

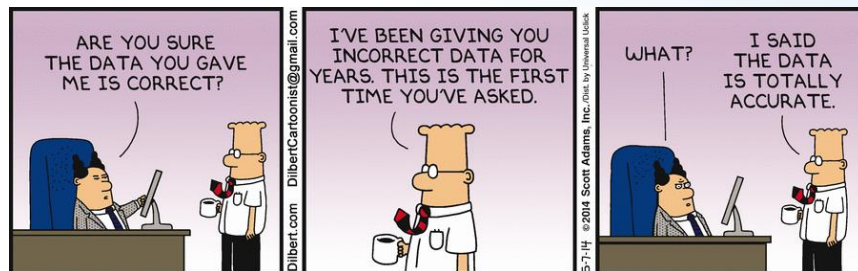
Text-Scraping und Textanalyse

Daria Liakhovets

Stefan Kostecky

Christoph Jungbauer

Katalin Feichtinger



Ziele



**FACHHOCHSCHULE
WIENER NEUSTADT**
Austrian Network for Higher Education

- Eigenschaften der Artikel zum Thema Klimawandel evaluieren: Lesbarkeitsindex, Sentiment, Anzahl Wörter etc.
- Zwei Medien anhand erfasster Merkmale miteinander vergleichen.

- Verschiedene Methoden und Werkzeuge zur Messung der (technischen) Lesbarkeit wie:
 - „Flesch Reading Ease“,
 - „Flesch-Kincaid Grade Ebene“,
 - „Wiener Sachtextformel“,
 - „Einfaches Maß für den Gobbledygook“ (SMOG)
 - „Gunning Fog-Index“ (FOG)
- Erlauben dabei eine **quantitative** Analyse.
- Viele dieser Methoden sind spezifisch auf eine Sprache ausgerichtet.
- Syntax und Semantik sind fundamental anders

Flesch Reading Ease



- Grundannahme, dass kurze Wörter und kurze Sätze für Leserinnen und Leser leichter verständlich sind.
- Es ergibt sich eine Skala von 0-100. Je höher der Wert, desto einfacher verständlich ist laut Flesch der Text, „Flesch Reading Ease“

Flesch-Reading-Ease		
Wert	Einstufung	Zielgruppeneinschätzung
0–30	Sehr schwer	Akademiker
30–50	Schwer	
50–60	Mittelschwer	
60–70	Mittel	13–15-jährige Schüler
70–80	Mittelleicht	
80–90	Leicht	
90–100	Sehr leicht	11-jährige Schüler

Flesch Reading Ease

Die Schritte zur Berechnung der Formel nach Flesch sind dabei folgende:

- Durchschnittliche Satzlengthen bestimmen (SA)
- Anzahl der Silben bestimmen und durch die Anzahl der Wörter dividieren. (SI)

Daraus ergibt sich die Formel für den Flesch Reading Ease:

$$FRE = 206,835 - 84,6 \times SA - 1,015 \times SA$$

Anpassung an deutsche Sprache:

$$FRE = 180 - 58,5 \times SA - SI$$

Bamberger und Vanecek entwickelten 1984 die Wiener Sachtextformel explizit für deutschsprachige Texte. Ähnlich dem Gunning-Fog-Index gibt sie an, für welche Schulstufe ein Sachtext geeignet ist. Dabei beginnt die Skala bei 4 und endet bei 15 (entsprechend den Schulstufen 4 bis 12 und darüber hinaus bis 15 an den Universitäten). Folgende Parameter sind zur Berechnung nötig:

- MS = prozentualer Anteil an Wörtern mit drei oder mehr Silben
- SL = durchschnittliche Satzlänge (gesamte Wort-Zahl durch gesamte Satz-Zahl)
- IW = prozentualer Anteil an Wörtern mit mehr als 6 Buchstaben
- ES = prozentualer Anteil an Wörtern mit nur einer Silbe

Die erste Wiener Sachtextformel lautet:

$$\text{WSTF1} = 0,1935 * \text{MS} + 0,1672 * \text{SL} + 0,1297 * \text{IW} - 0,0327 * \text{ES} - 0,875$$

Die zweite Wiener Sachtextformel lautet:

$$\text{WSTF2} = 0,2007 * \text{MS} + 0,1682 * \text{SL} + 0,1373 * \text{IW} - 2,779$$

Sie verzichtet hierbei auf das Abzählen von einsilbigen Wörtern, ohne dabei an Genauigkeit zu verlieren.

Die dritte Wiener Sachtextformel lautet:

$$\text{WSTF3} = 0,2963 * \text{MS} + 0,1905 * \text{SL} - 1,1144$$

Mit leichten Einbußen an Präzision kann zur Näherung auch nur mit mehrsilbigen Wörtern (MS) und der durchschnittlichen Satzlänge (SL) gearbeitet werden. Dies spart Arbeitsaufwand bei der Berechnung.

Als Ausgangsbasis unabhängig von Wort- oder Satzlängen wird ein Grundwortschatz benötigt. Der aktive Wortschatz eines durchschnittlichen Erwachsenen liegt bei 12.000 bis 16.000 Wörtern. Der Grundwortschatz ist dabei definiert als die Anzahl der Wörter, die nötig sind, um 85 % eines beliebigen Gespräches zu verstehen. Alan Pfeffer ermittelte hierzu 1975 eine Zahl von 1285 Wörtern, die nötig sind um diesen Schwellenwert zu erreichen. Diese Zahl bezieht sich auf einfache Gespräche. Je spezifischer die Anforderungen sind, desto mehr Fachbegriffe kommen hinzu um den Text verstehen zu können.

Textabdeckung durch Wortschatz				
Art der Kommunikation	Unter 1000 Wörter	Bis 2000 Wörter	Bis 3000 Wörter	Bis 4000 Wörter
Gesprochene Sprache	85,2	89,2	90,0	91,9
Bestseller	77,8	83,0	85,3	87,1
Abenteuerromane	73,2	79,0	81,9	83,6
Gesellschaftsromane	73,7	79,0	81,9	83,8
Anspruchsvolle Literatur	73,8	79,4	82,0	83,9
Belletristik (Durchschnitt)	74,2	80,0	82,7	84,5
Zeitungstexte	67,4	73,9	77,3	79,4
Universitätseinführungen	68,9	75,9	79,6	81,9
Fachzeitschriften	66,3	73,5	77,4	79,6
Fachtexte (Durchschnitt)	67,6	74,4	78,5	80,7

Text Mining ist ein Analyseverfahren zur Entdeckung von Strukturen und Daten im Text.

Data Mining mit Hilfe von Datenbanken (wie Kundendaten einer E-Commerce Plattform) prägen den Begriff KDD (Knowledge Discovery in Databases).

Text Mining durch automatisierte Verarbeitung bezeichnet Tan analog dazu als KDT (Knowledge Discovery from Text).

Pang und Lee beschreiben im Kontext von Text Mining zwei wesentliche Techniken für das automatisierte Lernen:

- Linguistische Methoden und
- Maschinelles Lernen

Kombination der Beiden: Sentiment Analyse

Die Sentiment Analyse kombiniert linguistische Methoden und maschinelles Lernen im Kontext von Text Mining und versucht mittels statistischer Methoden Stimmungen der Nutzerinnen zu extrahieren.
Eine Liste mit bewerteten Wörtern bietet hierfür beispielsweise die Hochschule Hof.

Sentiment Bewertung				
Phrase	Opinionwert	Standardabweichung	Standardfehler	Typ
einfach gut	0.93	0.19	0.01	a
großartig	0.95	0.23	0.01	a
sehr gut	0.90	0.22	0.00	a
nur Schrott	-0.85	0.52	0.10	n
nur schlecht	-0.97	0.16	0.01	a

- Theoretischer Hintergrund
- Daten extrahieren und aufbereiten
- Lesbarkeitsindex berechnen
- Textsentiment evaluieren (textblob_de package)
- Weitere Merkmale wie Länge eines Artikels, Anzahl Kommentare u.a. erfassen, Ähnlichkeit der Datenquellen evaluieren
- Ergebnisse visualisieren

Spiegel.de:

1000 Artikel, August 2015 – Dezember 2019

<https://www.spiegel.de/thema/klimawandel/>

Bild.de:

100 Artikel, November 2019 – Dezember 2019

<https://www.bild.de/themen/ereignisse/klimawandel/news-nachrichten-news-fotos-videos-16877746.bild.html>

Für den Vergleich wurde ein entsprechendes Sample der Spiegel-Artikel verwendet

Data Scraping

- **requests, BeautifulSoup4**
- Funktionen für das Extrahieren von relevanten Daten
- Beispiel:

```
def get_date_time(art):  
    return [i.strip('\r\n\t') for i in art.find('span', class_ = "article-function-date")\  
        .find('time').get_text().split('\xa0')]
```

- Generator für Bearbeitung aller Artikellinks auf der Website

Data Scraping



**FACHHOCHSCHULE
WIENER NEUSTADT**
Austrian Network for Higher Education

```
def scrape_page(url):
    page = get_page(url)
    article_links = find_and_complete_links(page)

    for link in article_links:
        art = get_page(link)
        headline, headline_intro, intro = get_headline_intro(art)
        date, time = get_date_time(art)
        text = get_clean_text(art)
        comments = get_comments(get_forum_link(art))

        article_dict = {'headline': headline,
                        'headline_intro': headline_intro,
                        'date': date,
                        'time': time,
                        'intro': intro,
                        'text': text,
                        'thread': comments}

        yield article_dict
```

Generator (*Code vereinfacht*)

Input: URL der Website

Output: Generator Objekt, Dictionaries mit
Artikeldaten

Dictionary:

<code>'headline' : str,</code>	<i>Titel des Artikels</i>
<code>'headline_intro' : str,</code>	<i>Zusätzliche Info zum Titel</i>
<code>'date' : str,</code>	<i>Datum</i>
<code>'time' : str,</code>	<i>Uhrzeit</i>
<code>'intro': str,</code>	<i>Kurzer Text, der das Thema des Artikels erläutert</i>
<code>'text' : str,</code>	<i>Der ganze Artikeltext</i>
<code>'thread' : list}</code>	<i>Kommentare (Liste)</i>

Speichern und einlesen mit **json** package:

with open('file.txt', 'w') as outfile:

 json.dump(data, outfile)

with open('file.txt') as json_file:

 data = json.load(json_file)


```
# first 5 comments  
article['thread'][:5]
```

```
['1.Ein Klima-Gipfel der Unfähigkeit und der peinlichen Pannen. Dieser Aufwand war echt für die Katz. Das ganze Desaster hat hoffentlich noch ein Nachspiel.'],  
['2.Boykott Nun wird es Zeit, Produkte aus Australien, Brasilien, USA, China etc zu boykottieren. Die EU ist ein großer Markt und hat daher Macht.'],  
['3.CO2-Bilanz Für diesen "Klima-Gipfel" wurde vermutlich mehr CO2 produziert als die Beschlüsse insgesamt an Einsparungen mit sich bringen. Das Ganze wird langsam immer lächerlicher.'],  
['4 Keine Sorge, das Nachspiel wird kommen.Wir als Spezies haben es uns redlich verdient.'],  
['5 Wie wäre denn ein entsprechend angepasster Konsum? Produkte vermeiden, die aus ignoranten Ländern wie China, Indien und den USA kommen (ausser, sie sind gezielt besonders nachhaltig produziert). Also z. B. kein Huawei Handy, sondern ein Sony oder Samsung. Die Liste kann man auf viele alltäglichen Entscheidungen ausdehnen. Politisch muss man sich die Frage stellen, ob man mit diesen Staaten die Wirtschaftsbeziehungen ausbauen will. Das würde uns Wohlstand kosten - aber ohne Reduktion von Investitionen und Konsum wird der fossile Energieverbrauch sich kaum deutlich reduzieren lassen.']]
```

Datenstruktur Beispiel

```
{
  'headline':      'Das Desaster von Madrid',
  'headline_intro': 'Klimakonferenz',
  'date':          'Sonntag, 15.12.2019 ',
  'time':          '18:04 Uhr',
  'intro':         'Fast nichts wurde beschlossen beim Marathongipfel in Madrid, denn Rechtspopulisten und Kohlefreunde lähmten den Klimaschutz. Nun sind die EU-Staaten gefordert, die CO2-Ziele zu erreichen. Allen voran Deutschland.',
  'text':          'Nun ist sie vorbei, die Klimakonferenz von Madrid. Die meisten der rund 26.000 Teilnehmer haben das Messegelände längst verlassen, sie mussten den Verlängerungs-marathon am Schluss nicht mitmachen. Rund um den Tagungsort, auf Plakaten und in den Gängen des U-Bahnhofs <...>',
  'thread':        [['1.Ein Klima-Gipfel der Unfähigkeit und der peinlichen Pannen. Dieser Aufwand war echt für die Katz. Das ganze Desaster hat hoffentlich noch ein Nachspiel.'],
                   ['2.Boycott Nun wird es Zeit, Produkte aus Australien, Brasilien, USA, China etc zu boykottieren. Die EU ist ein großer Markt und hat daher Macht.'],
                   ['3.CO2-Bilanz Für diesen "Klima-Gipfel" wurde vermutlich mehr CO2 produziert als die Beschlüsse insgesamt an Einsparungen mit sich bringen. Das Ganze wird langsam immer lächerlicher.'], <...>]
}
```



Herausforderungen

- Leere Datensätze (Websites mit anderer html-Struktur)
- Sonderzeichen (z.B. Texte in Fremdsprachen)
- Interaktive Grafiken im Text
- Unvollständige Datensätze (kostenpflichtiger Inhalt)
- Keine Kommentare und Intro bei Bild-Artikeln



Sentiment Analyse TextBlob

- TextBlob ist die Basis für natural language processing (NLP) mit Python
- für Python 2 als auch 3 verfügbar
- Ansätze, wie Erkennen von Wortarten, Extraktion von Substantiven, Stimmungsanalyse und auch Klassifizierungen möglich
- TextBlob-de ist die Erweiterung zur Untersuchung von deutschen Texten

TextBlob Allgemeines

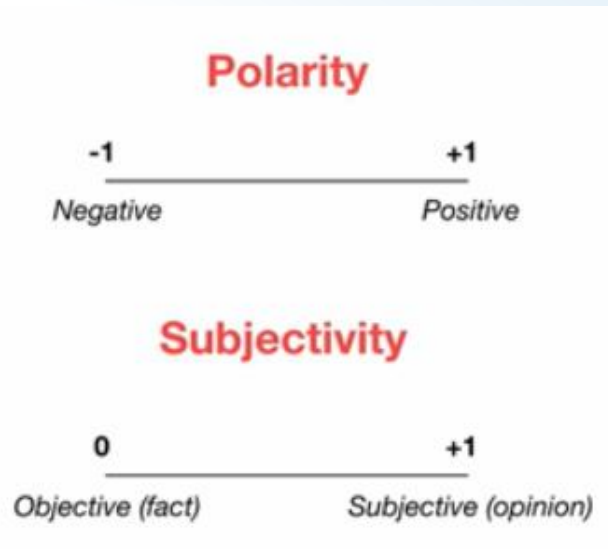
- TextBlob-Objekte können wie Python-Zeichenketten behandelt werden, die gelernt haben, wie man Natural Language Processing macht.
- Wesentliche TextBlob Objekte:
 - `blob.sentences` => Anzeige der einzelnen Sätze aus einem Text
 - `blob.tokens` => Anzeige der einzelnen Wörter
 - `blob.tags` => Anzeige der Wortarten
 - `blob.noun_phrases` => Zusammenfassung von zusammenhängenden Wörtern
 - `blob.sentiment` => ermittelt Polarity und Subjectivity

TextBlob Output Sentiment

TextBlob sentiment liefert ein namedtuple der Form Sentiment(polarity, subjectivity)

Polarity oder Sentiment score: gibt an wie positiv/negativ ein Wort/Text ist

Subjectivity gibt an wie meinungsbildend ein Wort/ Text ist



```
import nltk
from textblob_de import TextBlobDE as TextBlob
import numpy as np
```

```
from sa_class import Text_Sentiment # importiert die Klasse für Sentiment Analysis
```

```
import json
with open('article_2p9.txt') as json_file: # importiert Text Datei
    data = json.load(json_file)
```

```
ts = Text_Sentiment(data['text']) # unter Verwendung von der Klasse Text_Sentiment wird der Text des Arti-
kels ausgewählt, damit in weiterer Folge die Sentiments (polarity und subjectivity) ermittelt werden könne
n
```

```
ts.text_sentiment()
```

```
Sentiment(polarity=-0.0063888888888888892, subjectivity=0.06111111111111111)
```

Deskriptive Statistik - Übersichtstabelle SPIEGEL

	articles_counts	lesbarkeit	polarity	subjectivity	thread_counts	word_counts
count	625.0	625.000000	625.000000	625.000000	625.000000	625.000000
mean	1.0	50.216160	0.057845	0.084441	145.793600	620.521600
std	0.0	8.185391	0.107542	0.058467	150.558364	302.175037
min	1.0	12.200000	-0.370000	0.000000	10.000000	57.000000
25%	1.0	45.690000	-0.011785	0.042424	46.000000	381.000000
50%	1.0	50.300000	0.063953	0.080000	95.000000	578.000000
75%	1.0	55.780000	0.122500	0.120690	188.000000	825.000000
max	1.0	73.580000	0.385417	0.500000	1251.000000	1819.000000

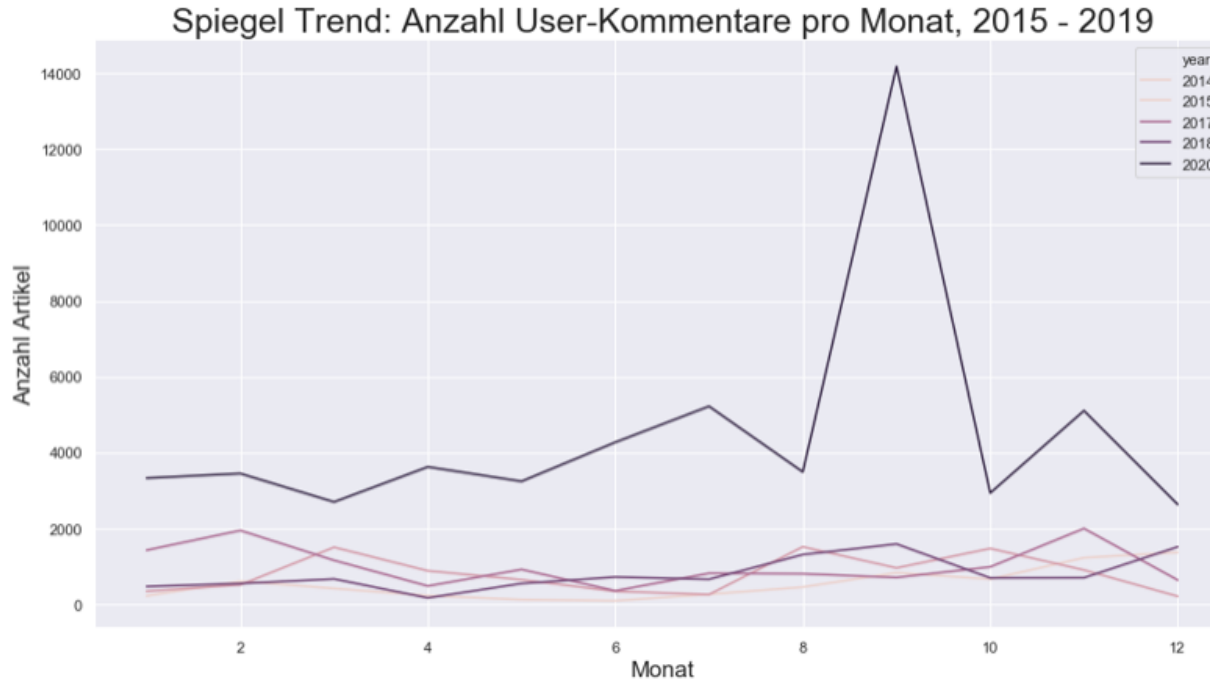
Heat-Map – Anzahl Artikel SPIEGEL nach Jahr und Monat

Ø 4,6 Ø 4,9 Ø 8,9 Ø 6,5 Ø 27



- Tendenzielle Bereitschaft in den Sommermonaten vermehrt Artikel zum Thema „Klimawandel“.
- Kontinuierliche Zunahme der Berichterstattung zu dieser Thematik
- Deutlicher Ausschlag im Vorjahr (2019) insb. im Monat 9 (Hitzerekord bis zu 40 Grad)

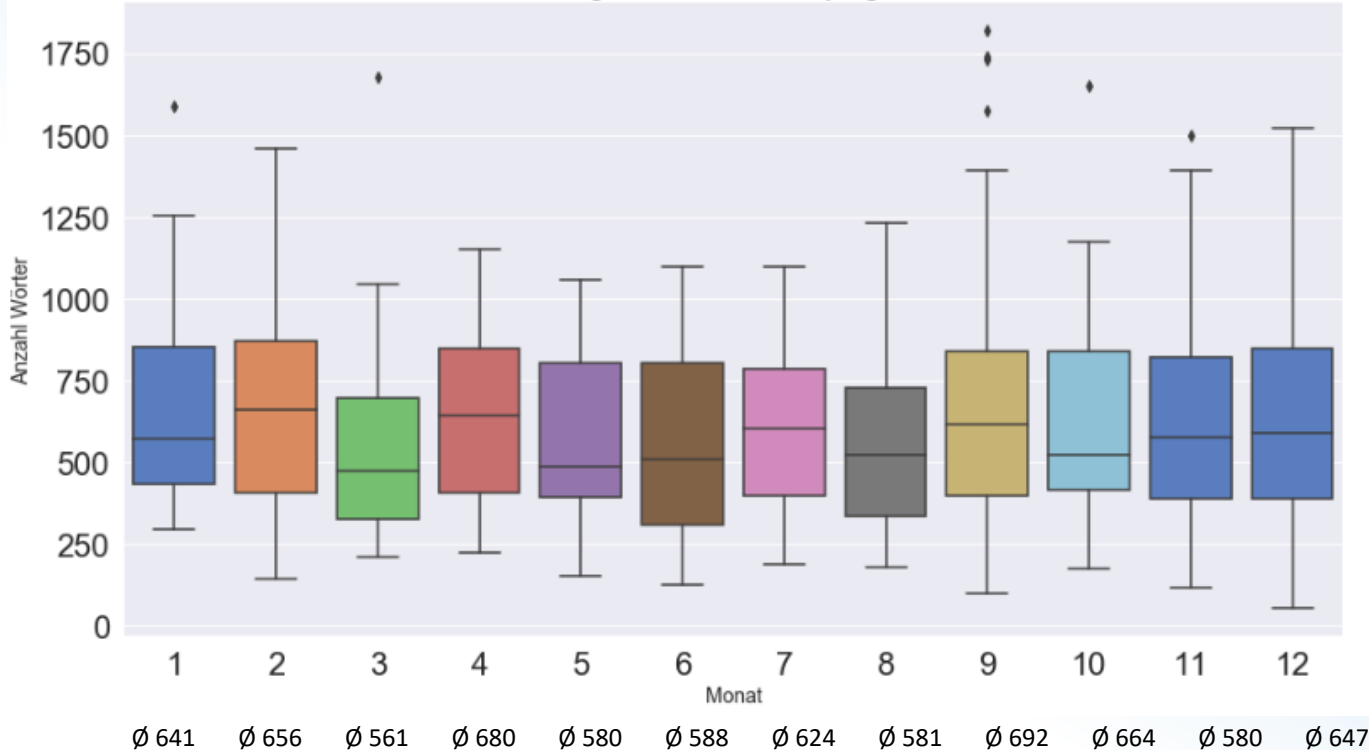
Line-Plot – Anzahl User-Kommentare SPIEGEL nach Jahr und Monat



- Synchrone Entwicklung der User-Kommentare mit der Anzahl der Artikel
- Rege Nutzer Aktivität in den Sommermonaten
- Spezielle Situation im Jahr 2019

Box-Plot – Verteilung der Anzahl Wörter SPIEGEL

Verteilung Anzahl Wörter Spiegel Artikel

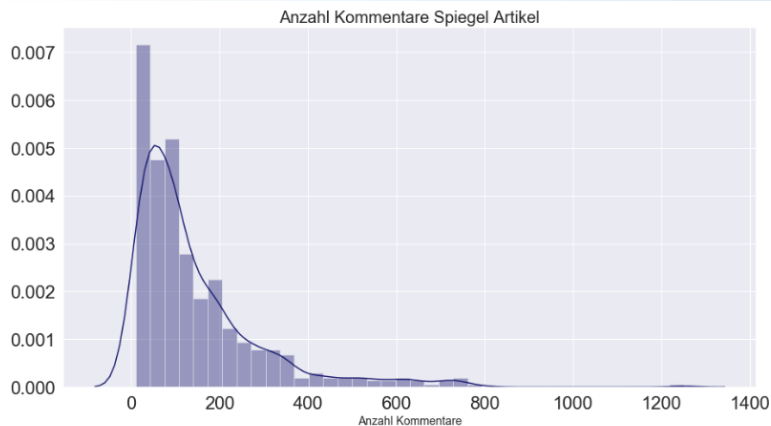
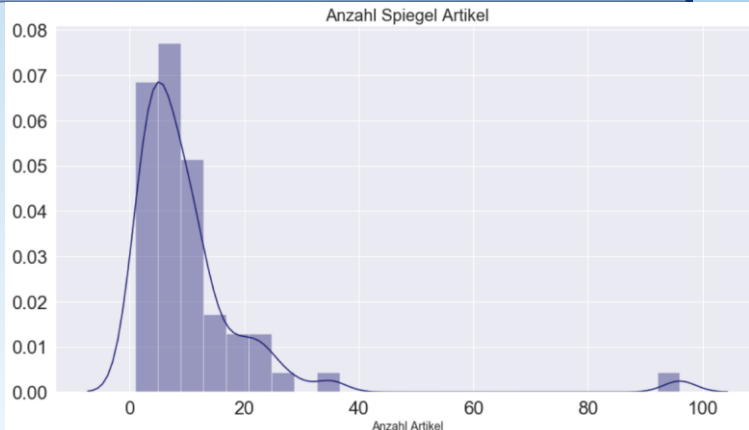
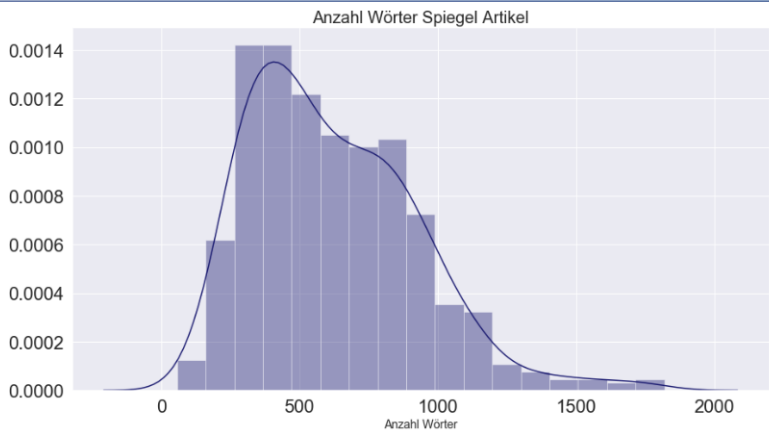


- Medianwerte unterliegen leichten Schwankungen
- Ausreißer in einigen Monaten insbesondere weist der Monat September Artikel mit sehr vielen Wörtern auf

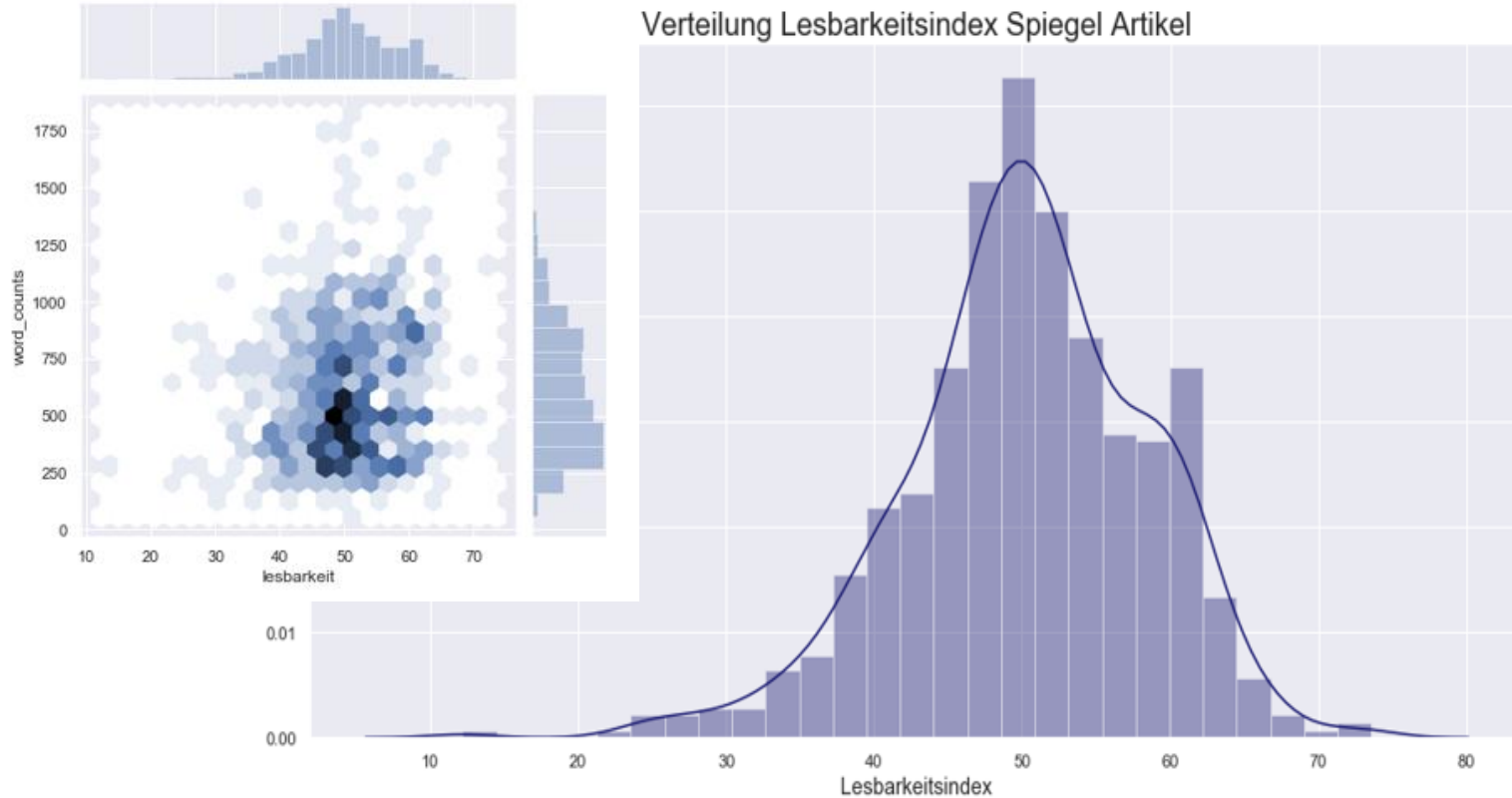
Verteilungskurven SPIEGEL Anzahl Wörter, Artikel, Kommentare



**FACHHOCHSCHULE
WIENER NEUSTADT**
Austrian Network for Higher Education



SPIEGEL Lesbarkeitsindex – Verteilung / Zusammenhang Anzahl Wörter

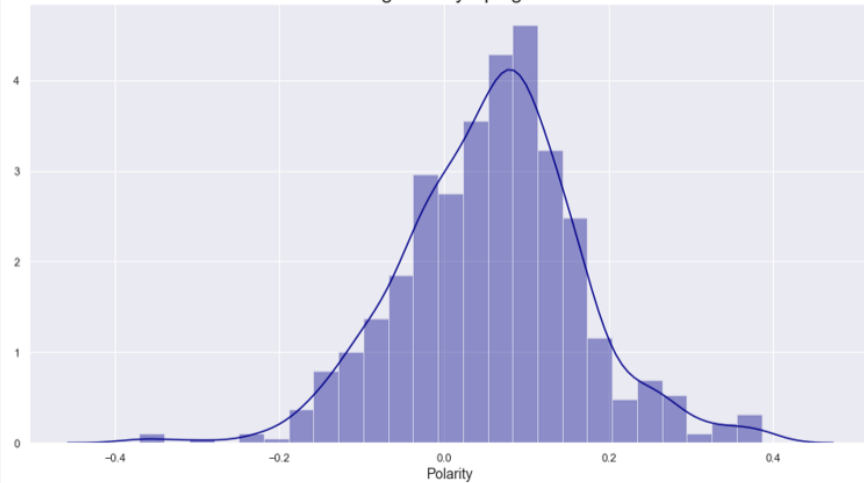


SPIEGEL Sentiments-Analyse Polarity / Subjectivity

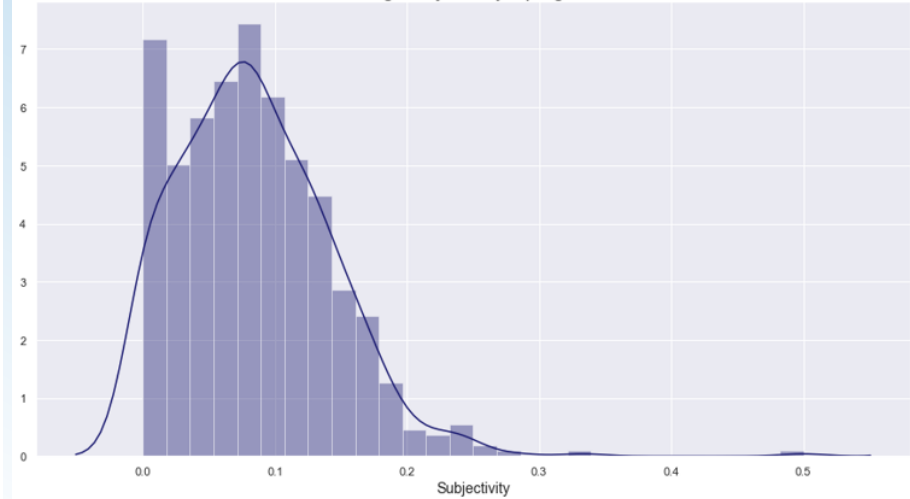


**FACHHOCHSCHULE
WIENER NEUSTADT**
Austrian Network for Higher Education

Verteilung Polarity Spiegel Artikel



Verteilung Subjectivity Spiegel Artikel

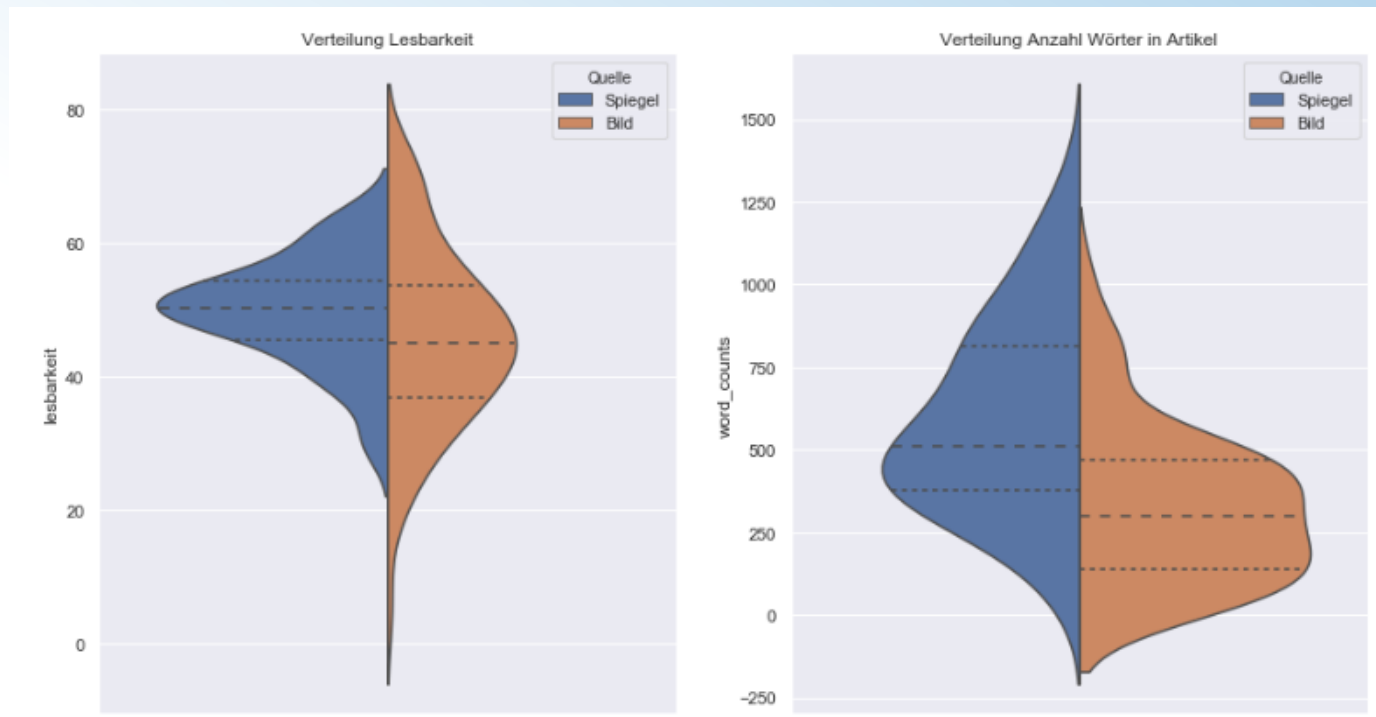


Vergleich BILD und SPIEGEL – Violin-Plot Lesbarkeit / Anzahl Wörter

Zeitraum 11 – 12 / 2019



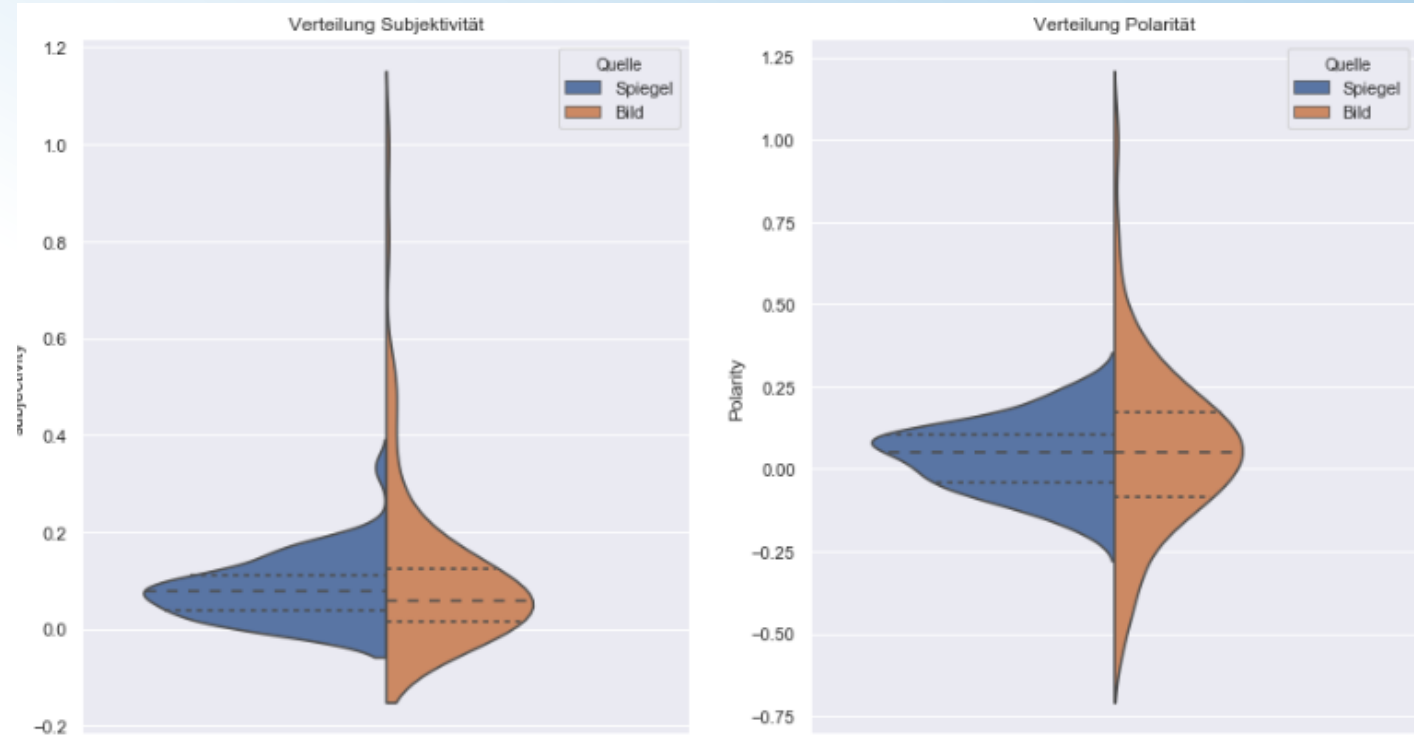
**FACHHOCHSCHULE
WIENER NEUSTADT**
Austrian Network for Higher Education



Vergleich BILD und SPIEGEL – Violin-Plot Sentiments-Analyse Zeitraum 11 – 12 / 2019



**FACHHOCHSCHULE
WIENER NEUSTADT**
Austrian Network for Higher Education



SPIEGEL

	word_counts	lesbarkeit	Polarity	subjectivity
count	53	53	53	53
mean	599.208	49.853	0.04	0.082
std	293.33	8.004	0.104	0.064
min	57	29.32	-0.184	0
25%	378	45.59	-0.04	0.04
50%	509	50.16	0.048	0.079
75%	815	54.32	0.103	0.111
max	1342	64	0.26	0.333

BILD

	word_counts	lesbarkeit	Polarity	subjectivity
count	72	72	72	72
mean	329.431	45.51	0.056	0.118
std	231.871	13.386	0.245	0.178
min	26	5.3	-0.5	0
25%	140.25	36.9	-0.084	0.014
50%	301.5	44.955	0.048	0.059
75%	469	53.717	0.172	0.125
max	1037	72.53	1	1

var_name	ttest_statistik	ttest_pvalue	mann_statistik	mann_pvalue
word_counts	5.741	0.000	895.0	0.000
lesbarkeit	2.101	0.038	1418.5	0.007
Polarity	-0.458	0.648	1823.0	0.336
subjectivity	-1.419	0.158	1827.5	0.344

Fragen?



FACHHOCHSCHULE
WIENER NEUSTADT

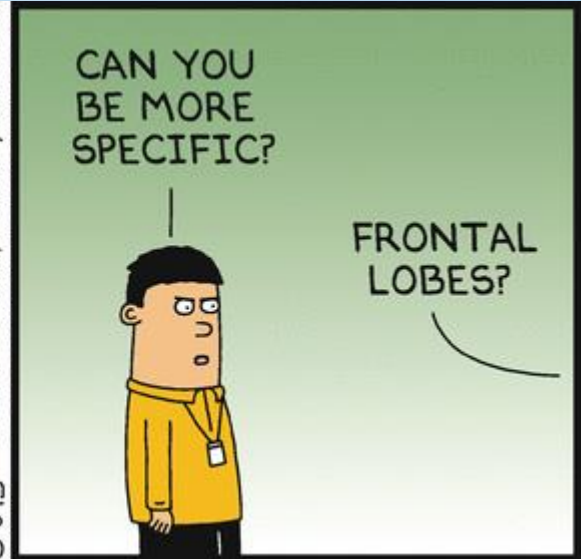


YES.

Dilbert.com DilbertCartoonist@gmail.com



8-6-15 © 2015 Scott Adams, Inc. /Dist. by Universal Uclick



FRONTAL LOBES?

- W. W. A. Initiative (WAI), „Easy-to-Read on the Web Online Symposium | Web Accessibility Initiative (WAI) | W3C“, W3C Web Accessibility Initiative (WAI). [Online]. Verfügbar unter: <https://www.w3.org/WAI/RD/2012/easy-to-read/>. [Zugegriffen: 28-Dez-2019].
- „G153: Making the text easier to read | Techniques for WCAG 2.0“. [Online]. Verfügbar unter: <https://www.w3.org/TR/2012/NOTE-WCAG20-TECHS-20121013/G153>. [Zugegriffen: 28-Dez-2019].
- „Europäische Regeln | Easy-to-Read“, 2017. [Online]. Verfügbar unter: <http://easy-to-read.eu/de/europaische-standards/>. [Zugegriffen: 28-Dez-2019].
- J. Philippi, *Einführung in die generative Grammatik*, Bd. 12. Vandenhoeck & Ruprecht, 2008.
- R. F. Flesch, *Art of readable writing*. New York: Harper & Row, 1949.
- T. Amstad, *Wie verständlich sind unsere Zeitungen?* Studenten-Schreib-Service, 1978.
- „fleschindex.de | Flesch-Index berechnen“. [Online]. Verfügbar unter: <http://fleschindex.de/berechnen>. [Zugegriffen: 28-Dez-2019].
- R. Gunning, *The technique of clear writing*. McGraw Hill Higher Education, 1952.
- A. Buttlar, *Grundzüge des Schulsystems der USA*. Wiss. Buchges., 1992.
- K. Fischer, *Satzstrukturen im Deutschen und Englischen: Typologie und Textrealisierung*, Bd. 1. Walter de Gruyter, 2013.
- R. Bamberger und E. Vanecek, *Lesen, verstehen, lernen, schreiben: die Schwierigkeitsstufen von Texten in deutscher Sprache*. Jugend und Volk, 1984.
- G. H. Mc Laughlin, „SMOG grading-a new readability formula“, *Journal of reading*, Bd. 12, Nr. 8, S. 639–646, 1969.
- C. H. Björnsson, *Läsbarhet*. Liber, 1968.
- W. Lenhard und A. Lenhard, „Berechnung des Lesbarkeitsindex LIX nach Björnson“. Unpublished, 2011.
- „Hohenheimer Verständlichkeitsindex: Klartext-Initiative“. [Online]. Verfügbar unter: <https://klartext.uni-hohenheim.de/hix>. [Zugegriffen: 28-Dez-2019].
- J. A. Pfeffer, „Grunddeutsch: Erarbeitung und Wertung dreier deutscher Korpora; ein Bericht aus dem Institut for Basic German, Pittsburgh“, 1975.
- R. Jones und E. Tschirner, *A frequency dictionary of German: Core vocabulary for learners*. Routledge, 2015.
- „Der Rat“. [Online]. Verfügbar unter: <http://www.rechtschreibrat.com/der-rat/>. [Zugegriffen: 28-Dez-2019].
- „Duden | Rechtschreibregeln“. [Online]. Verfügbar unter: <https://www.duden.de/sprachwissen/rechtschreibregeln>. [Zugegriffen: 28-Dez-2019].

- A.-H. Tan, „Text mining: The state of the art and the challenges“, gehalten auf der Proceedings of the PAKDD 1999 Workshop on Knowledge Discovery from Advanced Databases, 1999, Bd. 8, S. 65–70.
- A. Mehler und C. Wolff, „Einleitung: Perspektiven und Positionen des Text Mining“, gehalten auf der LDV-Forum, 2005, Bd. 20.
- B. Pang und L. Lee, „Opinion mining and sentiment analysis“, *Foundations and Trends® in Information Retrieval*, Bd. 2, Nr. 1–2, S. 1–135, 2008.
- T. M. Mitchell, *Machine learning*. WCB. McGrawHill, 1997.
- C. M. Bishop, *Neural networks for pattern recognition*. Oxford university press, 1995.
- W. Ertel, „Grundkurs Künstliche Intelligenz“, *Auflage*, Wiesbaden, 2009.
- „Human Brain Project Home“. [Online]. Verfügbar unter: <https://www.humanbrainproject.eu/en/>. [Zugegriffen: 28-Dez-2019].
- J. J. Hopfield, „Neural networks and physical systems with emergent collective computational abilities“, *Proceedings of the national academy of sciences*, Bd. 79, Nr. 8, S. 2554–2558, 1982.
- „Deep Learning“. [Online]. Verfügbar unter: <http://ufldl.stanford.edu/?papers>. [Zugegriffen: 28-Dez-2019].
- „Software links « Deep Learning“, 2017. [Online]. Verfügbar unter: http://deeplearning.net/software_links/. [Zugegriffen: 28-Dez-2019].
- „AlphaGo: using machine learning to master the ancient game of Go“, Google, 27-Jän-2016. [Online]. Verfügbar unter: <https://www.blog.google/topics/machine-learning/alphago-machine-learning-game-go/>. [Zugegriffen: 28-Dez-2019].
- „SePL (Sentiment Phrase List) / Opinion Mining - Opinion Mining“. [Online]. Verfügbar unter: <http://www.opinion-mining.org/SePL-Sentiment-Phrase-List>. [Zugegriffen: 28-Dez-2019].