

Themenspezifische Gruppierung deutscher Online-Zeitungen mit Natural Language Processing

Bachelorarbeit

Georg Donner
Matrikel-Nummer 553821

Betreuer	Prof. Dr. Gefei Zhang
Erstprüfer	Prof. Dr. Gefei Zhang
Zweitprüfer	Prof. Dr. Barne Kleinen

Inhaltsverzeichnis

Abbildungsverzeichnis	IV
1 Einleitung	1
2 Grundlagen	2
2.1 Natural Language Processing	2
2.1.1 Pipeline	3
2.2 Machine Learning	5
2.2.1 Feature Engineering	5
2.2.2 Dimensionsionalitätsreduktion	6
2.2.3 Klassifizierung	7
3 Datenverarbeitung	8
3.1 Verwendete Tools	8
3.1.1 Python	8
3.1.2 SpaCy	9
3.1.3 Scikit-learn	9
3.2 Datenselektion	10
3.2.1 Datensatz	10
3.2.2 Normalisierung	11
3.3 Textverarbeitung	12
3.4 Featuregenerierung	14
3.4.1 Lexikalisch	14
3.4.2 Syntaktisch	16
3.4.3 Featureset	18

Inhaltsverzeichnis

4	Datenauswertung	19
4.1	Überblick	19
4.2	Klassifizierung	19
4.2.1	Vorbereitung des Datensatzes	19
4.2.2	Feature Selection	20
4.2.3	Messung der Performance	20
4.2.4	Verwendete Verfahren	20
4.3	Clustering	21
4.3.1	Feature Extraction	21
5	Ergebnis	23
	Literaturverzeichnis	24

Abbildungsverzeichnis

2.1	Dependency Relations	4
3.1	Spiegel Online Struktur	11
3.2	Anteil dpa Artikel pro Zeitung	13

1 Einleitung

Natural Language Processing ist ein großes Feld, welches besonders in der letzten Zeit im Zuge der Digitalisierung viel an Aufmerksamkeit und Wichtigkeit gewonnen hat. Es ermöglicht uns Informationen schneller zu finden, Systeme durch gesprochene Sprache zu steuern oder ganze Texte zu generieren. Eine weitere Aufgabe ist es, eine große Menge an Texten in Kategorien einzuteilen, um die Daten auf eine gewünschte Teilmenge für eine spezifischere Suche oder Analyse zu reduzieren. Die Kategorisierung der Dokumente nach ihrem Inhalt ist hier der häufigste Anwendungsfall, es ist aber auch möglich Texte nach ihrem generellen Genre oder Schreibstil zu vergleichen.

Diese Arbeit wird am Beispiel deutscher Online-Zeitungen untersuchen, welche Möglichkeiten es gibt Texte unabhängig von ihrem Inhalt zu vergleichen. Dabei werden lexikalische, morphologische und syntaktische Merkmale, aber auch die Verwendung inhaltlich irrelevanter Wörter als Features verwendet. Es wird überprüft, inwiefern die Artikel gruppiert werden können und Rückschlüsse auf Unterschiede im Schreibstil ganzer Zeitungen statt nur einzelner Artikel zulassen.

Des Weiteren wird untersucht, ob und wie sich der Schreibstil einer Zeitung je nach Thema, wie z.B. Politik und Sport, unterscheidet.

2 Grundlagen

Für die Verarbeitung und Analyse von Zeitungsartikeln sind zwei Teilgebiete der Informatik besonders wichtig: Natural Language Processing und Machine Learning. Im Folgenden werden die für die Arbeit relevantesten Konzepte der beiden Bereiche genauer erklärt.

2.1 Natural Language Processing

Natural Language Processing ist ein Teilgebiet der Informatik, das Konzepte und Techniken der künstlichen Intelligenz und des Machine Learning verwendet, um natürliche Sprache zu verarbeiten. Der Begriff natürliche Sprache wird verwendet, um menschliche Sprachen zu beschreiben, die im Gegensatz zu künstlich entwickelten Plansprachen eine historische Entwicklung durchlebt haben. Diese Sprachen befinden sich in einem dauerhaften Entwicklungsprozess und sind häufig sehr variabel in ihrer Verwendung durch den Menschen. Die Analyse der Semantik eines Wortes oder Satzes ist für Computer besonders schwierig, da sich die Bedeutung häufig erst durch den Kontext ergibt.

Mit den schnellen Fortschritten im Bereich des Machine Learning in den letzten Jahrzehnten, eröffneten sich für die Verarbeitung natürlicher Sprache jedoch völlig neue Möglichkeiten. Die Erkennung von Syntax und Semantik wurde damit immer präziser und das Teilgebiet immer relevanter. Gegenwärtig basiert dies hauptsächlich auf Algorithmen des Supervised Learning, für die die Texte vorher manuell mit relevanten Markierungen versehen werden müssen. Ein bekanntes Beispiel für einen Korpus deutscher Sprache mit solchen Annotationen ist der TIGER Corpus ¹.

¹ <http://www.ims.uni-stuttgart.de/forschung/ressourcen/korpora/tiger.html>

TODO: Verwendung nicht annotierter Korpora

2.1.1 Pipeline

Bei der Analyse eines Textes werden in der Regel verschiedene Schritte abgearbeitet, die jeweils eigene Merkmale der Sprache untersuchen. Es entsteht eine sogenannte Pipeline, die je nach Anwendungsfall unterschiedlich aussieht. Das sequenzielle Ausführen dieser einzelnen Vorgänge ist notwendig, da beispielsweise die Analyse der Syntax voraussetzt, dass das Dokument bereits in Token zerlegt wurde. Im Folgenden werden die für diese Arbeit relevanten Schritte beschrieben.

TODO: eventuell noch andere Modelle beschreiben?

Tokenisierung

Tokenisierung beschreibt den Prozess, einen gegebenen Text in gleiche Einheiten, sogenannte Token, zu zerlegen. Meistens handelt es sich dabei um Wörter und Satzzeichen, je nach Anwendungsfall können es aber auch Wortgruppen oder Sätze sein. Die Tokenisierung ist häufig eine Voraussetzung einer tiefgehenderen Analyse des Textes, daher ist eine hohe Genauigkeit hier besonders wichtig. Eine Teilung an jedem Satzzeichen um eine Liste von Sätzen zu erhalten, ist zum Beispiel eine triviale Lösung, die bereits gute Ergebnisse erzielt. Es müssen jedoch viele Sonderregeln wie Abkürzungen und Zahlen beachtet werden, sodass das Problem wesentlich komplexer ist, als es zunächst scheint. Ein Ansatz um höhere Genauigkeit zu erreichen, ist ein Modell auf Basis bereits mit Annotationen versehener Korpora zu trainieren, welches die Regeln automatisch erlernt.

Part-of-speech-Tagging

Das Part-of-speech-Tagging, auch POS-Tagging, ist ein Verfahren, bei dem jedem Wort oder Satzzeichen die jeweilige Wortart zugeordnet wird. Bei der Analyse ist hierbei vor allem der Kontext in dem das Wort erscheint wichtig, da sich daraus häufig erst die Bedeutung ergibt. Die Informationen über die Wortart und oft auch weitere Details, geben sogenannte

2 Grundlagen

Tags, die meist aus einem festen Tagset stammen. In dieser Arbeit werden die Tags aus dem Stuttgart-Tübingen-Tagset (STTS) ² verwendet. Der Satz „Martin findet eine grüne Blechdose.“ sieht nach dem POS-Tagging beispielsweise so aus:

Martin/NE findet/VVFIN eine/ART grüne/ADJA Blechdose/NN ./.\$.

Die Tags geben Auskunft über mehr als nur die Wortart. Zum Beispiel sind die beiden Wörter *Martin* und *Blechdose* jeweils Nomen. Da jedoch beim POS-Tagging auch die Definition des Wortes überprüft wird, kann *Martin* korrekt als Eigennamen mit dem Tag NE identifiziert werden. Weiterhin wurde in diesem Satz die Verbform und der Adjektiv-Typ korrekt erkannt.

Dependency Parsing

Die Analyse der syntaktischen Struktur eines Satzes ist ein weiterer wichtiger Schritt in der Verarbeitung natürlicher Sprache. Beim Dependency Parsing wird zunächst jeder Satz nur auf seine Wörter reduziert, um dann die Beziehungen, sogenannte Dependency Relations, der Wörter innerhalb dieses Satzes zu bestimmen.

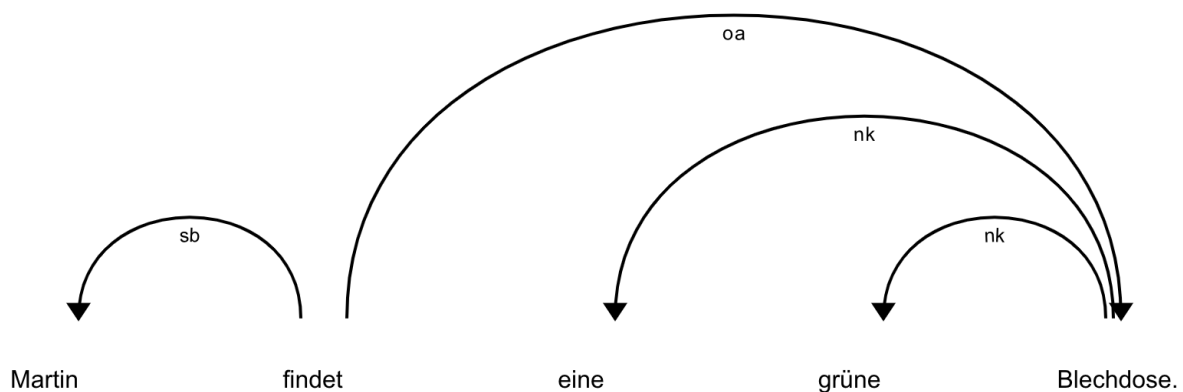


Abbildung 2.1: Visualisierung der Beziehungen mit Pfeilen

Diese Beziehungen ergeben letztendlich einen Baum, der navigiert werden kann und auch Aufschlüsse über die Komplexität eines Satzes zulässt. Die Ergebnisse des Dependency Parsing finden neben der Erkennung semantischer Beziehungen zwischen Wörtern auch noch

² <http://www.ims.uni-stuttgart.de/forschung/ressourcen/lexika/TagSets/stts-table.html>

2 Grundlagen

weitere Anwendungen. So werden sie z.B. von dem Natural Language Processing Tool SpaCy für die Erkennung der Satzenden verwendet [Exp19b].

Lemmatisierung

Für viele Analysen eines Textes ist es von Vorteil oder sogar notwendig, dass nicht jedes Wort welches unterschiedlich geschrieben wird, als ein anderes Wort behandelt wird. Um dies zu erreichen, werden alle Wörter auf ihre Grundform reduziert. Dieser Prozess heißt Lemmatisierung. Dies ist besonders hilfreich für die Feststellung der Häufigkeit eines Wortes in einem Dokument oder Korpus. So wäre die Frequenz der Wörter *findet*, *fand* und *finden* zunächst jeweils eins. Nach der Lemmatisierung gibt es nur noch das Lemma *finden* mit einer Frequenz von drei. Dies hilft dabei, das Rauschen innerhalb eines Textes zu reduzieren und das tatsächliche Vokabular genauer zu beurteilen.

2.2 Machine Learning

Die meisten der zuvor in Kapitel 2.1.1 erläuterten Schritte basieren auf Machine Learning und auch in dieser Arbeit werden einige Verfahren aus diesem Bereich angewandt, um den Schreibstil der Texte zu ermitteln und zu analysieren. Das Ziel des Machine Learning ist es, aus bestehenden Beobachtungen ein Modell zu trainieren, welches in diesen ein Muster oder Zusammenhänge erkennen soll und anschließend Aussagen über eine bisher unbekannte Beobachtung treffen kann. In diesem Kapitel werden einige für diese Arbeit relevante Themenbereiche aus dem Machine Learning genauer beschrieben.

2.2.1 Feature Engineering

Alle Beobachtungen aus welchen ein zu trainierendes Modell lernt, haben immer eine bestimmte Anzahl an Features, die jede Beobachtung so gut wie möglich repräsentieren sollen. Die Auswahl dieser Features bildet die Grundlage des Trainings und ist ausschlaggebend für die Performance des Modells. Die Verwendung voneinander unabhängiger Features, die mit der vorherzusagenden Klasse in Beziehung stehen, ist anzustreben um

2 Grundlagen

gute Resultaten zu erhalten [Dom12, Kap. 8]. Die Auswahl der richtigen Features ist also der Schlüssel für das erfolgreiche Lernen, aber es ist auch der schwierigste Teil. „Coming up with features is difficult, time-consuming, requires expert knowledge.“ [Ng13] Häufig werden viele Features generiert bzw. bestimmt, die möglicherweise nicht alle zielführend oder unabhängig sind. Hier gibt es verschiedene Prozesse, um die Anzahl an Features zu reduzieren und gleichzeitig kaum an Informationsgehalt zu verlieren. Die Auswahl von Features aus der großen Menge heißt *Feature Selection*. Die Kombination von Features, um daraus neue zu generieren, wird *Feature Extraction* genannt. Diese Reduktion ist oftmals erforderlich um das statistische Rauschen in den Daten zu minimieren und beim Training Zeit zu sparen. Besonders bei sehr rechenintensiven Algorithmen ist dies absolut notwendig.

2.2.2 Dimensionsionalitätsreduktion

Bei der Generierung von d verschiedenen Features bzw. Variablen für n Beobachtungen, entsteht ein Raum der Dimension d . Um so viele Informationen wie möglich zu jeder Beobachtung zu speichern und für z.B. eine Klassifizierung zu verwenden, werden sehr viele Features generiert, welche wiederum einen hoch-dimensionalen Raum ergeben. Dabei entsteht das Problem, dass die Datenverarbeitung immer mehr Zeit beansprucht, da auch wesentlich mehr Beobachtungen notwendig sind. Weiterhin führt es zu Schwierigkeiten bei der Bestimmung von Distanzen zwischen einzelnen Datenpunkten. In solch einem Raum ist es unvermeidbar, dass die Verteilung der Punkte sehr dünn ist. Darauf basierend wurde festgestellt, dass in einem hoch-dimensionalen Raum alle Punkte annähernd die gleiche Entfernung haben und die Suche nach dem nächsten Nachbarn nur sehr ungenau ist [HAK00, Kap. 1]. Eine häufige Annahme bei der Betrachtung dieser Räume ist, dass die Variablen nicht alle voneinander unabhängig sind und die Anzahl an Dimensionen reduziert werden kann, ohne dabei einen signifikanten Anteil an Informationen zu verlieren. Ein dafür häufig verwendetes Verfahren ist die *Hauptkomponentenanalyse* oder englisch *Principal Component Analysis (PCA)*. Hierbei werden iterativ, je nach gewünschter Anzahl an Features, Linearkombinationen der Variablen erstellt, die durch Maximierung der Varianz der Daten in die entsprechende Richtung bestimmt werden und so jeweils die höchste Aussagekraft besitzen. Bei einer Reduktion auf zwei oder drei Dimensionen können die Daten dann gut visualisiert werden, dabei gehen aber meistens zu viele Informationen verloren.

2.2.3 Klassifizierung

Um die Performance der Features messen zu können, wird in dieser Arbeit Klassifizierung verwendet. Dafür werden alle Beobachtungen zunächst in ein Trainings- und ein Testset unterteilt, um das trainierte Modell nach dem Lernen mit den Testdaten beurteilen zu können. Nach dem Training ist die Aufgabe des Modells, jeder neuen Beobachtung aus dem Testset eine Klasse aus einer Liste bereits bekannter Klassen zuzuordnen. Die Klassen der jeweiligen Beobachtungen aus dem Trainingsset sind bekannt.

Falls es nur zwei mögliche Klassen gibt, handelt es sich um eine *Binärer Klassifikation*. Eine häufige Herangehensweise bei mehr als zwei Klassen ist es, die Aufgabe in mehrere binäre Probleme zu zerlegen, bei der nach der *One-vs.-all* Strategie entschieden wird, welche Klasse zugewiesen wird. Dabei wird für jede mögliche Klasse eine Wahrscheinlichkeit bzw. Sicherheit berechnet und anschließend gewinnt die mit dem höchsten Wert.

Klassifizierung ist ein häufiges Problem, für das es aufgrund der vielen Anwendungsbeispiele und damit völlig unterschiedlichen Anforderungen keinen besten Algorithmus zur Bestimmung der Klasse bzw. der Wahrscheinlichkeit gibt. Die meisten dieser Algorithmen erstellen eine lineare Funktion, in die der Featurevektor anschließend eingesetzt wird. Jedes Feature bekommt hierbei eine eigene Gewichtung, die sich aus dem vorherigen Training ergibt. Dieses Training unterscheidet sich je nach Algorithmus und auch der zurückgegebene Score wird je nach Algorithmus unterschiedlich interpretiert.

3 Datenverarbeitung

Die Grundlage dieser Arbeit bildet ein großer Datensatz mit Artikeln verschiedener Online-Zeitungen. Damit diese Artikel verglichen werden können oder eine Klassifizierung durchgeführt werden kann, müssen diese zunächst jeweils in einen Featurevektor umgewandelt werden. Dieses Kapitel beschreibt die verwendeten Tools und nötigen Schritte um dieses Ziel zu erreichen.

3.1 Verwendete Tools

3.1.1 Python

Python ist eine sehr universelle Programmiersprache und findet in den unterschiedlichsten Bereichen Anwendung. Die Syntax ist einfach zu erlernen und legt besonderen Wert auf gute Lesbarkeit, was wiederum die Produktivität der Programmierer erhöht [Fou19]. Des Weiteren gibt es eine große Auswahl an Bibliotheken für viele verschiedene Anwendungsfälle. Besonders in der wissenschaftlichen Arbeit hat sich die Verwendung von Python zusammen mit den Packages SciPy, NumPy, Matplotlib und Pandas bewährt. Auf Basis dieser Kombination entwickelte sich die Anaconda Distribution¹, welche die wichtigsten Bibliotheken für Data Science und Machine Learning beinhaltet und die Verwaltung dieser pro Projekt erleichtert. Die Verwendung einer solchen Distribution ermöglicht es, out-of-the-box Optimierungen mathematischer Berechnungen durchzuführen, unter anderem auf Basis der IntelTM Math Kernel Library (MKL). Dadurch konnte in dieser Arbeit zum Beispiel die Analyse der Artikel mit spaCy um etwa 27% beschleunigt werden.

¹ <https://www.anaconda.com/what-is-anaconda/>

3.1.2 SpaCy

Auch für Natural Language Processing gibt es eine Vielzahl an Python Bibliotheken, die jedoch meist unterschiedliche Aufgaben erfüllen. *SpaCy* ist ein exzellentes Open-Source-Tool, welches die in Kapitel 2.1.1 beschriebenen Schritte der Pipeline vollständig abdeckt. Es legt besonderen Wert auf die Genauigkeit der Resultate und die Geschwindigkeit der Berechnungen, unter anderem durch die Verwendung und Optimierung des seiner Meinung nach besten Algorithmus für die jeweilige Aufgabe, statt mehrere Alternativen zur Verfügung zu stellen. SpaCy stellt trainierte Modelle für derzeit sieben verschiedene Sprachen zur Verfügung, weitere befinden sich bereits in der Entwicklung. Deutsch wurde im Jahr 2016 als erste Fremdsprache hinzugefügt und liefert trotz der reichhaltigeren Morphologie der Sprache sehr gute Ergebnisse [See16].

Eine häufig verwendete Alternative zu spaCy ist NLTK, kurz für Natural Language Toolkit. Es wurde in dieser Arbeit an einigen Stellen eingesetzt und für die Berechnung der linguistischen Merkmale in Betracht gezogen. Für die deutsche Sprache existiert hier jedoch kein fertiges Modell für das POS-Tagging, dieses müsste auf Basis eines annotierten Korpus selbst trainiert werden. Des Weiteren unterstützt NLTK kein Dependency Parsing, ein wichtiger Schritt der in dieser Arbeit verwendeten Pipeline. Die Verwendung und Ergebnisse dieser beiden Bibliotheken werden in Kapitel 3.3 genauer untersucht.

3.1.3 Scikit-learn

Scikit-learn ist eine Python Bibliothek, welche eine Vielzahl an effizienten Algorithmen aus dem Bereich des Machine Learning zur Verfügung stellt [PVG⁺11]. Sie baut auf den bereits in Kapitel 3.1.1 genannten Bibliotheken SciPy, NumPy und Matplotlib auf und bietet ebenfalls eine nahtlose Integration mit ihnen. Zudem gibt es eine sehr ausführliche und einsteigerfreundliche Dokumentation, die über die Beschreibung der in der Bibliothek enthaltenen Funktionen hinaus auch ausführliche Tutorials beinhaltet.

3.2 Datenselektion

Bevor die gegebenen Daten für das Extrahieren von Features verwendet werden, ist es ratsam, zunächst einen Überblick über die Struktur zu bekommen. Häufig befinden sich fehlerhafte Daten im Datensatz und auch das Format entspricht nicht immer den Erwartungen oder ist unregelmäßig. Die in dieser Arbeit getroffenen Maßnahmen, um den Datensatz bereit für die Textverarbeitung zu machen, werden in diesem Kapitel genauer erläutert.

3.2.1 Datensatz

Der in dieser Arbeit verwendete Datensatz wurde von Ole Wendt erstellt und enthält etwas mehr als zehn Millionen Artikel bis Juni 2018 für neun verschiedene Zeitungen. Tabelle 3.1 lässt jedoch erkennen, dass die Anzahl der Artikel je nach Zeitung sehr unterschiedlich ist. Für die Speicherung der Daten wurde das relationale Datenbanksystem *SQLite* gewählt.

Zeitung	Anzahl der Artikel
Ärzte-Zeitung	132753
Handelsblatt	942362
RP Online	1549537
SHZ	347335
Spiegel Online	620213
Der Tagesspiegel	975483
Westfälischer Anzeiger	318910
Welt	2980800
Zeit Online	2242198

Tabelle 3.1: Größe des Datensatzes

Die für diese Arbeit relevanten Datenfelder eines Artikels sind `content`, `url` und `date`, welche jeweils in Textform gespeichert sind. Es existieren weitere interessante Felder für z.B. den Autor oder die Zusammenfassung eines Artikels. Ein genauerer Blick in den Datensatz ergibt jedoch, dass diese bei dem Großteil der Einträge fehlen und somit nicht für die Analyse genutzt werden können.

3.2.2 Normalisierung

Für die Normalisierung der Daten werden zunächst alle Artikel entfernt, die weniger als 100 Zeichen haben. Dies ist notwendig, da ein kleiner Teil der Artikel entweder gar keinen Inhalt besitzt oder nur eine kurze Eilmeldung ist.

Eine weitere Auffälligkeit bei der Betrachtung des Datensatzes ist, dass das Datum je nach Zeitung in einem unterschiedlichen Format gespeichert ist. Zudem müssen Artikel ohne Datum herausgefiltert werden. Die Sortierung der Artikel nach Datum funktioniert trotzdem sehr zuverlässig, sodass die Artikel erst sortiert und anschließend das Datum formatiert werden kann. Obwohl das Datum später nicht direkt weiterverarbeitet wird, ist ein einheitliches Format, in dieser Arbeit wird der ISO 8601 Standard verwendet, dennoch wichtig für die Visualisierung der Daten. Für diesen Schritt ist die Python Funktion `datetime.strptime` sehr hilfreich, da damit auch Daten wie „1. März 2018, 9:11 Uhr“ einfach geparkt werden können.

Die Klassifizierung der Zeitungen soll *themenspezifisch* erfolgen, das heißt jedem Artikel soll idealerweise ein Thema bzw. eine Kategorie zugeordnet werden. Dafür werden im Rahmen dieser Arbeit sehr grobe Themen wie Politik, Wirtschaft oder Sport angestrebt. Für die Bestimmung dieser Themen bietet sich die URL des Artikels an, da die Webseiten der Zeitungen, wie Abbildung 3.1 verdeutlicht, meist nach solchen Themen strukturiert sind. Dies spiegelt sich dann im Format der URL wider.

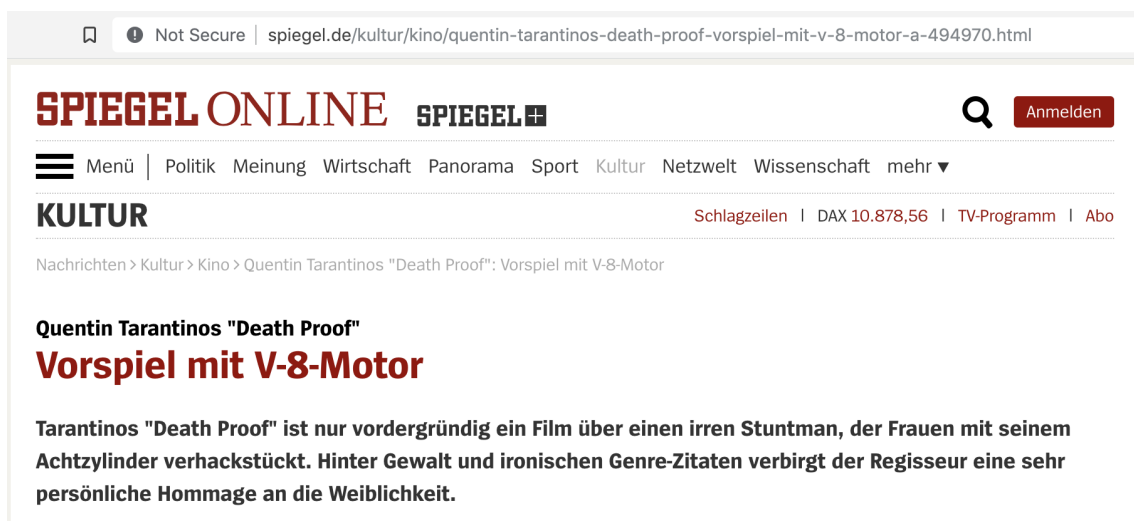


Abbildung 3.1: Die Struktur der Spiegel Online Webseite [Bor07]

Je nach Zeitung ergeben sich hierbei unterschiedliche Themen und nicht jedem Artikel kann auf diese Weise zuverlässig ein Thema zugeordnet werden. Zudem unterscheidet sich die Anzahl der Artikel je nach Thema stark, sodass einige aufgrund zu kleiner Repräsentation aussortiert werden.

3.3 Textverarbeitung

Bei der Textverarbeitung durchläuft eine Auswahl an Artikeln die in Kapitel 2.1.1 beschriebene Pipeline und die Ergebnisse werden anschließend für die weitere Verwendung zwischengespeichert.

Der erste Ansatz für die Verarbeitung der Texte war es, für jeden Schritt in der Pipeline ein explizit für dieses Verfahren entwickeltes Tool zu verwenden. Für die Tokenisierung wurde hier zunächst die populäre Bibliothek NLTK verwendet, welche bereits sehr gute Resultate bei annehmbarer Performanz lieferte. Für das POS-Tagging wurde der von Helmut Schmid entwickelte TreeTagger² verwendet, welcher zusätzlich das Lemma jedes Wortes bestimmt. Dieser unterstützt eine Vielzahl an Sprachen, da pro Sprache lediglich eine Parameterdatei notwendig ist. Die Verwendung des TreeTaggers bringt jedoch auch einige Nachteile mit sich. Er wird nur als Perl Skript zur Verfügung gestellt und kann so nur über Umwege in die Python Pipeline integriert werden. Zudem ist die Startzeit des Skriptes pro Artikel sehr hoch, sodass die Berechnung zu zeitaufwendig wird. Dies kann umgangen werden, indem die Artikel zuvor zu einem großen Text zusammengefügt und anschließend wieder getrennt werden und führte zu einer sehr hohen Performanz. Dadurch werden jedoch Ungenauigkeiten in den Ergebnissen riskiert und an einigen Stellen stoppte das Skript ohne einen Fehler zu werfen. Ein weiteres Problem dieses Ansatzes ist, dass durch die Ausführung der Pipeline in einzelnen Schritten die Ergebnisse jedes Schrittes zwischengespeichert werden müssen und eventuell später zu groß für den Arbeitsspeicher werden.

Im zweiten Ansatz wurde mehr Wert auf die Zuverlässigkeit der Ergebnisse statt der Optimierung der Geschwindigkeit des Prozesses gelegt. Hierfür wurden zunächst 1000 Artikel pro Zeitung pro Thema verwendet, insgesamt 43000 Artikel, und mit der Bibliothek spaCy analysiert. Die Erstellung und Ausführung einer Pipeline mit spaCy ist sehr intuitiv

² <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

3 Datenverarbeitung

und erfordert kaum Setup. Dennoch gibt es auch hier Möglichkeiten, die Performanz zu optimieren.

TODO: Überblick spaCy Geschwindigkeit und Optimierungen

Die Anzahl der analysierten Artikel wurde zunächst beschränkt, um die Größe der Dateien, welche die Ergebnisse der Textverarbeitung beinhalten, gering zu halten. Für 43000 Artikel werden hier etwa 3,4 GB benötigt, indem nur die für diese Arbeit relevanten linguistischen Merkmale der Token gespeichert werden. Dies ermöglicht einen schnellen, iterativen Prozess für die Generierung, Auswahl und Bewertung der in Kapitel 3.4 beschriebenen Features.

Nach der Analyse der Texte wurde eine weitere Filterung der Artikel vorgenommen, um den Anteil fehlerhafter Daten zu reduzieren. Es wurden hierfür jegliche Sätze mit weniger als vier Wörtern entfernt. Diese rigorose Filterung war notwendig da viele Artikel Wörter beinhalten, die nicht zum eigentlichen Inhalt gehören. Beispiele dafür sind Kürzel der Nachrichtenagentur, der Ort, der Autor, Quellen oder Verweise am Anfang oder Ende des Artikels.

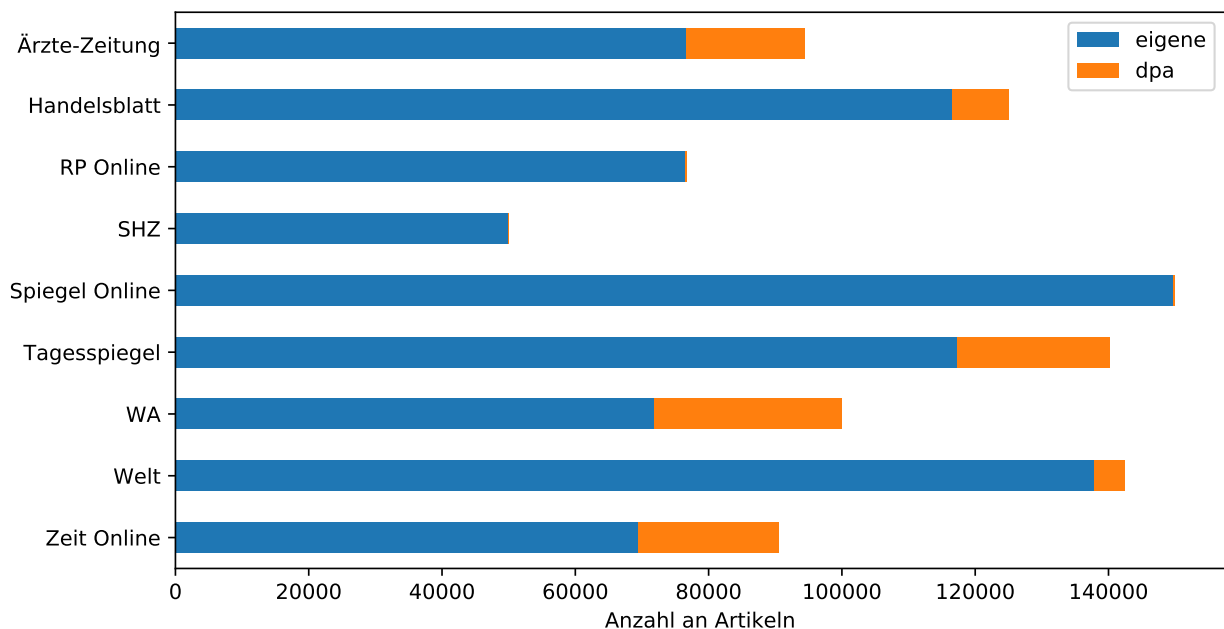


Abbildung 3.2: Anteil der dpa Artikel pro Zeitung

Zudem wurden Artikel der Nachrichtenagentur dpa herausgefiltert, um nur die Artikel

von eigenen Autoren der Zeitungen zu vergleichen. Wie in Abbildung 3.2 zu sehen, unterscheidet sich der Anteil je nach Zeitung sehr stark, ist allerdings auch abhängig von der Kategorie. So ist der Gesamtanteil an dpa Artikeln der *Zeit* etwa 25%, betrachtet man nur die Politik Artikel sind es nicht einmal 0,1%.

3.4 Featuregenerierung

Es gibt eine Vielzahl an Features, die für einen Text generiert bzw. ausgewählt werden können. In dieser Arbeit werden jedoch nur Features verwendet, welche unabhängig vom Inhalt der Texte sind. Auf diese Weise wird nur der Schreibstil einer Zeitung analysiert und die Auswahl der Features wird demnach für jeden weiteren Korpus einen sehr ähnlichen Einfluss besitzen. Besonders im Bereich der Erkennung der Autorschaft ist es ein häufiger Ansatz, die Texte anhand ihrer linguistischen Merkmale zu vergleichen. E. Stamatatos unterteilt diese weiter in lexikalische, syntaktische und semantische Features [Sta09]. Auch in der automatischen Erkennung von Textgenres, wie z.B. Literatur und Wissenschaft, wird geforscht, wie gut sich linguistische Merkmale für die Klassifizierung von Texten eignen [CWD⁺17]. Im weiteren Verlauf dieses Kapitels werden die in dieser Arbeit verwendeten Features genauer beschrieben und deren Verwendung begründet.

3.4.1 Lexikalisch

Für die Bestimmung lexikalischer Features ist lediglich die Tokenisierung eines Textes notwendig und erfordert somit einen vergleichsweise geringen Aufwand. Die einfachsten Merkmale sind hierbei die durchschnittliche Satz- oder Wortlänge und werden schon seit den Anfängen der Forschung im Bereich der Erkennung der Autorschaft oder des Genres eines Textes verwendet [SFK00, S. 473]. Obwohl sie nicht besonders allgemeingültig sind, spielen diese Merkmale doch immer noch eine Schlüsselrolle in der Unterscheidung von Textgenres [CWD⁺17] und sind daher auch in dieser Arbeit Teil des Featuresets.

Weiterhin existieren diverse Indizes, welche versuchen, die Lesbarkeit auf einer Skala von 0 (einfach) bis 100 (schwer) zu ermitteln oder dem Text eine Schulstufe zuzuordnen,

3 Datenverarbeitung

die abgeschlossen sein muss, um den Text verstehen zu können. Diese berechnen sich meist aus einer Kombination der Anzahl an Silben, Wörtern und Sätzen im Text. Einer der ältesten und populärsten Indizes ist der 1948 veröffentlichte Flesch-Grad [MP82]. Dieser wurde zwar ursprünglich für die englische Sprache entwickelt und ist damit nicht besonders gut geeignet für den in dieser Arbeit verwendeten Korpus, aber Toni Amstad konnte den Flesch-Grad später auf die deutsche Sprache übertragen [Ams78]. Dieser berechnet sich wie folgt:

$$FG = 180 - \frac{W}{T} - (58,5 \cdot \frac{S}{W})$$

S ist die Gesamtzahl der Silben

W ist die Gesamtzahl der Wörter

T ist die Gesamtzahl der Sätze

In dieser Arbeit soll überprüft werden, ob die Lesbarkeit eines Artikels zur Erkennung der Zeitung beiträgt und wird daher zunächst in das verwendete Featureset aufgenommen und später genauer analysiert.

Ein weiteres, sehr häufig gemessenes Merkmal eines Textes ist die Reichhaltigkeit des Vokabulars. Ein typisches Beispiel hierfür ist die Type-Token-Relation V/N , wobei V für die Gesamtzahl der einzigartigen Wörter und N für die Gesamtzahl aller Wörter steht [Sta09, S. 540]. Für ein genaueres Resultat bei vor allem kurzen Texten ist es empfehlenswert, die Lemmata der Wörter zu verwenden, da sich dadurch die Frequenzen der Wörter erhöhen und dabei helfen können, das Rauschen zu minimieren. Es zeigte sich jedoch, dass dieses und weitere Maße für die Reichhaltigkeit des Vokabulars abhängig von der Länge des Textes sind und besonders für kürzere Texte an Aussagekraft verlieren [TB98]. A. Cimoni et al. zeigten allerdings, dass die Type-Token-Relation für die Erkennung bestimmter Genres eine wichtige Rolle spielt [CWD⁺17, S. 5] und auch für die Erkennung der Autorschaft wird dieses Maß immer noch verwendet [Sta09, S. 540]. Aus diesen Gründen wird die Type-Token-Relation auch in dieser Arbeit in das Featureset aufgenommen.

Die Frequenzen der Wörter oder Lemmata eines Textes zu berechnen, ist ein sehr einfacher, aber doch sehr typischer Ansatz, um den gegebenen Text in einen Featurevektor umzuwandeln. Um diese miteinander vergleichen zu können, muss vorher bekannt sein, welche Wörter dafür insgesamt berücksichtigt werden. Eine sehr simple Methode hierfür ist das Bag-of-words Model, welches für jeden Text die Frequenz jedes im gesamten Korpus

3 Datenverarbeitung

vorkommende Wortes oder Lemmas berechnet. Kombiniert mit einem Index wie dem Tf-idf-Maß, liefert diese Darstellung eines Textes auch für die Erkennung der Autorschaft gute Ergebnisse [Seb02, S. 12]. Mit dieser Herangehensweise wird neben dem Schreibstil aber auch maßgeblich der Inhalt des Textes repräsentiert, weswegen diese Methode in dieser Arbeit nicht in Betracht gezogen wird. Hier ergibt es Sinn, nur für Synsemantika, Wörter welche keine inhaltliche Relevanz besitzen, die Frequenzen zu berechnen. Auf diese Weise ist die Analyse des Textes weniger vom Inhalt abhängig und die verwendeten Wörter sind repräsentativer, da sie von jedem Autor genutzt werden. Zudem wird argumentiert, dass die Verwendung von Synsemantika durch den Autor weniger kontrolliert ist und somit den Schreibstil stärker prägen [Kes14, S. 60-61]. Die Verwendung dieser Frequenzen als Features hat sich sowohl im Bereich der Erkennung der Autorschaft, als auch des Genres als sehr nützlich erwiesen [CWD⁺17, Sta09, ZLCH06, AL05]. Dafür werden in Anlehnung an [CWD⁺17] in dieser Arbeit die 100 häufigsten Lemmata des Korpus bestimmt und anschließend für jeden Artikel die Frequenzen errechnet.

Es gibt viele weitere Ansätze, Texte nur mithilfe von lexikalischen Merkmalen in einen Featurevektor umzuwandeln. Eine sehr bewährte Methode ist es, statt ganzen Wörtern oder Lemmata, N-Gramme der Buchstaben als Features zu verwenden [Sta09, S. 542]. So wird bei der Berechnung von Trigrammen zum Beispiel das Wort *Rausch* in *Rau*, *aus*, *usc* und *sch* zerlegt. Untersuchungen von E. Stamatatos ergaben, dass diese Methode selbst bei der Verwendung von Trainings- und Testkorpora zu einem jeweils unterschiedlichen Thema, noch gute Ergebnisse erzielt und N-Gramm-Frequenzen der Buchstaben als Features wesentlich robuster als Wortfrequenzen sind [Sta13]. Ein hiermit einhergehendes Problem ist jedoch die unglaublich hohe Dimensionalität der daraus resultierenden Featurevektoren. Für das Clustering oder die Gruppierung von Artikeln sind diese Ansätze damit weitaus schlechter geeignet als für die Klassifizierung.

3.4.2 Syntaktisch

Ein wesentlich aufwendigeres Verfahren um weitere Features zu generieren, ist die Analyse der Syntax. Es wird argumentiert, dass sich Autoren unbewusst syntaktische Muster beim Schreiben angewöhnen und die Syntax damit den Schreibstil zuverlässiger repräsentiert als lexikalische Merkmale [Sta09, S. 542]. Auch für die Erkennung von Genres wurde

3 Datenverarbeitung

festgestellt, dass syntaktische Merkmale eine Schlüsselrolle spielen und die Genauigkeit signifikant erhöhen [CWD⁺17].

Für die Bestimmung syntaktischer Merkmale eines Textes ist allerdings nicht nur die Tokenisierung notwendig, sondern auch das POS-Tagging. Da dieses Verfahren nicht komplett akkurat ist, wird durch die Verwendung darauf aufbauender Features möglicherweise weiteres Rauschen hinzugefügt, welches die Ergebnisse beeinträchtigen kann. Da in dieser Arbeit nur Zeitungsartikel analysiert werden, welche kaum syntaktische Fehler enthalten und die Genauigkeit für das POS-Tagging von spaCy derzeit bei 97.15% liegt [Exp19a], wird davon ausgegangen, dass die Verwendung syntaktischer Features keinen signifikanten negativen Einfluss hat. In dieser Arbeit wurden auch die Ergebnisse des Dependency Parsing mit in Betracht gezogen, da diese einen noch tieferen Einblick in die Struktur der Sätze erlauben. Das Risiko des Rauschens ist aufgrund der geringeren Genauigkeit von 89.75% von spaCy [Exp19a] jedoch hier wesentlich höher.

Ähnlich zu den Wortfrequenzen als lexikalisches Merkmal, werden in dieser Arbeit auch für die POS-Tags die Frequenzen für jeden Artikel berechnet, hierbei spricht man auch von POS-Monogrammen. Analysen von A. Cimoni ergaben, dass sich die Verwendung von Bigrammen der POS-Tags im Vergleich zu Monogrammen nicht signifikant auf die Genauigkeit der Klassifizierung auswirkt [CWD⁺17], daher wurde diese Möglichkeit in dieser Arbeit ebenfalls nicht weiter erforscht. Es ist wichtig zu erwähnen, dass die Frequenzen der einzelnen Wortarten für jeden Artikel teilweise nur sehr gering sind. Hier erhöht sich die Aussagekraft der Frequenzen als Features ebenfalls mit der Länge des Artikels.

Die lexikalische Dichte eines Textes ist sein Verhältnis von Autosemantika bzw. Inhaltswörtern zu der Gesamtzahl an Wörtern. Eine höhere lexikalische Dichte bedeutet, dass der Text die enthaltenen Informationen präziser und verständlicher wiedergibt. Für die Berechnung des Anteils der Autosemantika ist das POS-Tagging die Grundlage, da in der deutschen Sprache die Wortart bestimmt, ob ein Wort eine eigene lexikalische Bedeutung besitzt. Diese Wortarten sind Substantive, Verben, Adjektive und Adverbien.

Weiterhin werden die Monogramme der beim Dependency Parsing berechneten Dependency Relations als Features verwendet, indem die jeweiligen Frequenzen pro Artikel bestimmt werden. Auch hier ist es möglich Bigramme zu berechnen, allerdings würde das die

3 Datenverarbeitung

Dimensionalität der Featurevektoren beträchtlich erhöhen und die Frequenzen verringern. Damit spiele die Länge der Texte eine noch größere Rolle.

3.4.3 Featureset

In dieser Arbeit werden insgesamt 201 Features verwendet und untersucht. Im Folgenden wird ein Überblick über das Featureset gegeben:

1. Durchschnittliche Satzlänge
2. Flesch-Grad
3. Type-Token-Relation
4. Lexikalische Dichte
- 5-59. Frequenzen der Monogramme aller Wortarten (POS-Tags)
- 60-101. Frequenzen der Monogramme aller Dependency Relations
- 102-201. Frequenzen der 100 häufigsten Lemmata des Korpus

4 Datenauswertung

Wie können die Daten überhaupt ausgewertet werden? Es müssen gleichzeitig die Performance der Features, als auch des Algorithmus zur Klassifizierung analysiert und ausgewertet werden.

4.1 Überblick

Ein paar Plots zeigen, die die Korrelation ausgewählter Features zeigen?

4.2 Klassifizierung

Warum versuche ich die Klassifizierung? Weil wenn Klassifizierung nicht wirklich funktioniert und kein Model trainiert werden kann, welches die Zeitung auf Basis der Daten predicten kann, dann ist auch keine Gruppierung möglich.

Ist zudem eine sehr gute Methode dafür festzustellen, an welchen Features man die Zeitungen am besten voneinander unterscheiden kann und welche eventuell völlig unnötig sind und so bei der Gruppierung nur stören würden.

4.2.1 Vorbereitung des Datensatzes

Für das Training muss der Datensatz in Train/Testset aufgeteilt werden. Welche Aufteilung habe ich hier gewählt und warum?

4 Datenauswertung

Die Daten müssen vorher normalisiert oder standardisiert werden. Kommt allerdings auch auf die ausgewählten Features und den gewählten Algorithmus an.

4.2.2 Feature Selection

Welche Verfahren wurden zur Feature Selection verwendet? Warum? Warum war das überhaupt nötig. Vergleich zu bzw. Verweis auf bisherige(n) Studien zur Feature Selection zur Genre/ Author detection.

(evtl. welche Features sind pro Zeitung relevant gewesen? weil unterschiedliche thetas pro Kategorie bei One vs Rest Klassifizierung)

4.2.3 Messung der Performance

Welche verschiedenen Ansätze gibt es hier?

F-Maß

Wahrheitsmatrix

Beurteilung der Wahrscheinlichkeiten

Wahrheitsmatrix kann auch mit den probabilities gemacht werden. Sagt noch mehr über die Sicherheit der Vorhersagen aus.

4.2.4 Verwendete Verfahren

Warum müssen verschiedene Verfahren getestet werden?

Linear Discriminant Analysis

Logistische Regression

Random Forest

Support Vector Machines

4.3 Clustering

Zuerst wurde untersucht, ob beim Clustering der einzelnen Artikeln jeder Zeitung Cluster entstehen, welche die einzelnen Zeitungen repräsentieren. Weiterhin wird überprüft, ob sich dabei Cluster ergeben, die Artikel verschiedener Zeitung haben. Das wäre dann schon ein sehr guter Indikator dafür, dass zwei oder mehr Zeitungen einen ähnlichen Schreibstil haben.

Ein weiterer Ansatz ist es, den Durchschnitt jeder Zeitung zu berechnen und anschließend die Zeitungen zu clustern. Dies hat jedoch Nachteile: Der Durchschnitt einer Zeitung ist nicht besonders repräsentativ, besonders wenn die Standardabweichung hoch ist. Zudem gibt es hier je nach Kategorie nur etwa 8 oder weniger "Beobachtungen" die geclustert werden können.

4.3.1 Feature Extraction

Wie in den Grundlagen bereits beschrieben, ist es beim Clustering besonders wichtig, dass die Dimensionalität nicht hoch ist. Vor allem für die Visualisierung ist es notwendig, die Features auf zwei Dimensionen zu reduzieren.

PCA

Hauptkomponentenanalyse: welche Ergebnisse gibt es hier?

4 Datenauswertung

t-SNE

Ein Verfahren, was genau dafür gedacht ist visuell Cluster zu zeigen, bei denen die Distanz untereinander im Plot keine Aussagekraft hat. Wie sehen hier die Ergebnisse aus? Spoiler: KACKEE

5 Ergebnis

Literaturverzeichnis

- [AL05] Argamon, Shlomo; Levitan, Shlomo: Measuring the usefulness of function words for authorship attribution. In: *Proceedings of the 2005 ACH/ALLC Conference*, 2005
- [Ams78] Amstad, T: *Wie verständlich sind unsere Zeitungen?*[How readable are our newspapers?], Doctoral thesis, Universität Zürich, Switzerland, Diss., 1978
- [Bor07] Borcholte, Andreas: *Quentin Tarantinos "Death Proof": Vorspiel mit V-8-Motor.* <http://www.spiegel.de/kultur/kino/quentin-tarantinos-death-proof-vorspiel-mit-v-8-motor-a-494970.html>, Juli 2007. – letzter Zugriff: 17.01.2019
- [CWD⁺17] Cimino, Andrea; Wieling, Martijn; Dell’Orletta, Felice; Montemagni, Simonetta; Venturi, Giulia: Identifying Predictive Features for Textual Genre Classification: the Key Role of Syntax. In: *CLiC-it 2017 11-12 December 2017, Rome (2017)*, S. 107
- [Dom12] Domingos, Pedro: A few useful things to know about machine learning. In: *Communications of the ACM* 55 (2012), Nr. 10, S. 78–87
- [Exp19a] ExplosionAI: *German - spaCy Models Documentation.* <https://spacy.io/models/de>, 2019. – letzter Zugriff: 23.01.2019
- [Exp19b] ExplosionAI: *Linguistic Features - spaCy Usage Documentation.* <https://spacy.io/usage/linguistic-features#dependency-parse>, 2019. – letzter Zugriff: 12.01.2019
- [Fou19] Foundation, Python: *What Is Python? Executive Summary.* <https://www.python.org/doc/essays/blurb>, 2019. – letzter Zugriff: 15.01.2019

Literaturverzeichnis

- [HAK00] Hinneburg, Alexander; Aggarwal, Charu C.; Keim, Daniel A.: What is the nearest neighbor in high dimensional spaces? In: *26th Internat. Conference on Very Large Databases*, 2000, S. 506–515
- [Kes14] Kestemont, Mike: Function Words in Authorship Attribution. From Black Magic to Theory? In: *Proceedings of the 3rd Workshop on Computational Linguistics for Literature (CLFL)*, 2014, S. 59–66
- [MP82] McCallum, Douglas R.; Peterson, James L.: Computer-based readability indexes. In: *Proceedings of the ACM'82 Conference* ACM, 1982, S. 44–48
- [Ng13] Ng, Andrew: *Machine Learning and AI via Brain simulations*. 2013
- [PVG⁺11] Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.; Cournapeau, D.; Brucher, M.; Perrot, M.; Duchesnay, E.: Scikit-learn: Machine Learning in Python. In: *Journal of Machine Learning Research* 12 (2011), S. 2825–2830
- [Seb02] Sebastiani, Fabrizio: Machine learning in automated text categorization. In: *ACM computing surveys (CSUR)* 34 (2002), Nr. 1, S. 1–47
- [See16] Seeker, Wolfgang: *spaCy now speaks German*. <https://explosion.ai/blog/german-model>, 2016. – letzter Zugriff: 16.01.2019
- [SFK00] Stamatatos, Efstathios; Fakotakis, Nikos; Kokkinakis, George: Automatic text categorization in terms of genre and author. In: *Computational linguistics* 26 (2000), Nr. 4, S. 471–495
- [Sta09] Stamatatos, Efstathios: A survey of modern authorship attribution methods. In: *Journal of the American Society for information Science and Technology* 60 (2009), Nr. 3, S. 538–556
- [Sta13] Stamatatos, Efstathios: On the robustness of authorship attribution based on character n-gram features. In: *Journal of Law and Policy* 21 (2013), Nr. 2, S. 421–439
- [TB98] Tweedie, Fiona J.; Baayen, R H.: How variable may a constant be? Measures of lexical richness in perspective. In: *Computers and the Humanities* 32 (1998), Nr. 5, S. 323–352

Literaturverzeichnis

- [ZLCH06] Zheng, Rong; Li, Jiexun; Chen, Hsinchun; Huang, Zan: A framework for authorship identification of online messages: Writing-style features and classification techniques. In: *Journal of the American society for information science and technology* 57 (2006), Nr. 3, S. 378–393

Erklärung

Hiermit versichere ich, dass ich die vorliegende Arbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe, dass alle Stellen der Arbeit, die wörtlich oder sinngemäß aus anderen Quellen übernommen wurden, als solche kenntlich gemacht und dass die Arbeit in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegt wurde.

Ort, Datum

Unterschrift