Introduction to Digital History

https://wissen-ist-acht.github.io/digital history.intro/

Ina Serif

Last compiled on 02.05.2024

Table of contents

ı	W	elcome	4							
1	Was ist Digital History?									
2	Fors	schung und Lehre	8							
	2.1	Digitalisierte Quellen, digitale Quellen	9							
	2.2	Digitale Tools zur Analyse	10							
	2.3	Digitale Tools zur Kommunikation	12							
	2.4	Projekte und Ressourcen für Forschung und Lehre	13							
		2.4.1 Alte Geschichte	13							
		2.4.2 Mittelalter und Frühe Neuzeit	14							
		2.4.3 Moderne und Zeitgeschichte	15							
		2.4.4 Jüdische Geschichte	15							
		2.4.5 Geschichte Afrikas	16							
		2.4.6 Osteuropäische Geschichte	16							
		2.4.7 Epochen-/Areaübergreifend:	17							
3	Digital Literacy, Digital Criticism									
	3.1	Digital Literacy, Data Literacy	18							
	3.2	Digital Criticism, Data Criticism	19							
	3.3	Künstliche Intelligenz und Algorithmenkritik	22							
4	Datenerhebung, -aufbereitung und -analyse									
	4.1	Datenerhebung	26							
	4.2	Datenaufbereitung	30							
	4.3	Datenanalyse	32							
	4.4	Datensicherung	33							
5	FAIR, CARE und LOUD									
	5.1	FAIRe Daten	34							
	5.2	CARE-Prinzipien	35							
	5.3	TO(II)D	35							

П	Praxis	37							
6	Briefedition 'Der Sturm'	39							
7	Durch den Vordereingang								
8	Durch die Hintertür 8.1 Shell 101	43 44 45 49							
9	Ausblick	55							
Α _Ι	ppendices	60							
Α	Glossar	60							
В	Literatur, Tools, Tutorials B.1 Was ist digital history? B.2 Einführungen und Guides B.3 Digital history tools (free/open source) B.3.1 Allgemein B.3.2 Datenanalyse B.3.3 Karten B.3.4 Netzwerkanalyse B.3.5 Text-/Korpusanalyse B.3.6 Storytelling B.3.7 Visualisierung	62 62 63 63 63 64 64 64 64							
	B.4 Digital Literacy, Digital Criticism	65 65 66 66							

Part I Welcome

Der vorliegende Guide, erstellt im Herbstsemester 2022, begleitet die Einführungskurse im Fach Geschichte an der Universität Basel und soll einen ersten Einblick in den Bereich Digital History geben. Es ist ein *living document*, das regelmäßig aktualisiert wird und dabei auch die epochen- und areaspezifischen Inhalte der verschiedenen Einführungskurse berücksichtigt, indem Verweise auf verschiedene Digital-History-Projekte aus unterschiedlichen Bereichen mit der Zeit in den Guide einfließen. Für die Teilnehmer:innen der Einführungskurse wird der Guide von einer Präsenzsitzung begleitet, bietet aber hoffentlich auch unabhängig davon einen Mehrwert. Für Kommentare, Anregungen oder Beschwerden freue ich mich über eine Nachricht.

Der Guide ist in zwei Teile gegliedert: Die Kapitel 1–5 sollen eine erste Übersicht über Digital History bieten und den Blick auf Neuerungen und Veränderungen richten, die sich in den Geschichtswissenschaften aus der Nutzung digitaler Methoden ergeben. Der anschließende praktische Teil zeigt an einem konkreten Beispiel die Anwendung verschiedener Techniken auf, die sich (nicht nur) für Historiker:innen bei der Arbeit mit Quellenmaterial anbieten. Der Praxisteil verfolgt dabei zwei Ziele: Zum einen sollen Hemmungen bei der Arbeit mit dem Computer, die über die Nutzung als elektronische Schreibmaschine hinausgeht, abgebaut werden. Zum anderen soll ein grundlegendes Verständnis dafür hergestellt werden, welche Möglichkeiten computergestützte Analysen bieten und wie diese in der historischen Arbeit eingesetzt werden können.

Die Übersicht soll möglichst knapp gehalten werden – es gibt zahlreiche ausführliche Grundlagenwerke, weswegen viele Themen nur kurz angeschnitten, dafür aber mit weiterführenden Verweisen versehen werden. Dasselbe gilt für den Praxisteil: Weiterreichende Anleitungen, Tutorials oder Onlinekurse werden an entsprechender Stelle verlinkt. Vollständigkeit wird an keiner Stelle beansprucht; Hinweise auf weitere Online-Angebote nehme ich gerne auf.

1 Was ist Digital History?

Über die Antwort zur Frage, was Digital History ist oder umfasst, kann man ausgiebig diskutieren. Als Teilgebiet der Digital Humanities, der digitalen Geisteswissenschaften, kann die folgende aktuelle und pragmatische Definition von Blaney et al. (2021) hilfreich sein:

Digital humanities, in our view, is a question of approach: if you are actively and critically using digital tools to aid your work in researching, teaching or learning, you are probably doing digital humanities. We would encourage anyone to learn to program if they are interested in doing so, but we do not see it as a defining characteristic of work in digital humanities. [1, S. 6]

Dabei umfassen "digital tools" eine große Bandbreite – und es wird sich kaum eine: finden, der: die Studium, Forschung oder Lehre völlig ohne die Nutzung digitaler Techniken betreibt. Wir sind alle Historiker:innen im digitalen Zeitalter, und als solche müssen wir ohnehin neue Kompetenzen entwickeln. Wir nutzen ganz selbstverständlich Textverarbeitungsprogramme für das Abfassen von Seminararbeiten oder Aufsätzen, das Recherchieren in Suchmaschinen und Online-Bibliothekskatalogen oder -Lexika gehört zur täglichen Arbeit als Historiker:in. Wir können uns aber zudem dafür entscheiden, für ein Forschungsprojekt Methoden und Techniken einzusetzen, die über die traditionellen Werkzeuge der Geschichtswissenschaften hinausgehen – Analyse und Interpretation von Quellen durch deren genaue Lektüre, sogenanntes close reading -, und uns durch den Computer unterstützen lassen. Ob wir hierbei auf vorhandene Software zurückgreifen oder selbst Programme schreiben, um uns nicht nur als Historiker:innen im digitalen Zeitalter, sondern auch als digitale Historiker:innen zu verstehen, mögen manche als Glaubensfrage auffassen; eine inkludierende Haltung zu dieser Frage scheint mir dabei nur Vorteile zu haben. Entgegen einer häufig zitierten Aussage von Emmanuel Le Roy Ladurie (*1929), der Historiker von morgen werde Programmierer sein, oder er werde nicht sein: "L'historien de demain sera programmeur ou il ne sera pas." 5]

Für eine erste Idee dafür, wie man historische Fragestellungen mithilfe digitaler Methoden beantworten kann und wie unterschiedlich digital unterstützte Forschungsprojekte aussehen können, bietet sich unter anderem der Übersichtsartikel "State of the Field: Digital History" von Romein et al. (2020) an.[9] Eine anwachsende Liste an Beispielprojekten aus unterschiedlichen Epochen bzw. Themenbereichen findet sich unter Projekte und Ressourcen in Kapitel 2.

Um eine Annäherung an die aktive, kritische und reflektierte Nutzung digitaler Methoden in Forschung und Lehre mit einem Fokus auf deren Anwendung in den Geschichtswissenschaften

geht es im vorliegenden Guide. Weiterführende Texte zur Frage, was Digital History ist bzw. umfasst, finden sich unter Literatur, Tools, Tutorials

2 Forschung und Lehre

Die fortschreitende Digitalisierung in ganz unterschiedlichen Lebensbereichen zieht Veränderungen und Entwicklungen auch für die historische Arbeit nach sich, und dies auf mehreren Ebenen: in Bezug auf die Arbeit bzw. den Umgang mit Quellen, hinsichtlich des Einsatzes digitaler Methoden nicht nur zur Analyse von Forschungsergebnissen, sondern auch für deren Kommunikation, und schließlich für die Hochschullehre.



Figure 2.1: Randall Munroe, History Department, xkcd.com (17.12.2018).

2.1 Digitalisierte Quellen, digitale Quellen

Als Historiker:innen steht die Arbeit mit Quellen im Mittelpunkt unserer Analysen. Das bedeutet gleichzeitig, dass der Zugang bzw. die Verfügbarkeit von Dokumenten einen Einfluss darauf hat, welche Fragen wir beantworten oder welche Analysen wir vornehmen können. Zugangsbeschränkungen, die die Größe und Zusammensetzung unseres Untersuchungskorpus beeinflussen, können dabei von Gedächtnisinstitutionen – also Museen, Archiven, Bibliotheken – ausgehen, beispielsweise wenn bei zeitgenössischen Akten eine Schutzfrist festgesetzt wird oder wenn ein Objekt zu fragil für die Benutzung ist. Auch kann es aus finanziellen

und/oder organisatorischen Gründen schwierig sein, bestimmte Archive an weiter entfernten Orten aufzusuchen, um weitere Dokumente für die Untersuchung zu berücksichtigen. Groß angelegte Digitalisierungsprojekte in Bibliotheken und Archiven bergen damit die Möglichkeit, zusätzliche Quellen nicht nur über einen Eintrag im Bibliothekskatalog zu finden, sondern die entsprechenden Dokumente in digitaler Form auf den eigenen Rechner zu laden. Gerade auch für wertvolle historische Bestände – antike Papyri, Handschriften aus dem Frühmittelalter, einzelüberlieferte Frühdrucke usw. – entsteht hier die Möglichkeit, diese einem größeren Kreis verfügbar zu machen, ohne das Objekt zu großer Belastung durch häufige Benutzung auszusetzen, und ohne dass die Benutzer:innen lange Reisen auf sich nehmen müssten. Für mittelalterliche und frühneuzeitliche Handschriften und Drucke beispielsweise existieren mittlerweile mehrere (meist nationale) Portale, die eine zentrale Suche über alle Bestände ermöglichen; eine Auswahl findet sich unter Section 2.4.2.

Neben der Digitalisierung vorhandener Quellen (Retrodigitalisierung) steht die unaufhörliche Entstehung neuer Quellen in rein digitaler Form (born digital data). Der relativen Knappheit von Quellen – und damit Daten –, die Vormodernehistoriker:innen oftmals zu beklagen haben, steht eine Überfülle an zeitgenössischem Material gegenüber, und beide Situationen – zu wenig/zu unvollständige und zu viele/zu unübersichtliche Datenmengen – bergen methodische Probleme: Wie stellt man ein Korpus, also eine Sammlung von Quellen zusammen, das ausreichend Dokumente beinhaltet, um Fragestellungen zu beantworten, Thesen zu stützen, neue Erkenntnisse zu erhalten, das aber gleichzeitig in einem Forscher:innenleben bewältigbar bleibt? Historiker:innen müssen neue Kompetenzen erwerben, um mit solchen Fragen reflektiert umzugehen. Zur klassischen Quellenkritik kommt die digitale Quellenkritik, zur Fähigkeit, analoge Quellen zu lesen und zu verstehen, ein Äquivalent für den digitalen Bereich. Etwas ausführlicher geht es in Kapitel 3 um Digital Literacy und Digital Criticism.

2.2 Digitale Tools zur Analyse

Die hier bereits zitierte Definition, die aktive und kritische Nutzung digitaler Werkzeuge in Forschung, Lehre oder Studium sei es, was Digital Humanities ausmachten, wirft die Frage auf, was genau unter digitalen Werkzeugen, unter digital tools zu verstehen ist, und zu welchem Zweck man sie einsetzt. Allein schon das Lesen dieses Guides ist ohne digitale Hilfsmittel nicht möglich – es existiert kein gedrucktes Exemplar davon. Lesen am Bildschirm allein macht noch keinen digital humanist, aber man muss nicht erst eine Programmiersprache lernen, um den Computer für die eigene Arbeit zu nutzen und zu Ergebnissen zu kommen, die mit klassischen Methoden – im Bereich der Geschichtswissenschaften etwa papierbasiertes close reading von Quellen und Forschungsliteratur – nicht im selben Ausmaß erzielt werden könnten.

Untersuchungen, die digitale Methoden einsetzen, sind im Normalfall skalierbar – wenn man eine Software benutzt, die die Häufigkeit von Begriffen in einem Dokument zählt, sollte es keinen Unterschied in der Anwendung machen, ob man eines oder einhundert Dokumente auswerten will. Würde man dasselbe per Hand tun, wäre man analog zum Anwachsen der

Dokumente mit der Auszählung beschäftigt. Digitale Werkzeuge ermöglichen es also unter anderem, Untersuchungen auf größere Mengen von Dokumenten auszuweiten. Sie ermöglichen es auch, an ein so erweitertes Korpus andere Fragen zu stellen, als dies mit einer kleineren Quellen/Datengrundlage möglich wäre. Die vorherrschende Überlieferung historischer Quellen besteht aus Text, handgeschrieben, gemeißelt oder gedruckt – und durch die Möglichkeit, diesen mittels Texterkennung in computerlesbare Daten umzuwandeln, ergeben sich neue Perspektiven für die Arbeit von Historiker:innen: Wenn Texte als Daten verstanden werden, lassen sich aus Textquellen Datenbestände erstellen, die mithilfe quantitativer Methoden untersucht und ausgewertet werden können.¹

Für die Literaturwissenschaften beispielsweise ist ein wichtiges Anwendungsfeld die Überprüfung von Autor:innenschaft: Ob ein anonym überliefertes Werk einem:r namentlich bekannten Autor:in zugeschrieben werden kann, lässt sich entweder durch close reading von Literaturwissenschaftler:innen überprüfen, oder durch die Suche nach patterns, Mustern, nach quantifizierbaren Eigenschaften eines Textes, wie beispielsweise die Häufigkeit von Funktionswörtern, Partikeln, Satzzeichen usw. Der unter dem Pseudonym Robert Galbraith veröffentlichte Kriminalroman The Cuckoo's Calling konnte mit entsprechender Software Joanne K. Rowling zugeschrieben werden – damit dauerte die Untersuchung dreißig Minuten, was etwa dem Lesen von zwanzig Romanseiten entspricht. Zu einem Artikel, der diesen Fall thematisiert und in das Feld der linguistischen Forensik einbettet, die Straftäter:innen mithilfe quantitativer Textanalyse ermittelt, geht es hier. Ein Video zur Entwicklung und Anwendung von Software zur Zuschreibung von Autor:innenschaft finden Sie hier. Die genutzte Software, JGGAP, ² lässt sich offensichtlich auch für historische Analysen nutzen - man denke nur an Herrschaftssysteme, in denen strenge Zensur geübt wird/wurde und viele Autor:innen daher nicht unter ihrem Klarnamen publizier(t)en. Durch eine Identifikation anonymer Schreiber:innen lassen sich weitere Aspekte rund um die Thematik Zensur untersuchen – welche Akteur:innen waren öffentlich bekannt, wer publizierte gleichzeitig anonym und unter Klarnamen, welche Autor:innen schrieben aus dem Exil, welche Netzwerke lassen sich rekonstruieren usw. Dadurch, dass ein Programm durch quantitative Auswertungen die Kärrnerarbeit der Identifikation abnehmen kann – um einen reflektierten Umgang mit Daten und Algorithmen geht es in Kapitel 3 -, bleibt mehr Zeit für die qualitative Arbeit; gleichzeitig fußt die Analyse auf einem aussagekräftigen Datensatz, anstatt nur Einzelbeispiele beleuchten zu können.

Quantitative und qualitative Methoden sollen hier keinesfalls gegeneinander ausgespielt werden; vielmehr soll verdeutlicht werden, dass beide Herangehensweisen Vor- und Nachteile haben, und dass sie im besten Fall gewinnbringend miteinander kombiniert werden können – quantitative Auswertungen nur um ihrer selbst willen und ohne eine spezifische historische Fragestellung generieren kaum je einen Mehrwert.

Je nach Datengrundlage, Analysezweck und Forschungsfrage bieten sich unterschiedliche Tools

¹In der Zeitgeschichte ist dies einfacher möglich, weil viele Texte hier schon digital vorliegen; für Epochen bzw. Forschungsbereiche, die eher eine Quellenarmut zu beklagen haben, ist eine Generierung von Daten aus Texten umso interessanter.

²Java Graphical Authorship Attribution Program, http://evllabs.github.io/JGAAP/.

zur Nutzung an; für die meisten Forschungsvorhaben bis zum Ende des Studiums dürfte existierende Software ausreichen, sei es für die Akquise und Aufbereitung von Daten(-sätzen), für verschiedene Arten von Textanalysen, statistische Auswertungen, Netzwerkanalysen, Geomapping oder Visualisierungen. Eine Auswahl an Tools – alle kostenfrei/open source – für spezifische Analysen findet sich unter Literatur, Tools, Tutorials. Für gewisse Analysen bietet es sich an, Programmierkenntnisse zu erwerben – das Erstellen eigener Skripts, also kleiner Programme, beinhaltet die umfassende Kontrolle darüber, wie Daten eingelesen, aufbereitet, angereichert, analysiert und visualisiert werden; bei wiederkehrenden Prozessen, die händisch einige Arbeitszeit in Anspruch nehmen würden, lässt sich so zusätzlich Zeit sparen.

Für geisteswissenschaftliche Projekte werden zurzeit vor allem zwei Programmiersprachen genutzt, R und Python. Da sich beide großer Beliebtheit in den Humanities erfreuen, existieren mittlerweile zahlreiche Packages, die Data und Text Mining, also groß angelegte Daten- und Textanalysen, sehr einfach machen. Solche Packages für Programmiersprachen kann man sich wie Plug-Ins für Programme vorstellen, beispielsweise ein AdBlocker für den Browser. So etwas war von den Entwickler:innen ursprünglich nicht vorgesehen, aber jemand hatte Bedarf, Werbeanzeigen zu blockieren, hat hierzu ein Programm geschrieben und es der Allgemeinheit zur Verfügung gestellt. Der Unterschied zu einem Package ist, dass dieses verschiedene Funktionen zur Verfügung stellt – auswählen und ausführen müssen die Anwender:innen. Wer in Schule und Studium keine Berührungspunkte mit Programmieren hatte, wird zu Beginn vielleicht größere Berührungsängste haben – aber noch einmal: Sie müssen nicht programmieren können, um quantitativ zu arbeiten. Speziell an Historiker:innen ohne Programmier-Vorkenntnisse richtet sich das Projekt "The Programming Historian", das seit 2008 zahlreiche Tutorials veröffentlicht, um verschiedene Tools, Techniken und Workflows für die geschichtswissenschaftliche Forschung und Lehre vorzustellen.

2.3 Digitale Tools zur Kommunikation

Auch ohne quantitative/computergestützte Analysen bei der historischen Arbeit selbst bieten sich mithilfe digitaler Unterstützung verschiedene Möglichkeiten an, um Forschungsergebnisse zu kommunizieren. Dies betrifft zum einen unterschiedliche Formen der digitalen Publikation, zum anderen die Nutzung digitaler Tools zur Visualisierung von Ergebnissen. Gedruckte wissenschaftliche Zeitschriften und Buchverlage haben gewisse Kriterien bei der Annahme eines Textes zur Publikation – inhaltliche Qualität ist bei jeder wissenschaftlichen Arbeit wichtig; aber Formalia wie Länge/Umfang/Bebilderung sind für Onlineformate weniger zentral. So können erste Resultate aus einem neuen Forschungsprojekt oder auch Ergebnisse aus einer Seminararbeit einem interessierten Publikum in Form eines Blogbeitrags ohne viel organisatorischen Vorlauf und Rücksicht auf die Auslastung von Druckerpressen präsentiert werden.

Die Zahl der wissenschaftlichen Blogs ist im letzten Jahrzehnt kontinuierlich gestiegen, sodass für verschiedenste Fächer und Themengebiete eine Auswahl an passenden Publikationsorten besteht. Sie verfügen über eine Redaktion, also verantwortlich zeichnende

Wissenschaftler:innen, die eine Qualitätssicherung der Beiträge leistet. Auch viele Promotionsprojekte werden von Blogs begleitet – das liefert Sichtbarkeit für die Forschung der Doktorand:innen und bietet eine Plattform für Themen, die in der Dissertation keinen Platz finden, aber doch erwähnenswert sind.

Hypothèses hat sich als eine wichtige Plattform etabliert, die zahlreiche Wissenschaftsblogs hostet; ein Katalog aller Blogs, filterbar nach Themen und Epochen, findet sich hier.

Aus Lehrveranstaltungen und Forschungsprojekten am Departement Geschichte der Universität Basel sind u.a. folgende Blogs hervorgegangen:

- Exilland Schweiz. Kulturschaffende und Intellektuelle im Schweizer Exil, 2021–
- Global Health Africa. Circulating Knowledge and Innovations, 2019–
- Materialized Histories. Materielle Kultur und digitale Forschung, 2021-
- Materialized Identities. Objects, Affects and Effects in Early Modern Culture (1450–1750), 2016–2021.
- Stadt.Geschichte.Basel. Blog zum Forschungsprojekt, 2019-
- The Color Line. Race Relations in Schlüsseltexten amerikanischer Autor:innen (1881–1953), 2022–

Als relativ junges Format haben sich Data Stories etabliert – Narrative, die auf der Basis von (quantitativen oder qualitativen) Daten und Analysen einen Sachverhalt darstellen. Diese Art von (Daten-)Publikation wird u.a. im journalistischen Bereich genutzt, auch, um Interaktivität und Aktualität von Daten zu gewährleisten. Es gibt verschiedene Tools, um solche Stories zu erstellen, zum Teil bereits mit integrierter Publikationsmöglichkeit; eine Auswahl findet sich unter Literatur, Tools, Tutorials.

Ein Beispiel, das auch Input durch die Nutzer:innen ermöglicht, ist Darüber spricht der Bundestag, eine Data Story der "ZEIT", die alle Reden im deutschen Bundestag seit 1949 analysierbar macht. Ein Beispiel für eine kartenzentrierte Darstellung mit Einbettung verschiedener Medien ist die Reise von Arya aus Game of Thrones. Ein Beispiel für die Nutzung von Volkszählungsdaten stammt vom Schweizer Bundesamt für Statistik, Die Schweiz (er)zählen.

2.4 Projekte und Ressourcen für Forschung und Lehre

2.4.1 Alte Geschichte

Projekte:

• D-Scribes: Projekt zur Erkennung antiker Schreiber:innen griechischer und koptischer Papyri

Ressourcen/Portale:

- Papyprus Portal: digitale Papyrussammlung
- The Arabic Papyrology Database: Datenbank für vormoderne arabische Schriften auf Papyrus, Pergament oder Papier vom 7. bis zum 16. Jahrhundert

2.4.2 Mittelalter und Frühe Neuzeit

Projekte:

- Burchards Dekret Digital: Digitale Edition, die die handschriftliche Überlieferung des Decretum Burchardi erschließt.
- hallerNet: Editions- und Forschungsplattform mit mehreren Textsammlungen von Gelehrten aus der Übergangszeit zwischen 1700 und 1850, mit Strukturdaten verknüpft
- Printed Markets: Projekt zur Digitalisierung und Anreicherung der frühneuzeitlichen Basler Anzeigenzeitung "Avisblatt" (1729–1844)
- Repertorium Academicum: Projekt zur Erfassung europäischer Gelehrter zwischen 1250 und 1550

Ressourcen/Portale:

- dMGH: Monumenta Germaniae Historica online (Beta-Version)
- e-codices: Virtuelle Handschriftenbibliothek der Schweiz
- Fragmentarium: Laboratory for Medieval Manuscript Fragments
- Handschriftenportal: Zentraler nationaler Nachweis für Buchhandschriften in deutschen Bibliotheken und in deutscher Sprache (Entwicklungsstadium)
- Innovating Knowledge: Datenbank und digitale Edition zu Isidors von Sevilla "Etymologiae"
- e-manuscripta: Digitalisierte handschriftliche Quellen aus Schweizer Bibliotheken und Archiven
- e-rara: Plattform für digitalisierte Drucke aus Schweizer Institutionen
- Gallica: Digitalisierte Quellen aus französischen Bibliotheken
- Stapfer Enquête: Edition der Helvetischen Schulumfrage von 1799
- swisscollections: Suchplattform für historische Schweizer Bestände

• transcriptiones: Plattform zum Erstellen, Teilen und Nutzen von Transkriptionen historischer Manuskripte

2.4.3 Moderne und Zeitgeschichte

Projekte:

- impresso. Media Monitoring of the Past: Projekt zur Aufbereitung, semantischen Anreicherung, Repräsentation, Exploration und Recherche in historischen Medien (Zeitungen und Radio) über Zeit, Sprachen und Nationalgrenzen hinweg
- Living with Machines: Forschungsprojekt über die Auswirkungen der Mechanisierung der Arbeit während der Industrialisierung
- Refugee History: Wissenschaftliches Blog und interaktives Netzwerk zu aktuellen Debatten um das Thema "Flüchtlinge"

Ressourcen/Portale:

- Datenbank Bild + Ton zur Geschichte (Schweizer) sozialer Bewegungen
- Dodis: Wissenschaftliche Edition von Dokumenten zur Schweizer Außenpolitik
- Gallica: Digitalisierte Quellen aus französischen Bibliotheken
- e-newspaperarchives.ch: Schweizer Zeitungen online
- e-periodica: Schweizer Zeitschriften online
- Historische Statistik der Schweiz (HSSO)
- histat: Zeitreihen zur Historischen Statistik

2.4.4 Jüdische Geschichte

Projekte:

• Digital Jewish Studies Online, Stroum Center for Jewish Studies, University of Washington

Ressourcen/Portale:

• Blavatnik Archive: Archiv zur Erhaltung und Verfügbarmachung von Material zur (jüdischen) Geschichte des 20. Jahrhunderts mit Fokus auf die zwei Weltkriege und Sowjetrussland.

 Menny, Anna; Rürup, Miriam; Siegel, Björn: Jüdische Geschichte im deutschsprachigen Raum, in: Busse, Laura u. a. (Hg.): Clio-Guide. Ein Handbuch zu digitalen Ressourcen für die Geschichtswissenschaften, Berlin 2018, S. E.2-1–E.2-56. Online: https://doi.org/ 10.18452/19244.

2.4.5 Geschichte Afrikas

Projekte:

- Emandulo: Digitales Archiv, das archivalische/museale Sammlungen und Präsentationen über präkoloniale südafrikanische Geschichte zusammenführt und neu zusammenstellt.
- Legacies of British Slavery: Forschungsprojekt zum britischen Sklavenhandel und -besitz

Ressourcen/Portale:

- FHYA: Experimentelle digitale Forschungsplattform über präkoloniale südafrikanische Geschichte
- Legacies of British Slavery: Datenbank zum britischen Sklavenhandel und -besitz
- Slave Voyages: Datenbanken zum transatlantischen und interamerikanischen Sklavenhandel und Personendatenbank

2.4.6 Osteuropäische Geschichte

Projekte:

- Gulag: Many Days, Many Lives: Archiv und Präsentationsplattform zu den sowjetischen Gulags
- Gulag Online: Virtuelles Museum mit Präsentationen und Quellen zum Leben im Gulag, zu Personen und Objekten
- Seventeen Moments in Soviet History: Multimediales Online-Archiv mit ausgewählten Quellen zu Ereignissen in der sowjetischen Geschichte anhand 17 verschiedener Jahre zwischen 1917 und 1991
- The Imperiia Project: Eine räumliche Geschichte des Russischen Reichs

Ressourcen/Portale:

Blavatnik Archive: Archiv zur Erhaltung und Verfügbarmachung von Material zur (jüdischen) Geschichte des 20. Jahrhunderts mit Fokus auf die zwei Weltkriege und Sowjetrussland.

• The Other Side: Webarchiv von Interviews ehemaliger Ostarbeiter:innen, Kriegsgefangener und Insassen deutscher Lager; Publikationsplattform

2.4.7 Epochen-/Areaübergreifend:

Projekte:

• Lord of the Rings Project: Interaktive Analyse der Werke J. R. R. Tolkiens

Ressourcen/Portale:

- Around DH in 80 days: Portal zur Vorstellung achtzig verschiedener Digital-Humanities-Projekte weltweit und aus verschiedenen Disziplinen
- Internet Archive: digitale Bibliothek zur Archivierung von Büchern, Bildern, Filmen, Software, Musik und Webseiten

3 Digital Literacy, Digital Criticism

3.1 Digital Literacy, Data Literacy

Unter Data Literacy wird die Kompetenz verstanden, Daten zu sammeln, zu managen, zu evaluieren und zu nutzen, [7, S. 8] eine Kompetenz, die jede:r für den mittlerweile unvermeidlichen Umgang mit Daten verschiedenster Art im eigenen Alltag entwickeln sollte. Je nach Forschungsdisziplin ergeben sich weiter gewisse Spezifika, wobei Studierenden der Geisteswissenschaften ein Thema wie Algorithmenkritik nicht als erstes in den Sinn kommt, wenn es um die im Studium zu erwerbenden Kompetenzen geht. Aber auch ohne den Quellcode von machine-learning-Software im Detail zu verstehen, ermöglicht ein grundlegendes Verständnis von und ein Wissen über die Funktionsweisen solcher Anwendungen einen reflektierten Umgang mit diesen. Eine solche Art von Digital bzw. Data Literacy ist vor allem dann relevant, wenn es um die Interpretation von Ergebnissen geht, die scheinbar objektiv sind, bzw. scheinbar objektiv entstanden. Ein gutes Beispiel hierfür sind die Ergebnislisten bei Suchanfragen in einer Suchmaschine. Je nachdem, welchen Anbieter Sie nutzen, spielen verschiedene Umstände in die Generierung von Trefferlisten hinein, beispielsweise Ihre Suchhistorie, sodass search neutrality nicht mehr gewährleistet ist. 2

Gehen Sie auf die Bilder-Suche von Google und suchen Sie nach "historian". Was sehen Sie?

Wüsste ich nichts über Geschichtswissenschaftler:innen, würde ich aufgrund der Ergebnisse meiner Suche davon ausgehen, "a historian" wäre meist ein alter, weißer Mann mit Brille, Bart und einem großen Bücherregal; wenn Sie sich am Departement Geschichte der Uni Basel umsehen, dürfte ein etwas anderer Eindruck entstehen. Die Ergebnisse von Suchmaschinen, die für ihr Funktionieren Algorithmen anwenden, sind biased, verzerrt: Sie beruhen auf vorangegangenen Suchen, Vorlieben, geographischem Standort – und auf von Menschen eingegebenen Metadaten, also Daten mit Informationen über andere Daten. Ein Bewusstsein hierfür und das Hinterfragen von Datensätzen gehören also mit zur Arbeit in einer digitalisierten Welt.

¹Verschiedene Organisationen und Interessensverbände machen sich hierzu Gedanken; die deutsche "Gesellschaft für Informatik" hat ein Grundlagenpapier erstellt, in dem es um Data Literacy und Data Science Education als digitale Kompetenzen in der Hochschulausbildung geht: https://gi.de/fileadmin/GI/Hauptseite/Aktuelles/Aktionen/Data_Literacy/GI_DataScience_2018-04-20_FINAL.pdf.

²Alternative Suchmaschinen sind zum Beispiel Startpage oder DuckDuckGo.

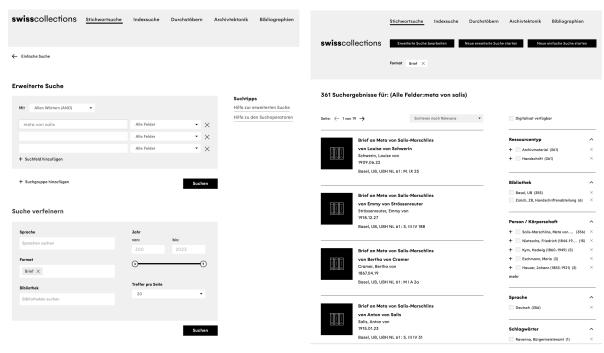
3.2 Digital Criticism, Data Criticism

Digitalisierte Quellen ebenso wie rein digitale erfordern eine erweiterte Art von Quellenkritik – im Einführungskurs an der Universität Basel lernen Sie die Grundlagen klassischer Quellenkritik:

Woher kommt eine Quelle, wer hat sie unter welchen Umständen und zu welchem Zweck erstellt? Welche Absichten können darin verborgen sein, und welche Verzerrungen können sich durch sie ergeben? Welche Tendenzen könnten sich in hochmittelalterlichen Herrscherchroniken verstecken, wenn der Verfasser in direkter Abhängigkeit des Auftraggebers stand? Wie sind Zeug:innenaussagen in Hexenprozessen zu bewerten, die unter Anwendung körperlicher Strafe entstanden sind? Mit wie viel Vorsicht sind die Inhalte eines Tagebuchs zu bewerten, das allem Anschein nach mit Blick auf eine spätere Veröffentlichung verfasst wurde?

Neben der inneren Kritik geht es bei der Arbeit mit Quellen immer auch um Fragen der Korpusbildung: Wie kann eine Quellengrundlage erstellt werden, die für Beantwortung einer spezifischen historischen Fragestellung belastbar und aussagekräftig genug ist und gleichzeitig in angemessener Zeit bearbeitet werden kann? Hinzu kommen Spezifika bei der Arbeit mit unterschiedlichen Quellenformen bzw. -formaten: Bei analogen Quellen, die auch in digitaler Form zur Verfügung stehen, besteht die Gefahr, dass ein Thema, ein Bereich, ein Aspekt vernachlässigt wird, wenn nur die unmittelbar verfügbaren, digitalisierten Bestände zur Korpusbidlung genutzt werden. Wenn Sie sich beispielsweise für die Schweizer Historikerin und Frauenrechtlerin Meta von Salis (1855–1929) und deren briefliche Korrespondenz – Friedrich Nietzsche war einer ihrer Brieffreunde – interessieren und über die Suchplattform für historische Schweizer Bestände, swisscollections, in nationalen Bibliotheken und Archiven nach entsprechenden Dokumenten suchen, erhalten Sie 361 Treffer:

Digital verfügbar waren hiervon im Oktober 2022 lediglich drei Einträge, wobei der erste ein Brief von Nietzsche an Meta von Salis ist, der zweite Eintrag umfasst sieben Briefe von Caroline Farner, und der dritte Eintrag ist weder an noch von Meta von Salis, sondern hat sie nur zum Thema:

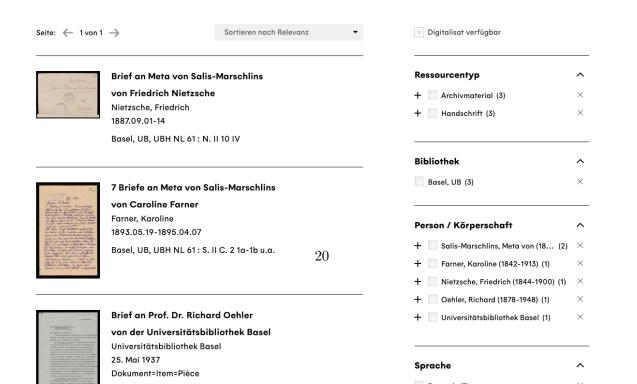


(a) Erweiterte Suchmaske von swisscollections

(a) Suchergebnisse für "Meta von Salis" + "Brief"



3 Suchergebnisse für: (Alle Felder:meta von salis)



Ihnen würde bei einer Korpuserstellung vom Schreibtisch aus, also nur mit den angezeigten Digitalisaten, also der Großteil der Überlieferung fehlen, und Ihre Untersuchungsergebnisse wären wohl sehr verzerrt, würden Sie statistische Aussagen treffen wollen: Meta von Salis unterhielt brieflichen Kontakt zu einem Mann und einer Frau, das Geschlechterverhältnis wäre also ausgeglichen; und Frauen schreiben im Schnitt mehr Briefe an Meta von Salis als Männer. Beim Blick auf alle Suchergebnisse würden sich Ihre Aussagen aber sehr ändern, und es würde sich lohnen, diese Verzerrung, diesen Bias aus Ihrer Datengrundlage zu entfernen.

Hinzu kommt natürlich immer das grundlegende Problem bei der Suche nach Quellen: swisscollections und ähnliche Portale können nur anzeigen, was die Kooperationspartner:innen zur Verfügung stellen. Hat eine Bibliothek Briefe von Meta von Salis in ihrem Bestand, diese aber noch nicht als Datensatz erfasst, wissen Sie im Gegensatz zum obigen Beispiel nicht einmal, dass Ihnen etwas entgehen würde, dass in Ihrem Korpus überhaupt ein Bias vorhanden ist.

Ähnliche Vorsicht zur Vermeidung von Verzerrungen in der Datengrundlage gilt bei der Arbeit mit rein digitalen Daten, beispielweise bei der Auswertung von Datensätzen aus Befragungen. Wenn Sie sich am 27.10.2022 vor die Universitätsbibliothek in Basel stellen und einen Tag lang mithilfe eines kurzen Fragebogens und einer Tabellendatei erfassen, wie zufrieden die befragten Personen mit dem Essen in der Unimensa sind, werden Sie am Ende einen Datensatz erhalten, in dem sich vermutlich über 80% der Befragten für besseres und nahezu 100% für günstigeres Essen in der Mensa aussprechen – eine gute Schlagzeile für die BZ, die sich auf die neuesten Ergebnisse einer wissenschaftlichen Studie berufen kann. Führen Sie die gleiche Umfrage eine Woche später, mitten während der Herbstmesse durch, werden die Ergebnisse wohl erheblich anders aussehen. Die Wahrscheinlichkeit, dass die Mensa infolge der BZ-Schlagzeile innerhalb weniger Tage den Menüplan überarbeitet und die Preise herabgesetzt hat, ist dabei wohl geringer als diejenige, dass sich Ihr Sample, die Auswahl an Datenpunkten, also befragten Personen, durch die Messe stark verändert hat: Im Umkreis der Bibliothek treffen Sie nun nicht mehr vor allem Studierende und andere Uni-Angehörige an, sondern auch Messebesucher:innen vom Petersplatz. Auch hier sind Verzerrungen entstanden, ähnlich wie beim vorherigen Beispiel mit den Briefen: Wenn aus einer Gesamtheit nur eine spezifische Untermenge beobachtet wird, die sich durch ein gemeinsames Merkmal von der Gesamtheit unterscheidet – digitalisierte Quelle oder Besucher:in der Universitätsbibliothek -, ist die Datengrundlage und damit die Untersuchungsergebnisse biased. Um bei Daten, die Sie nachnutzen, eventuell vorhandene Verzerrungen nicht weiterzutransportieren, ist das Üben von Datenkritik eine essentielle Kompetenz.

Zur Tatsache, dass Daten eben nicht "gegeben" sind (lat. dare, datum: geben, gegeben), sondern gemacht, und daher entsprechend interpretiert werden müssen, finden Sie ein gutes Interview von Roopika Risam (2020);[8] zur Zementierung von Klischees durch Übersetzungsalgorithmen gibt es einen Artikel in der Republik von Marie-José Kolly und Simon Schmid (2021);[4] und über die Macht von Data Science und dem Änderungspotential von Data Feminism haben Catherine D'Ignazio und Lauren F. Klein 2020 ein ganzes Buch veröffentlich.[3]

Zur Frage, wie sich die digitale Wende, der digital turn, auf die Quellenkritik auswirkt, sehen Sie sich dieses kurze Video des Projekts Ranke.2 – Quellenkritik im digitalen Zeitalter an:³

Eine Handreichung zum Umgang mit digitalisierten und digitalen Daten, das im selben Projekt erarbeitet wurde, finden Sie hier.

3.3 Künstliche Intelligenz und Algorithmenkritik

Mit dem Release von ChatGPT-3 der Firma OpenAI im Winter 2022 hat sich die Nutzung künstlicher Intelligenz (KI) auf einen Schlag für einen sehr großen Personenkreis geöffnet: Ein Chatbot steht, nach einmaliger Registrierung, über eine URL den unterschiedlichsten Fragen Rede und Antwort, und erklärt wahlweise Fermats letzten Satz, produziert Sonette im Stil Shakespeares über die Liebe einer Emmentalers zu einer Büroklammer oder fasst den in den 1980er-Jahren in Deutschland geführten Historiker:innenstreit über die Singularität des Holocaust zusammen. Technisch handelt es sich bei ChatGPT um ein generatives Sprachmodell (Generative Pre-trained Transformer), ein sog. Large Language Model (LLM), das mit einem sehr großen Datensatz trainiert wurde – etwa 1 Billion Wörter aus sechs Millionen Wikipedia-Artikeln, Millionen Webseiten und digitalisierten Büchern, insgesamt etwas 570 GB –, um auf dieser Grundlage neue, der Anfrage bzw. Eingabe der Nutzer:innen entsprechende Texte zu produzieren.⁴ Bittet man ChatGPT selbst, das Prinzip eines generativen Sprachmodells und seines Algorithmus zu erklären, erhält man folgende Antwort:

Hallo! Gerne erkläre ich dir das Prinzip hinter einem generativen Sprachmodell und meinem Algorithmus.

Ein generatives Sprachmodell ist ein künstlicher Intelligenz-Algorithmus, der darauf trainiert wird, Texte zu erzeugen, die einer bestimmten Sprache oder einem bestimmten Stil entsprechen. Das Modell nutzt dabei eine statistische Methode namens Wahrscheinlichkeitsverteilung, um die Wahrscheinlichkeit zu bestimmen, mit der ein bestimmtes Wort oder eine bestimmte Phrase in einem Text auftaucht. Das Ziel des Algorithmus besteht darin, Texte zu generieren, die natürliche Sprachstrukturen und -muster aufweisen, sodass sie für den Menschen verständlich und lesbar sind. Das generative Sprachmodell kann dabei verschiedene Ansätze verfolgen, wie zum Beispiel die Verwendung von rekurrenten neuronalen Netzen (RNN) oder Transformer-Modellen.

Als ChatGPT, basiere ich auf der GPT-3.5 Architektur und nutze ebenfalls ein Transformer-Modell. Mein Algorithmus nutzt eine große Menge an Textdaten, um Sprachmuster und -strukturen zu erlernen und kann anschließend Texte generieren,

³Eine Selbstlerneinheit zum Thema finden Sie hier: https://ranke2.uni.lu/u/archival-digital-turn/.

⁴Eine ausführliche und anschauliche Erklärung von Helmut Linde findet sich unter https://www.golem.de /news/kuenstliche-intelligenz-so-funktioniert-chatgpt-2302-171644.html, ein Einführungsvortrag mit besonderem Blick auf Konsequenzen für die Hochschullehre von Thomas Arnold unter https://www.youtube.com/watch?v=-c8ogAwX6KI.

die dieser Sprache ähnlich sind.

Um Texte zu erzeugen, nutzt mein Algorithmus eine sogenannte autoregressive Methode, bei der er das zuvor erzeugte Wort oder die zuvor erzeugte Phrase als Input verwendet, um das nächste Wort oder die nächste Phrase vorherzusagen. Auf diese Weise wird ein Text Schritt für Schritt erzeugt, wobei jeder Schritt auf dem vorherigen basiert.

Durch die ständige Anpassung der Gewichtungen im Modell, je nachdem wie gut es darin ist, den nächsten Textsatz vorherzusagen, wird das generative Sprachmodell im Laufe der Zeit immer besser darin, realistische und gut klingende Texte zu erzeugen.⁵

Das bedeutet, dass die Antworten besonders gut bzw. sinnvoll für Fragen zu Themen sind, für die sehr viel Trainingsmaterial vorlag. Das bedeutet auch, dass alles, was nach Abschluss des Trainings im November 2021 online gestellt wurde, nicht mehr Eingang in das Modell 3.5 gefunden hat; hier wird zwar aktuell nachgeholt, und je nach Release und Modell gibt es hier Unterschiede, aber Ereignisse wie den Ukraine-Krieg kennt ChatGPT (noch) nicht.⁶ Auch ist es nicht immer einfach, die Ausgewogenheit und Qualität der Antworten richtig abzuschätzen – bei der Erstellung von Programmierfunktionen lässt sich der Vorschlag eines Chatbots direkt testen. Bei Fragen, die etwas offener gestellt sind, bieten die LLMs sehr überzeugend klingende Antworten – die aber natürlich geprägt sind vom zugrundeliegenden Trainingsmaterial, Modell und Trainingsmethoden, die nicht immer transparent sind.⁷ Auf die Frage nach den Namen von "ten medieval philosophers" gibt ChatGPT-40 zehn Männer aus, sieben davon aus West-/Zentraleuropa. Erst auf den Hinweis, dass die Liste recht männerlastig und weiß sei, bekommt man eine größere Auswahl arabischer und persischer Philosophen und (westlicher) Frauen. Es braucht noch eine weitere Nachfrage, um auch indische und chinesische Namen zu erhalten. Der Großteil der Texte, die ins Training einflossen, sind englisch und damit größtenteils aus dem angloamerikanischen Raum. Ein Ungleichgewicht zugunsten Wissensproduktion und Meinungen aus diesem (Sprach-)Raum ist daher unvermeidlich und muss bei der Nutzung stets mitbedacht werden.

Mit Blick auf die obige, tadellos formulierte Antwort zum eigenen Funktionieren liegt es natürlich nahe, LLMs für alle Arten von Schreibarbeiten zur Hilfe zu ziehen, und die Aufregung

⁵ChatGPT-3.5-Antwort, Anfrage zu "Wie funktioniert ein generatives Sprachmodell und der dazugehörige Algorithmus?", 11. November 2022.

⁶Die neueste Version von ChatGPT, GPT-4.o, hat keine neueren Trainingsdaten als 3.5; perplexity gibt Juli 2023 als cut-off-Datum an, Googles Gemini den Zeitraum Ende 2023 bis Anfang 2024. Zwar fließen neueste Onlineartikel mit in Antworten ein, aber je nach Ereignis fehlen substantielle Hintergrundinformationen aus der Zeit zwischen dem Abschluss des Trainings und einer Online-Suche einige Monate später.

⁷Das "Open" im Firmennamen OpenAI, das hinter ChatGPT steht, ist also nicht Programm. Der vollständige technische Report, der viele grundlegende Parameter auslässt, ist hier abrufbar. Zur Kritik am Vorgehen von OpenAI siehe u.a. einen Artikel von Leigh Mc Gowran. Eine Open-source-Alternative mit offenem Trainingsset und Modell ist z.B. dolly. Das LLM perplexity wirbt damit, mit Fußnoten auf genutzte Onlinequellen zu verweisen und die Ergebnisse somit nachvollziehbarer zu machen.

⁸ChatGPT-40-Antwort, Anfrage zu "Name ten medieval philosophers", 3. September 2024.

im (Hoch-)Schulbereich war mitunter groß⁹ – auch, weil die Programme erfundene Fakten im Brustton der Überzeugung formulieren oder auf nicht existierende Referenzen verweisen. Bei unkritischer Übernahme der Antworten und Ausgabe des Textes als eigenes Produkt kommt es also nicht nur zur Täuschung in Bezug auf die eigenen Fähigkeiten, sondern es werden auch falsche Informationen übernommen und mitunter verbreitet, beispielsweise in Referaten. Zahlreiche Hochschulleitungen haben auf die neue Technik reagiert und Codes of Conduct erstellt; an der Universität Basel gibt es eine erweiterte Eigenständigkeitserklärung, in der versichert wird, mit KI erstellte Passagen als solche zu kennzeichnen – Software, die mithilfe KI erstellte Texte als solche identifiziert, gibt es natürlich auch. [Das Vizerektorat Lehre hat einen Leitfaden "Aus KI zitieren" erstellt.]

Überlegungen zum konstruktiven Einsatz generativer Sprachmodelle, als Lernassistent, ¹⁰ Schreibcoach oder Hilfe bei der Recherche, gibt es ebenso, und gelungene Beispiele, wie ChatGPT in der Lehre eingesetzt werden kann. ¹¹ Zentral für eine gewinnbringende und informierte Nutzung ist sicher ein ungefähres Wissen über die Funktionsweise des Algorithmus und die Bereitschaft, erhaltene Antworten kritisch zu prüfen und ggf. zu verwerfen.

⁹Siehe z.B. in Basel jetzt, 19.01.2023; NDR, 09.02.2023; bajour, 20.03.2023; 20 minuten, 24.03.2023.

¹⁰Siehe z.B. Pedagogical Educational Tutor (PET), entwickelt an der Universität Hohenheim, mit Einführungsvideo.

¹¹U.a von Andrea Klein oder Ethan Mollick.

4 Datenerhebung, -aufbereitung und -analyse

Jede Art von Forschung ist auf Daten angewiesen, seien sie mittels Personenbefragungen, medizinischer Messungen, Web Scraping oder interpretierender Analysen von Texten erhoben. Auf Grundlage von Daten können Forschungsfragen beantwortet, Thesen aufgestellt, Behauptungen widerlegt, Narrative untermauert werden. Analysen, die sich mit einem kleinen Set von Quellen bzw. Daten befassen, präsentieren Ergebnisse dabei oft in Form von Synthesen, die sich aus einer vorangehenden Interpretation der zugrundeliegenden Dokumente ergeben. Über das Quellenverzeichnis und entsprechende Anmerkungen im Text wird die Grundlage nachvollziehbar; dass ein bestimmter Abschnitt, ein Satz oder ein Wort auf eine gewisse Weise ausgelegt werden, wird aber auch durch die jeweiligen Forscher:innen selbst beeinflusst – eine Literaturwissenschaftlerin beispielsweise, die über Männerfiguren bei Joanne K. Rowling promoviert hat, wird bei der Diskussion um deren mögliche Autorschaft von The Cuckoo's Calling (siehe Section 2.2) diesen Text anders lesen und andere Argumente dafür oder dagegen aufwerfen als ein langjähriger Harry-Potter-Fan mit viel Leseerfahrung, aber anderer bzw. weniger formaler Ausbildung. Beide werden fundierte Aussagen treffen und Begründungen geben können, ob und wieso The Cuckoo's Calling von Rowling verfasst wurde oder nicht; beide werden auf ihre Erfahrung und gründliche Auseinandersetzung mit Rowlings Werk verweisen; und beide werden mit einzelnen Sätzen oder Passagen für eine Sichtweise argumentieren, die von einer dritten Person genau gegenteilig genutzt würde. Die Datengrundlage ist also dieselbe und nachvollziehbar, die Auswertung bzw. die Auswertungsstrategien hingegen sind es nicht mehr, und somit auch nicht die daraus gewonnenen Ergebnisse, die ja auch wieder Forschungsdaten darstellen.

Computergestützte Analysen haben den Anspruch, in allen Schritten nachvollziehbar zu sein und dadurch auch nachnutzbare Daten zu produzieren: Nicht nur die Quellengrundlage, also die Erhebung von Daten und die Erstellung eines Datensatzes, sondern auch alle Schritte von der Datenanreicherung und -verfeinerung über die genutzten Methoden bzw. Programme für die Auswertung bis hin zur Sicherung und Aufbewahrung sollen transparent, gut dokumentiert und nachvollziehbar sein. Zum einen, um die Resultate und die darauf fußenden Aussagen belastbar zu machen; zum anderen, um die gewonnenen Daten zur weiteren Nutzung kostenfrei und offen verfügbar zu machen. Zu den Prinzipien, die bei der Arbeit mit Daten berücksichtigt werden sollten, geht es nochmals in Kapitel 5. An dieser Stelle stehen die konkreten Arbeitsschritte bei der Datenerhebung und -aufbereitung, der Datenanalyse und -sicherung im Zentrum, die in Digital-History-Projekten häufig vorkommen.

4.1 Datenerhebung

Es gibt verschiedene Möglichkeiten, Daten für die historische Forschung zu erheben bzw. zu erstellen, von denen einige im Folgenden kurz angesprochen werden.

Für Zeiträume, in denen Quellen vergleichsweise knapp sind und keine seriellen Daten existieren, bietet sich die **Digitalisierung von Texten** und deren anschließende Analyse an. Digitalisierung beinhaltet dabei nicht nur die Transformation von einer physischen Quelle in ein digitales Bild, sondern auch die Anreicherung des Bilds mit Layout und Text: Erst durch eine Markierung von Bereichen, in denen Text vorkommt, ist es in einem zweiten Schritt möglich, diesen als solchen zu erkennen und damit maschinenlesbar und auswertbar zu machen. Eine solche Umwandlung vom Bild zum Text ist dabei sowohl für moderne Texte, die als Typoskript vorliegen, als auch für vormoderne Handschriften und Drucke möglich, in lateinischer ebenso wie in arabischer, chinesischer oder japanischer Schrift. Es gibt kostenpflichtige Programme wie den Abbyy FineReader, aber auch Open-Source-Tools mit und ohne Graphical User Interface (GUI). Weit verbreitet ist Transkribus, das viele Funktionalitäten bündelt; die Texterkennung ist ab einer gewissen Menge Seiten allerdings kostenpflichtig, wobei studentische Projekte auf Anfrage unterstützt werden können. Programme, die über die Kommandozeile laufen, gänzlich kostenfrei sind und ebenfalls zahlreiche Funktionalitäten bieten, sind beispielsweise Kraken, OCR4all, OCRopus oder Calamari.

Zur Extraktion von Daten aus digitalen/digitalisierten Texten existieren verschiedene Möglichkeiten mithilfe kleiner Kommandozeilenprogramme (eher mühsam und schwierig zu lesen) oder mit Packages für Programmiersprachen, für die Geisteswissenschaften vor allem R oder Python (siehe dazu auch Section 2.2). So können beispielsweise aus digitalisierten Telefonbüchern Entitäten, also Einheiten, wie Personen, Straßennamen oder Berufe oder aus alten Theaterprogrammheften gespielte Stücke, beteiligte Schauspieler:innen und verantwortliche Regisseurinnen extrahiert und als Datensätze weitergenutzt werden.¹

Der anfängliche Aufwand, der einer automatisierten Datenextraktion vorangeht, und die steile Lernkurve bei der Bedienung mancher Programme können abschreckend wirken. Wenn Sie nur ein Theaterprogramm detaillierter auswerten wollen, sind Sie sicher schneller, wenn Sie die entsprechenden Daten in eine Tabellensoftware abtippen. Wenn Sie aber einen größeren Quellenbestand zur Verfügung haben, der in sich ähnlich strukturiert ist, wie das bei Telefonbüchern oder einer Serie von Theaterprogrammheften der Fall sein dürfte, macht es kaum einen Unterschied mehr, ob Sie zehn oder hundert Theaterprogramme mithilfe eines Skripts analysieren möchten. Zudem können Sie Ihr erstelltes Skript, Ihr kleines Computerprogramm, anderen zur Verfügung stellen oder für ähnlich strukturierte Quellen in einem anderen Projekt nachnutzen.

¹Ein gut nachvollziehbares Tutorial zur Extraktion von Daten aus Telefonbüchern hat Lindsey Wieck für einen DH-Kurs an der St. Mary's University in San Antonio erstellt: https://lindseywieck.com/fall_2016_sf/gatheringdatatutorial.html. Derek Miller arbeitet zu Broadway-Vorstellungen, Visualizing Broadway, ein Projekt, das hier beschrieben wird; hier gibt es dazu ein Video in Vorlesungslänge.

Wenn Sie mit bereits digitalisierten Beständen aus öffentlichen Institutionen wie Galerien, Bibliotheken, Museen oder Archiven arbeiten wollen (sog. GLAMs: Galleries, Libraries, Archives, Museums), besteht oft die Möglichkeit, Daten über Schnittstellen herunterzuladen.² Solche Schnittstellen, engl. API (Application Programming Interface), ermöglichen eine Kommunikation zwischen zwei Computern, ohne dass hierfür der Umweg über eine graphische Oberfläche nötig ist. Anstatt also beispielsweise über die Suchmaske der Staatlichen Museen zu Berlin nach Objekten oder Dokumenten mithilfe verschiedener Schlagwörter zu suchen und die Ergebnisse dann einzeln herunterzuladen, kann Ihr Computer mit der Schnittstelle des Museums direkt kommunizieren und mit einfachen Befehlen ganze Ergebnislisten zur Weiterarbeit herunterladen. Für solche Abfragen können ein Kommandozeilenprogramm oder Programmiersprachen genutzt werden, die Abfrage besteht dabei im Wesentlichen aus einer Zeile, wie hier in der Programmiersprache R:

library(jsonlite)

cats <- fromJSON("https://smb.museum-digital.de/json/objects?&s=katze")</pre>

Wenn Sie die Schritte nachvollziehen möchten, können Sie R hier herunterladen. Wenn Sie das Programm öffnen, müssen Sie zuerst das Paket jsonlite installieren: install.packages("jsonlite")

Mit "Enter" wird das Paket installiert. Dann können Sie die zwei Zeilen oben eintippen und ebenfalls mit "Enter" ausführen. Die Ergebnisse Ihrer Suche können Sie sich mit

cats + "Enter" anzeigen lassen.

Das Ergebnis der Suchanfrage nach "katze" wird in der Variable cats gespeichert, und diese kann zur Weiterarbeit in ein Tabellenformat exportiert werden:

write.csv(cats, "docs/cats_smb.csv")

Die Funktion write.csv speichert den Inhalt der Variable cats als csv-Datei³ unter dem Dateipfad "docs/cats_smb.csv" auf der Festplatte.

²Unter openglam.ch finden sich Informationen zu Schweizer GLAM-Einrichtungen, die offene Daten anbieten.

³comma separated value ist ein Format in dem einzelne Werte andree über spezifische Tropper, meist commas

³comma separated value ist ein Format, in dem einzelne Werte, values, über spezifische Trenner, meist commas, eindeutig abgrenzbar sind und somit in einem Tabellenformat angezeigt werden können, wobei jeder Wert in einer separaten Zelle steht. Tabellensoftware wie Excel, Google Sheets oder Numbers kann csv-Dateien öffnen.

	Α	В	С	D	E	F	G	Н	1
1	(objekt_name	objekt_inventarnr	objekt_erfasst_am	institution_id	institution_name	image	total
2	1		Statuette der Göttin Bastet in Gestalt einer sitzenden Katze	ÄM 2598	2021-11-02 21:15:59	9	Ägyptisches Museum und Papyrussammlung	data/smb/resources/images/201806/200w_21081501496.jpg	134
3	2	52029	Shintô-Götterschrein mit heiliger Katze	I D 17657 a,b	2021-11-02 21:15:59	11	Ethnologisches Museum	data/smb/resources/images/201807/200w_06173101822.jpg	13
4	3	108895	Katze auf Rädern zum Ziehen	N (35 F) 917/1995,a	2021-11-02 21:15:59	2	Museum Europäischer Kulturen	data/smb/resources/images/201808/200w_04160418311.jpg	13
5	4		Einseitig bemaltes Ostrakon mit Darstellung einer Katze vor einem Opfertisch mit Gans (Tefnut-Legende)	ÄM 3317	2021-11-02 21:15:59	9	Ägyptisches Museum und Papyrussammlung	data/smb/resources/images/202009/200w_5f5f700958cb2.jpg	13
6	5		Figur der Göttin Bastet in Gestalt einer sitzenden Katze	ÄM 11385	2021-11-02 21:15:59	9	Ägyptisches Museum und Papyrussammlung	data/smb/resources/images/201806/200w_21081927984.jpg	13-
7	6		"Hälfte eines breiten Rings, darauf eine gelagerte säugende Hündin" (tatsächlich Teil eines ägyptischen Sistrums mit säugender Katze)	Misc. 8482	2021-11-02 21:15:59	10	Antikensammlung	data/smb/resources/images/201806/200w 27194138039.jpg	134
8	7	7012	Vierfüßiges Tier. Katze? (es handelt sich um einen geometrischen Löwen)	Misc. 7899	2021-11-02 21:15:59		Antikensammlung	data/smb/resources/images/201806/200w_27194155897.jpg	13-
9	8	230034	Amorette mit Katze	1928107	2021-11-02 21:15:59	14	Kunstgewerbemuseum	data/smb/resources/images/202009/200w_5f5ea0d3c91bf.jpg	13
10	9	63364	Mädchen mit Katze	60025-04.415	2021-11-02 21:15:59		Museum für Asiatische Kunst	data/smb/resources/images/201807/200w_15173931440.jpg	134
11	10			N (47 B) 3/2017,35	2021-11-02 21:15:59		Museum Europäischer Kulturen	data/smb/resources/images/201808/200w_04153954180.jpg	134
12	11	50739		I D 51881	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06165932553.jpg	13-
13	12		Katze mit Schellenbaum	I D 51948	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06171126839.jpg	13
14	13		Katze 猫	I D 50252	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06171645114.jpg	13
15	14	51441	Ema: Katze	I D 52073	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06171906607.jpg	13
16	15		Ema: Katze	I D 52049	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06171932133.jpg	13
17	16	51562	Ema: Katze	I D 52125	2021-11-02 21:15:59	11	Ethnologisches Museum	data/smb/resources/images/201807/200w_06172105136.jpg	13
18	17		Katzen	I D 52252 a,b	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06172711680.jpg	13
19	18		Katzen	I D 52251 a-c	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06172750863.jpg	13
20	19	51890		I D 52290	2021-11-02 21:15:59	11	Ethnologisches Museum	data/smb/resources/images/201807/200w_06172843317.jpg	13
21			Nikko "Nemuri-no-neko" "die schlafende Katze" vor dem Eingang zur Cryptomerien Allee, die zur						
	20		Grabstätte des lyeyasu führt. Nach Hidori Gingoro	VIII D 12502	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06172852239.jpg	13
22	21		Kauernde Katze 猫 Neko	I D 51748	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06165645371.jpg	13
23	22	50966		I D 51834	2021-11-02 21:15:59		Ethnologisches Museum	data/smb/resources/images/201807/200w_06170156251.jpg	13
24	23			KdZ 26603	2021-11-02 21:15:59		Kupferstichkabinett	data/smb/resources/images/201807/200w_28193302235.jpg	13
25	24	85574	St. Goar, Blick auf St. Goarshausen und Burg Katz	SM 9.13	2021-11-02 21:15:59	15	Kupferstichkabinett	data/smb/resources/images/201807/200w_28195922385.jpg	13

Figure 4.1: Beginn der Trefferliste für "katze" über die API der Staatlichen Museen zu Berlin

Um Abfragen zu vermeiden, die die Server überlasten, haben die meisten APIs entweder eine Authentifizierung oder eine maximale Trefferanzahl pro Abfrage eingebaut. Beim obigen Beispiel erhalten Sie dadurch nicht die gesamte Trefferanzahl (134, aus der Spalte "total" ersichtlich), sondern nur die ersten 24 – diese Einstellungen haben die Entwickler:innen der Schnittstelle gemacht. Um dennoch alle Treffer mit einer Abfrage zu erhalten, müssten Sie die Dokumentation der API lesen und die Abfrage etwas modifizieren. Wenn Sie das interessiert, finden Sie Details in der Fußnote.⁴

Wenn Webseiten keine Schnittstellen zur Verfügung stellen, besteht die Möglichkeit, mit Web Scraping an gewünschte Daten zu kommen. Je nach Webseite bzw. Inhalten ist die Rechtslage allerdings nicht ganz klar. Zum Download von Webseiten mit der Programmiersprache Python gibt es eine Lektion im Programming Historian von William J. Turkel und Adam Crymble. Ein weiteres Tutorial zur Datenakquise, von Zach Coble, Liz Rodrigues, Erin Pappas, Chelcie Rowell, und Yasmeen Shorish, findet sich hier.

⁴Die API aus dem Beispiel ist so konfiguriert, dass bei Abfragen mit Ergebnissen über 24 Treffern nur die ersten 24 ausgegeben werden; das ist etwas ungewöhnlich, aber wir können damit umgehen, indem wir die maximale Trefferausgabe pro Anfrage auf 10 setzen – diese Zahl ist nicht zu hoch, und wir können gut damit rechnen. Der Parameter für die maximale Trefferzahl kann mit &breitenat=10 eingestellt werden. Den Startpunkt der Ausgabe kann man mit dem Parameter &startwert= ändern. Um also alle Treffer für eine Abfrage zu erhalten, können wir die Ergebnisse in 10er-Schritten abfragen und anschließend zusammenfügen. Damit das nicht zu einer copy-paste-Aktion wird, müssen wir etwas ausführlicher formulieren bzw. mehrere Variablen verwenden. Das hat aber den Vorteil, dass man auf diese Art dann nach jedem Begriff suchen kann.

```
base_URL <- "https://smb.museum-digital.de/json/objects?&s=katze"</pre>
cats <- fromJSON(base URL)</pre>
start <- 0
breite <- 10
iterations <- catstotal[1]\%/\%10 + 1
endsize <- cats$total[1]-(iterations-1) * 10</pre>
cat_list <- data.frame()</pre>
for (i in 1:iterations){
  if(i < iterations){</pre>
cat_list <- rbind(cat_list, fromJSON(paste(base_URL, "&gbreitenat=10&startwert=", start , sep="")))</pre>
cat_list <- rbind(cat_list,fromJSON(paste(base_URL, "&gbreitenat=",</pre>
                                                   endsize, "&startwert=", start, sep="")))
  }
  start <- start + 10
  write.csv(cat_list, "Desktop/cat_list.csv")
     Zuerst machen wir den Code übersichtlicher und speichern den Großteil der URL in base URL:
   base_URL <- "https://smb.museum-digital.de/json/objects?&s=katze"
   Die Ergebnisse der Suchanfrage werden wieder im Objekt cats gespeichert:
   cats <- fromJSON(base_URL)</pre>
   Die Anzahl Durchgänge für eine Abfrage ergibt sich aus der Anzahl der totalen Treffer/10 + 1; die
   Anzahl der Treffer lässt sich aus der Spalte "total" im Objekt cats entnehmen. In R formuliert man das so:
   Im Katzenbeispiel sind es 134 Gesamttreffer, also (134/10 ohne Rest)+1, also 14 Durchgänge:
   iterations <- cats$total[1]%/%10 + 1</pre>
   Dann setzt man den Startwert auf 0:
   start <- 0
   Und die Maximaltreffer auf 10:
   breite <- 10
   Die letzte Iteration muss dabei nicht die nächsten 10 Treffer abfragen, sondern nur noch 4 (die letzten 4
   nach 130):
   endsize <- cats$total[1]-(iterations-1) * 10</pre>
   Dann erstellen wir eine leere Tabelle, einen data frame, die wir mit unseren Anfragen nach und nach
   befüllen. (Bei kleinen Datenmengen kann die Funktion rbind zur Verbindung von Einzeltabellen genutzt
   werden; bei größeren Datenmengen ist das iterative Verlängern von data frames nicht empfohlen):
   cat_list <- data.frame()</pre>
   Wenn wir diese Variablen festgelegt haben, können wir einen Loop, eine Schleife bauen, die unter bestimmten
   Bedingungen verschiedene Aktionen ausführt:
   for (i in 1:iterations){
   Falls die letzte Iteration noch nicht erreicht ist, wird die Abfrage in 10er-Schritten durchgeführt, wobei der
   Startwert bei jedem Durchgang um 10 verschoben wird und die Ergebnisse hintereinander in cat_list
   geschrieben werden.
   if(i < iterations){</pre>
   cat_list <- rbind(cat_list, fromJSON(paste(base_URL, "&gbreitenat=10&startwert=", start ,</pre>
   sep="")))
   } else {
   Sobald die letzte Iteration erreicht ist, werden nicht mehr die nächsten 10, sondern so viel Treffer, wie in
   endsize gespeichert, abgefragt, in unserem Beispiel 4:
   cat_list <- rbind(cat_list,fromJSON(paste(base_URL, "&gbreitenat=", endsize, "&startwert=",</pre>
```

start, sep="")))

4.2 Datenaufbereitung⁵

Bei der Arbeit mit Datensätzen, seien sie selbst erhoben oder von Dritten übernommen, ist es häufig der Fall, dass Informationen fehlen oder uneinheitlich erhoben wurden, was eine spätere Analyse erschwert.

Wenn in einer Umfrage unter Studierenden das Studienfach mit in eine Tabelle aufgenommen wurde, ohne zuvor Werte für diese Kategorie zu definieren, finden sich für "Geschichte" und "Deutsch" vielleicht auch folgende Varianten: "Gesch.", "Geschichtswissenschaft", "Geschichtswissenschaften", "Geschihcte", "Germanistik", "Dt.", "Germ.". Werten für zwei Studienfächer gibt es neun – ohne, dass sich das Fächerspektrum erweitert hätte. Im besten Fall werden solche Varianten schon bei der Erhebung der Daten vermieden, indem eine feste Liste an Werten erstellt wird. Erhält man jedoch einen Datensatz mit verschiedenen Varianten für ein und denselben Wert, muss man diese zusammenführen, um eine saubere Datengrundlage zu erhalten. Sie können entweder mit Strg-R versuchen, verschiedene Schreibweisen zu finden und zu ersetzen; in Tabellenprogrammen wie Excel, Open Office oder Google Sheets können Sie sich einzigartige Werte einzelner Spalten anzeigen lassen und zusammengehörende Varianten zu einem Grundwert zusammenführen; am hilfreichsten, recht voraussetzungslos zu bedienen und dabei auch für große Datensätze nutzbar ist die Software OpenRefine, mit der Sie Daten extrahieren,⁶ säubern/vereinheitlichen⁷ und anreichern⁸ können, um eine für Ihre Forschungsfrage und dafür notwendige Analysen sinnvolle Datengrundlage zu erhalten.

Für Textdaten sind verschiedene Schritte zur Aufbereitung notwendig, je nachdem, welche Methode bzw. Software Sie nutzen möchten. Für die meisten Analysen ist es sinnvoll, mit sogenannten Stopword-Listen zu arbeiten. Stopwords sind Wörter, die vor einer Analyse aus einem Korpus entfernt werden, um aussagekräftigere Ergebnisse zu erhalten, gerade, wenn es um rein quantitative Methoden zur inhaltlichen Erschließung geht. Stopwords sind Wörter mit grammatikalischen Funktionen, die in großer Zahl in Dokumenten vorkommen, jedoch wenig Bedeutung tragen. Wenn man den unbearbeiteten Text dieses Guides nach Worthäufigkeiten auswertet, hier mit Voyant-Tools, lässt sich nur schwerlich erahnen, worum es geht – "digital"

[}]start <- start + 10
Zum Schluss, in diesem Fall nach 14 Iterationen, wird die Tabelle in eine Datei geschrieben:
write.csv(cat_list, "Desktop/cat_list.csv")
}</pre>

⁵Eine häufige Aussage ist, zur Datenvorbereitung/Preprocessing würde 80% der Arbeitszeit verwendet, zur Analyse und Interpretation blieben nur 20%. In einem Blogartikel von 2020 geht Leigh Dodds diesen Zahlen nach – ganz so dramatisch ist das Verhältnis in Wahrheit wohl nicht.

⁶Evan Peter Williamson: Fetching and Parsing Data from the Web with OpenRefine, Programming Historian 6 (2017), https://doi.org/10.46430/phen0065.

⁷Seth van Hooland, Ruben Verborgh, Max De Wilde: Cleaning Data with OpenRefine, Programming Historian 2 (2013), https://doi.org/10.46430/phen0023.

⁸Karen Li-Lun Hwang: Enriching Reconciled Data with OpenRefine, The Bytegeist Blog 2018, https://medium.com/the-bytegeist-blog/enriching-reconciled-data-with-openrefine-89b885dcadbb

steht auf Platz 12, viel häufiger sind Artikel und Präpositionen. Mit Hilfe einer Stopword-Liste, die die häufigsten nicht-sinntragenden Wörter aus dem Text entfernt, wird der Inhalt klarer:

	Term	Count
1	und	191
2	die	170
3	https	150
4	in	119
5	der	101
6	sie	96
7	für	89
8	von	87
9	zu	83
10	mit	74
11	ist	72
12	digital	66
13	sich	64
14	data	61
15	oder	50
16	zur	49
17	eine	49
18	daten	49
19	ein	47
20	das	46
21	es	44
22	werden	42
23	den	37
24	auf	37
25	um	36

	Term	Count
1	https	150
2	digital	66
3	data	61
4	daten	49
5	history	35
6	wiki	32
7	doi.org	29
8	tools	24
9	online	23
10	en.wikipedia.org	23
11	quellen	21
12	command	21
13	chapter	21
14	shell	18
15	digitale	18
16	text	17
17	arbeit	17
18	analyse	16
19	terminal	15
20	line	15
21	interface	14
22	geschichte	14
23	forschung	14
24	ressourcen	13
25	literacy	13

(a) Worthäufigkeiten roher Text

(a) Worthäufigkeiten ohne Stopwords

Weitere Schritte beinhalten oft eine Tokenisierung, also die Segmentierung in Einheiten der Wortebene, und eine Lemmatisierung, also die Rückführung von verschiedenene Formen eines Worts auf eine Grundform – aus "ist", "war" und "sind" wird "sein". Wie bei den Schreibvarianten der Studienfächer haben die verschiedenen Flexionsformen für die meisten Forschungsfragen keinen Mehrwert und können zur weiteren Analyse zusammengeführt werden. Für solche vorbereitenden Schritte gibt es existierende Software und Packages für Programmiersprachen, sodass hier das Rad nicht neu erfunden werden muss, vor allem für moderne, weit verbreitete Sprachen, siehe auch Section B.3. Schwieriger wird es für nicht-standardisierte Sprachen

bzw. Sprachformen, also dialektal geprägte oder vormoderne Texte. Zwar gibt es auch hierfür Programme, die tatsächlich erreichte Präzision muss dabei jedoch je nach Quelle beurteilt werden.

4.3 Datenanalyse

Wenn Sie einen Datensatz zur Analyse zur Verfügung haben, aus selbst erhobenen Daten oder durch Nachnutzung eines vorhandenen, und für Ihre Zwecke aufbereitet haben, folgt (endlich) auch die Analyse. Welche Software oder Methoden Sie verwenden, hängt dabei nicht nur von der Art und Menge der Daten, sondern auch dem Datenformat und vor allem auch Ihrer Forschungsfrage ab. Wenn Sie eine Personendatenbank haben, in der Briefschreiber:innen und Empfänger:innen aufgenommen sind und der Wohnort der Personen bekannt ist, Sie es jedoch versäumt haben, die Datierungen der Einzelbriefe zu verzeichnen, können Sie nur eine räumliche Verteilung, keine raum-zeitliche Entwicklung eines Briefschreiber:innennetzwerks darstellen. Wenn Sie aber nur an der örtlichen Verteilung weiblicher und männlicher Verfasser:innen interessiert sind und die zeitliche Komponente für Sie keine Rolle spielt, erübrigt sich auch ein raum-zeitliche Analyse. Bevor Sie sich also für eine Methode entscheiden, sollten Sie sich fragen, zu welchem Zweck Sie Ihren Datensatz nutzen wollen und welche Frage(n) er beantworten soll.

In einem nächsten Schritt sollte über die konkrete Art der Analyse nachgedacht werden, die mit den vorhandenen Daten möglich ist. Unter den zahlreichen Möglichkeiten für die Arbeit mit strukturellen Daten sind für die Geschichtswissenschaften u.a. die Netzwerkanalyse oder die Regressionsanalyse häufig genutzte Methoden. Für textuelle Daten bieten sich ebenfalls verschiedene Arten der Analyse an, darunter beispielsweise Auszählungen von Worthäufigkeiten als Teil der Stylometrie/Zuschreibung von Autor:innenschaft (siehe Section 2.2), Topic Modeling als statistische Methode zur Identifizierung wiederkehrender Themen in größeren Textbeständen, oder Sentimentanalyse, um Stimmungen, Gefühle, Bewertungen aus Textpassagen zu extrahieren. Wenn Sie über georeferenzierte Daten verfügen, können Sie verschiedene Analysen mithilfe von GIS (Geographic Information System) durchführen und visualisieren.

Ob Sie für Topic Modeling ein eigenes Skript schreiben oder vorhandene Software nutzen, ob Sie Regressionsanalysen selbst durchführen oder auf Webseiten durchführen lassen, ist dabei Ihre Entscheidung; oftmals ist das Nutzen vorhandener Webangebote für erste kurze Analysen sinnvoll, um zu überlegen, ob die vorgesehene Methode überhaupt sinnvolle Ergebnisse liefern kann. Für größere Projekte, in denen komplexere Analysen über einen längeren

⁹Ein Großprojekt an der Universität Stanford, "Mapping the Republic of Letters", hat für das 18. Jahrhundert das Briefnetzwerk europäischer Gelehrter modelliert. Ein Fallbeispiel ist das Netzwerk Voltaires, in verschiedenen Visualisierungen: http://republicofletters.stanford.edu/publications/voltaire/letters/. Dan Edelstein. Interactive Visualization for Voltaire's Correspondence Network. Letters in Voltaire's Network [Created using Palladio, http://hdlab.stanford.edu/palladio].

Zeitraum durchgeführt werden sollen, bietet sich die Arbeit mit Programmiersprachen schon allein deswegen an, weil so ein sehr hohes Maß an Anpassungen von vorhandenen Funktionen für die eigenen Zwecke und die völlige Kontrolle über die eigenen Daten ermöglicht wird. Eine Auflistung häufig genutzter Tools für die historische Arbeit findet sich in Section B.3.

4.4 Datensicherung

In Kapitel 5 wird es um Fragen zur nachhaltigen Speicherung von Forschungsdaten gehen; an dieser Stelle sei darauf hingewiesen, dass die Sicherung von Daten am besten auch mit einer Versionierung und mit einer Dokumentation einhergeht. Datenversionierung hat den Vorteil, dass Schritte wieder rückgängig gemacht werden können, Datensätze in unterschiedlichen Stadien gespeichert und für eine spätere Weiterarbeit genutzt werden können und einzelne Schritte einzelnen Projektmitarbeiter:innen zugeschrieben werden können. Zusätzliche Versionierung geht dabei über die Funktionalitäten von Backup-Programmen oder Cloudspeichern wie Dropbox oder Switchdrive hinaus, und für Einzelprojekte wie auch für kollaboratives Arbeiten hat sich in der Wissenschaft wie in der Wirtschaft git etabliert, häufig in Kombination mit Daten-/Coderepositorien auf GitHub. Die meisten von Ihnen werden vermutlich keine eigenen GitHub-Repositorien anlegen, aber das System dennoch irgendwann nutzen, am ehesten durch den Download von dort zur Verfügung gestellten Daten – die Textdaten für diesen Guide liegen auch in einem GitHub-Repositorium. Die **Dokumentation** von gespeicherten Daten schließlich beinhaltet Informationen zur Entstehung des Datensatzes: Wie und von wem wurden die Daten erhoben? Wie wurden sie annotiert? In welchem Format sind die Daten vorhanden? Welche Software wurde an welcher Stelle benutzt? Was stellen die Daten dar? Die Sicherung von Daten an mehreren Orten, bspw. auf der lokalen Festplatte, in einem Cloudspeicher und auf einem USB-Stick, schützt sicher vor Datenverlust. Eine Dokumentation und die Sicherung in einem Repositorium, einem Langzeitspeicher für Daten, sorgt zusätzlich für Sichtbarkeit und die Möglichkeit zur Nachnutzung von Ergebnissen. Als Fachrepositorien für die Geisteswissenschaften existieren beispielsweise DARIAH-DE oder das DaSCH, es gibt spezialisiertere Repositorien wie AMAD (Mittelalter), oder für alle Disziplinen offene wie Zenodo (fächerübergreifend, betrieben durch das CERN). Sie können Ihre Forschungsdaten dort kostenfrei ablegen, Ihre Urheberschaft nachweisen und die Daten/Publikation mit einem Digital Object Identifier (DOI), also einem eindeutigen und dauerhaften digitalen Identifikator, nachhaltig zitierbar machen.

5 FAIR, CARE und LOUD

Bereits beim Beginn eines Projekts, sei es eine Proseminararbeit oder ein kollaboratives Großprojekt, sollten Fragen nach Sicherung, Austauschbarkeit und Nachnutzbarkeit von Forschungsdaten gestellt werden. Denn oftmals enden Projekte, ohne dass erstellte Daten für anschließende Forschungen verfügbar gemacht werden, sei es, weil nicht rechtzeitig nach Lösungen zur langfristigen Speicherung gesucht wurde, sei es, weil Daten in einer Form erhoben und gespeichert wurden, die eine Nachnutzung erschwert oder auch unmöglich macht. Spätestens beim ersten Gang ins Archiv – sei es analog oder digital –, bei dem Sie Quellen transkribieren, werden Sie sich vermutlich fragen, ob das nicht entweder schon von einer anderen Person erledigt wurde und Sie sich die Arbeit sparen könnten, oder aber ob Sie in Zukunft anderen die von Ihnen erstellten Transkriptionen zur Verfügung stellen wollen. Die Frage ist nur, wo und wie.

Zu Beginn des Studiums sind solche Überlegungen wohl noch nicht zentral; dennoch sollen einige Fragen rund um Speicherung, Aufbewahrung und Nachnutzbarkeit von Daten/Datenformaten hier kurz thematisiert werden, um dafür zu sensibilisieren; auch, weil sie den Prozess der Datenerhebung beeinflussen.

5.1 FAIRe Daten

Die Prinzipien FAIRer Daten wurden 2016 von einem Konsortium aus Wissenschaftler:innen und Organisationen wie folgt definiert:[10] Findability, Accessibility, Interoperability, Reuse of digital assets.

Daten sollen also auffindbar und zugänglich sein, zudem interoperabel, also mit verschiedenen Systemen nutzbar, und wiederverwendbar. Wenn Sie für eine Proseminararbeit zehn Testamente aus dem 18. Jahrhundert im Staatsarchiv Basel fotografieren, anschließend transkribieren, die vererbten Gegenstände identifizieren, zwischen den Erblasser:innen vergleichen und Ihre Ergebnisse ausgedruckt bei dem:r Dozierenden einreichen, sind Ihre Daten das genaue Gegenteil: Niemand weiß, dass Sie die Daten erhoben haben, sie sind über gängige Suchmethoden nicht auffindbar und nur über persönliche Kontakte zugänglich; und wenn Ihr:e Dozent:in Ihre Ergebnisse anderen Studierenden zur Verfügung stellen will, um weitere Forschung anzuregen, geht dies nur in Form von Kopien Ihrer gedruckten Arbeit; Papierkopien sind dabei weder interoperabel noch sind Ihre Daten vernünftig wiederverwendbar

– sie müssten via Abtippen erst wieder maschinenlesebar gemacht werden, um damit weiterarbeiten zu können. Wenn Sie Ihre transkribierten Texte und die identifizierten Objekte in Standardformaten und mit offener Lizenz auf einem Repositorium veröffentlichen, machen Sie nicht nur wichtige Teile Ihrer eigenen Arbeit sichtbar, sondern erleichtern so auch anschließende Forschungen.¹ Zudem kann so vermieden werden, dass geleistete Arbeit wie beispielsweise Transkriptionen doppelt gemacht wird.²

5.2 CARE-Prinzipien

Anschließend an die FAIR-Prinzipien wurden 2019 von der Global Indigenous Data Alliance die CARE-Prinzipien für den Umgang mit indigenen Daten formuliert:[2] Collective Benefit, Authority to Control, Responsibility, Ethics.

Das Augenmerk liegt dabei darauf, nicht einfach offene Daten und Datenaustausch zu propagieren, sondern auch die Menschen und den Zweck zu berücksichtigen, um bestehende Machtunterschiede zwischen verschiedenen Akteur:innen nicht zu verstärken. Indigene Daten sollen dem kollektiven Nutzen dienen, ein Recht auf Kontrolle soll gegeben sein, Verantwortung für die Datennutzung übernommen und Ethische Prinzipien beachtet werden.³ Auch wenn diese Richtlinien speziell für die Arbeit mit indigenen Daten ausgearbeitet wurden, ergänzen sie den datenzentrierten Ansatz der FAIR-Prinzipien um eine Dimension, die den Entstehungskontext der Daten mitberücksichtigt und zur Reflexion über die (Weiter-)Arbeit mit Daten anregt.

5.3 LO(U)D

Tim Berners-Lee, der Erfinder des World Wide Web, propagierte früh die Verlinkbarkeit von standardisierten digitalen Daten und dadurch die Entstehung eines Semantic Web, in dem Daten mensch- und maschinenlesbar austauschbar sind:

The Semantic Web isn't just about putting data on the web. It is about making links, so that a person or machine can explore the web of data. With linked data, when you have some of it, you can find other, related, data.⁴

¹Zu verschiedenen Repositorien, also digitalen Speicherorten für Daten, vgl. das Unterkapitel Section 4.4.

²Die Plattform transcriptiones ermöglicht niederschwelliges Teilen von Transkriptionen historischer Dokumente.

³Eine auf deutsch übersetzte Fassung der CARE-Prinzipien findet sich hier.

⁴Tim Berners-Lee: Linked Data, 2009. Online: https://www.w3.org/DesignIssues/LinkedData.html.

Um andere Daten nicht nur finden, sondern auch weiternutzen und kombinieren zu können, wurden auch für Linked Open Data Prinzipien formuliert, die bei der Erstellung bzw. Veröffentlichung von Daten beachtet werden sollten; bekannte Beispiele für LOD-Datensätze sind Wikibase, Wikidata oder GeoNames.

Berners-Lee hat ein Fünf-Sterne-System zur Klassifizierung offener Datensätze vorgeschlagen, wobei 5 Sterne Linked Open Data entspricht:

- 1 Stern: Ein Datensatz ist offen verfügbar, in einem beliebigen Format, etwa als PDF-Datei.
- 2 Sterne: Ein Datensatz ist offen verfügbar in einem strukturierten Format, etwa im Dateiformat von Microsoft Excel (.xls).
- 3 Sterne: Ein Datensatz ist offen verfügbar in einem nicht-proprietären strukturierten Format, etwa als Comma-separated values (.csv).
- 4 Sterne: Ein Datensatz folgt Standards des World Wide Web Konsortiums (W3C), wie die Nutzung des Resource Description Frameworks (RDF) und die Vergabe von Uniform Resource Identifiern (URIs).
- 5 Sterne: Ein Datensatz erfüllt alle obigen Bedingungen und enthält zudem Links zu anderer Linked Open Data.

Über das Fünf-Sterne-System hinaus wurden Grundsätze für Linked Open Usable Data entwickelt, die darauf abzielen, Daten nicht nur in offener und verlinkter Form zur Verfügung zu stellen, sondern sie auch verständlich zu strukturieren und zu dokumentieren, um ihre Nutzbarkeit zu verbessern.

Über Dinge wie das W3C, RDF und URIs müssen Sie sich im Verlauf Ihres Studiums vermutlich keine Gedanken machen, aber behalten Sie im Auge, dass Ihre Forschung(sdaten) sichtbarer und nachhaltig nachnutzbarer sind, wenn Sie nicht bereits bei der Erstellung einer Datei ein Format wählen, das nur auf einem bestimmten Betriebssystem läuft oder den Erwerb eines kostenpflichtigen Programms erfordert.

Part II

Praxis

Im Praxisteil sollen verschiedene Schritte rund um die **Datenerhebung** – wie komme ich von Quellen zu (strukturierten) Daten, und was ist das überhaupt? –, die **Datenaufbereitung** – wie kann bzw. muss ich die Daten für meine Zwecke bearbeiten – und die **Datenanalyse** – wozu sind die strukturierten Daten da, und was mache ich damit? – an einem kleinen Beispiel durchgeführt werden. Dabei werden viele Praktiken und Konzepte nur angeschnitten, die bei Interesse mithilfe weiterführender Literatur und Tutorials vertieft werden können. Manch eine:r wird Inhalte vermissen – Lücken sind unvermeidbar, aber Anregungen sind herzlich willkommen.

Als Quellenbeispiel dient eine Briefsammlung, die im Rahmen der digitalen Edition "Der Sturm" an der Akademie der Wissenschaften in Mainz erstellt wurde.⁵ Das Projekt, das Briefe von Personen der internationalen Avantgarde rund um die Zeitschrift "Der Sturm" aufbereitet, hatte dabei unterschiedliche Nutzer:innengruppen im Blick und stellt die Quellen bzw. die erstellten Daten auf verschiedene Weise zur Verfügung: Auf der Webseite des Projekts können die bisher edierten Briefe am Bildschirm gelesen werden, und es gibt ein Register der in den Texten genannten Entitäten, also Einheiten, hier der Personen, Orte und Werke; darüberhinaus gibt es die Möglichkeit, über eine Schnittstelle die Quellen oder die Registerdaten herunterzuladen. Wir können uns den Briefen also über den Vordereingang, die Webseite, oder durch die Hintertür, mittels Kommandozeile, nähern, und beide Herangehensweisen vergleichen. Grundlegende Konzepte für die Arbeit mit Daten, bzw. für die Schritte von der Quelle zum Datensatz, werden angesprochen.

Briefe sind eine gängige Quellengattung in den verschiedenen Epochen/Areas, und sie können sowohl für Textanalysen als auch zur Gewinnung von Strukturdaten genutzt werden. Das Beispiel zielt also auch auf eine größtmögliche Nähe zur tatsächlichen historischen Arbeit.

⁵DER STURM. Digitale Quellenedition zur Geschichte der internationalen Avantgarde, erarbeitet und herausgegeben von Marjam Trautmann und Torsten Schrade. Mainz, Akademie der Wissenschaften und der Literatur, 2018. Online: https://sturm-edition.de/, Lizenz: CC-BY-4.0.

6 Briefedition 'Der Sturm'

Die Webseite https://sturm-edition.de/ dient als Portal für die Arbeit mit Quellen zum avantgardistischen STURM-Unternehmen, das mit der Gründung der gleichnamigen Zeitschrift in Berlin im Jahr 1910 durch Herwarth Walden begann und mit den nachfolgenden Gründungen einer Galerie, einer Bühne und eines Verlags internationale Bedeutung erlangte. Neben Walden waren weitere Akteur:innen an diesem Unternehmen beteiligt, und bisher (Stand Herbst 2024) sind 179 Briefe von drei Künstler:innen an Walden im Portal verfügbar. Zudem wurden die Zeitschrift, Ausstellungskataloge, Jahrbücher, Verlagsschriften und weitere Materialien wie Plakate, Fotografien oder Einladungskarten digitalisiert.

Gehen Sie auf die Startseite und lesen Sie die Kurzbeschreibung zu den Briefen, Personen, Orten und Werken, um eine erste Idee vom Material zu bekommen. Klicken Sie dann die Briefabteilung an und schauen sich Brief Nummer 8 von Franz Marc an Herwarth Walden an.

Wie Sie sehen können, wurde der Brief nicht nur digitalisiert, also in ein digitales Bild umgewandelt, das Sie über einen externen Viewer betrachten können, sondern auch historisch-kritisch ediert und nach den Richtlinien der TEI P5 in XML codiert.

Was bedeutet das genau, und wieso ist es wichtig für unsere Arbeit als Historiker:innen?

Bei einer historisch-kritischen Edition wird keine reine Leseversion eines Textes erstellt, wie man es beispielsweise von der Textdarstellung in einem Roman kennt, sondern zusätzliches Material zum besseren Verständnis des Texts herangezogen und als Information in einem sog. kritischen Apparat zur Verfügung gestellt. Beispielsweise werden Quellen, die ein:e Verfasser:in für bestimmte Textpassagen als Vorlage genutzt hat, genannt, oder auf Ereignisse während der Entstehungszeit, die Einfluss auf den Text hatten, verwiesen. Entitäten wie Personen, Orte, Werke usw. werden erklärt oder kommentiert. Der Text selbst wird so quellengetreu wie möglich dargestellt, orthographische oder grammatikalische Fehler werden nicht korrigiert, und meist werden auch extratextuelle Elemente wie Durchstreichungen oder Hervorhebungen im Lesetext dargestellt.

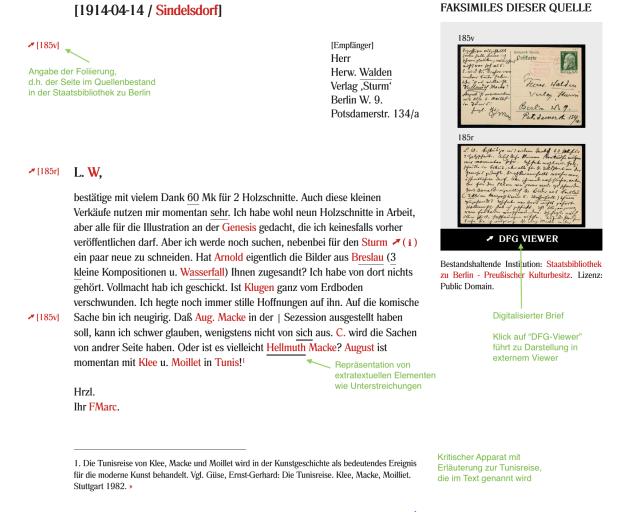


Figure 6.1: Brief Nr. 8 von Franz Marc an Herwarth Walden. Rote Einfärbung kennzeichnet eine Verlinkung, meist auf Entitäten wie Personen oder Orte.

Wie genau die Editor:innen jeweils vorgegangen sind – dies unterscheidet sich von Edition zu Edition –, wird in den jeweiligen Editionsrichtlinien vermerkt. Für die Quellen des STURM-Projekts gibt es verschiedene Richtlinien, je nach Quellenart. Die der Briefedition finden Sie hier.

Eine Codierung in XML bedeutet, dass eine Textdatei mit der Extensible Markup Language ausgezeichnet wurde, d.h. Strukturen im Text werden durch festgelegte Zeichen so markiert, dass sie sowohl von Menschen als auch von Computern interpretiert werden können. Etwas

 $^{^1}$ Online unter https://sturm-edition.de/quellen/briefe/fma/Q.01.19140414.FMA.01.html

geläufiger ist Ihnen vermutlich HTML, **H**ypertext **M**arkup **L**anguage, eine Sprache zur Auszeichnung elektronischer Dokumente. Die Prinzipien sind dabei ähnlich: Festgelegte Elemente werden durch einen öffnenden und einen schließenden Tag ausgezeichnet, z.B.

<salute>Hrzl.</salute>,

um die Grußformel in einem Brief (wie in Brief Nummer 8, siehe oben) als solche zu markieren. Das hat den Vorteil, dass bei der Suche nach Grußformeln nicht konkrete Begriffe formuliert werden müssen, sondern nach dem entsprechenden Element <salute> gesucht werden kann, und zwar sowohl vom Menschen als auch von der Maschine. Darauf kommen wir nochmals zurück. Zunächst betrachten wir die Briefe auf der Webseite.

7 Durch den Vordereingang

Stellen Sie sich vor, Sie würden für ein Forschungsprojekt unter anderem die Briefe von Franz Marc an Herwarth Walden untersuchen wollen. Sie wären bei Ihrer Untersuchung daran interessiert, welche Grußformeln Marc in seinen Briefen an Walden benutzte, um aus eventuellen Änderungen Rückschlüsse auf das Verhältnis der beiden ziehen zu können. Wie würden Sie vorgehen, um sich eine Übersicht zu verschaffen? Würden Sie eine Downloadmöglichkeit der Webseite nutzen oder die Briefe online lesen? Welche Schritte würden Sie durchführen, um die Texte aller Briefe von Marc und Walden auf Ihrem Computer zu speichern? Notieren Sie Ihr Vorgehen zu den Fragen in einem Dokument.

Es gibt zahlreiche Wege, die zum Ziel führen, und keiner ist dabei besser oder schlechter; aber manche sind möglicherweise effizienter als andere – das heißt, Sie sparen Zeit, die Sie für andere Dinge verwenden können, sei es fürs Studium oder in der Badi.

Egal, wie Sie vorgegangen sind, erledigen Sie folgende Aufgaben:

- 1. Erstellen Sie eine Übersicht über alle Grußformeln in den 45 Briefen von Franz Marc an Herwarth Walden. Welches Format analog oder digital bzw. welchen Dateityp Sie wählen, bleibt Ihnen überlassen.
- 2. Nehmen Sie Brief Nr. 8 von Marc an Walden und unterteilen Sie ihn in strukturelle Elemente

Ob Sie den Brief hierfür ausdrucken und Strukturelemente mit einem Stift markieren oder ob Sie am Computer arbeiten, können Sie wählen. Was Sie als Strukturelement verstehen, ist Teil der Aufgabe.

(Bei einem Gedicht wären mögliche Strukturelemente etwa Titel und Strophe.)

3. Nehmen Sie denselben Brief und markieren Sie Entitäten, also Einheiten wie Personen, Orte, etc. Erstellen Sie hierfür ein Tabellendokument (mit Excel, Google Sheets, Open Office o.Ä.), in das Sie die Entitäten aufnehmen.

Nachdem Sie die Aufgaben erledigt haben, lesen Sie "The Ten Commandments of Inputting Data" in Kapitel 3 im Buch "Quantitative Methods in the Humanities".[6, S. 57–60] Das Kapitel finden Sie im ADAM-Workspace Ihres Einführungskurses – es lohnt sich als Ganzes, aber die "Ten Commandments" reichen auch. Eine eher praxisorientierte Onlinefassung finden Sie hier im das Buch begleitenden Blog.

Würden Sie nach der Lektüre bei Aufgabe 3 anders vorgehen? Notieren Sie sich etwaige Änderungen bzw. Erkenntnisse und bringen Sie Ihre Resultate zur begleitenden Sitzung mit, in analoger oder digitaler Form.

8 Durch die Hintertür

Bevor wir uns näher mit dem Beispielkorpus der Briefedition befassen, werfen wir einen kurzen Blick auf die Interaktionsmöglichkeiten mit dem Computer und wie wir diese für unsere Arbeit als Historiker:innen nutzen können, sei es für die Erhebung, die Aufbereitung oder die Analyse von Daten.

Es gibt zwei Arten, um mit einem Computer zu interagieren bzw. ihn zu nutzen: über ein Graphical User Interface (GUI), also vor allem mit der Maus und durch das Anklicken von Objekten, oder, etwas direkter, über die Kommandozeile.¹ Um via GUI eine Datei "Brief1.txt" im Ordner "Briefe" zu löschen, öffnet man den Finder (Mac), den Explorer (Windows) oder den Filebrowser der Wahl (Linux), klickt sich zum Ordner "Briefe", macht einen Rechtsklick auf die zu löschende Datei "Brief1.txt", klickt "In den Papierkorb legen" oder zieht die Datei mit der Maus direkt dorthin. Dieselbe Aktion kann man als Kommando schreiben: Man öffnet das Terminal (Linux oder Mac; den Finder öffnen und im Suchfenster "Terminal" eingeben und Programm öffnen) oder eine PowerShell (Windows; mit der rechten Maustaste auf das Startsymbol klicken, dann "Windows PowerShell" auswählen), navigiert im sich öffnenden Fenster mit Texteingabe zum entsprechenden Ordner, bspw. cd Documents/Briefe + 'Enter' (Mac und Linux) bzw. cd ./Documents/Briefe (Windows) und gibt dort das Kommando rm Brief1.txt ein, das mit der Entertaste ausgeführt wird.

(base) serina00@dg-19-mac-02 ~ % cd Documents/Briefe

(base) serina00@dg-19-mac-02 Bilder % rm Brief1.txt

Vorgehen in der Kommandozeile bzw. im Terminal auf MacOS

Die beiden Vorgehensweisen unterscheiden sich dabei in drei Punkten:

1. Das Kommando rm ist endgültig, die Datei ist ohne Übergangszeit im Papierkorb gelöscht.

¹Kommandozeile/Command Line, Bash, Shell oder Prompt finden sich oft als synonym genutzte Begriffe für Command Line Interfaces. Auf UNIX-basierten Betriebssystemen wie Mac OS und Linux ist das Terminal als Interface weit verbreitet; für Details: https://en.wikipedia.org/wiki/Command-line_interface#History. Windowsnutzer:innen kommen mit der PowerShell ganz gut zurecht, es empfiehlt sich eventuell die Installation von Cygwin oder MinGW, um mit einem UNIX-basierten Interface arbeiten zu können.

- 2. Das Kommando lässt sich relativ simpel auf eine Vielzahl von Dokumenten anwenden, wobei ganz unterschiedliche Bedingungen beachtet werden können, und es lässt sich mit anderen Kommandos verbinden.
- 3. Terminal sight k3wl aus.

Bevor wir den zweiten – und für unsere Arbeit hilfreichsten – Unterschied genauer anschauen, kurz zur Kommandozeile.

8.1 Shell 101

In einem Terminal/einer Shell – zur Unterscheidung siehe Fußnote 1 – können Kommandos bzw. Programme ausgeführt werden, die auf der *Struktur*ebene stattfinden – wie beispielsweise das Löschen einer Datei, rm Dateiname.xyz (rm für remove), oder das Erstellen eines Ordners, mkdir NeuerOrdner (mkdir für make directory). Ebenso möglich sind Operationen auf der *Inhalts*ebene – wie beispielsweise das Suchen eines bestimmten Begriffs in einer Textdatei, grep 'Begriff' Textdatei.txt (Mac/Linux) bzw. Select-String -Path Textdatei.txt -Pattern 'Begriff' (Windows), oder das Auszählen mehrerer Begriffe und das Speichern des Ergebnisses in einer neuen Datei, grep -Ec '(Begriff1|Begriff2)' Textdatei.txt | wc -1 > Ergebnisse.txt (Mac/Linux) bzw. (Select-String -Path Textdatei.txt -Pattern '(Begriff1|Begriff2)'.Matches.Count > Ergebnisse.txt (Windows) – die Kommandos werden weiter unten nochmals einzeln erklärt.

Woher weiss Ihre Shell aber, was sie ausführen soll, wenn Sie rm oder grep/String-Select eintippen? Es gibt zahlreiche Shell-Programme, die bereits auf Ihrem System vorinstalliert sind, und mit denen Sie vieles tun können – öffnen Sie Ihre Shell, tippen Sie date ein und drücken 'Enter': Das aktuelle Datum mit Uhrzeit erscheint. (Ihre Shell sucht nach dem ersten Argument, dem Befehl date, im Filesystem des Computers, und wenn sie fündig wird, führt sie eine Aktion mit den entsprechenden Parametern aus.)

Note

tmi: Wenn Sie echo \$PATH im Terminal (Mac/Linux) bzw. \$env:PATH (Windows) eingeben, sehen Sie eine Auflistung der Orte, an denen nach Befehlen gesucht wird. Tippen Sie which date ein und drücken Sie 'Enter', um zu sehen, wo das Programm "date" in Ihrem Computer liegt.

Falls Sie einen Befehl eintippen, den es nicht gibt bzw. für den kein installiertes Programm auf Ihrem Computer existiert, bekommen Sie eine simple Fehlermeldung – kaputtgehen kann dabei nichts:

(base) serina00@dg-19-mac-02 ~ % nonsense

command not found: nonsense

Der Output auf Windows ist etwas ausführlicher:

nonsense: The term 'nonsense' is not recognized as a name of a cmdlet, function, script file Check the spelling of the name, or if a path was included, verify that the path is correct as

Das aktuelle Datum wird Ihnen wahrscheinlich auch in Ihrer Toolbar angezeigt, und einen neuen Ordner können Sie per Rechtsklick erstellen, dazu brauchen Sie das Terminal nicht unbedingt. Um einen Begriff in einem Textdokument zu finden und alle Vorkommen zu zählen, können Sie das Dokument öffnen, Strg-F drücken, den Begriff eingeben und das Ergebnis sehen. Wenn Sie nach mehreren Begriffen suchen wollen, müssen Sie dieselbe Aktion zweimal ausführen: Strg-F, Begriff 2. Und wenn Sie mehrere Dateien durchsuchen möchten, beispielsweise um herauszufinden, wie oft die Grußformel "Mit herzlichem Gruß" in einer Briefsammlung vorkommt, müssen Sie die Suche in jeder Datei einzeln ausführen. Wenn Sie dann noch nach der Variante "Mit herzlichen Grüßen" oder gar "Herzl. Gruß" suchen wollen, vervielfacht sich Ihre Arbeit.

Sie können dasselbe auch mit dem Terminal machen und einige der Built-in-Programme nutzen, um sich Zeit und Arbeit zu sparen.

8.2 Strg-F 2.0

Wir arbeiten wie im vorangegangenen Kapitel mit einem Teilkorpus der Quellenedition "Der Sturm", nämlich mit allen Briefen, die von Franz Marc verfasst wurden. Um die folgenden Schritte nachzuvollziehen, laden Sie sich den Ordner "letters_Der_Sturm" herunter. Sie können dazu entweder das vollständige GitHub-Repository zu diesem Guide als Zip-Datei herunterladen und entpacken, und im Ordner "docs" befindet sich der Ordner "letters_Der_Sturm".

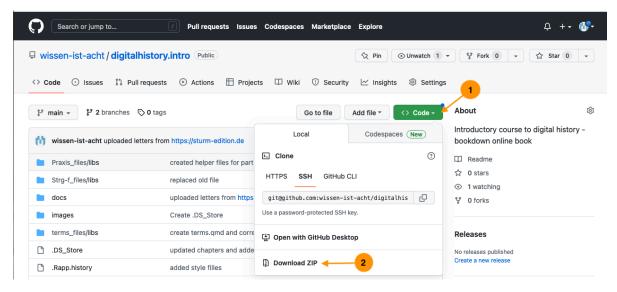


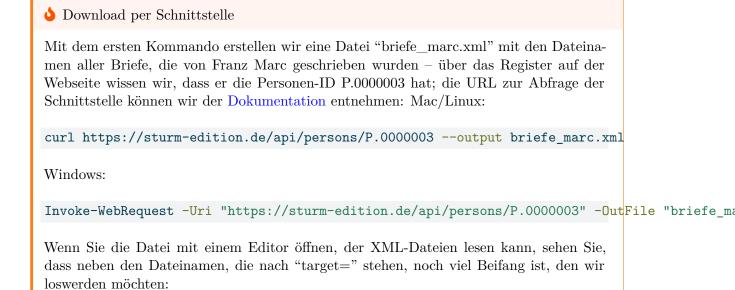
Figure 8.1: GitHub-Repositorium mit Quellcode, Download via "Code" und "Download ZIP".

Sie können das Repositorium auch über die Kommandozeile klonen

```
git clone https://github.com/wissen-ist-acht/digitalhistory.intro.git
```

oder als bequeme Variante diesen Direktlink nutzen.

Wenn Sie die Schnittstelle der Webseite ausprobieren wollen, kommen Sie mit wenigen Kommandos an die Dateien.



Denn eigentlich brauchen wir nur die Dateinamen, um die Dateien mit einem entsprechenden Befehl herunterladen zu können. Auch sehen wir, dass nicht nur Briefe mit dem Kürzel "FMA" für Franz Marc gelistet sind, sondern auch neun mit "JVH", Jacoba van Heemskerck. Beim Blick in die betreffenden Dateien zeigt sich, dass in diesen Franz Marc jeweils erwähnt und im TEI-XML mit persName key="P.0000003" getaggt wurde, daher Teil des Abfrageergebnisses ist. Mit dem zweiten Kommando erstellen wir eine neue Datei, in der die einzelnen extrahierten Dateinamen ohne diejenigen von Jacoba van Heemskerck mit dem Download-Kommando curl kombiniert und um die entsprechende URL zum Download ergänzt werden:

Mac/Linux:

```
curl -o Q.01.19160128.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19160128.FMA.01.xm curl -o Q.01.19160205.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19160205.FMA.01.xm curl -o Q.01.19160302.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19160302.FMA.01.xm curl -o Q.01.19160101.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19160101.FMA.01.xm curl -o Q.01.19160122.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19160122.FMA.01.xm curl -o Q.01.19160115.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19150703.FMA.01.xm curl -o Q.01.19150703.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19150417.FMA.01.xm curl -o Q.01.19150417.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19151106.FMA.01.xm curl -o Q.01.19150918.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19150918.FMA.01.xm curl -o Q.01.19150918.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19150918.FMA.01.xm curl -o Q.01.19150827.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19150827.FMA.01.xm curl -o Q.01.19150303.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19150827.FMA.01.xm curl -o Q.01.19150303.FMA.01.xml https://sturm-edition.de/api/files/Q.01.19150303.FMA.01.xm curl -o Q.01.19150303.
```

cat briefe_marc.xml (Mac/Linux)/-Path "briefe_marc.xml" (Windows) gibt den Inhalt der Datei ins Terminal; grep -o 'Q.*FMA.01.xml\b' (Mac/Linux) bzw. Select-String -Pattern 'Q.*FMA.*\.xml\b' (Windows) findet in diesem Inhalt alle Zeichenketten zwischen "Q" und "FMA.01.xml", wobei nach "xml" durch das Hinzufügen von \b das Zeichenende angezeigt ist; die 45 gefundenen Zeichenketten werden in je eine neue Zeile geschrieben, wobei mit curl -o \$_ der Befehl curl -o und mit \$_ (Mac/Linux) bzw. \$match (Windows) als Platzhalter die Zeichenkette (also der Dateiname) geschrieben wird, gefolgt von https://sturm-edition.de/api/files/\$_ (Mac/Linux) bzw. https://sturm-edition.de/api/files/\$match und mit \$_ bzw. \$match wieder die Zeichenkette (also wieder der Dateiname). Mit einem dritten Kommando, bash, führen wir die erstellte Datei aus, d.h. die in ihr stehenden Kommandos werden ausgeführt - also via curl (Client URL) die Briefe heruntergeladen.

```
bash dateinamen_briefe_marc.txt
```

Egal, wie Sie die Dateien heruntergeladen haben, sollten Sie 45 Briefe im xml-Format vorfinden. Öffnen Sie dann das Terminal (Mac/Linux) bzw. die PowerShell (Windows) und bewegen sich mit cd, also change directory, in den Ordner (directory), in dem Ihre Textdateien liegen. In meinem Fall ist das unter Documents/GitHub/digital_history_intro/docs/letters_Der_Sturm.

```
(base) serina00@dg-19-mac-02 ~ % cd Documents/GitHub/digital_history_intro/docs/letters_Der_s
```

Bei den meisten von Ihnen ist das vermutlich unter "Downloads" – probieren Sie es aus.

(Um zu prüfen, was in einem Ordner liegt, können Sie im Terminal 1s (für 1ist) eingeben, bzw. in der PowerShell dir (für directory)):

ls

```
Q.01.19140115.FMA.01.xml
                            Q.01.19150315.FMA.02.xml
Q.01.19140119.FMA.01.xml
                            Q.01.19150327.FMA.01.xml
Q.01.19140121.FMA.01.xml
                            Q.01.19150417.FMA.01.xml
Q.01.19140124.FMA.01.xml
                            Q.01.19150501.FMA.01.xml
Q.01.19140125.FMA.01.xml
                            Q.01.19150615.FMA.01.xml
Q.01.19140125.FMA.02.xml
                            Q.01.19150703.FMA.01.xml
Q.01.19140409.FMA.01.xml
                            Q.01.19150710.FMA.01.xml
Q.01.19140414.FMA.01.xml
                          Q.01.19140421.FMA.01.xml
Q.01.19150827.FMA.01.xml
                          Q.01.19140507.FMA.01.xml
```

8.3 Erste Schritte

Wenn Sie in den Ordner navigiert sind, in dem die Briefdateien liegen, können Sie mit einem einzeiligen Kommando die Suchvorgänge nach einer Grußformel, die Sie in einem Texteditor nacheinander mit Strg-F mit jeder einzelnen Datei ausführen würden, im Terminal/in der PowerShell mit dem Programm grep (Global Regular Expression Print, Mac/Linux) bzw. Select.String (Windows) für alle Briefe in diesem Ordner auf einmal vornehmen, indem Sie alle Dateien, die auf ".xml" enden, in die Suche aufnehmen. Die Ergebnisse – bei dieser Suche nach den Grußformeln "Mit herzlichem Gruß" oder "Mit herzlichen Grüßen" ein Treffer in einem Brief – können Sie sich im Terminal/in der PowerShell anschauen:

Mac/Linux:

```
grep -E -i '(Mit herzlichem Gruß|Mit herzlichen Grüßen)' *.xml
```

Windows:

```
Select-String -Path *.xml -Pattern "(Mit herzlichem Gruß|Mit herzlichen Grüßen)"
```

Output:

```
Q.01.19160115.FMA.01.xml: <salute>Mit herzlichen Grüßen für Sie beide</salute>
```

Die Formulierung "Mit herzlichen Grüßen" kommt also einmal im Korpus vor, und zwar im Dokument Q.01.19160115.FMA.01.xml.

Sie können auch mit wc -1 (Mac/Linux) den Wordcount, die Anzahl der gefundenen Treffer auf Zeilenebene, -1 zählen, bzw. mit Matches.Count (Windows), und mit > in eine neue Datei schreiben (die während der Ausführung des Kommandos erstellt wird):

Mac/Linux:

Windows:

```
(Select-String -Path *.xml -Pattern "(Mit herzlichem Gruß|Mit herzlichen Grüßen)").Matches.Co
```

Wenn Sie die neu erstellte Datei count_greetings.txt öffnen, die sich im selben Ordner wie die Briefe befindet, sollte dort "1" stehen, weil unsere Suche einen Treffer ergeben hat.

Das Kommando grep (Mac/Linux) hat im obigen Befehl den Zusatzparameter E bekommen, das Kommando Select.String (Windows) den Parameter -Pattern, d.h. wir suchen nicht eine exakte Zeichenkette, sondern nutzen Möglichkeiten zur Mustersuche, zur Suche nach Patterns. Diese werden formuliert als sog. Extended Regular Expressions (von hier kommt das E), als reguläre Ausdrücke. Wir haben in unserer Suchabfrage nämlich nicht nur nach "Mit herzlichem Gruß" gesucht, sondern auch nach "Mit herzlichen Grüßen", formuliert mit dem Zeichen "|", hier als "oder" zu lesen. Mithilfe Regular Expressions können wir unsere Suche weiter ausbauen und nach verschiedenen Varianten/Schreibweisen auf einmal suchen.

Note

Regular Expressions haben verschiedene flavours – je nach Programmiersprache werden Dinge etwas anders formuliert, und manche Defaulteinstellungen unterscheiden sich. In unserem Fall benötigt grep noch den Parameter –i, um Groß- und Kleinschreibung zu ignorieren. Select.String ignoriert dies by default und braucht keinen zusätzlichen Parameter. Solche Feinheiten sind bei der Arbeit mit Regular Expressions wichtig zu wissen, aber das lernt man on the go.

Mac/Linux:

```
grep -E -i '(Mit herzlichem Gru(ß|ss)|Mit herzlichen Grü(ß|ss)en|H(e|.?)rzl. Gru(ß|ss))' *.x
```

Windows:

```
(Select-String -Path *.xml -Pattern "(Mit herzlichem Gru(fiss)|Mit herzlichen Grü(fiss)en|H(
```

So formuliert finden wir 17 Treffer für eine Grußformel, mit den möglichen Schreibweisen "Mit herzlichem Gruß", "Mit herzlichem Grüßen", "Mit herzlichen Grüßen", "Mit herzlichen Grüßen", "Herzl. Gruß", "Herzl. Gruß", "Hrzl. Gruß", "Hrzl. Gruß".

Wenn wir herausfinden möchten, ob Grüße mal herzlich, mal hrzl. oder freundlich verschickt wurden, können wir die Suche und die Art der Ausgabe modifizieren:

Mac/Linux:

```
grep -E -i 'Gr(u|ü)(ß|ss)' *.xml
```

Windows:

```
Select-String -Path *.xml -Pattern "Gr(u|ü)(ß|ss)"
```

Output:

```
Q.01.19140115.FMA.01.xml:
                                               stets sofort antworte; es muß verloren gegangen
Q.01.19140119.FMA.01.xml:
                                               <salute>Hrzl. Gruß</salute> <signed>Ihr <persNa</pre>
Q.01.19140125.FMA.02.xml:
                                               <salute>Hrzl. Gruß</salute>
Q.01.19140421.FMA.01.xml:
                                                   <closer>Gute Besserung persName key="P.000
                                               <salute>besten Gruß</salute>
Q.01.19140507.FMA.01.xml:
Q.01.19140831.FMA.01.xml:
                                               <salute>Hrzl. Gruß von Eurem Freund in Waffen/:
                                               <salute>Hrzl. Gruß 1 x 2</salute> <signed>Ihr <</pre>
Q.01.19141113.FMA.01.xml:
                                               <salute>Hrzl. Gruß Ihnen beiden</salute>
Q.01.19150112.FMA.01.xml:
Q.01.19150116.FMA.01.xml:
                                               <salute>Mit herzl. Gruß Ihnen beiden</salute> <</pre>
                                               <salute>Herzl. Gruß</salute> <signed>Ihr <persN</pre>
Q.01.19150121.FMA.01.xml:
```

Mit diesem Kommando durchsuchen wir also den Text nach dem Muster $Gr(u|\ddot{u})(\beta|ss)$, also Beginn mit Gr oder gr, dann folgt entweder ein u oder ein \ddot{u} , dann entweder ein β oder ss. Weil wir kein Wortende markiert haben (das ginge mit \b), werden auch "Grüße" oder "Grüssen" gefunden.

Wenn Sie sich während der Lektüre des vorangegangenen Kapitels auf der Webseite durch die Briefe geklickt haben, werden Sie festgestellt haben, dass ein Brief nicht immer mit "Gruß" oder "Grüßen" endet. Beim Output der Suchanfragen im Terminal sehen Sie, dass alle Grußformeln von einem Tag-Paar umgeben sind: <salute> kennzeichnet den Beginn des Grußes, </salute> das Ende. Öffnen Sie eine der Briefdateien und suchen Sie nach "salute". (Wenn Sie keinen XML-fähigen Editor auf dem Computer haben, öffnen Sie die Datei einfach mit einem Browser.)

```
112
                      <opener>
                          <salute>Lieber <persName key="P.0000001" ref="http://d-</pre>
113
                          nb.info/gnd/118770950">Walden</persName>,</salute>
114
                      ich vermute, daß <persName key="P.0000058" ref="http://d-
115 ▼
                      nb.info/gnd/118870645">Cassirer</persName> es auf den Proceß ankommen läßt; ich
                      rate nicht dazu:
116
                          Sie wi<hi rend="underline">sse</hi>n ja meine Bedingung, daß ich auf gar
117
                          Fall irgend welche Kosten von der Sache haben möchte. Das ist mir weder
                          <persName key="P.0000058" ref="http://d-</pre>
118
                          nb.info/gnd/118870645">Cassirer</persName> noch <term type="journal"</pre>
                          key="W.0000070" ref="http://d-nb.info/gnd/4127687-5">die
119
                              Aktion</term> wert.
120
                      Von <persName key="P.0000057" ref="http://d-
                      nb.info/gnd/143669230">Reiche</persName> resp. <persName key="P.0000059"
                      ref="http://d-nb.info/gnd/2097840-6">Arnold</persName> hab ich illustr. Katalog
                      erhalten; Sie werden ihn wohl
121
                         auch inzwischen gesehen haben; sonst kann ich Ihnen den meinen schicken.
                      Von
122 ▼
                          <persName key="P.0000033" ref="http://d-</pre>
123
                          nb.info/gnd/118969161">Filla</persName> hab ich die <term type="artwork"</pre>
                          key="W.0000011">
124
                              <hi rend="underline">Häringe</hi>
125
                          </term> gewählt, von
                              <persName key="P.0000060" ref="P.0000060">Beneé</persName> ein größeres
126
                              Aquarell, ich glaube mit blau (hell,
127
                          geometrisch) kein, hochformat.<pb xml:id="S.178v.02" n="178v"</pre>
                          facs="http://resolver.staatsbibliothek-berlin.de/SBB0000DAA400000001"/> Auf
                          beide Bilder haben sie
128
                          auf der Rückseite mit Blei "Marc" geschrieben. Sie werden es schon finden.
                          Sobald ich die Sachen habe, sende ich 2 Aquarelle; geben Sie mir bitte die
129
130
                          Adressen.
131 ▼
                          <salute>Hrzl. Gruß</salute> <signed>Ihr <persName key="P.00000003"</pre>
132
                          ref="http://d-nb.info/gnd/11857745X">F. Marc</persName>
133
134
                      </closer>
```

Figure 8.2: Ausschnitt aus Brief Nr. 1 von Franz Marc an Herwarth Walden

Wie Sie sehen, gibt es das Tag-Paar <salute>-</salute> zweimal, einmal umrahmt vom Tag-Paar <opener>-</opener>, einmal von <closer>-</closer>. Die Anrede ist mit dem ersten, die Grußformel mit dem zweiten Tag-Paar markiert. Wir können also, wenn wir mit Dokumenten arbeiten, die nach festgelegten Richtlinien ausgezeichnet wurden, nach dem Element Grußformel suchen, ohne erst einen Blick in die Texte werfen zu müssen, um verschiedene Suchabfragen zu formulieren. Wir formulieren unsere Suchabfrage um und suchen nun nach einer Abfolge von Zeichen mit dem Beginn <closer>, gefolgt von keinem bis zu beliebig vielen (.*) Zeichen der Klasse cntrl, also nicht sichtbare Zeichen wie Tabs, Seiten- oder Zeilenumbruch. Danach folgt <salute>, wiederum gefolgt von keinem bis zu beliebig vielen (.*) Zeichen, keinem bis zu beliebig vielen (.*) Zeichen der Klasse cntrl und nochmal keinem bis zu beliebig vielen (.*) Zeichen, bis der Beginn des Schlusstags zu </salute> kommt. Damit werden die verschiedenen Fälle in den Briefen abgedeckt, dass zwischen <closer> und <salute> Text oder ein Zeilenumbruch stehen kann oder auch nicht, und dass zwischen <salute> und </salute> Text, kein Text oder ein Zeilenumbruch kommen kann.

Mac/Linux:

```
grep -E -zo '<closer>[[:cntrl:]].*<salute>.*[[:cntrl:]].*<' *.xml</pre>
```

Output:

```
Q.01.19140115.FMA.01.xml:<closer>
                     <salute>Hrzl.
Q.01.19140119.FMA.01.xml:<closer>
                     <salute>Hrzl. Gruß</salute> <signed>Ihr <persName key="P.0000003" ref="h</pre>
Q.01.19140121.FMA.01.xml:<closer>
                     <salute>Hrzl.</salute> <signed>Ihr <persName key="P.0000003" ref="http://</pre>
Q.01.19140125.FMA.02.xml:<closer>
                     <salute>Hrzl. Gruß
Q.01.19140409.FMA.01.xml:<closer>
                     <salute>Herzl.
Q.01.19140414.FMA.01.xml:<closer>
                     <salute>Hrzl.</salute> <signed>Ihr <persName key="P.0000003" ref="http://</pre>
Q.01.19140507.FMA.01.xml:<closer>
                     <salute>besten Gruß
Q.01.19140512.FMA.01.xml:<closer>
                     <salute>Hrzl.</salute> <signed>Ihr <persName key="P.0000003" ref="http://</pre>
Q.01.19140606.FMA.01.xml:<closer>
                     <salute>Hrzl.</salute> <signed>Ihr <persName key="P.0000003" ref="http://</pre>
Q.01.19140608.FMA.01.xml:<closer>
                     <salute>hrzl.</salute> <signed>Ihr <persName key="P.0000003" ref="http://</pre>
```

Wenn wir die Ergebnisse direkt in eine Datei schreiben wollen, können wir das natürlich auch tun:

Mac/Linux:

```
grep -E -zo '<closer>[[:cntrl:]].*<salute>.*[[:cntrl:]].*<' *.xml > Grussformeln.txt
```

Spätestens jetzt wäre es aber an der Zeit, das Instrumentarium zu wechseln: Mit dem Terminal bzw. der Shell kann man verschiedenste Operationen durchführen, und es gibt zahlreiche kleine Programme, die man zusätzlich installieren kann – zum Parsen, also Zerlegen von XML-Dateien, zur Bearbeitung von Bilddateien oder zum Download von YouTube-Videos. Die Übersichtlichkeit ist allerdings recht begrenzt, und gerade für die Analyse von Struktur- und Textdaten gibt es weitaus geeignetere Programmiersprachen wie R oder Python, wie in Section 2.2 bereits erwähnt.

Um beispielsweise mit dem Volltext der Briefe weiterzuarbeiten, sei es für ein $close\ reading$ oder für quantitative Textanalysen, können für eine bessere Lesbarkeit und einfachere Auswertungen die Auszeichnungen des TEI-XML entfernt werden.

9 Ausblick

Wenn Sie die von Ihnen erstellten Strukturierungen von Brief Nr. 8 von Franz Marc an Herwarth Walden und die identifizierten Entitäten mit denjenigen in den XML-Dateien der STURM-Editor:innen vergleichen, ergeben sich vermutlich einige Unterschiede. Ein zentraler ist sicher, dass sich die Herausgeber:innen bei ihrer Strukturierung an ein Schema gehalten haben, das im Bereich der Texteditionen Standard ist, TEI XML, und das kann zahlreiche Vorteile haben. So können Sie beispielsweise die im vorherigen Kapitel gezeigten Abfragen nach der Grußformel durchführen, ohne sich Gedanken darüber machen zu müssen, ob sich die Benennung des Tags auf halber Strecke ändert. Und wenn Sie sich während der Forschungsarbeit dafür entscheiden würden, anstatt Grußformeln besser die Anrede zu untersuchen, oder aber das Korpus auf die Briefe von Jacoba van Heemskerck auszuweiten, könnten wir dies mit wenigen Änderungen in unseren Abfragen machen, weil auch die Anrede mit einem einheitlichen Tag codiert ist; würden wir die Quellen nur über die Webseite lesen und unsere Auszählungen von Hand machen, würde unsere Arbeit mit dem Hinzufügen neuer Dokumente von vorne beginnen.

Das Erstellen von standardisierten Daten mithilfe eines Schemas bzw. bestimmter Richtlinien ermöglicht es auch, verschiedene Datensätze miteinander zu kombinieren oder mit weiteren Daten anzureichern. Bei einem Blick in die Briefe, auf der Webseite oder in der XML-Datei, wird ersichtlich, dass Entitäten wie Personen oder Orte nicht nur als solche markiert und projektintern verlinkt, sondern auch mit weiteren Normdaten verbunden wurden, beispielsweise mit dem dazugehörigen Eintrag in der GND, der Gemeinsamen Normdatei der Deutschen Nationalbibliothek, oder in Geonames, einer Datenbank für geographische Daten.

```
<salute>Lieber <persName key="P.0000001" ref="http://d-
                          nb.info/gnd/118770950">Walden</persName>,</salute:
114
                      </opener>
                      ich habe nicht das geringste vom Anwalt erhalten; Sie kennen mich ja, daß ich
115
                         stets sofort antworte; es muß verloren gegangen sein. Grüßen Sie bitte D<hi rend="super">
                              <hi rend="underline">r</hi>
                          </hi>
118
                          119
                      In einer niederrheinischen Zeitung soll von einer bevorstehenden großen
                          ref="http://sws.geonames.org/2952539">Barmen</placeName> von <persName key="P.0000009"
                         ref="http://d-nb.info/gnd/118559737">Kandinsky</persName> mir u.s.w.
berichtet worden sein; ich kann mir nur denken, daß die <persName key="P.0000057"
                          ref="http://d-nb.info/gnd/143669230">Reiche</persName>-<placeName key="0.0000032"
                          ref="http://sws.geonames.org/2935022">Dresden</placeName> Collection dahin
                              kommt;<pb xml:id="S.177v.02" n="177v" facs="http://resolver.staatsbibliothek-
                              berlin.de/SBB0000DAA300000001"/> ich
                         bin natürlich einverstanden, nur soll die frei<hi rend="underline">e
Rück</hi>fracht <placeName key="0.0000045"</pre>
                          ref="http://sws.geonames.org/2952539">Barmen</placeName> - <placeName key="0.00000002"
                          ref="http://sws.geonames.org/2950159">Berlin</placeName> gesichert
127
                              sein.<note>Im Jahr 1913 fand in der Stadt Barmen eine STURM-Ausstellung mit
128
                                  Werken der Künstlergruppe "Der Blaue Reiter" statt. In der Literatur gilt
                              diese Ausstellung als nicht ausreichend nachgewiesen; dieser Brief gibt
129
                              endgültigen Aufschluss über das tatsächliche Stattfinden der Ausstellung.
                              Vgl. Enders, Rainer: Ausstellungen außerhalb der Berliner Galerie. URL: <ref
                              thistoricum.net/themen/portale/sturm/ausstellungen/</ref>
                              (Aufruf 06.04.2017).</note> Aber es ist natürlich nur Vermutung von mir. Von
                              <persName key="P.0000057" ref="http://d-nb.info/gnd/143669230">Reiche</persName> höre ich
                              gar nichts.
                          <salute>Hrzl.</salute>
                          <signed>Ihr <persName key="P.0000003" ref="http://d-nb.info/gnd/11857745X">F.
136
137
                                  Marc</persName>
                          </signed>
```

Figure 9.1: Ausschnitt aus Brief Nr. 1 von Franz Marc an Herwarth Walden, Normdaten gelb hervorgehoben.

Wenn Sie den Link zu Kandinski oder zu Berlin anklicken, erhalten Sie auf den Seiten der GND bzw. Geonames zahlreiche zusätzliche Informationen zur Person bzw. zum Ort, unter anderem Lebensdaten bzw. Geokoordinaten. Wenn Sie nun beispielsweise wissen wollten, welche Orte in den Briefen Franz Marcs genannt werden, könnten Sie diese nicht nur mithilfe des Tag-Paars <placeName>-</placeName> extrahieren, sondern mit den dazugehörigen Geokoordinaten anreichern und sich auf einer Karte anzeigen lassen.

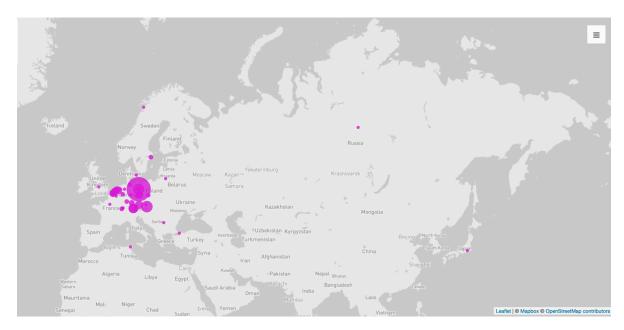


Figure 9.2: Orte, die im Briefkorpus von Franz Marc an Herwarth Walden genannt werden (Visualisierung erstellt mit Palladio).



Figure 9.3: Orte, die im Briefkorpus von Franz Marc an Herwarth Walden genannt werden, Detail (Visualisierung erstellt mit Palladio).

Für solche Vorgänge reicht ein kurzes (aber nicht unbedingt schnell erstelltes ...) Skript, an

das sich, je nach Ziel bzw. gewünschter Visualisierung, noch weitere Schritte zur Datenverfeinerung bzw. -anreicherung anschließen lassen. 1 Diese Abfolge lässt sich recht unaufwendig auf weitere Dokumente ausweiten – ob Sie nur die Orte in den Briefen Franz Marcs oder auch in denen Jacoba van Heemskercks extrahieren und visualisieren wollen, spielt mit Blick auf die Rechenzeit des Skripts keine Rolle.

Das Erstellen von Datensätzen nach bestimmten Richtlinien, einerseits formal, andererseits auch mit Blick auf die FAIR-Prinzipien, bietet also viele Vorteile für die eigene Arbeit – so müssen beispielsweise Schemata zur Klassifizierung nicht von Neuem erfunden werden – ebenso wie für die Arbeit anderer – grundlegende Informationen können übernommen werden und es bleibt mehr Zeit für die inhaltliche Forschung.

Ein anschauliches Beispiel für die Weiternutzung von Daten ist das von Studierenden erstellte Projekt quoteSalute, eine Webseite, auf der Sie historische Grußformeln generieren lassen können, falls Ihnen beim Briefeschreiben der Standardgruß zu langweilig geworden ist. Das Projekt hat hierfür mehrere XML-codierte Briefkorpora, die alle als offen nutzbare Daten online verfügbar sind, kombiniert, die Grußformeln extrahiert und angereichert. Die genaue Projektbeschreibung finden Sie hier.

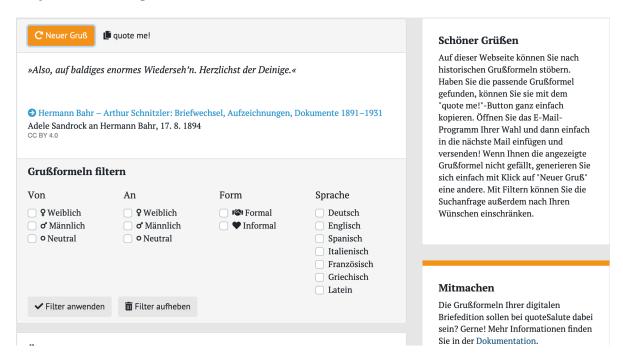


Figure 9.4: Startseite von quoteSalute

¹In diesem Fall wurde die mithilfe des verlinkten Skripts (in der Programmiersprache R) erstellte csv-Datei mit OpenRefine weiter bearbeitet, um den Adressen Geokoordinaten hinzuzufügen. Die so angereicherte Tabelle wurde in das Onlinetool Palladio geladen, das Koordinaten auf einer Karte darstellen kann.

Wie Sie sehen, können Sie sich Grüße nicht nur generieren lassen, sondern auch filtern, beispielsweise nach Geschlecht der Absender:innen bzw. der Adressat:innen oder nach Sprache. Sie können dies deswegen, weil die vorhandenen Codierungen der Textdateien im Projekt ausgeweitet und zum Beispiel um das Geschlecht der beteiligten Personen, sofern ersichtlich, ergänzt wurden. Der Code für das Projekt ist dabei offen auf einem GitHub-Repositorium zugänglich, und neue Korpora werden gerne aufgenommen.

Computergestütztes und computerbasiertes Arbeiten, das sollte dieser Guide vermitteln, vereinfacht, beschleunigt oder ermöglicht gar Prozesse, die wir für unsere historische Forschung nutzen können. Dabei gibt es mittlerweile zahlreiche Programme, die über eine graphische Oberfläche bedient werden können und die meisten Anforderungen aus geisteswissenschaftlicher Sicht erfüllen. Manche Analysen hingegen erfordern sehr spezifische Schritte oder viel Rechenkapazitäten, für die sich das Erlernen einer Programmiersprache lohnen kann. Gemeinsam ist beidem, dass es sowohl für das Formulieren einer interessanten Fragestellung, für die Auswahl der Datenbasis als auch für die Interpretation der Ergebnisse, die die Maschine liefert, informierte Expertise braucht – ein Programm zur Bestimmung der Autor:innenschaft eines unbekannten Textes rechnet unendlich lange, wenn es alle existierenden Texte mit einem unbekannten vergleichen muss; hier braucht es eine Eingrenzung auf eine Auswahl durch Literaturwissenschaftler:innen. Historische Forschung ist immer auch Analyse des Einzelfalls, des Partikularen, eines close readings; die Möglichkeit, den Blick mithilfe des Computers auszuweiten, kann in den meisten Fällen gewinnbringend in Betracht gezogen werden.

A Glossar

API Application Programming Interface: a facility offered by a web

resource which allows search queries independent of a GUI, often

performed using scripts

bash default program that runs in the **command line**

bias systematic error that results from an unbalanced sample

big data huge amount of data, identifiable through repeated freezing of your

standard program when opening a file

born digital data data which originated in a digital form

CLI Command Line Interface, text interface that allows interaction with

the computer; see also bash

close reading careful and attentive interpretation of a text

CMS Content Management System

Console See CLI

Crowdsourcing projects that include the active participation of the public to

generate content, transcribe sources etc.

csv comma separated values, a structured text format, using commas as

separators between columns

distant reading quantitative approach to huge amounts of texts, using computational

methods to search for interpretable patterns

GUI Graphical User Interface

HTML Hypertext Markup Language, a structured text format, like the

format this guide is written in, to render documents in a browser

Jupyter notebook web application/interactive coding environment that runs in a

browser; let's you create and share code (https://jupyter.org)

machine learning umbrella term for different methods that use data to do a task in a

specific way, using data to learn and to improve the results

machine readable transformation of, for example, text into a data format that is

processable by a computer

OCR Optical Character Recognition, process of transforming text on an

image into a data format

OS Operating System

open source freely available source code that can be used, modified and

redistributed without limitations

OSS Open Source Software

Regular Expression	syntax for search and replace text using patterns (instead of exact
	matches)
terminal	See CLI
web scraping	extracting data from websites

B Literatur, Tools, Tutorials

B.1 Was ist digital history?

- Antenhofer, Christina, Kühberger, Christoph, Strohmeyer, Arno (Hrsg.): Digital Humanities in den Geschichtswissenschaften. Wien 2023.
- Brennan, Sheila A.: Digital History, in: The Inclusive Historian's Handbook, https://inclusivehistorian.com/digital-history/, 04.06.2019.
- Hohls, Rüdiger: Digital Humanities und digitale Geschichtswissenschaften, in: Busse, Laura u. a. (Hg.): Clio-Guide. Ein Handbuch zu digitalen Ressourcen für die Geschichtswissenschaften, Berlin 2018, S. A.1-1–B.1-34. Online: https://doi.org/10.18452/19244.
- Romein, C. Annemieke u. a.: State of the Field: Digital History, in: History 105 (365), 04.2020, S. 291—312. Online: https://doi.org/10.1111/1468-229X.12969.
- Winters, Jane: Digital History, in: Tamm, Marek; Burke, Peter (Hg.): Debating New Approaches to History, London 2019, S. 277–300.
- Art. "Digital history", in: Wikipedia, 07.09.2022. Online: https://en.wikipedia.org/w/index.php?title=Digital history&oldid=1220662831, Stand: 25.04.2024.

B.2 Einführungen und Guides

- Battershill, Claire; Ross, Shawna: Using Digital Humanities in the Classroom. A Practical Introduction for Teachers, Lecturers, and Students, London u.a. 2022.
- Blaney, Jonathan u. a.: Doing Digital History. A Beginner's Guide to Working with Text as Data, Manchester 2021.
- Cohen, Daniel J.; Rosenzweig, Roy: Digital History. A Guide to Gathering, Preserving, and Presenting the Past on the Web, Philadelphia 2006. Online: https://chnm.gmu.edu/digitalhistory/.
- Döring, Karoline u. a. (Hg.): Digital History. Konzepte, Methoden und Kritiken Digitaler Geschichtswissenschaft, Boston 2022, Online: https://doi.org/10.1515/9783110757101.

- Dougherty, Jack; Nawrotzki, Kristen (Hg.): Writing History in the Digital Age, 2013. Online: https://doi.org/10.3998/dh.12230987.0001.001.
- Graham, Shawn u. a.: Exploring Big Historical Data. The Historian's Macroscope, 2022. Online: https://doi.org/10.1142/12435.
- Lemercier, Claire; Zalc, Claire: Quantitative Methods in the Humanities. An Introduction, Charlottesville 2019.
- Peroni, Silvio: The CTP Book. A book for teaching Computational Thinking and Programming skills to people with a background in the Humanities.

B.3 Digital history tools (free/open source)

B.3.1 Allgemein

• Programming Historian: Tutorials zu verschiedenen Tools und Methoden für historische Forschung und Lehre

B.3.2 Datenanalyse

- Datasette: interaktives Tool zur Datenexploration und -visualisierung
- Heurist: webbasierte Forschungsumgebung zur Datenerhebung, -exploration und -visualisierung
- Nodegoat: webbasierte Forschungsumgebung zur Datenerhebung, -exploration und visualisierung

B.3.3 Karten

- Dariah GeoBrowser: webbasiertes Tool zur Darstellung raumzeitlicher Daten
- QGis: Software zum Erstellen, Bearbeiten, Anzeigen und Analysieren räumlicher Information
- Map in Seconds: webbasiertes Tool zur Erstellung von Karten (optimiert für Chrome Browser)
- Palladio: webbasiertes Tool zur Visualisierung strukturierter Daten, u.a. Geodaten

B.3.4 Netzwerkanalyse

- Gephi: Programm zur Analyse und Visualisierung von Graphen und Netzwerken
- Nodegoat: webbasierte Forschungsumgebung zur Datenerhebung, -exploration und visualisierung, u.a. auch Netzwerkanalyse
- Palladio: webbasiertes Tool zur Visualisierung strukturierter Daten, u.a. auch Netzwerkanalyse

B.3.5 Text-/Korpusanalyse

- AntConc: Korpusanalyse-Toolkit
- Lemmatisierung: Sammlung der FID Romanistik
- Natural Language Toolkit, Package für Python zur Tokenisierung, Lemmatisierung usw.: NLTK
- Tokenisierung: Tutorial von fortext zu NLTK
- Voyant-Tools: Sammlung von Tools zur Textanalyse, browserbasiert oder standalone

B.3.6 Storytelling

- Omeka: webbasierte Plattform zur Kuration digitaler Sammlungen, für die Veröffentlichung und Ausstellung von Bildern/Objekten
- StoryLine JS: Tool zur Erstellung von interaktiven Diagrammen
- StoryMap JS: Tool zur Erstellung von interaktiven Karten/Bildern mit verschiedenen Elementen
- Timeline JS: Tool zur Erstellung von interaktiven Timelines

B.3.7 Visualisierung

- Bostock, Michael; Heer, Jeffrey; Ogievetsky, Vadim: A Tour through the Visualization Zoo. A Survey of Powerful Visualization Techniques, from the Obvious to the Obscure, in: Queue 8, Nr. 5 (2010). Online: https://queue.acm.org/detail.cfm?id=1805128
- Data Visualisation Catalogue: Guide zur Auswahl von Visualisierungsformen
- Datasette: interaktives Tool zur Datenexploration und -visualisierung
- FID Romanistik: Sammlung von Tools zur Datenvisualisierung

- Palladio: webbasiertes Tool zur Visualisierung strukturierter Daten
- RAWGraphs: Tool zur Datenvisualisierung von tabularen Daten (.tsv-, .csv-, .dsv- oder .json-Dateien)

B.4 Digital Literacy, Digital Criticism

- Ekström, Andreas: The Moral Bias behind your Search Results, TED talk 7.12.2015 (9:18), Online: https://www.youtube.com/watch?v=_vBggxCNNno.
- Gibbs, Frederick W.: New Forms of History: Critiquing Data and Its Representations, in: The American Historian, February 2016. Online: http://tah.oah.org/february-2016/new-forms-of-history-critiquing-data-and-its-representations/.
- Tavani, Herman; Zimmer, Michael: Search Engines and Ethics, in: Edward N. Zalta (Hg.): The Stanford Encyclopedia of Philosophy (Fall 2020 Edition), Online: https://plato.stanford.edu/archives/fall2020/entries/ethics-search/, Kap. 3.1.

B.5 Programmiersprachen & Editoren

- Python:
 - Althage, Melanie u. a.: Python für Historiker:innen. Ein anwendungsorientierter und interaktiver Einstieg, 20.07.2022, Jupyter Book, (v1.0)
 - Karsdorp, Folgert: Python Programming for the Humanities. Interactive Tutorial and Introduction to Programming with Python for the Humanities.
 - Mattingly, William: PythonHumanities.com. A Free Course for the Digital Humanities.
 - Editor z.B. PyCharm
- R:
 - Mullen, Lincoln A.: Computational Historical Thinking With Applications in R
 - Phillips, Nathaniel D.: YaRrr! The Pirate's Guide to R.
 - Wickham, Hadley; Grolemund, Garrett: R for Data Science.
 - Editor z.B. RStudio
- Editoren: Anaconda; Sublime Text; Visual Studio Code

B.6 Regular Expressions

- Knox, Doug: Understanding Regular Expressions, Programming Historian 2 (2013), ht tps://doi.org/10.46430/phen0033. (self-learning lesson)
- RegexOne: Learn Regular Expressions with simple, interactive exercises. (interactive self-learning tutorial)

B.7 Terminal/Command Line/Shell

- Dawson, Ted: Introduction to the Windows Command Line with PowerShell, Programming Historian 5 (2016), https://doi.org/10.46430/phen0054. (self-learning lesson)
- MIT Computer Science Department: 1-hour-lecture on the Shell (video)
- Milligan, Ian; Baker, James: Introduction to the Bash Command Line, Programming Historian 3 (2014), https://doi.org/10.46430/phen0037. (self-learning lesson)
- datacamp course:Introduction to Shell (interactive self-learning lesson)
- Jeroen Janssens: Data Science at the command line (book)

B.8 XML

• Latex Ninja Blog: A shamelessly short intro to XML for DH beginners (includes TEI) (blog post)

References

- [1] Jonathan Blaney et al. *Doing digital history: a beginner's guide to working with text as data.* eng. IHR research guides. Manchester: Manchester University Press, 2021. ISBN: 978-1-5261-3268-0.
- [2] Stephanie Russo Carroll et al. "The CARE Principles for Indigenous Data Governance". en. In: *Data Science Journal* 19 (Nov. 2020), p. 43. ISSN: 1683-1470. DOI: 10.5334/dsj-2020-043. URL: http://datascience.codata.org/articles/10.5334/dsj-2020-043/ (visited on 11/28/2022).
- [3] Catherine D'Ignazio and Lauren F. Klein. *Data feminism*. English. OCLC: 1157171756. 2020. ISBN: 978-0-262-35852-1. URL: https://direct.mit.edu/books/book/4660/Data-Feminism.
- [4] Marie-José Kolly and Simon Schmid. "Sie ist hübsch. Er ist stark. Er ist Lehrer. Sie ist Kindergärtnerin". de. In: Republik (Apr. 2021). URL: https://www.republik.ch/2021/04/19/sie-ist-huebsch-er-ist-stark-er-ist-lehrer-sie-ist-kindergaertnerin (visited on 08/23/2022).
- [5] Emmanuel Le Roy Ladurie. "La fin des érudits". In: Le Nouvel Observateur (Aug. 1968).
- [6] Claire Lemercier and Claire Zalc. Quantitative Methods in the Humanities. An Introduction. Charlottesville: University of Virginia Press, 2019. ISBN: 978-0-8139-4270-4.
- [7] Chantel Ridsdale et al. Strategies and Best Practices for Data Literacy Education Knowledge Synthesis Report. Jan. 2015. DOI: 10.13140/RG.2.1.1922.5044.
- [8] Roopika Risam. "It's Data, Not Reality": On Situated Data With Jill Walker Rettberg. en. June 2020. URL: https://medium.com/nightingale/its-data-not-reality-on-situated-data-with-jill-walker-rettberg-d27c71b0b451 (visited on 08/16/2022).
- [9] C. Annemieke Romein et al. "State of the Field: Digital History". en. In: History~105.365~ (Apr. 2020), pp. 291–312. ISSN: 0018-2648,~1468-229X. DOI: 10.1111/1468-229X.12969. URL: https://onlinelibrary.wiley.com/doi/10.1111/1468-229X.12969 (visited on 09/15/2022).
- [10] Mark D. Wilkinson et al. "The FAIR Guiding Principles for scientific data management and stewardship". en. In: *Scientific Data* 3.1 (Mar. 2016). Number: 1 Publisher: Nature Publishing Group, p. 160018. ISSN: 2052-4463. DOI: 10.1038/sdata.2016.18. URL: https://www.nature.com/articles/sdata201618 (visited on 11/09/2022).