

# Syntaktische Profile für Interpretationen jenseits der Textoberfläche

## Andresen, Melanie

melanie.andresen@uni-hamburg.de  
Universität Hamburg, Deutschland

## Begerow, Anke

anke.begerow@haw-hamburg.de  
Hochschule für Angewandte Wissenschaften Hamburg

## Franken, Lina

lina.franken@uni-hamburg.de  
Universität Hamburg, Deutschland

## Gaidys, Uta

Uta.Gaidys@haw-hamburg.de  
Hochschule für Angewandte Wissenschaften Hamburg

## Koch, Gertraud

gertraud.koch@uni-hamburg.de  
Universität Hamburg, Deutschland

## Zinsmeister, Heike

heike.zinsmeister@uni-hamburg.de  
Universität Hamburg, Deutschland

Viele Verfahren des Text Mining und Distant Reading beschränken sich auf eine wortbasierte Auswertung von Texten. Auch wenn auf Basis der Wortformen und ihrer linearen Abfolge bereits neue Perspektiven auf Texte gewonnen werden (z. B. mittels der Voyant Tools, Sinclair / Rockwell 2016), schöpfen diese Methoden das Potential von Texten bei Weitem nicht aus. Insbesondere, wenn die Auswertungsergebnisse Gegenstand weiterführender Interpretationen werden, z. B. um soziale Phänomene zu beschreiben, sehen wir einen Mehrwert in der Auswertung zusätzlicher sprachlicher Strukturen. Konkret verwenden wir syntaktische Annotationen, die präzisere Informationen zu Wortkombinationen liefern können, etwa „X ist Subjekt von Y“ anstelle von „X steht im Kontext von Y“. Zudem bestehen viele syntaktische Relationen über eine längere Distanz an der Oberfläche hinweg und können deshalb nur durch eine syntaktische Perspektive erfasst werden (z. B. Andresen / Zinsmeister 2017). Dies gilt für unterschiedliche Sprachen in unterschiedlichem Ausmaß. Das Deutsche verfügt über deutlich mehr Distanzstrukturen als das Englische, für das die meisten Analyseverfahren ursprünglich entwickelt wurden.

In diesem Beitrag vergleichen wir zwei Ansätze zur Berechnung von Kollokationen, einen oberflächenorientierten Ansatz und einen auf Dependenzannotationen basierten. An zwei Fallstudien aus den Fächern Kulturanthropologie und Pflegewissenschaft wird demonstriert, wie die beiden Ansätze eine qualitative Interpretation von Textdaten in Hinblick auf gesellschaftliche bzw. soziale Phänomene unterstützen können. Die Erstellung eines eindeutigen Goldstandards, der eine formale Evaluation erlauben würde, ist bei dieser Art Fragestellung nicht möglich. Stattdessen wird auf qualitative Weise das Potential dieser Analyse zur Bearbeitung der geistes- und sozialwissenschaftlichen Fragestellungen beschrieben.

## Syntaktische Profile: Forschungsstand

Die Nutzung syntaktischer Informationen zur Charakterisierung von Wortverwendungen ist vor allem in der Lexikografie betrieben worden. Populär wurde das Konzept unter dem Namen „word sketch“ besonders durch die korpuslinguistische Software SketchEngine (Kilgarriff et al. 2004, Kilgarriff et al. 2014). Für ein gegebenes Suchwort wird hier angegeben, welche anderen Wörter besonders häufig in spezifischen syntaktischen Relationen zum Suchwort stehen, z. B. *glimpse* als frequentes Objekt von *catch* (Kilgarriff et al. 2014: 9). Im Digitalen Wörterbuch der deutschen Sprache (DWDS) kann eine entsprechende Darstellung als sog. DWDS-Wortprofil abgerufen werden (Geyken 2011).

Weitere Anwendungen gibt es in der Literaturwissenschaft und Linguistik: Googasian / Heuser (2019) vergleichen die syntaktischen Kontexte von Menschen und Tieren in einem Korpus sog. „wild animal stories“. Andresen (2018) nutzt syntaktische n-Gramme für einen Vergleich der Wissenschaftssprachen in den Fächern Literaturwissenschaft und Linguistik. Eine Anwendung auf sozialwissenschaftliche Fragestellungen erfolgt vor allem in den Politikwissenschaften: Anhand syntaktischer Muster wie z. B. Argumentstrukturen oder Mustern der Redewiedergabe werden hier politische Akteure und ihre Positionen identifiziert, miteinander in Relation gesetzt und so Diskurse charakterisiert (van Atteveldt et al. 2008, Kleinnijenhuis / van Atteveldt 2014, Wüest et al. 2011, Blessing et al. 2013). In den Fächern Pflegewissenschaft und Kulturanthropologie steht die Exploration des Mehrwertes syntaktischer Analysen noch aus.

## Daten und Fragestellungen der Fallstudien

Das Potential syntaktisch definierter Kollokationen wird an zwei Fallstudien mit komplementärer Datenlage erprobt. Die erste nutzt ein großes Korpus geschriebener Sprache, das dadurch methodisch eine sehr sichere Grundlage bietet. Die zweite Fallstudie basiert auf einem eher kleinen Korpus mit gesprochener Sprache, was mehr methodische Herausforderungen erwarten lässt.

Die kulturalanthropologische Fallstudie befasst sich mit dem Themenkomplex der Telemedizin und insbesondere der Frage nach der (Nicht-)Akzeptanz telemedizinischer Anwendungen durch unterschiedliche gesellschaftliche Akteursgruppen. Das hierfür erstellte Korpus umfasst 8.784 Texte mit insgesamt 14,8 Mio. Token und basiert auf einem Webcrawling (Adelmann / Franken 2020). Dafür wurden Webseiten von Krankenkassen, Ärzte- und Patientenverbänden als Ausgangspunkt genutzt und dann Links zu Seiten verfolgt, die mindestens eines von mehreren Wörtern aus einem Wortfeld zur Telemedizin enthalten (Koch / Franken im Druck).

Grundlage der pflegewissenschaftlichen Fallstudie ist ein Korpus aus 31 Dialogen, die mit schwerkranken und sterbenden Menschen in palliativer Versorgung geführt wurden. Es handelt sich um Transkripte gesprochener Sprache im Umfang von gut 100.000 Token. Gegenstand der Studie sind die Deutungen von Entscheidungen hinsichtlich der gesundheitlichen Versorgung der Betroffenen.

## Methode

Die Texte beider Korpora werden mithilfe des Parsers MATE (Bohnet 2010), trainiert auf der Hamburger Dependency Treebank (Foth et al. 2014), mit Lemmata, Wortarten und syntaktischen Abhängigkeiten annotiert. Unter Kollokationen verstehen wir „a combination of two words that exhibit a tendency to occur near each other in natural language“ (Evert 2008: 1214). Bei der Operationalisierung von „near each other“ können Kriterien an der Textoberfläche oder syntaktische Kriterien angesetzt werden: Für den einfachen Ansatz ohne Annotationen betrachten wir Wörter in einem Kontextfenster von  $\pm 3$  Wörtern als benachbart, für den syntaktischen Ansatz Wörter mit einer direkten Abhängigkeitsrelation. In beiden Fällen wird mithilfe des Log-Likelihood-Ratios (LLR, Dunning 1993) berechnet, welche Kombinationen häufiger im Korpus vorkommen, als basierend auf den Einzelfrequenzen der Wörter zu erwarten wäre. Im Falle der syntaktischen Kollokationen werden dafür die Einzelfrequenzen in der spezifischen syntaktischen Relation genutzt. Die Ergebnisse basieren auf den Lemmata und werden nach Schlüsselwörtern gefiltert, die für die jeweiligen Fragestellungen als

bedeutsam ausgewählt wurden (*T/telemed* bzw. *E/entscheid*). Das hierfür verwendete Analyseskript steht auf GitHub zur Verfügung.<sup>1</sup> Für die Interpretation werden die Top 10 beider Listen verglichen und nach Bedarf weitere Einträge gesichtet.

## Ergebnisse der Fallstudien

### Kulturalanthropologie

Tabelle 1 zeigt die oberflächenbasierten Kollokationen zu Lemmata mit *T/telemed* im kulturalanthropologischen Korpus mit den höchsten LLR-Werten. *Telemedizin* ist sehr stark mit dem verwandten Wort *Telematik* assoziiert, was die enge Verknüpfung der Bereiche anzeigt. Manche Wortpaare sind Bestandteil mehrteiliger Eigennamen (*Bayerische TelemedAllianz*, [*Zentrum für Telematik und Telemedizin GmbH*]), die für die Interpretation einen eingeschränkten Mehrwert haben, aber doch für den Diskurs potentiell relevante und ggf. bisher unbekannte Akteure sichtbar machen. Mit dem *Tag der Telemedizin* wird ein Fachkongress als wichtiger Begegnungspunkt dieser Akteure aufgeführt. Außerdem liegen allgemeine Konzepte wie *telemedizinisch* und *Anwendung* hoch im Ranking.

Tabelle 1: Top 10 der oberflächenbasierten Kollokationen zu Lemmata mit *T/ telemed* im kulturalanthropologischen Korpus

Wort 1	Wort 2	LLR	abs. Frequenz
Telematik	Telemedizin	2044,88	468
telemedizinisch	Anwendung	1753,90	465
bayerisch	TelemedAllianz	1497,94	204
Telemedizin	GmbH	1007,39	340
Tag	Telemedizin	845,74	274
telemedizinisch	Betreuung	841,88	212
bayerisch	Telemedallianz	731,35	97
der	Telemedizin	644,28	2533
Gesellschaft	Telemedizin	632,95	241
telemedizinisch	Zentrum	585,73	165

Tabelle 2: Top 10 der syntaxbasierten Kollokationen zu Lemmata mit *T/telemed* im kulturalanthropologischen Korpus

Wort 1	Relation	Wort 2	LLR	abs. Frequenz
GmbH	ist Apposition von	Telemedizin	2261,71	353
telemedizinisch	ist Attribut von	Anwendung	2236,42	456
bayerisch	ist Attribut von	TelemedAllianz	2002,48	204
Telemedizin	ist Genitivattribut von	Tag	1904,26	274
telemedizinisch	ist Attribut von	Betreuung	981,87	200
bayerisch	ist Attribut von	Telemedallianz	967,93	98
DGTelemed	ist Apposition von	V.	899,83	84
telemedizinisch	ist Attribut von	Zentrum	772,92	163
Telemedizin	ist Apposition von	Fachkongress	727,37	64
Telemedizin	ist Genitivattribut von	Möglichkeit	720,75	144

Die syntaktischen Informationen in Tabelle 2 machen den Zusammenhang zwischen den Bestandteilen der Eigennamen in der Relation der Apposition explizit und bieten damit mehr Informationen zur Einordnung dieser Datenpunkte. Die Annotationen ermöglichen außerdem, die Gesamtliste nach bestimmten syntaktischen Relationen zu filtern. Die genannten Appositionen beispielsweise können anhand des Relationslabels ausgeblendet werden. Auch hier gibt es sehr allgemeine Konzepte wie *telemedizinisch* als Attribut von *Anwendung*, die zwar frequent, aber nicht sehr informativ sind. *Telemedizin* als Genitivattribut von *Möglichkeit* weist daraufhin, dass eben deren Möglichkeiten noch Gegenstand des Diskurses sind. In der Durchsicht der Kollokationen jenseits der Top 10 finden sich verwandte Themen des Potentials und der Projekthaftigkeit (*Potential der Telemedizin*, *Telemedizin-typisches Potential*, *evaluiertes Telemedizinprojekt*, *vielversprechendes Telemedizinprojekt*), die anzeigen, dass sich die Umsetzung der Telemedizin in einer frühen Phase befindet und ihre Akzeptanz als Regelversorgung noch nicht abschließend verhandelt ist.

Auch zur Kollokation *telemedizinisch* als Attribut von *Betreuung* finden sich weiter unten im Ranking ähnliche Verwendungen zum Thema Betreuung (*telemedizinisch betreuen*, *telemedizinisch betreut* ...) und Unterstützung (*telemedizinisch unterstützt*, *telemedizinisch-unterstützte* (sic!) *Versorgung*, *Telematikunterstützung* ...). Dies weist auf die (bisher) eher ergänzende Rolle der Telemedizin im Verhältnis zur medizinischen Regelversorgung hin.

## Pflegewissenschaft

Tabelle 3 zeigt die zehn ersten oberflächenbasierten Kollokationen des Dialogkorpus zu Lemmata mit *E/entscheid*. Hier werden zunächst Probleme in den Daten deutlich: Mit *Finan/* liegt ein für gesprochene Sprache typischer Abbruch eines Wortes (vermutlich: *Finanzentscheidung*) vor. Zudem ist *Entscheidungsvariant* eine fehlerhafte Lemmaform zu *Entscheidungsvarianten*. Insgesamt sind die Frequenzen aufgrund der geringen Korpusgröße klein, lassen aber trotzdem hilfreiche Schlüsse für die Analyse zu. Im syntaxbasierten Gegenstück in Tabelle 4 sind zusätzliche Probleme erkennbar, die durch die automatische Verarbeitung gesprochener Sprache entstehen. Die Relation zwischen *hab* und *entscheiden* ist fehlerhaft als adverbial (korrekt: auxiliär) bezeichnet. Allerdings wird ein direkter Zusammenhang zwischen diesen Wörtern erst durch die syntaktischen Annotationen überhaupt erkennbar, da sie im Satz häufig nicht benachbart stehen. Zusätzlich gibt es vollständig falsche Analysen wie die Relation zwischen *entscheidend* und *Puh*.

Tabelle 3: Top 10 der oberflächenbasierten Kollokationen zu Lemmata mit *E/entscheid* im pflegewissenschaftlichen Korpus

Wort 1	Wort 2	LLR	abs. Frequenz
Entscheidung	treffen	33,41	7
richtig	Entscheidung	23,06	7
Tablettenform	entscheiden	21,28	4
der	Entscheidung	17,54	28
dieser	Entscheidung	13,02	7
selbst	entscheiden	12,88	4
Entscheidung	überlassen	10,87	2
Entscheidungsvarianten	nein	10,38	1
Finan/	Versorgungsentscheidung	10,38	1
entschieden	Abraten	10,38	1

Auch für die pflegewissenschaftliche Interpretation bieten die Kollokationen mit den höchsten LLR-Werten erste Anhaltspunkte, die dann durch eine Sichtung der weiteren Rangplätze ergänzt werden können. Die häufigsten Kollokatoren von Entscheidungen stehen für eine Realisierung eigener Entscheidungen der Betroffenen. Die Kollokation *richtig* macht Bewertungen der Entscheidungen sichtbar. Es zeigen sich zudem gegensätzliche Dimensionen des Phänomens, wie „selber entscheiden“ vs. „Entscheidung abgeben“, die hinter

der Kollokation mit *überlassen* stehen. Insbesondere die Subjekt- und Objektrelationen (*Entscheidung treffen*, *Entschluss entstehen*, *Entscheidung überlassen*) sind durch die syntaktische Analyse adäquater und theoretisch fundierter abgebildet. Dieser Nutzen wird jedoch durch Fehler in der automatischen Annotation eingeschränkt. Zudem werden diese Relationen in der gesprochenen Sprache mit kürzeren Sätzen möglicherweise auch durch die oberflächenbasierte Analyse besser erfasst als in anderen sprachlichen Registern. Insgesamt betrachtet werden durch den quantitativen Zugang Verwendungszusammenhänge des Phänomens „Entscheidung“ transparent, die wiederum auf wichtige Handlungskontexte in der Versorgungsrealität von schwerkranken und sterbenden Menschen verweisen.

Tabelle 4: Top 10 der syntaxbasierten Kollokationen zu Lemmata mit *E/ entscheid* im pflegewissenschaftlichen Korpus

Wort 1	Relation	Wort 2	LLR	abs. Frequenz
entscheiden	ist Adverbial von	hab	33,25	9
richtig	ist Attribut von	Entscheidung	24,54	6
Entscheidung	ist Akkusativobjekt von	treffen	23,47	4
Entschluss	ist Subjekt von	entstehen	21,93	2
selbst	ist Adverbial von	entscheiden	18,74	4
entscheidend	ist Adverbial von	Puh	16,42	1
Tablettenform	ist Subjekt von	entscheiden	16,13	2
Entscheidung	ist Akkusativobjekt von	überlassen	15,59	2
Entscheidung	ist Akkusativobjekt von	treff	13,52	2
für	ist Präposition zu	entscheiden	13,46	7

## Schlussfolgerungen

Es hat sich gezeigt, dass die Berechnung von Kollokationen auf der Grundlage der sprachlichen Oberfläche bzw. der Syntax für qualitative Fragestellungen

informativ sein kann. Für die Auswertung einer spezifischen Fragestellung ist die Assoziationsstärke allein jedoch nicht immer das entscheidende Kriterium. Die Kollokationen geben Hinweise auf Zusammenhänge innerhalb des Korpus, die neue Fragestellungen und Perspektiven generieren können. Gleichzeitig werden durch die Relationsannotationen bereits kleine Datenmengen in erweiterter Form auswertbar.

Die beispielhaften Analysen haben gezeigt, dass die syntaktischen Annotationen eine für die Interpretation hilfreiche Differenzierung bieten, indem präziser angegeben wird, in welcher Relation zwei Wörter stehen. Das ermöglicht auch das Filtern nach interessanten Relationstypen. Zudem werden durch den Einbezug der Syntax Relationen zwischen Wörtern in Distanzstellung sichtbar, was insbesondere vom Verb abhängige Satzteile besser sichtbar macht. Andererseits erfordern die syntaktischen Annotationen eine aufwendigere Vorverarbeitung, die mehr Zeit und technische Fähigkeiten erfordert. Außerdem stellen sie eine zusätzliche Fehlerquelle dar. Dies gilt besonders für die gesprochensprachlichen Daten. Eine systematische Überprüfung und Rückbindung an konkrete Korpusbelege ist deshalb wichtig und verbessert die Interpretationsmöglichkeiten aus qualitativer Sicht.

Anschließend an diese Arbeiten ist geplant, stärker Kontexte zu aggregieren: In grammatischer Hinsicht wird das durch Koreferenzannotationen erfolgen, die für das pflegewissenschaftliche Korpus bereits vorliegen. Auf semantischer Ebene verfolgen wir den Ansatz, Wortgruppen zu Konzepten zusammenzufassen, z. B. können *Ärztin*, *Arzt*, *Hausarzt*, *Onkologin* usw. auf ein gemeinsames Konzept *ÄRZT\*INNEN* abgebildet werden (vgl. den Ansatz von Wüest et al. 2011). Wir sehen außerdem weitere Anwendungsfälle über die Fächergrenzen hinweg, etwa zur literaturwissenschaftlichen Beschreibung von Geschlechterzuschreibungen, indem die sprachlichen Kontexte weiblicher und männlicher Vornamen verglichen werden.

## Fußnoten

1. <https://github.com/melandresen/DHd2020>

## Bibliographie

**Adelmann, Benedikt / Franken, Lina** (2020): Thematic web crawling and scraping as a way to form focussed web archives, in: *Engaging with Web Archives Conference Book of Abstracts*. To be published at <https://ewaconference.com/>.

**Andresen, Melanie** (2018): Sprachliche Variation in der Germanistik: eine n-Gramm-basierte Stilanalyse, in: *Book of Abstracts of DHd 2018*. Köln, Deutschland, 311–15.

**Andresen, Melanie / Zinsmeister, Heike** (2017): The benefit of syntactic vs. linear n-grams for linguistic description, in: *Proceedings of the 4th International Conference on Dependency Linguistics (Depling 2017)*. Pisa, Italy, 4–14 <http://aclweb.org/anthology/W17-6503>.

**van Atteveldt, Wouter / Kleinnijenhuis, Jan / Ruigrok, Nel** (2008): Parsing, Semantic Networks, and Political Authority Using Syntactic Analysis to Extract Semantic Relations from Dutch Newspaper Articles, in: *Political Analysis*, 16(4): 428–46 doi:10.1093/pan/mpn006.

**Blessing, Andre / Sonntag, Jonathan / Kliche, Fritz / Heid, Ulrich / Kuhn, Jonas / Stede, Manfred** (2013): Towards a Tool for Interactive Concept Building for Large Scale Analysis in the Humanities, in: *Proceedings of the 7th Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*. 55–64 <https://www.aclweb.org/anthology/W13-2708>.

**Bohnet, Bernd** (2010): Very High Accuracy and Fast Dependency Parsing is not a Contradiction, in: *Proceedings of the 23rd International Conference on Computational Linguistics (COLING 2010)*. Beijing, China, 89–97. <https://www.aclweb.org/anthology/C10-1011>

**Dunning, Ted** (1993): Accurate Methods for the Statistics of Surprise and Coincidence, in: *Computational Linguistics*, 19(1): 61–74. <https://www.aclweb.org/anthology/J93-1003>

**Evert, Stefan** (2008): Corpora and collocations, in: Lüdeling, Anke / Kytö, Merja (Hg.), *Corpus linguistics: an International Handbook*, Vol. 2. (Handbücher zur Sprach- und Kommunikationswissenschaft 29). Berlin, Boston: De Gruyter, 1212–1248.

**Foth, Kilian A. / Köhn, Arne / Beuck, Niels / Menzel, Wolfgang** (2014): Because Size Does Matter: The Hamburg Dependency Treebank, in: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, 2326–2333. Reykjavik, Iceland. [http://www.lrec-conf.org/proceedings/lrec2014/pdf/860\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2014/pdf/860_Paper.pdf).

**Gaidys, Uta / Gius, Evelyn / Jarchow, Margarete / Koch, Gertraud / Menzel, Wolfgang / Orth, Dominik / Zinsmeister, Heike** (2017): hermA: Automated modelling of hermeneutic processes, in: *Hamburger Journal für Kulturanthropologie*(7): 119–23.

**Geyken, Alexander** (2011): Statistische Wortprofile zur schnellen Analyse der Syntagmatik in Textkorpora, in: Abel, Andrea / Zanin, Renata (Hg.), *Korpora in Lehre und Forschung*. Bozen: Bolzano Univ. Press, 129–54.

**Googasian, Victoria / Heuser, Ryan J.** (2019): Digital Animal Studies: Modeling Anthropomorphism in Animal Writing, 1870–1930, in: *Book of Abstracts of DH 2019* <https://dev.clariah.nl/files/dh2019/boa/0458.html>.

**Kilgariff, Adam / Baisa, Vít / Bušta, Jan / Jakubíček, Miloš / Kovář, Vojtěch / Michelfeit, Jan, Rychlý / Pavel. / Suchomel, Vít** (2014): The Sketch Engine: ten years on, in: *Lexicography*, 1(1): 7–36 doi:10.1007/s40607-014-0009-9.

**Kilgariff, Adam / Rychlý, Pavel / Smrz, Pavel / Tugwell, David** (2004): The Sketch Engine: in:

*Proceedings of the 11th EURALEX International Congress*. 105–15 <https://euralex.org/publications/the-sketech-engine/>.

**Kleinnijenhuis, Jan / van Atteveldt, Wouter** (2014): Positions of Parties and Political Cleavages between Parties in Texts, in: Kaal, Bertie / Maks, Isa / van Elfrinkhof, Annemarie (Hg.), *Discourse Approaches to Politics, Society and Culture*, Vol. 55. Amsterdam: Benjamins, 1–20 doi:10.1075/dapsac.55.01kle.

**Koch, Gertraud / Franken, Lina** (im Druck): Automatisierungspotenziale in der qualitativen Diskursanalyse. Das Prinzip des Filterns, in: Schilling, Samuel / Klimczak, Peter (Hg.): *Die Gesellschaft im Spiegellabyrinth sozialer Medien*. Wiesbaden.

**Sinclair, Stéfan / Rockwell, Geoffrey** (2016): Voyant Tools. Web. <http://voyant-tools.org/>.

**Wüest, Bruno / Clematide, Simon / Bünzli, Alexandra / Laupper, Daniel** (2011): Semi-Automatic Core Sentence Analysis: Improving Content Analysis for Electoral Campaign Research, in: *International Relations Online Working Paper*(1). [https://www.sowi.uni-stuttgart.de/dokumente/forschung/irowp/IROWP\\_Series\\_2011\\_1\\_Wueest\\_Clematide\\_Buenzli\\_Laupper\\_Content](https://www.sowi.uni-stuttgart.de/dokumente/forschung/irowp/IROWP_Series_2011_1_Wueest_Clematide_Buenzli_Laupper_Content)