened?

- We represent texts as "bags of words"

- Create a large frequency table:

  - each column = text

  - each row = word

  - each cell = relative frequency

- Check out `table_with_frequencies.txt`

| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 1 | bre_charette | bre_graal | bre_graalcon | bre_graalcon | bre_tristanbe | bre_tristanth | bre_yvain | fra_aliscans | fra_asp |
| 2 | le | 6,21921182 | 7,19456215 | 6,55339806 | 6,69576325 | 6,70977187 | 5,52029697 | 5,1766025 | 7,7381 |
| 3 | il | 6,052076 | 5,45550847 | 4,66221683 | 5,00319353 | 4,22105648 | 5,00538858 | 5,2139787 | 3,8564 |
| 4 | en | 4,934905 | 4,12252825 | 3,65088997 | 3,93868427 | 5,27022081 | 5,69991618 | 4,84956083 | 4,0917 |
| 5 | et | 4,97888811 | 4,90819209 | 4,55097087 | 4,84351714 | 1,93973405 | 3,59238415 | 4,56923939 | 3,2935 |
| 6 | de | 2,33990148 | 2,65713277 | 3,00364078 | 3,21481797 | 2,41551787 | 3,41276494 | 3,07419174 | 3,6716 |
| 7 | que | 3,80893737 | 4,05190678 | 3,09466019 | 2,89546519 | 2,8059046 | 3,46066339 | 4,35432629 | 1,7896 |
| 8 | a | 1,77691766 | 1,65077684 | 2,05299353 | 2,00127741 | 2,48871538 | 1,98778589 | 1,78471314 | 2,4281 |
| 9 | je | 2,95566502 | 2,83368644 | 1,69902913 | 1,4477326 | 2,90350128 | 3,04155191 | 2,99009531 | 1,6131 |
| 10 | avoir | 1,77691766 | 1,40360169 | 2,18446602 | 2,05450287 | 2,18372575 | 2,74218656 | 1,53242385 | 2,5793 |
| 11 | i | 2,7885292 | 2,6924435 | 2,52831715 | 0,52160954 | 0,59777968 | 0,37121303 | 2,89665483 | 0,6805 |
| 12 | estre | 0,94123856 | 1,32415254 | 1,83050162 | 1,52224824 | 1,43955106 | 0,76637528 | 1,20538217 | 1,688 |
| 13 | si | 1,76812104 | 1,43008475 | 1,56755663 | 1,69256973 | 1,06136391 | 1,02981679 | 1,9342179 | 1,0082 |
| 14 | son | 0,89725545 | 1,0240113 | 0,95064725 | 1,45837769 | 1,1589606 | 2,02370974 | 0,97178098 | 1,5291 |
| 15 | qui | 1,18754398 | 1,63312147 | 1,77993528 | 1,34128167 | 1,18335977 | 0,56280685 | 1,5043917 | 1,0082 |
| 16 | ester | 1,16995074 | 1,14759887 | 1,53721683 | 1,52224824 | 1,09796267 | 1,36510598 | 1,36423098 | 1,176 |
| 17 | vos | 1,25791696 | 1,35063559 | 1,12257282 | 1,0112838 | 1,17116018 | 1,43695366 | 1,23341432 | 0,7981 |
| 18 | dire | 1,12596763 | 1,34180791 | 1,0315534 | 0,88354269 | 0,97596682 | 1,05376602 | 0,76621192 | 0,8738 |
| 19 | se | 1,12596763 | 1,14759887 | 1,09223301 | 1,20289547 | 1,11016225 | 1,44892827 | 1,36423098 | 0,8149 |
| 20 | un | 1,14356087 | 1,06814972 | 1,15291262 | 0,89418778 | 1,02476516 | 0,50293378 | 0,77555597 | 1,2939 |
| 21 | par | 0,80928923 | 0,71504237 | 1,00121359 | 0,68128593 | 1,28095645 | 0,85019758 | 0,72883573 | 1,0082 |
| 22 | tot | 0,82688248 | 0,76800847 | 0,82928803 | 1,28805621 | 0,76857387 | 0,50293378 | 0,99046907 | 0,9326 |

# Bag of words?

- We **ignore** word order, position of word in document, syntax, …

- Only use word counts

- Relative frequencies

# Only use 3,000 words

- Most Frequent Words: `MFW`

- Better for statistics

- Check out `wordlist.txt`

- What kind of words are most frequent?

| MFW SETTINGS: | Minimum | Maximum |
|---|---|---|
| | 3000 | 3000 |
| CULLING: | Minimum | Maximum |
| | 0 | 0 |

# Distance matrix

| Dist() | Text1 | Text2 | Text3 |
|--------|-------|-------|-------|
| Text1 | 0.0 | Dist(Text1, Text2) | Dist(Text1, Text3) |
| Text2 | Dist(Text2, Text1) | 0.0 | Dist(Text2, Text3) |
| Text3 | Dist(Text3, Text1) | Dist(Text3, Text2) | 0.0 |

# Build tree

- Now we build a tree bottom-up

- First, join 2 texts that are most similar

- Combine them in a new node

- Work you way up the three

- Until all texts are joined

- Horizontal axis reflects (dis)similarity

# Do it yourself (1)

1. Try out different parameters:

   - Vary the number of `MFW` (under `features` tab): 30, 50, 1000, 5000, … (Always update `Minimum` and `Maximum` simultaneously!)

   - Vary the `distance metric` (under `statistics` tab)

   - Do you get different results? "Better" results?

- *Graal*, *Yvain* and *Charette* always cluster together. Can you think of an explanation why?

# Do it yourself (2)

1. Under the `sampling` tab, select `Normal sampling` and insert 3,000 under `Sample size`.

2. Run the analysis again. There are much data points now: can you guess what happened?

3. Set the `Sample size` at an absurd size: e.g. 20,000. Do you get an error? Why?

# Unstability

- Cluster Analyses can be unstable (cf. 30 > 31 `MFW`)

- Very different results for small change in parameters

- Rerun experiment with for `MFW: Minimum=50, Maximum=3000, Increment=50`

- We now iteratively run cluster analyses for different frequency bands: `50-100 MFW, 100-150 MFW, 150-200, …, 2900-2950 MFW, 2950-3000 MFW.`

- Do you see the tree change in each picture?

# Bootstrap Consensus Trees

- Bootstrap Consensus Trees (BCT)

- Gives "summary" of different cluster analyses

- Only visualises nodes on which there is a consensus among the trees (50% majority vote)

- Rerun analysis, but select `Consensus Tree` (under `statistics`), but leave `Consensus strength` to `0.5`

**genres**
**Bootstrap Consensus Tree**

50-3000 MFW  Culled @ 0%
Classic Delta distance Consensus 0.5

# Do it yourself (+)

- We have seen that the cluster analyses easily distinguish Jean Bodel's *matière de Bretaigne* and *matière de France* without supervision. But what about the *matière de Rome*? Add the `rom_*` texts under data to the `corpus` folder.

- Rerun various cluster analyses on this expanded data set and experiment with the BCT. Experiment with different MFWs and sample sizes. What is the result? Do you get pretty clusters? How do you interpret this? Which two Arthurian texts behave strangely?

genres
Bootstrap Consensus Tree

50-3000 MFW  Culled @ 0%
Classic Delta distance Consensus 0.5

# Text selection

- Sometimes you don't want to analyse all texts under `corpus`

- Under `features`, tick `Select files manually`

- You will get a dialogue window:

  - (De)select individual texts using `Control+Click`

  - Select a range of texts using `Shift+Click`

- Try to run an analysis using only the `bre_*` and `rom_*` texts

# Do it yourself (1)

- I downloaded the entire oeuvre by Dante Alighieri (1265-1321) from <u>danteonline.it</u>

- (I don't know anything about Dante, and I don't speak Italian)

- Still, analyse his oeuvre: "Distant" Reading!

- Type `setwd("pitt17/data/dante")` in R to navigate to the correct directory

# Do it yourself (2)

- Run various (normal) cluster analyses on Dante's work: try different `MFW`s. (Don't use sampling yet: `No sampling`) Do you see a clear clustering of texts?

- Analyze these two clusters using `oppose()`. Don't forget to create the folders necessary for this: divide the texts in a `primary` and `secondary set`. Result? Silly me! Can you too find out why these two clusters are there?

- Add cluster labels followed by "_" in the file names under corpus to sort our the colouring of the cluster plots. Each file should get a title = `clustername_title.txt`

**dante**
**Bootstrap Consensus Tree**

Convivio 7
Convivio 8
Convivio 6
Convivio 3
Convivio 4
Convivio 5
Convivio 9
Convivio 1
Convivio 2
Convivio 8
Convivio 14
Convivio 13
Convivio 12
Convivio 11
Convivio 8
Convivio 9
Convivio 10

Fiore 1
Fiore 3
Fiore 4
Fiore 2

Rime 3
Rime 2
Rime 1

VitaNuova 1
VitaNuova 3
VitaNuova 2

Paradiso 2
Paradiso 1
Paradiso 6
Paradiso 5
Paradiso 4
Paradiso 3

Inferno 4
Inferno 5
Inferno 6
Inferno 2
Inferno 3
Inferno 1

Purgatorio 2
Purgatorio 1
Purgatorio 5
Purgatorio 4
Purgatorio 3
Purgatorio 6

100-2000 MFW  Culled @ 0%
Classic Delta distance Consensus 0.5

# Do it yourself (3)

- Now analyse <u>only</u> the Italian works using `stylo()`.

- Now run Bootstrap Analysis Trees for various `MFWs` (adjust `Minimum`, `Maximum` and `Increment`).

- Try out different sample sizes (e.g. 5,000). You can leave out *DettoDAmore*, which is too short. Do you see clusters here? Can you explain them using the internet?

- Which two parts of the *Commedia* are closest to each other?

- Use `oppose()` to find out which words are typical of `Paradiso` (in comparison to the other parts).

# Spelling variation

- No printing press: manual copying
- Scribes, copyists
- No standard language, spelling
- Regional, personal preferences
- Especially vernacular texts
- Each copy unique

# Recognizable?

| | |
|---|---|
| D | Ter stont ende ter seluer vren |
| E | Tier stont ende ter seluer vren |
| F | TIere stont enter seluer vren |
| G | Tottien stonden en ter uren |
| H | TEn stonden ende ter seluer vren |
| I | Tjerst stont ende tier veren |
| J | Tyer stont ende tier seluer vren |
| N | TJer stont tier seluer vre |

# Huge issue

- Issue for computational text analysis

- Lemmatize, part-of-speech tag

- Often seen as problem…

- E.g. stemmatology: reconstruction

- But also interesting!

- Study scribal behaviour

# Angus McIntosh

- Middle English philology

- *Linguistic Atlas of Late Medieval English*

- Scribal language

- Interested in modelling scribal behaviour

# Hypothesis

- Each scribe has unique 'profile'
- Combination of:
  - Graphetic profile (*handwriting*)
  - Linguistic profile (*language*)
- Today focus on language:
  - *alt* vs. *olt* (dialect)
  - *tijt* vs. *tyt* (spelling)

# 3. Chaucer

- Scribal profile in 4 MSS

- Chaucer, *Canterbury Tales*

- Well-studied scribes

- Parallel copies of 1 tale

- *The Man of Law*

- Data courtesy of J. Thaisen

# Parallel content:
## focus on linguistic differences

hateful harmN condiciounN of povert

with thrist with cold with hungR so counfoundid

ohatefull harme condicyouN of pouert

with thurste witħ colde witħ hungR so coNfounded

O hate full harme condiciouN of pouerte

wt thrust wt colde and honger so confounded

# Principal Components Analysis

- (My favourite)

- enter `setwd("~/Desktop/pitt17/data/chaucer")`

- Check out `corpus` folder

- launch `stylo()`:
    - Use `MFW=500`
    - Set method to `PCA (corr.)` under `statistics`
    - `Normal sampling; size=500 (Sampling)`

**Stylometry with R: enter analysis parameters**

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

STATISTICS: Cluster Analysis    MDS    PCA (cov.)    PCA (corr.)    tSNE

Consensus Tree    Consensus strength

0.5

DISTANCES:   Classic Delta    Argamon's Delta    Eder's Delta   Eder's Simple

Manhattan    Canberra    Euclidean

OK

# What do you see?
## (Does this make sense?)

# Do it yourself

- Select a number of different manuscript *pairs* and *triples* tick `Select files manually` and use `Control+Click`. Can you describe what you see? Where are the samples positioned?

- Use 2 manuscripts. Set `PCA flavour=Symbols` and steadily decrease the sample size (500, 300, …, 50,). How small can samples get before the plot gets fuzzy? What does this tell us?

# Loadings

- Extremely helpful feature of PCA

- Tells on which specific word differences the PCA is based

- Use 2 manuscripts. Set `PCA flavour= Loadings.` The loadings will be plotted in dark; the samples in lightgrey. (If difficult to read, lower the `MFW=100`)

- What is there results? Inspect the original files: do the loadings make sense?

**chaucer**
**Principal Components Analysis**

PC2 (8%)

PC1 (20.2%)
150 MFW  Culled @ 0%
Correlation matrix

PCA flavour=

Technical

# Character n-grams

- Words are not always used in stylometry

- Also character n-grams

- Under `features` **tab:**
  - `features = chars`
  - `ngram size = 3`

- Make sure to set `PCA flavour=Loadings`

- Can you guess what character n-grams are?

**chaucer**
**Principal Components Analysis**

150 MFC 3-grams Culled @ 0%
Correlation matrix

PC1 (11%)
PC2 (9.7%)

Quartz 2 [*]

# Hildegard of Bingen

- Influential women writer

- 1098–1179

- Germany

- Divine visions

- "Sybil of the Rhine"



[Wiesbaden, Landesbibliothek, 1, fol. 1r.]

# Varied oeuvre

- Visions

- Music

- Scientific texts

- Recipes

- Medical treatises

- Letters (pope, emperor, ...)



[Dendermonde, St.-Pieters & Paulusabdij, Ms. Cod. 9]

# Early 2012

- Sara Moens

- Jeroen Deploige

- Dept. History, UGhent

- Editing two texts

- Collaborate?

# Secretaries

- Wrote in Latin...

- But was bad at it!

- No formal training as woman

- Assisted by male secretaries

- Gender issues...

- Dictated



[Hildegard and her 1st secretary Volmar]

# Correction grammatical mistakes

(Only form, not content!)

# Two shorter texts...

*Visio ad Guibertum missa* &

*Visio de sancto Martino*

- "Attributed" to Hildegard

- *Opera omnia...*

- But style not typical of her

- Doubts authorship?

- Last secretary....

[MS Brussels, Royal Library, 5527-34, fol. 141v. ]

# Guibert of Gembloux

- Monk from Brabant

- Hildegard's last secretary

- Fascination St Martin

- Very elaborate style

- "Pushy"

[MS Brussels, Royal Library, 5527-34, fol. 141v. ]

# Stylometry?

When you correct [this text], keep to this rule: that [...] you apply your skill only to make corrections where the order or the rules of correct Latin are violated. Or if you prefer – and this is something I have conceded in this letter beyond my normal practice – you need not hesitate to clothe the whole sequence of the vision in a more becoming garment of speech, preserving the true sense

[*Visio de St. Martino*, trans. Newman, 1987, p. 23]

# Corpus

- *Corpus Christianorum* (Brepols)

- Complete materials

- *Epistolaria*

- Hildegard, Guibert

- Bernard of Clairvaux (1090-1153)

- 3x +100k tokens

# Do it yourself

- Check out folder `pitt17/data/hildegard`:

  - `B_ep.txt` = Letters from Bernard of Clairvaux

  - `B_Mart.txt` = Sermon about St. Martin by Bernard

  - `D_Mart.txt` = Dubious

  - `D_Missa.txt` = Dubious

  - `G_ep.txt` = Letters by Guibert

  - `H_epG.txt` = Letters by Hildegard, with Guibert

  - `H_epNG.txt` = Letters by Hildegard, before Guibert

- All texts lemmatised

# Wordlist

- Restrictive wordlist `wordlist_master.txt`

- Non-function words removed via hashtag (#)

- *Copy* `wordlist_master.txt` and rename copy to `wordlist.txt`

- Restrict analysis: tick `Use existing wordlist`

- Stylo will look for `wordlist.txt and use only these words`

# Run PCA

- Sample size = 10,000
- MFW = 65
- **Select PCA**
- PCA Flavour = Technical
- **Select** B_ep.txt, G_ep.txt, H_epNG.txt
- Existing wordlist + Select Texts Manually
- **Same plot?**

# Test PCA



[ss=10,000; 65 MFW; content words 'culled']

# Play with `Sample size`



[ss=5,000]    [ss=1,000]

# Boxplot

- **Plot differences in** `MFW`

- **Use** `oppose():`
  - `primary_set = G_ep.txt`
  - `secondary_set = B_ep.txt + H_ep.txt`

- Also try:
  - *in* for Hildegard
  - *non* for Bernard



**Boxplot for "et"**

(y-axis) Absolute frequency per slice (2000 words)

Primary (62/62)    Secondary (124/124)

(Wilcoxon rank sum: p < 0.05)

# "Anonymous" text?



[ss=3,706; Bernard's *Sermo in festo sancti Martini* as "anonymous test case"; ]

# Add text

- Bernard's *Sermo in festo sancti Martini*

- "Anonymous" test case

- Add `B_Mart.txt`

- Set `Sample size` to **3,500** (length of `B_Mart.txt`)

- Attribution?

# Bigger picture



## Principal Components Analysis

# Concluding experiment

- Use all texts

- Sample size = 3,000

- Don't forget: Existing wordlist

- PCA Flavour = Classic

- Stable? Try out different settings!

# Synergy Hypothesis

- Pennebaker (e.g. 2011)

- *The Secret Life of Pronouns*

- Federalist papers and Beatles songs

- Collaborative writing style?

- "unlike either of one of the styles that the collaborating authors would produce on their own"

- Practical ànd theoretical relevance

# "Hollywood version"? Online documentary
## vimeo.com/70881172

# References

•Argamon S. (2008). Interpreting Burrows's Delta: Geometric and Probabilistic Foundations, Literary and Linguistic Computing, 23(2): 131–47.

•Burrows, J. (2002). 'Delta': A Measure of Stylistic Difference and a Guide to Likely Authorship, Literary and Linguistic Computing, 17(3): 267–87.

•Chrupała, G. et al. (2008). Learning Morphology with Morfette. In Proceedings of LREC 2008. Marrakech, Morocco: pp. 2362–67.

•Eder, M., Kestemont, M. and Rybicki, J. (2013). Stylometry with R: A Suite of Tools. In Digital Humanities 2013. Conference Abstracts. University of Nebraska-Lincoln, pp. 487–89.

•Ferrante, J. (1998). Scribe quae vides et audis. Hildegard, Her Language, and Her Secretaries. In Townsend, D. et al. (eds), The Tongue of the Fathers. Gender and Ideology in Twelfth-Century Latin. Philadelphia: University of Pennsylvania Press, pp. 102–35.

•Kestemont, M., Daelemans, W. and De Pauw, G. (2010). Weigh your Words – Memory-Based Lemmatization for Middle Dutch, Literary and Linguistic Computing, 25(3): 287–301.

•Newman, B. (1987). Sister of Wisdom. St. Hildegard's Theology of the Feminine. LA: University of California Press.

•Passarotti, M. and Dell'Orletta, F. (2010). Improvements in Parsing the Index Thomisticus Treebank. Revision, Combination and a Feature Model for Medieval Latin. In Calzolari, N. et al. (eds), Proceedings of LREC 2010. Valetta, Malta, pp. 1694-71.

•Pennebaker, J. (2011). The Secret Life of Pronouns. What our Words Say About Us. NY: Bloomsbury.

•Petrie, K., Pennebaker, J. and Sivertsen, B. (2008). Things We Said Today: A Linguistic Analysis of the Beatles, Psychology of Aesthetics, Creativity, and the Arts, 2(4), 197–202.