

Fernuniversität Hagen

Modul 63458 - Fachpraktikum

Natural Language Processing, Information Extraction
und Retrieval

WiSe 2022/23

Praktikumsbericht

**Automatische Erstellung
einer Wissensrepräsentation
aus einem medizinischen Text**

eingereicht von

Anne Koch, Clara Jansen, Dietrich Tönnies

Betreuer: Prof. Dr. Hemmje
Dr. Christian Nawroth
Stephanie Heidepriem, M.Sc.

Hagen, 25. Januar 2023

INHALTSVERZEICHNIS

1	EINLEITUNG	1
1.1	Motivation	1
1.2	Problembeschreibung	2
1.3	Forschungsfragen	3
1.4	Forschungsmethode	3
1.5	Forschungsziele	5
1.6	Lösungsansatz	6
2	STAND DER WISSENSCHAFT	10
2.1	FZ 1.1. Aufbau und Struktur medizinischer Texte	10
2.2	FZ 2.1. Überblick über depressive Erkrankungen	10
2.2.1	Allgemein	10
2.2.2	Symptome	10
2.2.3	Formen der Depression	10
2.2.4	Ursachen	11
2.2.5	Behandlung	11
2.3	FZ 2.2. Automatisierung durch NLP	11
2.3.1	Natural Language Processing	11
2.3.2	Named Entity Recognition	11
2.3.3	Entity Linking	12
2.3.4	spaCy	12
2.3.5	pySBD	12
2.3.6	Weitere Python-Bibliotheken für Machine Learning (ML) und NER	12
2.4	FZ 3.1. Wissensrepräsentation mittels RDF	13
2.4.1	Resource Description Framework Schema	13
2.4.2	Dublin Core Metadata Initiative	14
2.4.3	Serialisierung von RDF-Graphen	14
2.4.4	SPARQL	14
3	KONZEPTUELLE MODELLIERUNG UND ENTWURF	16
3.1	User Centered System Design	16
3.1.1	Anwendungskontext und Anwendungsfall	16
3.1.2	Aktivitätsmodell	17
3.1.3	Informationsmodell	17
3.1.4	Architektur- und Komponentenmodell	18
3.2	FZ 1.2 Theoriebildung zur Vorverarbeitung medizinischer Texte	18
3.3	FZ 2.3 Theoriebildung zur Überführung medizinischen Fach- vokabulars in maschinenlesbare Form zur weiteren Verarbei- tung durch NLP	19
3.4	FZ 3.2 Theoriebildung zur Wissensrepräsentation	20
3.5	Zusammenfassung	21
4	PROOF-OF-CONCEPT IMPLEMENTIERUNG	23

4.1	FZ 1.3 Implementierung zur Vorverarbeitung medizinischer Texte	23
4.2	FZ 2.4 Implementierung zur Überführung medizinischen Fachvokabulars in maschinenlesbare Form zur weiteren Verarbeitung durch NLP	23
4.3	FZ 3.3 Implementierung zur Wissensrepräsentation	23
5	EVALUIERUNG	24
5.1	FZ 1.4 Evaluierung zur Vorverarbeitung medizinischer Texte	24
5.1.1	Fähigkeiten	24
5.1.2	Limitierungen	24
5.2	FZ 2.5 Evaluierung zur Überführung medizinischen Fachvokabulars in maschinenlesbare Form zur weiteren Verarbeitung durch NLP	24
5.3	FZ 3.4 Evaluierung zur Wissensrepräsentation	24
5.4	Zusammenfassung	24
6	ZUSAMMENFASSUNG UND AUSBLICK	25
6.1	Zusammenfassung	25
6.2	Ausblick	25

DCMI	Dublin Core Metadata Initiative
MENHIR	Mental health monitoring through interactive conversations
MeSH	Medical Subject Headings
NER	Named Entity Recognition
NLM	United States National Library of Medicine
NLP	Natural Language Processing
RDF	Resource Description Framework
SPARQL	SPARQL Protocol And RDF Query Language
STOP	STop Obesity Platform
UMLS	Unified Medical Language System
W3C	World Wide Web Consortium

EINLEITUNG

Diese Praktikumsarbeit behandelt die automatische Erstellung einer Wissensrepräsentation aus einem medizinischen Text. Eine prototypische Softwarelösung wird entwickelt, die für einen vorgegebenen Beispieltext einen Wissensgraphen ermittelt.

Gemäß einem Artikel des *Economist* von 2017 [[the_economist_worlds_2017](#)], der in Verbindung mit den aktuellen technischen Entwicklungen, Künstlicher Intelligenz und Data Science viel zitiert wird, ist nicht mehr Öl die wertvollste Ressource, sondern Daten. Ebenso wie Öl müssen die Daten jedoch erst aufbereitet werden, um tatsächlich von Nutzen zu sein. Auch Informationen, die in Form von für den Menschen lesbaren Texten zur Verfügung stehen, sind nicht direkt für Computer auswertbar. Daher wurden mit den Methoden des *Natural Language Processing (NLP)* Verfahren entwickelt, um automatisch Informationen aus Texten extrahieren zu können. Die allgemeinsprachlichen NLP-Verfahren bedürfen jedoch noch einer Spezialisierung, um auch für Fachtexte wie zum Beispiel aus dem Bereich der Medizin nutzbringend angewendet werden zu können.

Aus diesen Rahmenbedingungen ergibt sich die Motivation und die Problembeschreibung für diese Praktikumsaufgabe, die in den beiden nachfolgenden Abschnitten beschrieben werden.

1.1 MOTIVATION

Die spezielle Herausforderung der Praktikumsaufgabe liegt darin, automatisiert Zusammenhänge zwischen Entitäten medizinischer Texte zu finden und diese Zusammenhänge zu repräsentieren. Dem zugrunde liegt das Ziel des Dissertationsprojektes von Stephanie Heidepriem, zu erforschen, wie ein textbasierter Chatbot zur Unterstützung psychisch kranker Menschen entwickelt werden kann.

Das Dissertationsprojekt ist unter anderem an das Projekt *Mental health monitoring through interactive conversations (MENHIR)* angelehnt [[noauthor_menhir_2022](#)]. Dieses Projekt dient der Erforschung von Konversationstechnologien, die psychisch kranke Menschen unterstützen sollen. In diesem Zusammenhang wird auch ein MENHIR-Chatbot entwickelt, der personalisierte Unterstützung und hilfreiche Bewältigungsstrategien bieten soll.

Das Forschungsprojekt *STop Obesity Platform (STOP)* beschäftigt sich mit der Extraktion von Wissen unter anderem aus Chatbots, das anschließend aufbereitet und mit weiteren Informationen kombiniert werden soll. Dieses Wissen soll medizinischem Fachpersonal zur Verfügung gestellt werden und außerdem Menschen mit Adipositas dabei helfen, eine gesündere Ernährung einzuhalten [[noauthor_stop_2022](#)].

Es ist aber auch denkbar, dass die Problemstellung der Praktikumsaufgabe auch darüber hinaus für weitere Anwendungen interessant ist und mögliche Ergebnisse in verschiedenen Gebieten eingesetzt werden können, in denen Informationen in (medizinischen) Texten zueinander in Zusammenhang gebracht werden sollen.

1.2 PROBLEMBESCHREIBUNG

Es soll eine Konsolenapplikationen entwickelt werden, die englischsprachige medizinische Texte analysiert und das darin enthaltene Wissen strukturiert in einem RDF-Graphen hinterlegt. Die Konsolenapplikation soll in Python programmiert werden, wobei für die Textanalyse Klassen und Methoden der Open-Source-Bibliothek *spaCy* o.ä. genutzt werden sollen.

Die medizinischen Texte enthalten z.B. Aussagen zu Krankheiten (z.B. Depression) und zählen deren Symptome (z.B. Motivationsverlust) auf. Aufgabe der Konsolenapplikation ist es, Schlüsselbegriffe im Text zu finden und miteinander in Bezug zu setzen, indem z.B. Symptome den ihnen zugrunde liegenden Krankheiten zugeordnet werden. Die von *spaCy* oder anderen NLP-Bibliotheken zur Verfügung gestellte Funktionalität besteht aus einer Pipeline von Analyse-Komponenten, die nacheinander auf den zu analysierenden Text angewendet werden. Diese Komponenten dienen dazu, Texte in einzelne Sätze zu zerlegen und die Sätze grammatikalisch zu analysieren.

PROBLEM 1: VORVERARBEITUNG EINES TEXTES Medizinische und wissenschaftliche Texte bestehen häufig aus Aufzählungen und Zwischenüberschriften. Der Text sollte so aufbereitet werden, dass anschließende Analyse-Algorithmen möglichst erfolgreich arbeiten können. Es soll erprobt werden, wie dies am besten bewerkstelligt werden kann, z.B. durch Umwandlung von Aufzählungen in mehrere Aussagesätze. Auch sollte automatisch erkannt werden, welche Sätze Aussagen (Wissen) formulieren etwa durch Unterscheidung von Indikativ und Konjunktiv oder durch Ignorieren von Fragesätzen.

PROBLEM 2: TRAINING DER NER - KOMPONENTE Neben der grammatikalischen Analyse besteht eine wesentliche Aufgabe von *spaCy* oder anderen NLP-Bibliotheken darin, Schlüsselbegriffe zu finden und nach Möglichkeit zu kategorisieren. Diese Art der Analyse wird als *Named Entity Recognition (NER)* bezeichnet. Bei *spaCy* übernimmt diese Aufgabe der regelbasierte *Entity-Ruler* oder der auf statistischen Modellen basierende *Entity-Recognizer*. Eine weitere Komponente ist der *Entity-Linker*, mit dem Begriffe eindeutig den in einer Wissensbasis gespeicherten Entitäten zugeordnet werden können. Auch der *Entity-Linker* basiert auf statistischen Modellen und wird mit Beispielsätzen trainiert.

Die Standard-NER-Funktionalität von *spaCy* erkennt medizinische Begriffe nur unzureichend. Es ist daher notwendig, eine Datenbank mit

medizinischen Begriffen und Kategorien zusammenzustellen, die für das Training der NER-Komponente verwendet werden kann.

Die *United States National Library of Medicine (NLM)* stellt mit dem *Unified Medical Language System (UMLS)* ein mächtiges Werkzeug für die Textanalyse zur Verfügung. Teil von UMLS ist der sogenannte Metathesaurus, der aus einer Vielzahl von Thesauri unterschiedlicher Organisationen zusammengestellt wird. Zu diesen Thesauri gehört u.a. der von der NLM entwickelte Thesaurus *Medical Subject Headings (MeSH)*. *MetaMap* und das weniger umfangreiche *MetaMapLite* sind eigene *Entity-Recognition*-Werkzeuge der *National Library of Medicine*. Die zugrundeliegende Datenbasis dieser Werkzeuge sollen im Rahmen dieses Praktikums dafür verwendet werden, die *Named Entity Recognition*-Komponenten von *spaCy* zu trainieren. Aus den von der NLM zur Verfügung gestellten Begriffslisten soll eine geeignete Auswahl erfolgen, die für das Training der NER-Komponente verwendet werden kann.

PROBLEM 3: ZUORDNUNG DER ENTITÄTEN Basierend auf den gefundenen Entitäten soll das Python-Programm in der Lage sein, Begriffe wie Krankheiten und Symptome richtig zuzuordnen. Hierzu muss die Struktur von Aussagesätzen, Fragen, Aufzählungen und Überschriften analysiert werden und so aufbereitet werden, dass eine automatische Analyse durch NER und *Entity-Linker* möglichst erfolgreich ist. Die Ausgabe der Konsolenapplikation erfolgt im RDF/XML-Format.

1.3 FORSCHUNGSFRAGEN

Dieser Abschnitt beschreibt die Forschungsfragen, die sich aus dem in Abschnitt 1.2 angeführten Problemen ableiten.

- FF 1 Wie kann ein medizinischer Text vorverarbeitet werden, so dass relevante Aussagen (Überschriften, Aufzählungen, Leerzeilen und Dialoge) automatisch erkannt werden?
- FF 2 Wie kann ein medizinisches Fachvokabular über depressive Erkrankungen automatisiert in eine maschinenlesbare Form überführt werden?
- FF 3 Wie kann eine Wissensrepräsentation über einen gegebenen Text maschinenlesbar erstellt werden?

1.4 FORSCHUNGSMETHODE

In dieser Arbeit kommt die in der Forschung zu Informationssystemen etablierte Methode von Nunamaker, Chen und Purdin [**nunamaker_systems_1990-1**] zum Einsatz, die einen strukturierten Rahmen zur Durchführung von Forschungsarbeiten bereitstellt.

Die Forschungsmethode umfasst die folgenden in Abbildung 1.1 dargestellten Phasen:



Abbildung 1.1: Forschungsmethode nach [nunamaker_systems_1990-1]

BEOBSACHTUNG: Insbesondere wenn der Untersuchungsgegenstand relativ unbekannt ist, werden Fallstudien, Feldversuche oder Umfragen durchgeführt, um ein „Gefühl“ für den Aufwand zu erhalten. Auf dieser Grundlage können dann konkrete Hypothesen erstellt werden, die durch Experimente geprüft werden.

In dieser Arbeit findet aufgrund der Art des Untersuchungsgegenstandes in der Beobachtungsphase hauptsächlich die Literaturrecherche und Ermittlung des Standes von Wissenschaft und Technik statt.

THEORIEBILDUNG: In dieser Phase werden neue Ideen, Konzepte Methoden oder Modelle entwickelt. Diese Theorien beschreiben das Systemverhalten allgemein, haben jedoch wenig praktische Bedeutung für die Zieldomäne. Sie können aber genutzt werden, um Forschungshypothesen aufzustellen, die Planung von Experimenten zu unterstützen und systematische Beobachtungen durchzuführen.

SYSTEMENTWICKLUNG: Diese Phase besteht aus fünf Teilen:

- Konzeptentwurf
- Erstellung der Systemarchitektur
- Erstellung von Prototypen
- Produktentwicklung
- Technologietransfer

EXPERIMENT: In dieser Phase werden die gefundenen Theorien und entwickelten Systeme evaluiert. Die Ergebnisse der Experimente können

genutzt werden, um die Theorien weiterzuentwickeln und die Systeme zu verbessern.

Obwohl die Phasen in der Methode keine vorgegebene Reihenfolge haben, sondern sich gegenseitig beeinflussen, wird in der vorliegenden Arbeit die oben angeführte Abfolge verwendet.

1.5 FORSCHUNGSZIELE

Dieser Abschnitt beschreibt die Forschungsziele, die sich aus den in Abschnitt 1.3 angeführten Problemen ableiten.

FZ 1 Wie kann ein medizinischer Text vorverarbeitet werden, so dass relevante Aussagen (Überschriften, Aufzählungen, Leerzeilen und Dialoge) automatisch erkannt werden?

FZ 1.1 Aufbau und Struktur medizinischer Texte (Observation)

FZ 1.2 Theoriebildung zur Vorverarbeitung medizinischer Texte

FZ 1.3 Implementierung zur Vorverarbeitung medizinischer Texte

FZ 1.4 Evaluation zur Vorverarbeitung medizinischer Texte

FZ 2 Wie kann ein medizinisches Fachvokabular über depressive Erkrankungen automatisiert in eine maschinenlesbare Form überführt werden?

FZ 2.1 Überblick über depressive Erkrankungen (Observation)

FZ 2.2 Automatisierung durch NLP (Observation)

FZ 2.3 Theoriebildung zur Überführung medizinischen Fachvokabulars in maschinenlesbare Form zur weiteren Verarbeitung durch NLP

FZ 2.4 Implementierung zur Überführung medizinischen Fachvokabulars in maschinenlesbare Form zur weiteren Verarbeitung durch NLP

FZ 2.5 Evaluation zur Überführung medizinischen Fachvokabulars in maschinenlesbare Form zur weiteren Verarbeitung durch NLP

FZ 3 Wie kann eine Wissensrepräsentation über einen gegebenen Text maschinenlesbar erstellt werden?

FZ 3.1 Wissensrepräsentation mittels RDF (Observation)

FZ 3.2 Theoriebildung zur Wissensrepräsentation

FZ 3.3 Implementierung zur Wissensrepräsentation

FZ 3.4 Evaluation zur Wissensrepräsentation

1.6 LÖSUNGSANSATZ

Die Wissensrepräsentation medizinischer Texte soll Entitäten unterschiedlicher Kategorie miteinander in Beziehung setzen. Naheliegender ist es z.B., automatisch in einem Text Symptome und Krankheiten zu identifizieren, diese miteinander in Beziehung zu setzen und eine Wissensrepräsentation der Art

loss of appetite - is a symptom of - depression

zu erzeugen. Dazu muss spaCy befähigt werden, nicht nur Fachbegriffe zu erkennen, sondern diese auch richtig zu kategorisieren, als z.B. Krankheit oder Symptom. Folgender Lösungsansatz ist angedacht:

MetaMapLite verwendet eine auf UMLS basierende Datenbasis mit medizinischen Fachbegriffen. Der Umfang der Datenbasis von MetaMapLite ist etwas reduziert und umfasst z.B. nur englische Fachbegriffe. Auf der MetaMapLite-Website kann das zip-Archiv

public_mm_data_lite_usabase_2022aa.zip

heruntergeladen werden. In diesem Archiv befindet sich die Datei *postings* in dem Unterordner

\public_mm_lite\data\ivf\2022AA\USAbase\indices\cuisourceinfo.

Die Datei *postings* enthält ca. 11 Millionen englischsprachige Einträge und ist mit einer Größe von etwa 739 MB etwas handlicher als die Datenbestände von UMLS. Um zu prüfen, ob diese Datenbasis geeignet ist, ist ein Auszug von Einträgen, die für den zu analysierenden Beispieltext über Depressionen relevant sind, in der folgenden Tabelle aufgelistet.

CUI bezeichnet den Concept Unique Identifier, mit dem Synonyme gefunden werden können. Der CUI besteht nur aus dem mit dem Buchstaben 'C' beginnenden Teil und ist hier ergänzt durch einen sogenannten Term-Type. Der Term-Type 'PT' kennzeichnet z.B. bevorzugte Einträge. SUI bezeichnet den String Unique Identifier, mit dem gleichlautende (und damit redundante) Bezeichnungen gefunden werden können.

Die Tabelle listet nur Einträge auf, in denen der Begriff, z.B. 'depression', einer in Klammern hinter dem Begriff stehenden Kategorie zugeordnet ist, hier z.B. 'disease'. Über den CUI können in der Datei *postings* Synonyme gefunden werden, wobei nur der mit 'C' beginnende Teil relevant ist. Diese Synonyme enthalten häufig keine Kategorien in Klammern. Die Kategorien können aber über den CUI den Synonymen zugeordnet werden. Die sehr zahlreichen Synonyme bzw. alternativen Ausdrücke sind in der Tabelle nicht aufgeführt.

CUI	SUI	Item	Source
FNC0232933	S3225525	Abnormal menstrual cycle (finding)	SNOMEDCT_US
PTC0424569	S3221759	Circumstances interfere with sleep (disorder)	SNOMEDCT_US
YC0009806	S3235964	Constipation (finding)	SNOMEDCT_US
LAC3845528	S14560529	Depressed mood (e.g., feeling sad, tearful)	LNC
SYC0344315	S3252744	Depressed mood (finding)	SNOMEDCT_US
GTC0011570	S1431189	depression (disease)	AOD
FNC2939186	S3264511	Disturbance in mood (finding)	SNOMEDCT_US
FNC1288289	S3312713	Fearful mood (finding)	SNOMEDCT_US
FNC0150041	S3313279	Feeling hopeless (finding)	SNOMEDCT_US
FNC0022107	S3313282	Feeling irritable (finding)	SNOMEDCT_US
FNC0424000	S3313310	Feeling suicidal (finding)	SNOMEDCT_US
ETC0917801	S3373158	Insomnia (disorder)	SNOMEDCT_US
PTC0015672	S3386372	Lack of energy (finding)	SNOMEDCT_US
FNC2981158	S3386389	Lack of libido (finding)	SNOMEDCT_US
FNC1971624	S3397609	Loss of appetite (finding)	SNOMEDCT_US
FNC0178417	S3397620	Loss of capacity for enjoyment (finding)	SNOMEDCT_US
FNC0456814	S3397668	Loss of motivation (finding)	SNOMEDCT_US
FNC0679136	S3398077	Low self-esteem (finding)	SNOMEDCT_US
PTC5444612	S20749480	mood (physical finding)	MTH
FNC2945580	S3485453	Poor self-esteem (finding)	SNOMEDCT_US
FNC0235160	S3513783	Restless sleep (finding)	SNOMEDCT_US
PTC0233481	S3620195	Worried (finding)	SNOMEDCT_US

In der Tabelle finden sich drei Kategorien, 'disease', 'disorder' und 'finding' (Die *postings*-Datei enthält noch weitere wie z.B. 'situation' oder 'procedure'). Hier bietet sich an, die Begriffe und Kategorien dieser Datenbasis für das Trainieren des Entity-Recognizers oder Entity-Rulers von spaCy zu verwenden. Die Zuordnung 'finding' bedeutet eigentlich Befund, kann hier aber auch als Symptom verstanden werden. 'disorder', also Störung, kennzeichnet in den meisten Fällen Krankheiten. Der Begriff 'insomnia' findet sich nicht im zu analysierenden Text, dort ist stattdessen von 'disturbed sleep' die Rede ist. Hier muss versucht werden, über Synonyme eine richtige Kategorisierung zu erreichen. Die meisten Einträge stammen von der Datenbank SNOMED CT (Systematized Nomenclature of Medicine Clinical Terms), die seit 2003 Teil des UMLS Metathesaurus ist.

Der Inhalt der Datei *postings* muss aufbereitet (und ggf. auch gekürzt) werden, so dass die Daten für das Training des Entity-Recognizers verwendet

werden können. Es muss untersucht werden, wie basierend auf der Wörterliste geeignete Trainingsdaten erzeugt werden können, etwa durch automatisch erzeugte Beispielsätze. Alternativ kann der Entity-Ruler zum Einsatz kommen, dem die Wörterliste einfach übergeben werden kann und dessen Funktion vorhersagbarer ist als die des Entity-Recognizers.

Der zu analysierende Text muss durch ein Python-Programm aufbereitet werden, mit dem Ziel, dass die NLP-Komponenten der Prozesspipeline des NLP-Frameworks, möglichst erfolgreich arbeiten. Ein Vorgehen könnte beispielsweise sein, Aufzählungen in mehrere vollständige Sätze zu zerlegen, so dass möglichst sinnvolle vollständige Sätze entstehen, die eine Krankheit und ein Symptom enthalten, so dass die automatische Textanalyse nicht überfordert wird, Krankheit und Symptome zusammenzubringen. Hierbei werden die von *spaCy* zur Verfügung gestellten Pipeline-Komponenten wie der Lemmatizer und Parser genutzt.

Es ist zu prüfen, inwieweit der *Dependency Parser* von *spaCy* genutzt werden kann, um Abhängigkeiten und Subjekt-Prädikat-Objekt-Beziehungen in Sätzen zu erkennen, die für die Erstellung des RDF-Graphen genutzt werden können.

Idealerweise entsteht jedoch bereits durch die NER ein Text mit Aussagesätzen, die jeweils eine Krankheit und ein oder mehrere Symptome enthalten. Für die Wissensrepräsentation soll der Entity-Linker eingesetzt werden. Normalerweise wird der Entity-Linker verwendet, um (mehrdeutige) Entitäten, die vom Entity-Recognizer oder Entity-Ruler gefunden werden, einer eindeutigen Entität zuzuordnen (z.B. einem Wikipedia-Eintrag). Der Entity-Linker wird trainiert, so dass die Zuordnung kontextbasiert erfolgt.

Für die Erstellung der Wissensrepräsentation kann der Entity-Linker auf unorthodoxe Weise verwendet und mit dieser Komponente die Wissensrepräsentation aufgebaut werden. Wesentlich ist hierbei die vorhandene Datenstruktur der Wissensbasis des Entity-Linkers. In dieser werden durch das Python-Programm als zusammengehörend erkannte Krankheiten und Symptome gespeichert. Dabei können Symptome mehreren Krankheiten zugeordnet werden. Die nach Analyse eines Textes in der Wissensbasis gespeicherten Daten werden dann als XML-Datei oder RDF-Modell ausgegeben. Alternativ zum Entity-Linker kann eine eigene Datenstruktur zum Zwischenspeichern der Wissensrepräsentation genutzt werden.

Ein Vorteil des Entity-Linkers ist, dass er mit Sätzen aus den zu analysierenden Texten trainiert werden kann. Eine kontinuierlich wachsende Wissensbasis (durch zahlreiche automatische Textanalysen) vorausgesetzt, erlaubt dem Entity-Linker, im Laufe der Zeit in beliebigen Texten kontextbasiert einem Symptom die richtige Krankheit zuzuordnen. Analysiert man einen Text dann mit dem Entity-Recognizer und danach mit dem Entity-Linker, dann entsteht automatisch die Wissensrepräsentation, wobei etwa die Entität 'lack of energy' vom Entity-Recognizer gefunden wird und als 'Symptom' kategorisiert wird und anschließend vom Entity-Linker der Krankheit 'depression' zugeordnet wird, sofern sich dies aus dem Kontext ergibt.

Arbeitspakete:

- Manuelle Erstellung einer Wissensrepräsentation basierend auf dem zu analysierenden Text.
- Identifizierung eines medizinischen Vokabulars, das zum Training der NER-Komponente verwendet werden kann.
- Untersuchung, auf welche Weise die Trainingsdaten für möglichst gute Ergebnisse der NER-Komponente aufbereitet werden müssen (z.B. durch automatisch erzeugte Beispielsätze). Alternativ Verwendung des Entity-Rulers.
- Entwicklung einer Programm-Komponente, die die zu analysierenden Texte vorverarbeitet, so dass die NLP-Pipeline möglichst effizient arbeitet.
- Entwicklung einer Programm-Komponente, die basierend auf den vom NER gefundenen Entitäten zusammengehörende Entitäten (etwa Symptom und Krankheit) identifiziert und eine Wissensbasis aufbaut, ggf. unter Ausnutzung der Wissensbasis des Entity-Linkers.
- Entwicklung einer Programm-Komponente, die die Wissensbasis als RDF/XML-Datei ausgibt.

STAND DER WISSENSCHAFT

2.1 FZ 1.1. AUFBAU UND STRUKTUR MEDIZINISCHER TEXTE

Medizinische Texte existieren in verschiedenen Formen mit unterschiedlichen Zielsetzungen. Beispielsweise gibt es Lehrbuch- und Enzyklopädietexte, die überblicksweise über medizinische Sachverhalte informieren, wissenschaftliche Studien, in denen neue wissenschaftliche Erkenntnisse vorgestellt werden und außerdem Arztberichte, in denen über individuelle Patienten und deren Krankheitsverlauf informiert wird. Die für diese Arbeit genutzten Analysetexte stellen enzyklopädische Texte dar, in denen insbesondere Symptome und Merkmale psychologischer Erkrankungen beschrieben werden. Charakterisiert sind diese Texte dadurch, dass sie durch Zwischenüberschriften gegliedert sind und neben Fließtext auch Aufzählungen enthalten.

2.2 FZ 2.1. ÜBERBLICK ÜBER DEPRESSIVE ERKRANKUNGEN

Als Quellen für diesen Abschnitt dienten **[mpg_depression]**, **[who_depression]** und **[psychrembel_depression]**.

2.2.1 Allgemein

Depression ist eine psychische Krankheit, die geschätzt 3,8% der Weltbevölkerung und 5% der Erwachsenen betrifft. Sie kann dazu führen, dass betroffene Menschen Schwierigkeiten haben im Arbeits- und Familienleben zurecht zu kommen.

2.2.2 Symptome

Übliche Symptome einer Depression sind eine traurige Grundstimmung, Konzentrationsprobleme, Hoffnungslosigkeit, Müdigkeit und Reizbarkeit. Auch Schlafstörungen, eine Libidostörung, verminderter oder gesteigerter Appetit und ein vermindertes Selbstwertgefühl oder Selbstbewusstsein sind Symptome. Im schlimmsten Fall treten Selbsttötungsgedanken auf.

2.2.3 Formen der Depression

Die Schwere einer typischen Depression wird durch leichte, mittelgradige und schwere Episoden unterschieden. Dabei sind die Übergänge fließend. Besondere Formen der Depression sind die chronische Depression, bei der trotz therapeutischer Maßnahmen nur wenig Besserung eintritt. Des Weiteren

ren gibt es die manisch-depressive Depression, bei der sich depressive und manische Phasen abwechseln. Außerdem gibt es auch kurze, akute depressive Verstimmungen, die nur zwischen einem Tag und zwei Wochen dauern.

2.2.4 Ursachen

Depressionen entstehen durch ein komplexes Zusammenspiel von sozialen, psychologischen und biologischen Faktoren. Dabei wird angenommen, dass für eine typische Depression die Genetik 50% der Ursachen ausmacht. Konkret können Kindheitserfahrungen, Verluste und Arbeitslosigkeit eine Depression begünstigen.

2.2.5 Behandlung

Je nach Schweregrad der Krankheit werden zur Behandlung von Depressionen psychologische Behandlungen angewandt, und/oder den betroffenen Personen Antidepressiva verschrieben.

2.3 FZ 2.2. AUTOMATISIERUNG DURCH NLP

2.3.1 Natural Language Processing

Natural Language Processing (NLP) ist ein Teil der Künstlichen Intelligenz. Er befasst sich mit der Aufgabe Maschinen den Umgang mit der menschlichen Sprache beizubringen, also das Verstehen der Sprache und dem richtigen Antworten darauf. Mithilfe von NLP-Methoden werden zum Beispiel Sätze analysiert, die Bedeutung von Texten erfasst, Chatbots erstellt und sogar ganze Geschichten geschrieben. Die grundlegendsten Schritte, die beim Erarbeiten eines NLP-Modells (oder allgemeiner eines Machine-Learning-Modells) ausgeführt werden, sind folgende:

1. Aufbereitung von Daten
2. Wählen eines geeigneten Algorithmus
3. Trainieren des Modells mit Trainingsdaten
4. Testen des Modells mit Testdaten
5. Modell anwenden / Vorhersagen treffen

2.3.2 Named Entity Recognition

Named Entity Recognition (NER) bezeichnet einen Teilbereich des Natural Language Processing, bei dem es darum geht wichtige Entitäten, wie zum Beispiel Personen, Orte oder Institutionen, in einem Text zu erkennen. Methoden zur Bewältigung dieser Aufgabe werden seit circa 30 Jahren entwickelt.

Darunter fallen grammatikbasierte und statistische Methoden, sowie auch Methoden des Machine Learning.

2.3.3 Entity Linking

Als Entity Linking wird im Bereich des NLP die Aufgabe beschrieben, die den Entitäten (z.B. Personen, Orte) in einem Text das korrekte Äquivalent in einer Wissensbasis zuordnen soll. In [shen_entity_2021] wird Entity Linking folgendermaßen definiert (übersetzt):

Definition 2.3.1 (Entity Linking) *Gegeben sei ein Dokument D , welches die erkannten Entitäten $M = \{m_1, m_2, \dots, m_{|M|}\}$ enthält, sowie eine Ziel-Wissensbasis KB , welches die Entitäten $E = \{e_1, e_2, \dots, e_{|E|}\}$ enthält. Das Ziel ist es jede Entität m_i in M seinem korrekten Äquivalent e_i in E zuzuordnen.*

2.3.4 spaCy

spaCy ist eine Python-Bibliothek, die die erforderlichen Daten und Algorithmen enthält, die zum verarbeiten von Texten mit natürlicher Sprache benötigt werden. SpaCy enthält vortrainierte Modelle für über 70 verschiedene Sprachen, unter anderem Spanisch, Englisch, Griechisch und Deutsch. Zudem ermöglicht spaCy das Einbinden von Modellen, die mithilfe anderer Python-Bibliotheken (zum Beispiel PyTorch oder TensorFlow) trainiert wurden. Anders als einige andere NLP-Bibliotheken konzentriert sich spaCy darauf, Software für den Produktionseinsatz zur Verfügung zu stellen. Erstmals wurde spaCy 2015 von Matthew Honnibal veröffentlicht. ([vasiliev2020natural], [github_spacy], [spacy])

2.3.5 pySBD

[sadvilkar_pysbd_2020] stellen die Python-Bibliothek *pySBD* vor, die Satzgrenzen in Texten ermittelt. Dieser Schritt der Datenaufbereitung ist wichtig, da viele weiterführende NLP-Schritte auf Satzebene durchgeführt werden und daher die korrekte Bestimmung der Satzgrenzen Auswirkungen auf den Erfolg der weiteren Verarbeitung hat. Die Autoren der Studie vergleichen die Genauigkeit ihrer Bibliothek mit anderen Bibliotheken anhand der sogenannten *Golden Rules*, die zahlreiche Beispiele für Sonderfälle von Satzgrenzen in mehreren Sprachen enthalten. Laut ihren Angaben übertrifft die Genauigkeit der Satzgrenzenerkennung von pySBD die Genauigkeit der Satzgrenzenerkennung von spaCy und anderen Bibliotheken.

2.3.6 Weitere Python-Bibliotheken für Machine Learning (ML) und NER

Eine der wichtigsten Python-Bibliotheken für NLP ist NLTK (Natural Language Toolkit) ([bird2006nltk]). Weitere Python-Bibliotheken, die für NLP-

Aufgaben genutzt werden können, sind unter Anderem Gensim ([vrehuuvrek2011gensim]), Pattern ([de2012pattern], [github_pattern]), scikit-learn und PyTorch.

2.4 FZ 3.1. WISSENSREPRÄSENTATION MITTELS RDF

Das *Resource Description Framework (RDF)* [w3c_all_2022] ist ein Framework zur Darstellung von Informationen im Semantischen Web, das von der *RDF Working Group* des *World Wide Web Consortium (W3C)* erstellt wurde. Das RDF-Modell besteht aus einem Datenmodell, mit dem Aussagen über Ressourcen in Form eines Graphen dargestellt werden. Die Informationen werden als Tripel von Subjekt, Prädikat und Objekt gespeichert und ermöglichen auf diese Weise eine maschinenlesbare Bereitstellung semantischer Informationen.

Die Subjekte und Objekte sind dabei die Knoten des Graphen und die Prädikate die Kanten zwischen diesen Knoten.

Abbildung 2.1 zeigt einen einfachen RDF-Graphen, in dem das Subjekt „depression“ mit dem Prädikat „hasSymptom“ mit zwei Symptomen verbunden und über das Prädikat „isA“ mit dem Objekt „mentalDisorder“ verbunden ist.

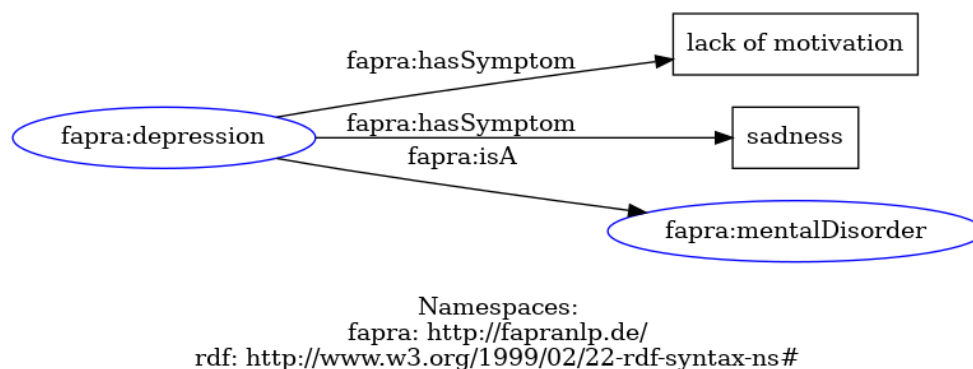


Abbildung 2.1: einfacher RDF-Graph (visualisiert mit <https://www.ldf.fi/service/rdf-grapher>)

2.4.1 Resource Description Framework Schema

Zusätzlich zu RDF stellt das W3C auch eine Empfehlung für ein standardisiertes Datenmodellierungsvokabular für RDF-Daten bereit: das RDF-Schema. Dies stellt eine Erweiterung des grundlegenden RDF-Vokabulars dar und bietet eine Reihe von vordefinierten Klassen (*rdfs:Class*), Ressourcen (*rdfs:Class*) und Eigenschaften (*rdf:Property*), die zur Definition von Klassen und Eigenschaften eines RDF-Modells verwendet werden.

Klassen können in einer Hierarchie organisiert werden, wobei Unterklassen Merkmale von ihren Oberklassen erben. Eigenschaften können als zu einer bestimmten Domäne (*rdfs:domain*) und einem bestimmten Bereich (*rdfs:range*) gehörend definiert werden, womit sich präzise Beziehungen zwischen Ressourcen erstellen lassen.

2.4.2 Dublin Core Metadata Initiative

Die Dublin Core Metadata Initiative (DCMI)[[DCMI_dcmi_2022](#)] ist ein Konsortium von Organisationen, die Standards für die Beschreibung von Online-Ressourcen entwickeln und pflegen. Die Metadatenelemente des *Dublin Core* umfassen 15 standardisierte Kernfelder, die zur Beschreibung der Merkmale von Online-Ressourcen dienen.

Die Dublin-Core-Metadatenelemente sind leicht zu verwenden, aber dennoch flexibel genug, um zahlreiche Ressourcentypen und -kontexte beschreiben zu können. Die Elemente umfassen grundlegende beschreibende Informationen wie Titel, Ersteller und Thema sowie speziellere Elemente zur Identifizierung der Art der Ressource, ihrer Sprache und ihrer Beziehung zu anderen Ressourcen. Außerdem umfasst das DCMI auch mehrere Codierschemata für Wortschätze bestimmter Anwendungsgebiete, unter anderem für MeSH.

2.4.3 Serialisierung von RDF-Graphen

Für die Serialisierung von RDF-Graphen stehen mehrere Formate zur Verfügung. So wird der MeSH-Datensatz im *N-Triples*-Format bereitgestellt und in diesem Praktikum erfolgt die Ausgabe der Wissensrepräsentation im RDF/XML-Format.

Ein Beispiel für die Serialisierung des in Abbildung 2.1 dargestellten Graphen im RDF/XML-Format ist in Listing 2.1 angeführt.

Listing 2.1: RDF/XML

```
<?xml version="1.0" encoding="utf-8"?>
<rdf:RDF
  xmlns:fapra="http://fapranlp.de/"
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
>
  <rdf:Description rdf:about="http://fapranlp.de/depression">
    <fapra:hasSymptom>lack of motivation</fapra:hasSymptom>
    <fapra:hasSymptom>sadness</fapra:hasSymptom>
    <fapra:isA rdf:resource="http://fapranlp.de/mentalDisorder"/>
  </rdf:Description>
</rdf:RDF>
```

Für Python steht mit `RDFLib` [[rdflib_team_rdflib_2022](#)] eine Bibliothek zur Arbeit mit RDF-Graphen bereit.

2.4.4 SPARQL

SPARQL Protocol And RDF Query Language (SPARQL) ist eine Abfragesprache für das semantische Web, mit der Daten, die im RDF-Format vorliegen, abgerufen und bearbeitet werden können. Sie ermöglicht die Abfrage mehre-

rer Datenquellen und -typen und enthält integrierte Funktionen und Operatoren zur Datentransformation. SPARQL kann auch zum Aktualisieren und Erstellen von RDF-Daten verwendet werden und enthält Befehle für Datenänderungen und Integritätsbedingungen.

KONZEPTUELLE MODELLIERUNG UND ENTWURF

In diesem Kapitel wird der Anwendungskontext und konkrete Anwendungsfälle nach UCSD sowie Modelle der technischen Umsetzung einer Automatisierungsunterstützung (Informationsmodell, Komponenten-/Dienstemodelle, Architekturmodell) entsprechend RUP modelliert.

3.1 USER CENTERED SYSTEM DESIGN

Der Anwendungskontext und die Anwendungsfälle werden mithilfe des Ansatzes des User Centered System Design nach [norman1986user] ermittelt. Auf der Grundlage der Anwendungsfälle wird das Informations- und Datenmodell entwickelt und daraus das Komponenten- und Architekturmodell.

Gemäß Aufgabenstellung umfasst die Entwicklung keine graphische Benutzeroberfläche.

3.1.1 Anwendungskontext und Anwendungsfall

Der Anwendungskontext ist eine Unterstützung bei der Behandlung und Betreuung psychisch erkrankter Personen durch Chatbots. Um den Chatbots den erforderlichen Kontext und das entsprechende Hintergrundwissen zu vermitteln, wird eine Wissensrepräsentation über die psychischen Erkrankungen benötigt. Die zu entwickelnde Konsolenapplikation *MedExtractor* ist dazu gedacht, Forschenden und später auch Chatbots zu ermöglichen, aus einem Input in Form einer reinen Textdatei die Wissensrepräsentation zu erstellen, siehe Abbildung 3.1. Dazu soll die Applikation mit dem Namen des Textes als Parameter aufgerufen werden. Das Ergebnis der Applikationsausführung ist eine RDF-Datei.

Als zusätzliche Funktionalität ist denkbar, die Wissensbasis in Form der trainierten spaCy-EntityLinker-Repräsentation auszugeben.

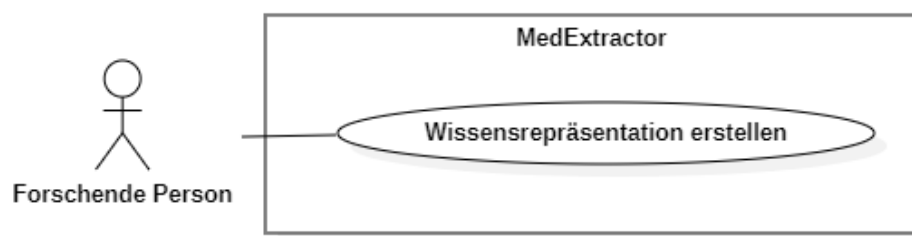


Abbildung 3.1: Anwendungsfall der Konsolenapplikation

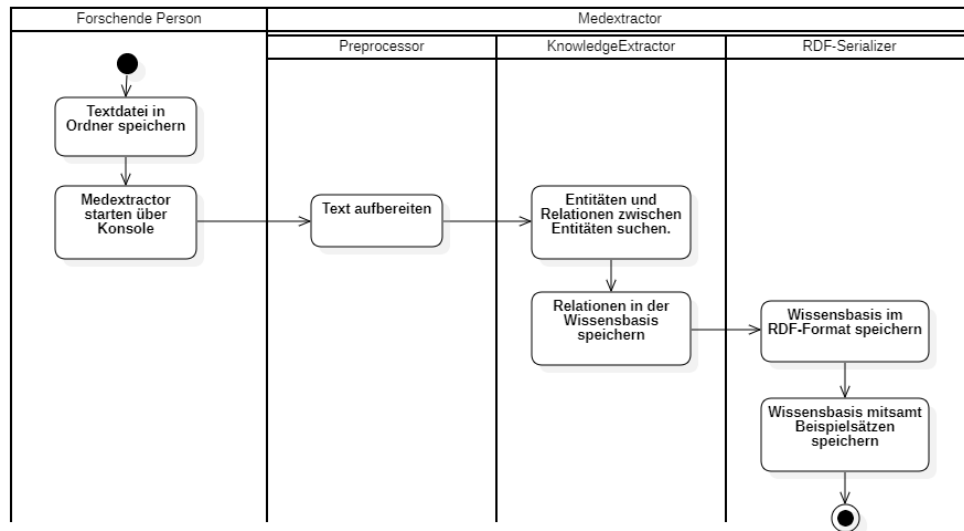


Abbildung 3.2: Aktivitätsdiagramm der Konsolenapplikation

3.1.2 Aktivitätsmodell

Für die Modellierung des Anwendungsfalls sind die in Abbildung 3.2 dargestellten Aktivitäten erforderlich.

3.1.3 Informationsmodell

Für den Aufbau der Wissensbasis sind erstens die gefundenen Entitäten und zweitens die Beziehungen zwischen diesen von Bedeutung. Diese werden in Form von semantischen Beziehungen in der Wissensbasis gespeichert, wie dies in Abbildung 3.3 des Informationsmodells (ER-Diagramm) der Applikation dargestellt ist.

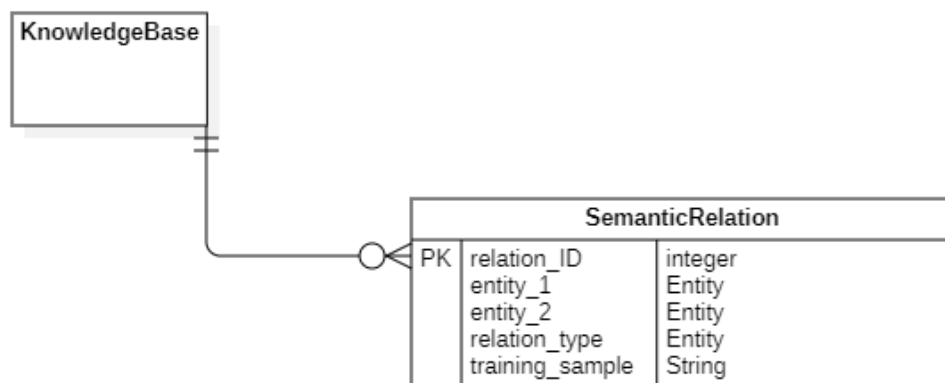


Abbildung 3.3: Informationsmodell der Konsolenapplikation

3.1.4 Architektur- und Komponentenmodell

Aus dem Anwendungsfall und dem Aktivitätsmodell ergibt sich die grobe Gesamtarchitektur der Applikation wie in Abbildung 3.4 dargestellt.

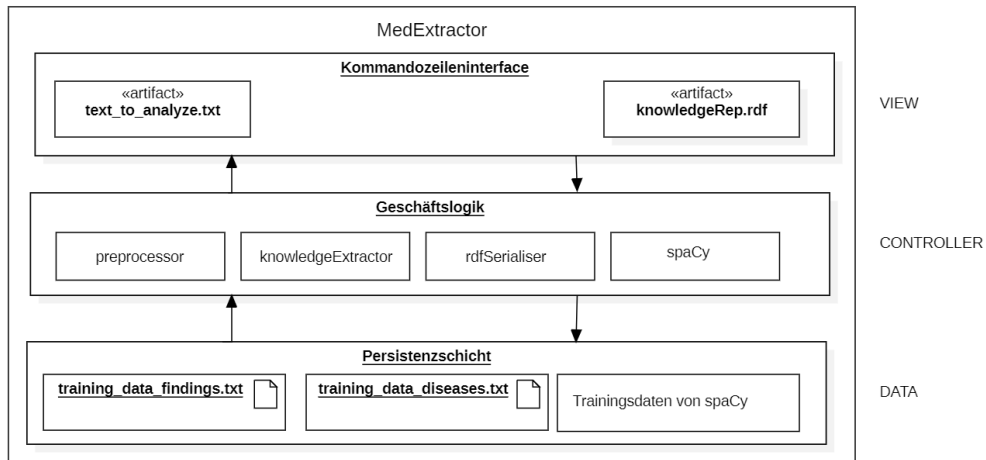


Abbildung 3.4: Architekturmodell der Konsolenapplikation

Abbildung 3.5 stellt das Komponentenmodell der Applikation dar.

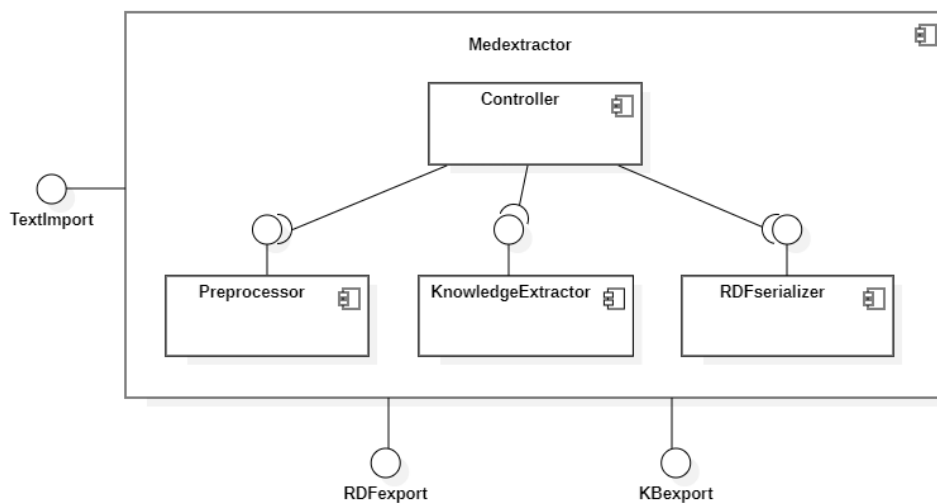


Abbildung 3.5: Komponentenmodell der Konsolenapplikation

3.2 FZ 1.2 THEORIEBILDUNG ZUR VORVERARBEITUNG MEDIZINISCHER TEXTE

Um die Vorverarbeitung medizinischer Texte durchführen zu können, wird dem Präprozessor (Abbildung 3.6) der Originaltext übergeben. Die Umwandlung und Ausgabe des Textes erfolgt über die Methode `get_preprocessed_text()`. Der Text wird in vollständige Sätze umgewandelt, die für sich alleine stehen

können. D.h. Aufzählungen werden in Sätze umgewandelt und Pronomen durch konkrete Substantive ersetzt.

Für die Erkennung der Satzeinheiten kann die Bibliothek *pySBD* - *python Sentence Boundary Disambiguation* verwendet werden, die regelbasiert auf Grundlage von 100 „Goldenen Regeln“ arbeitet.

Falls möglich, soll auch eine Indexstruktur (Kapitel-/Abschnittsüberschriften) ermittelt werden.

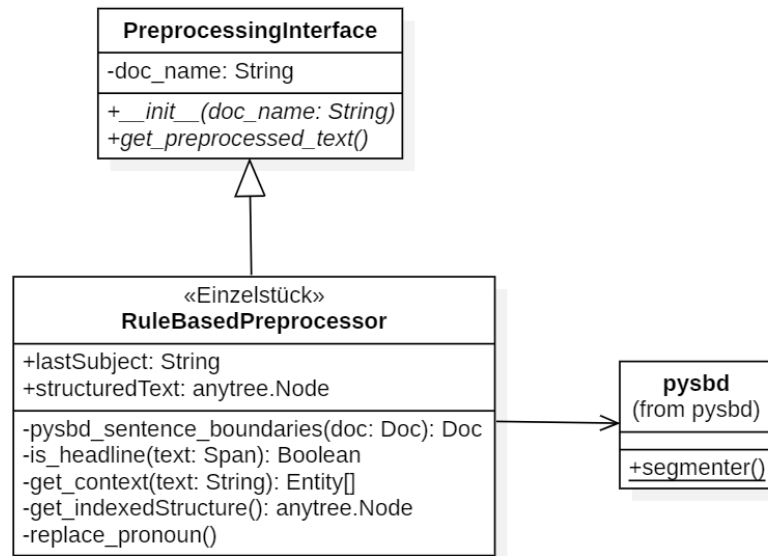


Abbildung 3.6: Klassendiagramm des Präprozessors der Konsolenapplikation

3.3 FZ 2.3 THEORIEBILDUNG ZUR ÜBERFÜHRUNG MEDIZINISCHEN FACHVOKABULARS IN MASCHINENLESBARE FORM ZUR WEITEREN VERARBEITUNG DURCH NLP

Beim Starten des MedExtractors wird eine KnowledgeExtractor-Instanz erzeugt (Abbildung 3.7). Zur Instanziierung wird der Dateiname der persistent gespeicherten Wissensbasis übergeben. Das KnowledgeExtractor-Objekt verwaltet das KnowledgeBase-Objekt. Der EntityRuler wird der spaCy-Pipeline hinzugefügt.

Dieser Instanz kann durch den Aufruf KnowledgeExtractor(text:String) ein vom Preprocessor vorbereiteter Abschnitt oder Satz zur Analyse übergeben werden. Daher muss KnowledgeExtractor die Methode `__call__(text:String)` implementieren. Zusätzlich kann über die Methode `set_context(Entity[])` dem KnowledgeExtractor mitgeteilt werden, um welches Thema es in dem Abschnitt oder Satz geht. Das Thema ergibt sich aus Entitäten, die in Kapitel- oder Abschnittsüberschriften gefunden werden.

`is_related()` ist eine private Methode, die vom KnowledgeExtractor-Objekt verwendet wird, um zu prüfen, ob zwei beliebige vom EntityRuler gefundene Entitäten in Beziehung zueinander stehen. Die SemanticRelation-Objekte

werden auf Basis eines wiederholten Aufrufs von `is_related()` mit unterschiedlichen Kombinationen gefundener Entitäten erzeugt.

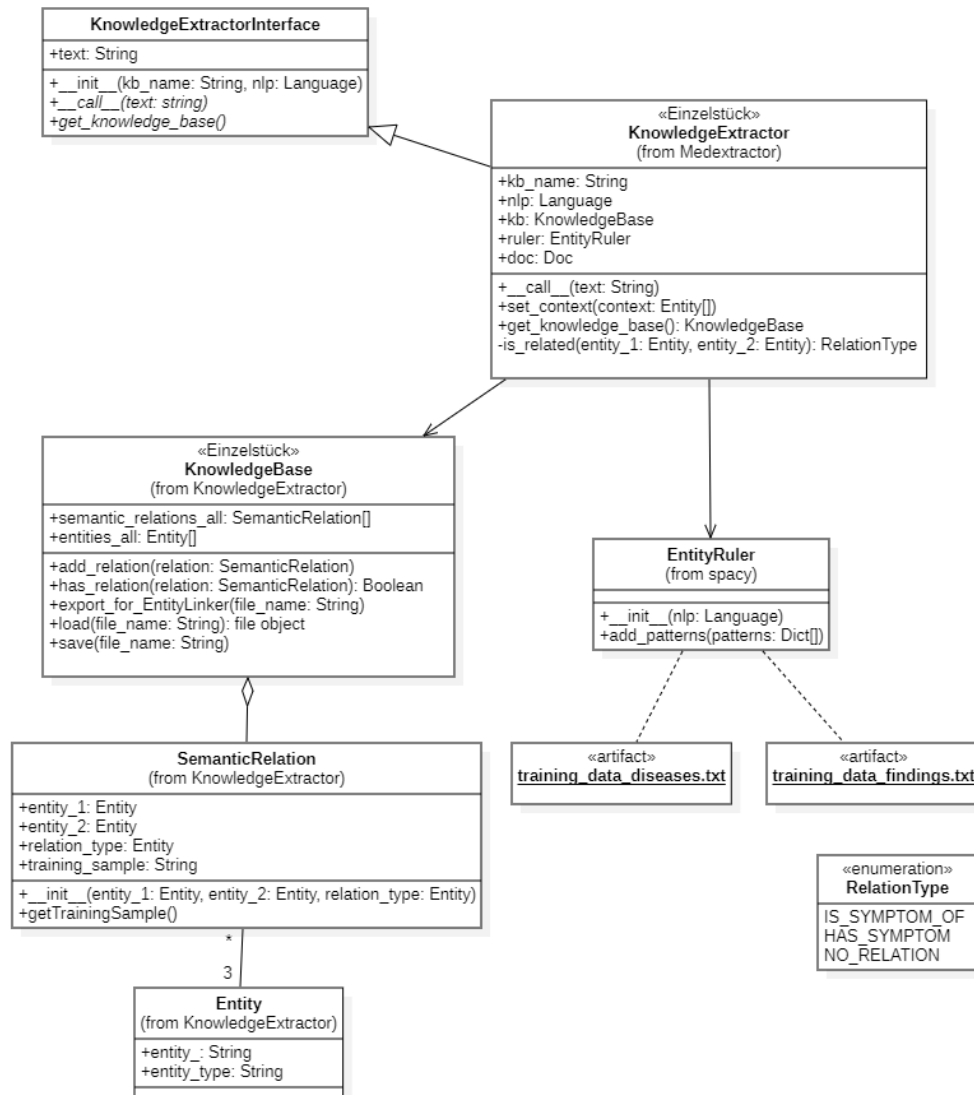


Abbildung 3.7: Klassendiagramm des KnowledgeExtractors der Konsolenapplikation

3.4 FZ 3.2 THEORIEBILDUNG ZUR WISSENSREPRÄSENTATION

Der RDFSerialiser (Abbildung 3.8) serialisiert die Wissensrepräsentation in eine RDF/XML-Datei. Dabei arbeitet der RDFSerialiser mit der Python-Bibliothek `rdflib`. Zunächst soll ein Graph ohne Inhalt erstellt werden mit der Funktion `create_graph`, dies geschieht mit der Klasse `GraphManager`. Anschließend sollen die Inhalte der übergebenen `KnowledgeBase` in den Graphen übertragen werden (`knowledgebase_to_graph`), dies geschieht mit den für `GraphManager` implementierten Methoden. Der so entstandene Graph wird dann mit der Methode `serialise_graph` in das in `serialisation_format`

spezifizierte Format übertragen und anschließend die so entstandene Wissensrepräsentation gespeichert.

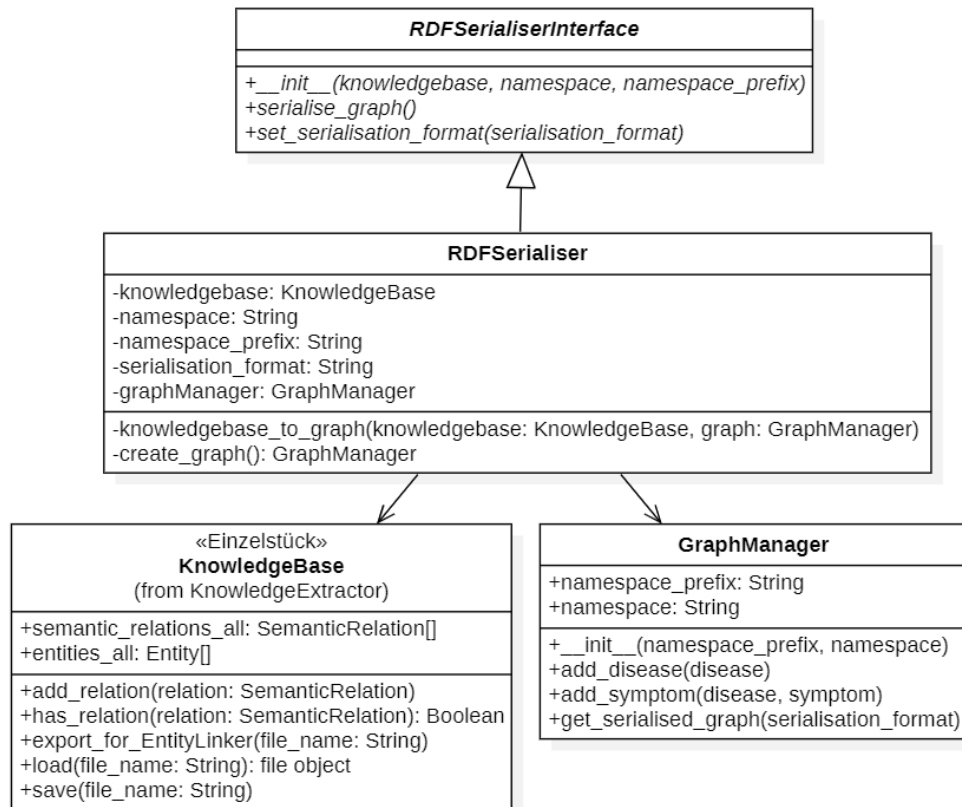


Abbildung 3.8: RDFSerializer der Konsolenapplikation

3.5 ZUSAMMENFASSUNG

Abbildung 3.9 stellt das Hauptklassendiagramm der Applikation dar.

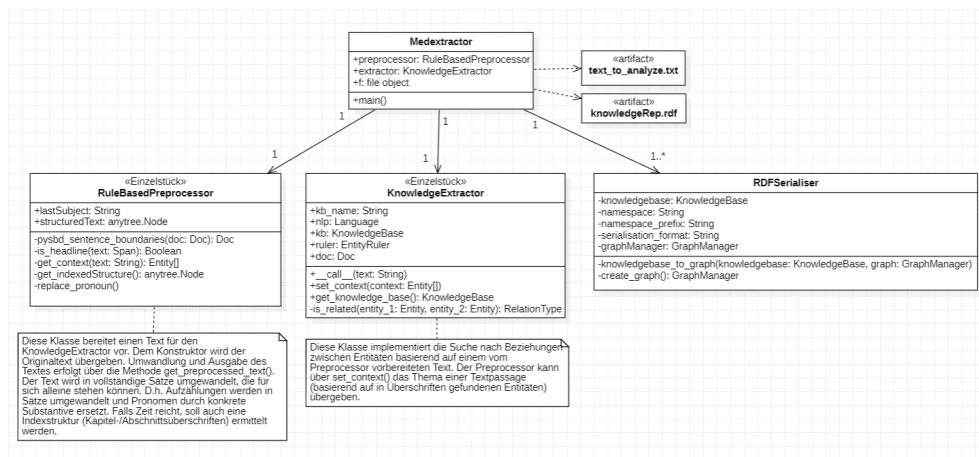


Abbildung 3.9: Hauptklassendiagramm der Konsolenapplikation

PROOF-OF-CONCEPT IMPLEMENTIERUNG

4.1 FZ 1.3 IMPLEMENTIERUNG ZUR VORVERARBEITUNG MEDIZINISCHER TEXTE

Da die zu verarbeitenden medizinischen Texte nicht nur Fließtext enthalten, sondern auch Überschriften und Aufzählungen, ist als erster Schritt eine sinnvolle Unterteilung in Satzeinheiten erforderlich. Dazu kommt die Bibliothek pySBD ([sadvilkar_pysbd_2020]) zum Einsatz. Anschließend werden die Satzeinheiten noch so bearbeitet und zusammengefügt, dass möglichst sinnvolle Sätze für die weitere Verarbeitung entstehen. Ist eine erkannte Satzeinheit zum Beispiel ein leerer String, so wird diese nicht weiter verarbeitet. Nachgestellte Leerzeichen werden gelöscht. Außerdem werden Aufzählungen möglichst in einzelne Sätze verwandelt. So wird beispielsweise aus der Auflistung

The psychological symptoms of depression include:
- continuous low mood or sadness
- feeling hopeless and helpless
- having low self-esteem

eine Reihe von Sätzen:

The psychological symptoms of depression include continuous low mood or sadness.
The psychological symptoms of depression include feeling hopeless and helpless.
The psychological symptoms of depression include having low self-esteem.

4.2 FZ 2.4 IMPLEMENTIERUNG ZUR ÜBERFÜHRUNG MEDIZINISCHEN FACHVOKABULARS IN MASCHINENLESBARE FORM ZUR WEITEREN VERARBEITUNG DURCH NLP

4.3 FZ 3.3 IMPLEMENTIERUNG ZUR WISSENSREPRÄSENTATION

EVALUIERUNG

5.1 FZ 1.4 EVALUIERUNG ZUR VORVERARBEITUNG MEDIZINISCHER TEXTE

5.1.1 Fähigkeiten

5.1.2 Limitierungen

Bei einigen Formatierungen von Sätzen / Informationsauflistungen hat die aktuelle Implementierung des Preprocessors noch Probleme. So kann zum Beispiel der folgende (verkürzte) Abschnitt einer Informationsseite zum Thema *Drug-Dependent People* ([**drug_dependent_people**]) noch nicht in sinnvolle Sätze transformiert werden:

All drug dependency exhibits similar characteristics— • Chronic, progressive, and relapsing disease • Denial • Lying and deceit • Changes in normal behavior

5.2 FZ 2.5 EVALUIERUNG ZUR ÜBERFÜHRUNG MEDIZINISCHEN FACHVOKABULARS IN MASCHINENLESBARE FORM ZUR WEITEREN VERARBEITUNG DURCH NLP

5.3 FZ 3.4 EVALUIERUNG ZUR WISSENSREPRÄSENTATION

5.4 ZUSAMMENFASSUNG

ZUSAMMENFASSUNG UND AUSBLICK

6.1 ZUSAMMENFASSUNG

6.2 AUSBLICK