

Quantitative Textanalyse frühneuzeitlicher Reiseberichte

1. Einleitung

Die Leitfrage dieses Projekts war, ob und wie sich eine vergleichende quantitative Textanalyse gewinnbringend für die Untersuchung eines historischen Textkorpus einsetzen lässt. Wir wollten herausfinden, ob sich die Texte im Hinblick auf Autorschaft, Entstehungszeit und Entstehungskontext in ihrer Wortwahl und Ausdrucksweise signifikant unterscheiden. In Bezug auf das Genre der frühneuzeitlichen Reiseberichte stand hier insbesondere die Frage im Zentrum, wie die Stationen der beschriebenen Reisen – die von den Autoren besuchten Städte und Ortschaften – sowie deren Zweck – sind die Autoren Forschungsreisende, politische Gesandte oder Abenteurer? – die Wortwahl beeinflussten. Darauf aufbauend stellten wir uns in methodologischer Hinsicht die Frage, was die Ergebnisse einer solchen Untersuchung für eine weitergehende qualitative Analyse eines historischen Textkorpus bedeuten könnten. Lassen sich mit den angewandten Methoden bereits sinnvolle Aussagen über den Inhalt der einzelnen Texte treffen? Oder dient eine solche «Vogelschau» mittels *Distant Reading* eher dem Zweck, gewisse Anhaltspunkte in den Texten zu finden, auf Grundlage derer eine Vertiefung stattfinden könnte? Zur Beantwortung dieser Fragen stellten wir in einem ersten Schritt ein Textkorpus zusammen, welches wir dann für einen maschinellen Vergleich aufbereiteten. Im zweiten Schritt folgte die Untersuchung des Korpus mittels verschiedener Methoden der quantitativen Textanalyse.

2. Textgrundlage

Als Ausgangskorpus wählten wir die Reiseberichte des Augsburger Kaufherrn Philipp Hainhofer (1578 – 1647). Diese sind in handschriftlicher Form überliefert und werden im Rahmen eines digitalen Editionsprojektes der *Herzog August Bibliothek* digitalisiert und transkribiert. Die Transkripte werden auf der Webseite des Projekts als TXT-

Dateien frei zur Verfügung gestellt.¹ Philipp Hainhofer war Kunstsammler, Kunsthändler und Agent und arbeitete in dieser Rolle für verschiedene Herzogen. Er unternahm zahlreiche Reisen, zumeist in diplomatischem Auftrag, die ihn hauptsächlich nach Süddeutschland, aber auch nach Österreich führten.² Für die vorliegende Arbeit wurden elf Einzeltexte analysiert, die während seiner Reisen zwischen den Jahren 1603 – 1636 entstanden sind. Im Konkreten handelt es sich um fünf Reiseberichte aus München, zwei aus Neuburg jeweils einen aus Nürnberg und Regensburg und einen, der von Wildbad über Heidelberg nach Durlach führt. Die Sammlung enthält noch weitere Texte, die bis ins Jahr 1594 zurückgehen. Diese werden noch bearbeitet und schrittweise der Öffentlichkeit zugänglich gemacht, standen zum Zeitpunkt der Projektarbeit jedoch noch nicht zur Verfügung.

Neben einem Vergleich der einzelnen Texte von Hainhofer sollte ein Vergleich mit weiteren Texten erfolgen, um die quantitativen sowie qualitativen Ergebnisse der angewandten Methoden besser einordnen zu können. Daher wurden im weiteren Prozess die Reiseberichte von drei zusätzlichen Autoren ausgewählt und die Texte Hainhofers gesammelt als vierter Vergleichskorpus betrachtet. Es handelt sich dabei um die Reiseberichte von Georg Christoph von Neitzschitz (ca. 1600 – 1636), der von 1630 bis 1636 als Forschungsreisender den Orient bereist hatte, und kurz nach seiner Rückkehr verstarb³, dem Zürcher Kunsttheoretiker und Pädagogen Johann Georg Sulzer (1720 – 1779)⁴, und dem Magdeburger Schriftsteller Friedrich Schulz (1762 – 1798).⁵ Im

¹ Philipp Hainhofer: Reiseberichte und Sammlungsbeschreibungen 1594–1636. Edition und Datensammlung zur Kunst- und Kulturgeschichte der ersten Hälfte des 17. Jahrhunderts [Wolfenbütteler Digitale Editionen, Nr. 4], hrsg. und eingeleitet von Michael Wenzel, Transkription und Kommentar von Ursula Timann und Michael Wenzel, Wolfenbüttel: Herzog August Bibliothek 2020ff. Die Links zum Download finden sich unter: <https://hainhofer.hab.de/informationen-zur-edition/downloads>, zuletzt abgerufen am 13.09.2024.

² Vgl. dazu: <https://hainhofer.hab.de/informationen-zur-edition/ueber-philipp-hainhofer>, zuletzt abgerufen am 12.09.2024.

³ Neitzschitz, Georg Christoph von: Sieben-Jährige und gefährliche WeltBeschauung Durch die vornehmsten Drey Theil der Welt Europa/ Asia und Africa. Bautzen, 1666. In: Deutsches Textarchiv <https://www.deutschestextarchiv.de/neitschitz_reise_1666>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-20878-1.

⁴ Sulzer, Johann Georg: Tagebuch einer von Berlin nach den mittäglichen Ländern von Europa in den Jahren 1775 und 1776 gethanen Reise und Rückreise. Leipzig, 1780. In: Deutsches Textarchiv <https://www.deutschestextarchiv.de/sulzer_reise_1780>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-25231-5.

⁵ Schulz, Friedrich: Reise eines Liefländers. 3 Bde. Berlin, 1795.
Bd. 1, H. 1: <https://www.deutschestextarchiv.de/schulz_reise0101_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-200905197358.
Bd. 1, H. 2: <https://www.deutschestextarchiv.de/schulz_reise0102_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-17297-4.

Gegensatz zu den Handschriften Hainhofers handelt es sich hier um Werke, die im Druck erschienen waren. Sie liegen im *Deutschen Textarchiv (DTA)* in digitalisierter Form sowie als TXT-Dateien vor. Neben dem Genre und der digitalen Verfügbarkeit dienten als Auswahlkriterien erstens die Länge der Texte, da uns ein vergleichbarer Datenumfang angesichts einer quantitativen Datenanalyse wichtig erschien, um Verzerrungen in den Ergebnissen weitestgehend zu vermeiden. Zweitens sollte eine zeitliche Distanz des Entstehungszeitraumes einen historischen Vergleich ermöglichen. Mit den Texten / Textsammlungen von Hainhofer und Neitzschitz, sowie von Sulzer und Schulz bestand das Korpus schliesslich aus vier Quellen, von denen jeweils zwei bezüglich ihrer Entstehungszeit relativ nahe beieinander liegen.

3. Vorgehen

Die Analyse wurde mit der Programmiersprache *Python* programmiert. Zur Ausführung der Einzelnen Analyseschritte und gleichzeitigen Präsentation der Ergebnisse haben wir mit einem *Jupyter-Notebook* gearbeitet, das neben diesem Arbeitsbericht ebenfalls im Github-Repository zur Verfügung steht.⁶ Im Notebook ist der gesamte verwendete Code inklusive Outputs, Visualisierungen, ergänzenden Kommentaren sowie zusammenfassenden Beschreibungen der Ergebnisse übersichtlich dargestellt. Daher werden im Folgenden nur ausgewählte Ergebnisse exemplarisch aufgeführt und für eine detaillierte Betrachtung auf das Notebook bzw. die daraus generierte PDF-Datei verwiesen⁷. Neben der Aufbereitung der Texte für maschinelle Vergleichsmöglichkeiten werden die verwendeten Methoden sowie durchgeführten Analysen beschrieben, die sich wie folgt aufbauen: Vereinheitlichung bzw. Normalisierung der Korpora, Vergleich absoluter Worthäufigkeiten vs. statistischer Gewichtung durch Tf-idf, die Notwendigkeit von Stoppwörtern, weiterführende lexikalische Analysen mittels *Natural Language*

Bd. 2, H. 3: <https://www.deutschestextarchiv.de/schulz_reise0201_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-200905197368.

Bd. 2, H. 4: <https://www.deutschestextarchiv.de/schulz_reise0202_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-17298-9.

Bd. 3, H. 4 u. 5: <https://www.deutschestextarchiv.de/schulz_reise03_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-200905197374.

⁶ https://github.com/heinen-st/python_travelogues.

⁷ Das entsprechende PDF ist auf GitHub zu finden unter: Heinen, Stefan: grand_tour_corpus_analysis In: python_travelogues, https://github.com/heinen-st/python_travelogues/tree/main/documentation, zuletzt abgerufen am 15.09.2024. Für eine bessere Lesbarkeit wird im Folgenden auf dieses Dokument als *Notebook* verwiesen.

Processing (NLP) wie Kategorisierung von Wortarten und Ähnlichkeitsmessung verwendeter Lemmata.

Als erste Analysemethode wählten wir die Kalkulation des Tf-idf-Masses für die einzelnen Wörter im Korpus. Tf-idf steht für «*term frequency* („Vorkommenshäufigkeit“) – *inverse document frequency* („inverse Dokumenthäufigkeit“)». Mit dieser Methode werden Wortvorkommnisse innerhalb eines Textkorpus statistisch gewichtet. So kann die relative Häufigkeit von Wörtern in Dokumenten errechnet werden. Als Ergebnis einer Tf-idf-Analyse erhält man einen Wert pro Wort, der angibt, wie spezifisch dieses Wort für ein Dokument ist, verglichen mit anderen Dokumenten im Korpus. Kommt also ein Wort in einem Dokument sehr häufig vor, im ganzen Korpus jedoch nur sehr selten, hat es einen hohen Tf-idf-Wert im jeweiligen Dokument. Kommt ein Wort in einem Dokument zwar sehr häufig vor, aber auch in allen anderen Dokumenten des Korpus, dann hat es einen tieferen Tf-idf-Wert. Somit kann herausgefunden werden, welche Wörter charakteristisch für ein Dokument sind, während häufig vorkommende, nicht-bedeutungstragende Wörter, wie Funktionswörter oder Partikel, nicht ins Gewicht fallen. Zur Vorbereitung haben wir als Erstes den Blogbeitrag zu Tf-idf von *Programming Historian* gelesen und uns mit dem dort zur Verfügung gestellten Code vertraut gemacht.⁸ Nachdem wir *Jupyter Notebook* und *Python* installiert hatten, haben wir ein Programm erstellt, mit dem alle Quellen einheitlich in ein TXT-Format gebracht und in einer CSV-Datei zusammengeführt wurden. Auf diese Weise sind Analysen mit Python-Modulen leicht durchführbar. Zunächst haben wir uns auf die Hainhofer-Texte konzentriert und die Tf-idf-Werte für jedes Wort von allen elf Texten des Hainhofer-Korpus errechnet. Als Ergebnis erhielten wir eine Liste/CSV-Tabelle je Dokument, welche die Wörter nach Tf-idf-Wert geordnet aufgelistet hat. Da den elf Texten elf unterschiedlichen Reisen entsprechen, die sich hinsichtlich des Zeitpunkts, des Anlasses und z.T. der Destinationen unterscheiden, wurde erwartet, dass sich diese Differenzen in unterschiedlichen Tf-idf-Werten ausdrücken würden. In der Tat wurden für alle Reiseberichte verschiedene Wörter als besonders spezifisch für den jeweiligen Text identifiziert, wodurch sich ein eigenes Wörter- bzw. Tf-idf-Profil pro Dokument darstellte. Es gab jedoch auch einzelne Wörter, welche nicht aussagekräftig waren wie «euer» oder «bei». Bei der Tf-idf-Methode wird die Relevanz von Wörtern so gewichtet, dass häufig auftretende, wenig bedeutungstragende Wörter einen niedrigeren Wert erhalten

⁸ <http://programminghistorian.org/en/lessons/analyzing-documents-with-tfidf>, zuletzt abgerufen am 12.09.2024.

und somit als wenig relevant eingestuft werden. Da nun dennoch solche bedeutungsschwachen Wörter z.T. hohe Werte erhalten haben, könnten diese mit manuell definierten Stoppwort -Listen herausgefiltert werden, damit sie in der Berechnung nicht berücksichtigt werden. Dies wurde für die weiterführende Vergleichsanalyse der vier verschiedenen Autoren berücksichtigt und angewendet.

Wir haben uns nun etwas genauer mit den Quellen selbst und deren Inhalt auseinandergesetzt, also das *Distant Reading* durch ein *Close Reading* ergänzt. Hierfür wurden für jedes Dokument die 20 Begriffe mit den höchsten Tf-idf-Werten extrahiert und mit der online zur Verfügung stehenden Kurzzusammenfassung des jeweiligen Reiseberichtes von Hainhofer verglichen. Wir wollten herausfinden, wie repräsentativ die errechneten Tf-idf-Werte und die ausgegebene Spezifität der Begriffe für jedes Dokument sind. Dabei konnte festgestellt werden, dass die ermittelten Begriffe durchaus als semantisch richtungsweisend anzusehen sind. Im Reisebericht von München 1613 beispielsweise geht Hainhofer nach München, um eine Hochzeit zu besuchen. Unter den ersten sieben Wörtern befinden sich «ehestand, braut, bräutigam, breutigam, hochzeiterin» Auch wenn die Liste die Zusammenfassung nicht ersetzt, wird das Thema also in der Tf-idf-Liste semantisch gut fassbar. Ein solches Profil ist bei allen Wortlisten erkennbar, auch wenn die Qualität unterschiedlich zu sein scheint. Zum Teil tauchen Wörter auf, die wenig bedeutungstragend sind, was die Aussagekraft beeinträchtigt.

In einem nächsten Schritt haben wir die drei weiteren Reiseberichte hinzugefügt, um Vergleichswerte ermitteln zu können. Da sich bei den Texten aus dem Deutschen Textarchiv die Repräsentation der Druckschrift in den publizierten TXT-Dateien wesentlich von derjenigen der Handschriften Hainhofers unterschied, mussten diese zuerst angeglichen werden, um einen sinnvollen maschinellen Vergleich zu ermöglichen (vgl. Notebook, 2.1). Weiter wurde das Textmaterial von Hainhofer und Sulzer, das aus mehreren Dateien bestand, jeweils zu einer TXT-Datei zusammengefasst. Mithilfe des *Count-Vectorizers* vom *Natural Language Toolkit (NLTK)* wurden Stopp-Wort-Listen errechnet, um bedeutungsschwache Wörter vor den weiterführenden Analysen gezielt zu entfernen. Zunächst wurden rein quantitativ die häufigsten Wörter für jedes Korpus ermittelt (vgl. Notebook, 2.3) und in Balkendiagrammen visualisiert, bevor anschließend die statistisch gewichtete Tf-idf-Analyse durchgeführt wurde (vgl. Notebook, 2.4). Beim Vergleich der beiden Ergebnisse fällt auf, dass es durchaus Ähnlichkeiten in den

Wortlisten der beiden Methoden gibt, dass aber die Tf-idf-Analyse spezifischere Wortlisten hervorbringt, bei denen reisespezifische Aspekte hervortreten (Erwähnung von Ortschaften oder «fremder Völker», wie z.B. «turcken» oder «mohren» bei Neitzschitz). Es wird deutlich, dass die Listen für die neueren Texte (Schulz 1795 und Sulzer 1780) bedeutungsstärkere Wörter beinhalten als für die älteren Texte (Hainhofer 1615 und Neitzschitz 1666). Hier werden die Grenzen der Analyse-Methoden wie Tf-idf deutlich, die für historische Texte aus verschiedenen Epochen nicht auf die gleiche Weise geeignet sind.

Für weiterführende semantisch-lexikalische Untersuchungen wurden Methoden des *Natural Language Processing* (NLP) angewendet. Hierfür nutzten wir das Python-Paket *spaCy* mit dem deutschen Sprachmodell *de_core_news_lg*, um die Wörter der Korpora zu lemmatisieren (vgl. Notebook 3.1). Dies bedeutet, dass die Wörter auf ihre Grundform reduziert und semantisch als Einheit betrachtet werden, unabhängig von ihren Flexionsformen (z.B. (ich) *sehe* und (er) *sieht* entsprechen beide dem Lemma *sehen*). Durch anschliessendes *Part-of-speech tagging* (POST) wurden die Lemmas nach Wortarten kategorisiert und visualisiert (vgl. Notebook 3.2). Wir konzentrierten uns auf Nomen und Adjektive, da sie uns zur Beschreibung von Reisedestinationen als besonders relevant erschienen. Die Analyse mit *spaCy* bestätigte die vorherigen Beobachtungen, dass in allen vier Korpora ähnliches Vokabular verwendet wurde. Vereinzelt lassen sich Schwerpunkte der Texte erahnen. Insbesondere bei Hainhofer wird der politische Aspekt seiner Reise deutlich, da Begriffe wie «fürstlich» und «kayserlich» häufig vorkommen. Schliesslich wurden mit *spaCy* Ähnlichkeitswerte des verwendeten Vokabulars der vier Korpora berechnet und in Balkendiagrammen miteinander verglichen (vgl. Notebook 3.3). Hier zeigte sich sowohl für wortartenspezifische Vergleiche (Nomen und Adjektive) als auch bei der Betrachtung der Gesamtkorpora, dass die beiden neueren Texte (Schulz 1795 und Sulzer 1780) bei den Adjektiven und in der Gesamtbetrachtung die höchsten Ähnlichkeitswerte aufweisen, während die zeitlich weit auseinanderliegenden Korpora (Hainhofer 1615 und Neitzschitz 1666 gegenüber Schulz und Sulzer) deutlich geringere Ähnlichkeitswerte hervorbringen. Dies kann als Hinweis gedeutet werden, dass epochenabhängige Unterschiede im Schreibstil und verwendeten Vokabular identifiziert werden. Gleichzeitig zeigen auch Hainhofer und Neitzschitz trotz zeitlicher Nähe eher geringe Ähnlichkeitswerte. Eine Ursache könnten die unterschiedlichen Reisemotive und Reiseziele sein. Hier wird deutlich,

dass eine rein quantitative Betrachtungsweise wertvolle Interpretationsansätze bieten kann, aber weiterführende Analysen weiterhin erforderlich bleiben.

4. Schwierigkeiten

Während der Projektarbeit haben sich uns verschiedene Herausforderungen ergeben, welche aber größtenteils lösbar waren. Der erste Themenkomplex war die Aufarbeitung der Texte. Während die Normalisierung gut gelang, scheint die Lemmatisierung trotz mehrerer Überarbeitungen nicht vollumfänglich zu greifen. Insbesondere bei Hainhofer werden dieselben Lemmata mit unterschiedlicher Flexionsform als Einzelwörter behandelt (vgl. Notebook, 3.2, Säulendiagramm). Ein anderes Problem dieser Kategorie war, dass auch viele irrelevante Wörter wie «einem» oder «bej» im ersten Teil der Tf-idf-Listen vorkamen. Um diesem Problem entgegenzuwirken, haben wir Stopppwortlisten erstellt, die je nach Fragestellung noch weiter angepasst werden können. Besonders deutlich wurde die unzureichende Spezifität von den zur Verfügung stehenden Modellen bezüglich historischer Quellen. Das Sprachmodell von spaCy wurde an modernem Deutsch trainiert (mit Schwerpunkt auf Nachrichten), weshalb es für die historischen Texte nur bedingt geeignet ist, da es verwendete Ausdrücke und Schreibweisen häufig nicht korrekt erfassen kann.

5. Fazit

Es stehen zahlreiche Methoden zur Verfügung, die auf verhältnismässig einfache Weise quantitative Textanalysen in grossem Umfang erlauben. Die erforderlichen *Python*-Kompetenzen sind jedoch nicht zu unterschätzen, da korpuspezifische Anpassungen vorgenommen werden müssen, um tatsächlich saubere Daten und darauf aufbauende Ergebnisse zu erhalten. Die in der vorliegenden Arbeit angewendeten Methoden wie Tf-idf oder Lemmatisierung ermöglichten es, verschiedene Korpora hinsichtlich ihrer zentralen Themen zu überblicken und erste Differenzierungen vorzunehmen. Dieses *Distant Reading* bietet wertvolle Einblicke und erste aussagekräftige Visualisierungen, die als Ausgangspunkt für weiterführende Hypothesen und Analysen dienen können. Potenzielle Verzerrungen bzw. mangelnde Spezifität der verwendeten Modelle muss stets mitgedacht werden, bevor Schlussfolgerungen gezogen werden. Um tatsächlich zu gehaltvollen Resultaten zu kommen ist eine Kontextualisierung und

Beachtung epochenspezifischer Eigenheiten von Textquellen unerlässlich. Ausserdem scheint es naheliegend, nicht nur einzelne Wörter und Lemmata zu betrachten, sondern «tiefer» in die Textstruktur einzutauchen. Dies ist einerseits ergänzend durch klassisches *Close Reading* möglich. Andererseits gibt es weitere automatisierte Analysemöglichkeiten zur Erfassung grösserer Spracheinheiten. Hier könnte im Weiteren der Kontext von bestimmten Zielwörtern betrachtet werden, um beispielsweise die Einbettung bzw. Beschreibung von Reisedestinationen qualitativ besser zu erfassen (z.B. durch N-Gramme oder *Noun chunks*).

6. Bibliografie

5.1 Quellen

Hainhofer, Philipp: Reiseberichte und Sammlungsbeschreibungen 1594–1636. Edition und Datensammlung zur Kunst- und Kulturgeschichte der ersten Hälfte des 17. Jahrhunderts [Wolfenbütteler Digitale Editionen, Nr. 4], hrsg. und eingeleitet von Michael Wenzel, Transkription und Kommentar von Ursula Timann und Michael Wenzel, Wolfenbüttel: Herzog August Bibliothek 2020ff.

Neitzschitz, Georg Christoph von: Sieben-Jährige und gefährliche WeltBeschauung Durch die vornehmsten Drey Theil der Welt Europa/ Asia und Africa. Bautzen, 1666. In: Deutsches Textarchiv <https://www.deutschestextarchiv.de/neitschitz_reise_1666>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-20878-1.

Schulz, Friedrich: Reise eines Liefländers. 3 Bde. Berlin, 1795.

Bd. 1, H. 1: <https://www.deutschestextarchiv.de/schulz_reise0101_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-200905197358.

Bd. 1, H. 2: <https://www.deutschestextarchiv.de/schulz_reise0102_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-17297-4.

Bd. 2, H. 3: <https://www.deutschestextarchiv.de/schulz_reise0201_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-200905197368.

Bd. 2, H. 4: <https://www.deutschestextarchiv.de/schulz_reise0202_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-17298-9.

Bd. 3, H. 4 u. 5: <https://www.deutschestextarchiv.de/schulz_reise03_1795>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-200905197374.

Sulzer, Johann Georg: Tagebuch einer von Berlin nach den mittäglichen Ländern von Europa in den Jahren 1775 und 1776 gethanen Reise und Rückreise. Leipzig, 1780. In: Deutsches Textarchiv <https://www.deutschestextarchiv.de/sulzer_reise_1780>, abgerufen am 19.07.2024, URN: urn:nbn:de:kobv:b4-25231-5.

5.2 Tutorials

Matthew J. Lavin: "Analyzing Documents with TF-IDF," *Programming Historian* 8 (2019), <https://doi.org/10.46430/phen0082>.

Die zahlreichen Beiträge auf *Stack Overflow* (<https://stackoverflow.com/>) die massgeblich zum Gelingen dieses Projektes beigetragen haben, können an dieser Stelle leider nicht einzeln gelistet werden. Unser Dank geht an die Python-Community!