

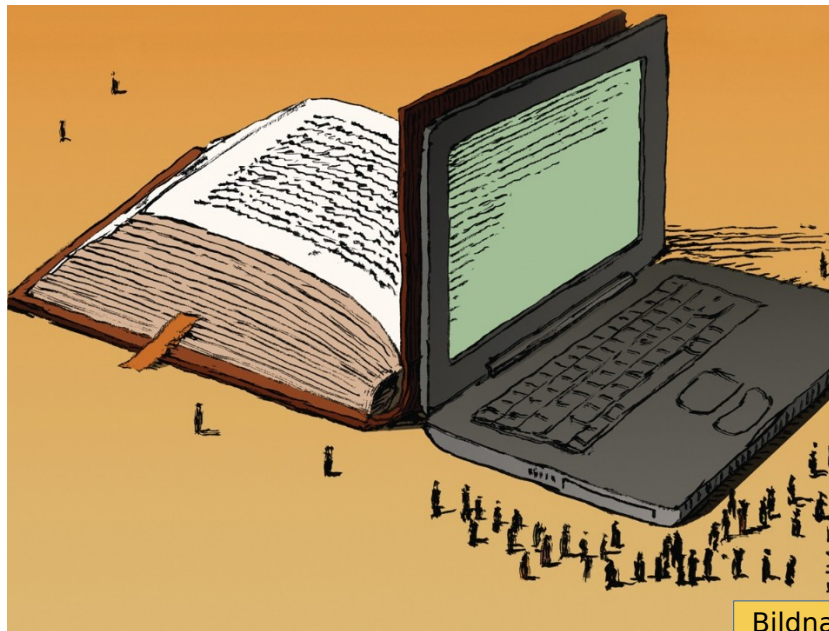
Eröffnung ZfdW
23. November 2018
Fritz Kliche

Übersicht

- Seminar „Wahlkampf in (a)sozialen Medien“
- Methoden der quantitativen Textanalyse
 - Termextraktion
 - Burrows' Zeta

Seminar „Wahlkampf in (a)sozialen Medien“

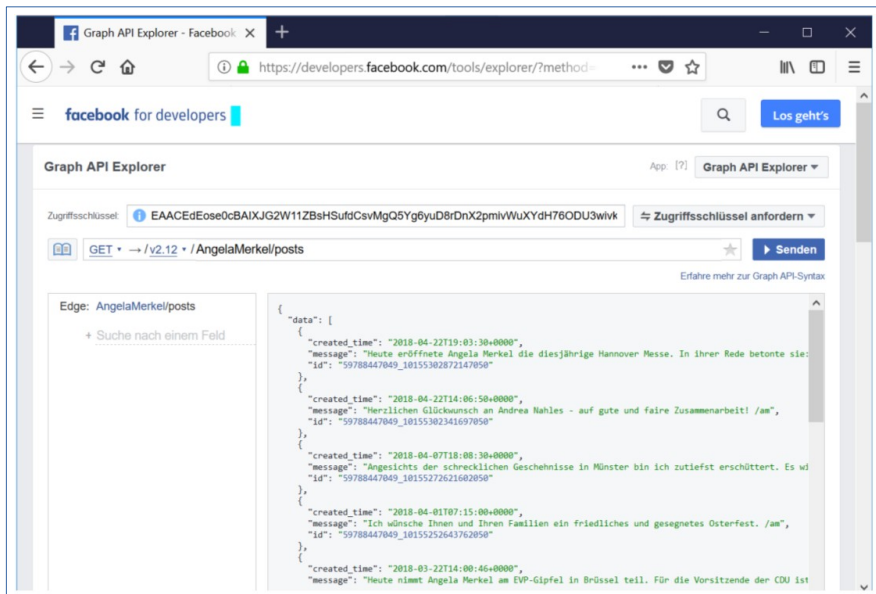
- Seminar in Anlehnung an laufendes Forschungsprojekt
- Kooperation Politikwissenschaft und Computerlinguistik
- Digitaler Wandel:
 - Computerlinguistische Methoden für politikwissenschaftliche Fragen
 - Konzept der Digital Humanities



Bildnachweis:

<https://mitpress.mit.edu/books/war-learning>

Seminar Wa(s)M: Das Facebook-Korpus



Seminar Wa(s)M: Das Facebook-Korpus

```
<doc id=294 subcorpus="AngelaMerkel" db_id="10152198616677050_1394383917241010" LOC="0"
PER="PER" ORG="0" MISC="0" date=201612" likecount=2 messagetype="comment" language="de"
length=143 cross="-" frequer=0 ptv=1 neg=1 sentdiff=0 sharecount=2>
<s>      0
Bitte  NN      Bitte  0
Frau   NN      Frau   0
Kanzlerin NN      Kanzlerin  0
.      $.      .      0
Hilft  VVFIN    helfen  0
evtl.  NE       evtl.  0
Herrn  NN       Herr   0
Sch<C3><A4>uble NE    Sch<C3><A4>uble I-PER
aber   ADV      aber   0
hat    VAFIN     haben  0
leider ADV      leider 0
die    ART     die    0
Linke  NN       Linke  0
auch   ADV      auch   0
schon  ADV      schon  0
festgestellt VVPP    feststellen  0
```

Begriffe

- Korpus
- Sub-/Teilkorpus
- „Messages“ (Begriff von FB)
- Metadaten

Annotationen

Im Beispiel verwendete Annotationen:

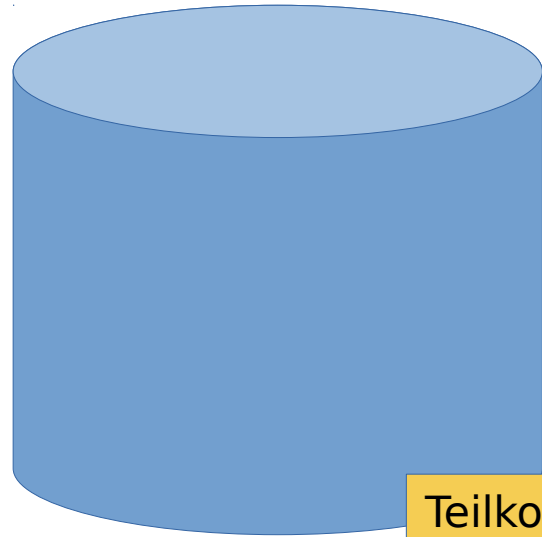
- Wortart
- Lemma
- Eigenname

Seminar Wa(s)M: Das Facebook-Korpus

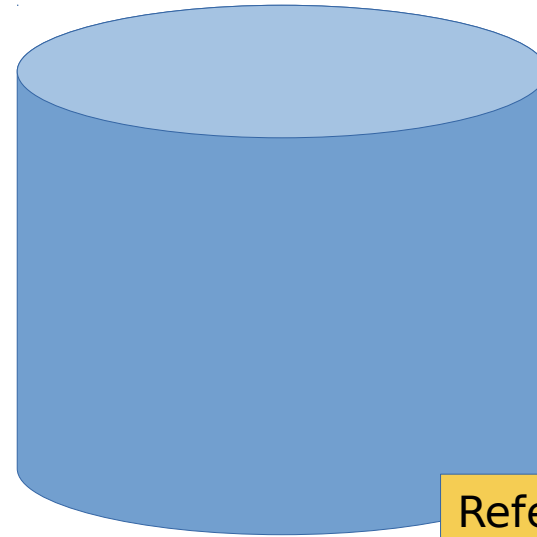
- AngelaMerkel	530.520	- linkspartei	120.629
- CDU	204.026	- sahrwagenknecht	333.454
- CSU	246.238	- DietmarBartschMdB	19.537
- joachimherrmannscsu	3.053	- FDP	100.273
- B90DieGruenen	127.847	- lindnerchristian	211.783
- GoeringEckardt	31.016	- SPD	156.911
- Cem	106.837	- martinschulz	336.652
		- alternativefuerde	761.036
		- aliceweidel	121.634

- Messages insgesamt: 3.411.448
- Posts: 9.081
- Kommentare: 1.851.945
- Subkommentare: 1.550.420
- Tokens (d. h.: Wörter und Satzzeichen): 113.639.602

Erste Idee: Termextraktion



Teilkorpus



Referenzkorpus

$$\textit{Weirdness Ratio}(x) = \frac{\frac{\text{Anzahl aller Vorkommen von "Einkommen" im Teilkorpus}}{\text{Anzahl aller Vorkommen von [N.*] im Teilkorpus}}}{\frac{\text{Anzahl aller Vorkommen von "Einkommen" im Vergleichskorpus}}{\text{Anzahl aller Vorkommen von [N.*] im Vergleichskorpus}}}$$

Schäfer et al, 2015

Erste Idee: Termextraktion

Partei und Spitzenkandidaten, mit cross="-", Pattern: [ADJ.*] [N.*]

##_linke	##_spd	##_gruene
1665 sozial Gerechtigkeit [#2]	1956 sozial Gerechtigkeit [#1]	433 deutsch Politiker [#2086]
1038 falsch Partei [#75882-#7]	1450 letzt Jahr [#143913-#14]	365 türkisch Politik [#106378]
698 deutsch Volk [#44491-#4]	1205 lieb Herr [#148238-#149]	352 deutsch Politik [#20504-#2]
648 letzt Jahr [#159153-#159]	983 geehrter Herr [#76467-#7]	306 deutsch Volk [#22539-#2]
598 einzig Partei [#64531-#6]	883 groß Koalition [#94347-#9]	302 letzt Jahr [#72953-#732]
512 bedingungslos Grundein	835 herzlich Glückwunsch [#]	255 froh Weihnachten [#393]
481 herzlich Glückwunsch [#]	721 leer Versprechung [#141]	236 lieb Grüne [#74528-#747]
476 geehrt Frau [#90929-#91]	611 klein Mann [#131940-#1]	228 eigen Land [#27787-#28]
461 groß Koalition [#111038-#11]	568 lieb Martin [#149678-#15]	210 cem özdemir [#15807-#1]
460 frau Wagenknecht [#812]	552 gut Gruß [#103743-#104]	203 ganz Welt [#41513-#417]
459 groß Teil [#114035-#114]	527 deutsch Volk [#39442-#3]	201 herzlich Glückwunsch [#]
430 ganz Welt [#89593-#900]	484 gut Mann [#106276-#106]	172 lieb Herr [#74787-#7495]
410 nah Wahl [#179289-#179]	437 jung Mensch [#125852-#125]	171 grün Partei [#52911-#53]
409 nah Jahr [#177971-#178]	409 heiß Luft [#113673-#114]	170 grün Politik [#53168-#53]
401 lieb Frau [#161729-#162]	404 bedingungslos Grundein	170 türkisch Volk [#107504-#107]

##_fdp	##_cdu	##_csu	##_afd
1446 herzlich Glückwunsch [3743 herzlich Glückwunsch [440 herzlich Glückwunsch [#]	2386 deutsch Volk [#77816-#8]
822 lieb Herr [#106501-#107]	2567 deutsch Volk [#64946-#6]	339 heiß Luft [#32936-#3327]	1606 herzlich Glückwunsch [#]
681 frei Demokrat [#48743-#48]	1043 frau merkel [#106590-#106]	292 letzt Jahr [#41976-#4220]	879 eigen Volk [#101466-#102]
491 letzt Jahr [#102332-#102]	1042 letzt Jahr [#201187-#201]	252 deutsch Volk [#12878-#128]	834 letzt Jahr [#231798-#232]
472 geehrter Herr [#54460-#54]	1012 lieb Frau [#205231-#206]	245 nah Jahr [#46318-#4656]	824 eigen Land [#97439-#982]
400 gut Mann [#75151-#7555]	975 geehrt Frau [#120534-#120]	224 nah Wahl [#46796-#4701]	784 etabliert Partei [#118547-#118]
396 gut Idee [#73825-#74220]	917 ganz Welt [#118491-#118]	192 deutsch Autofahrer [#10]	739 nah Jahr [#255737-#2564]
353 wichtig Thema [#169202]	890 eigen Volk [#83599-#844]	165 inner Sicherheit [#35839]	666 arm Deutschland [#35544]
324 groß Koalition [#66079-#66]	803 gut Frau [#150057-#1500]	158 eigen Land [#15911-#160]	664 deutsch Sprache [#74936]
291 jung Mensch [#90519-#905]	733 froh Weihnachten [#1118]	146 doppelt Staatsbürgersch	654 AfD Wähler [#7409-#8062]
272 inner Sicherheit [#87229]	691 neu Jahr [#226067-#226]	143 lieb Herr [#43156-#4329]	576 ganz Welt [#141080-#141]
224 nah Jahr [#114206-#114]	656 eigen Land [#80430-#810]	134 arm Deutschland [#5060]	563 deutsch Bürger [#65401-#65]
222 gut Abend [#69779-#700]	596 besorgt Bürger [#40965-#40]	131 eigen Volk [#16666-#167]	461 einzig Partei [#109059-#109]
211 liberal Partei [#104975-#104]	595 nah Jahr [#219154-#219]	118 deutsch Bürger [#11269]	458 groß Teil [#169839-#1702]
201 erst Mal [#41833-#4203]	576 inner Sicherheit [#17548]	110 ganz Welt [#23629-#237]	453 deutsch Pass [#71958-#719]
199 groß Problem [#67054-#67]	469 groß Teil [#144509-#144]	104 deutsch Pass [#12055-#120]	453 eigen Meinung [#98755-#98]
192 richtig Weg [#134846-#134]	433 WARNING Frau [#16986]	103 lieb CSU [#42794-#4289]	434 gut Nacht [#180318-#180]
187 kalt Progression [#9139]	433 deutsch Bürger [#56728]	100 klar Wort [#38061-#3816]	433 nah Wahl [#257807-#2582]
186 sozial Marktwirtschaft [394 jung Mann [#182215-#182]	98 geehrter Herr [#24056-#240]	408 jung Mann [#208270-#208]
179 gut Bildung [#71467-#714]	382 froh Ostern [#111352-#111]	98 leer Versprechung [#414]	387 deutsch Kultur [#69626-#69]
177 erst Linie [#41647-#4182]	361 gut Kanzlerin [#152542-#152]	83 wichtig Thema [#66944-#669]	387 dumm Mensch [#87893-#878]
177 sozial Gerechtigkeit [#1]	347 offen Grenze [#233116-#233]	81 etabliert Partei [#19679-#196]	379 gut Beispiel [#174606-#174]

Zweite Idee: Burrows' Zeta

- Ausgangspunkt: Stilometrie
 - Die Stilometrie bezeichnet „die Anwendung quantitativer Methoden zur Erfassung und Klassifizierung stilistischer Merkmale von Texten“.
- Eine Anwendung ist *authorship attribution*, d. h. die Zuordnung eines Textes, dessen Autor unbekannt oder unsicher ist, zu einem Autor mit ähnlichen stilistischen Merkmalen.
- Ein Maß für den kontrastiven Vergleich zweier Korpora/Texte mit dem Ziel einer stilometrischen Analyse ist Burrows' Zeta.

Viehhauser, 2015

Schöch, 2018

Zweite Idee: Burrows' Zeta

- Die zu vergleichenden Korpora/Texte werden – im Gegensatz zum vorgestellten Verfahren zur Termextraktion – zunächst in Einheiten segmentiert, beispielsweise in Abschnitte einheitlicher Länge (in Tokens).



Zweite Idee: Burrows' Zeta

- Die häufigsten Tokens der Textdaten werden ermittelt.
- Für jedes dieser Tokens wird für beide zu vergleichenden Korpora/Texte bestimmt, in wie vielen Segmenten sie mindestens einmal vorkommen.

$$dp_i(Z) = \frac{df_i(Z)}{n(Z)} \quad \text{bzw.} \quad dp_i(V) = \frac{df_i(V)}{n(V)}$$

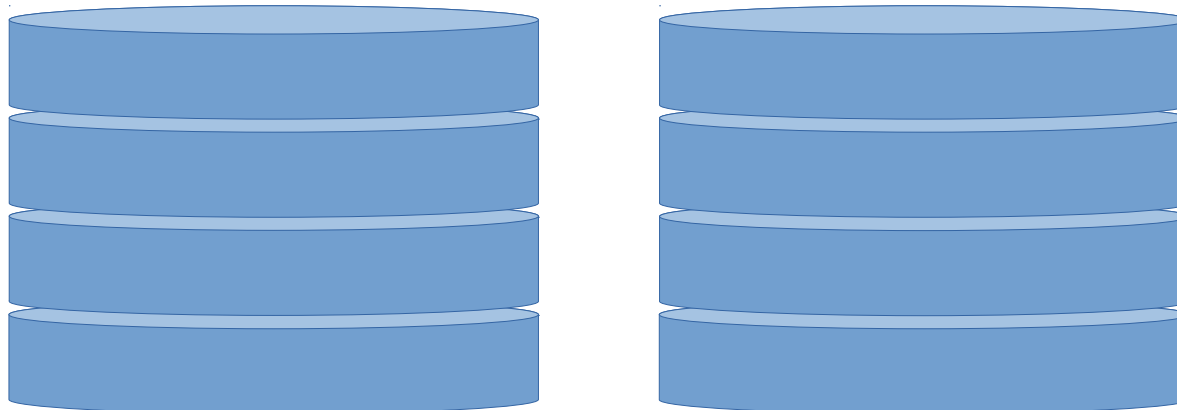
wobei

$dp_i(Z)$: die *document proportions* des Tokens i im Korpus/Text Z .

$df_i(Z)$: die *document frequency* des Tokens i im Korpus/Text Z .

$n(Z)$: die Anzahl der Segmente in Z

V : das Vergleichskorpus/der Vergleichstext



Zweite Idee: Burrows' Zeta

- Damit ergibt sich Zeta:

$$Zeta_i = dp_i(Z) - dp_i(V)$$

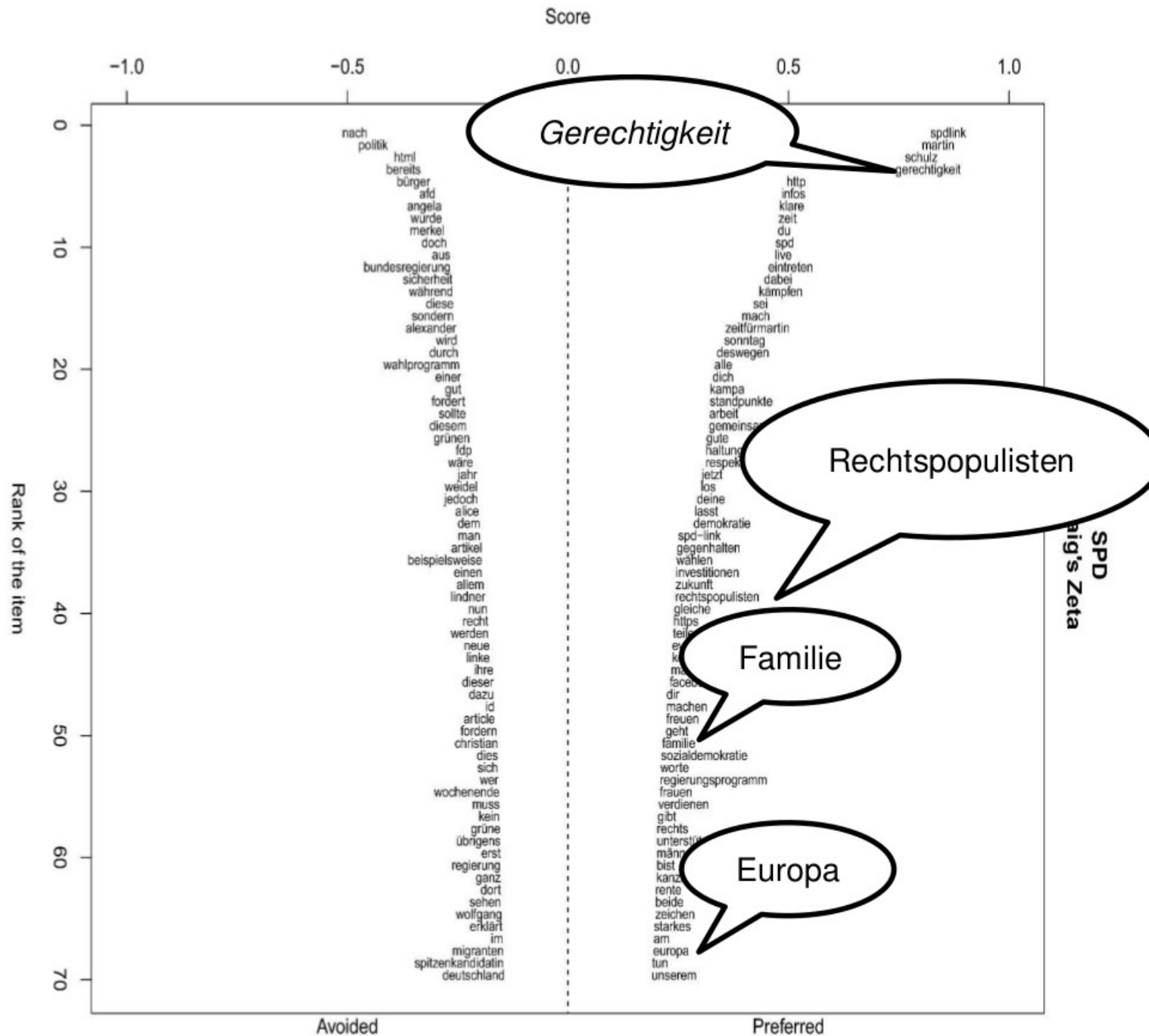
- Zeta liefert Ergebnisse von -1 bis 1.
- In den resultierenden Wert fließt vor allem ein, wie konsistent ein Token verwendet wird, und nicht allein, wie häufig es vorkommt.

Zweite Idee: Burrows' Zeta

- Umsetzung mit STYLO
 - R-Paket zur stilometrischen Analyse.
 - Bedingung:
 - Die Texte liegen in *plain text* vor.
 - Die Texte sind lemmatisiert.
- Hinweis: *stopword removal* ist nicht erforderlich.

Eder et al., 2016

Einige Ergebnisse



Einige Ergebnisse

Vergleich der Ergebnisse
Termextraktion - Burrows' Zeta

##	afd
2386	deutsch Volk [#77816-#8
1606	herzlich Glückwunsch [#
879	eigen Volk [#101466-#102
834	letzt Jahr [#231798-#232
824	eigen Land [#97439-#982
784	etabliert Partei [#118547-
739	nah Jahr [#255737-#2564
666	arm Deutschland [#35544
664	deutsch Sprache [#74936
654	AfD Wähler [#7409-#8062
576	ganz Welt [#141080-#141
563	deutsch Bürger [#65401-#
461	einzig Partei [#109059-#1
458	groß Teil [#169839-#1702
453	deutsch Pass [#71958-#7
453	eigen Meinung [#98755-#
434	gut Nacht [#180318-#180
433	nah Wahl [#257807-#2582
408	jung Mann [#208270-#208
387	deutsch Kultur [#69626-#
387	dumm Mensch [#87893-#
379	gut Beispiel [#174606-#1

gutmenschen
altparteien
pack
petry
höcke
arsch
weidel
gauland
lügenpresse
scheiss
idioten
hetze
nazi
meuthen
frau
alice
etablierten
muslime
fresse
scheiße
deppen
moslems
blau
unfassbar
ausländer
scheiß
nazis
dummheit
gesindel
islam
kotzen
heimat
islamisierung
volksverräter
dreck

Fragen

- Erweiterung der Merkmale:
 - Statt Tokens:
 - N-grams?
 - POS-Patterns?
 - Weitere Merkmale?

- Anwendung computerlinguistischer Methoden in Seminaren
 - Umgang mit R
 - Umgang mit der Kommandozeile
 - Verwendung bestehender Tools

Übersicht

- Eder, M., Rybicki, J. and Kestemont, M. (2016). Stylometry with R: A package for computational text analysis. *R Journal*, 8(1): 107-121.
- Schäfer, J., Rösiger, I., Heid, U. und Dorna, M. (2015). Evaluating noise reduction strategies for terminology extraction. In den *Proceedings der TIA 2015*. Granada.
- Schöch, C. (2018). Zeta für die kontrastive Analyse literarischer Texte. Theorie, Implementierung, Fallstudie. In T. Bernhart u. a. (Hrsg.): *Quantitative Ansätze in den Literatur- und Geisteswissenschaften. Systematische und historische Perspektiven* (S. 77-94). Berlin: De Gruyter.
- Viehhauser, G. (2015). Historische Stilometrie? Methodische Vorschläge für eine Annäherung textanalytischer Zugänge an die mediävistische Textualitätsdebatte. In C. Baum und T. Stäcker (Hrsg.): *Grenzen und Möglichkeiten der Digital Humanities. Sonderband der Zeitschrift für digitale Geisteswissenschaften*.