

Fernuniversität Hagen

Modul 63458 - Fachpraktikum

Natural Language Processing, Information Extraction und Retrieval  
WiSe 2022/23

Exposé

**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, 3. November 2022

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>3</b>
<b>2</b>	<b>Motivation</b>	<b>4</b>
<b>3</b>	<b>Problembeschreibung</b>	<b>5</b>
<b>4</b>	<b>Forschungsfragen und Ziele</b>	<b>6</b>
<b>5</b>	<b>Lösungsansatz</b>	<b>7</b>
<b>6</b>	<b>Stand der Wissenschaft</b>	<b>8</b>
<b>7</b>	<b>Arbeits- und Zeitplan</b>	<b>9</b>

# 1 Einleitung

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

## 2 Motivation

Die spezielle Herausforderung unserer Praktikumsaufgabe liegt darin, automatisiert Zusammenhänge zwischen Entitäten medizinischer Texte zu finden und diese Zusammenhänge zu repräsentieren. Dem zugrunde liegt das Ziel des Dissertationsprojektes von Stefanie Heidepriem, zu erforschen, wie ein textbasierter Chatbot zur Unterstützung psychisch kranker Menschen entwickelt werden kann. Es existieren bereits Frameworks, die in der Lage sind Entitäten, also wichtige Schlüsselbegriffe, automatisiert aus Texten zu extrahieren. Dies nennt man Named Entity Recognition.

Named Entity Recognition ... [TODO, einige Worte zu NER, aktuelle Entwicklungen, Erkenntnisse im medizinischen Bereich]

Das übergeordnete Dissertationsprojekt von Stefanie Heidepriem ist unter anderem an das Projekt MENHIR angelehnt. Dieses Projekt dient der Erforschung von Konversationstechnologien, die psychisch kranke Menschen unterstützen sollen. In diesem Zusammenhang wird auch ein MENHIR-Chatbot entwickelt, der psychisch kranken Menschen personalisierte Unterstützung und hilfreiche Bewältigungsstrategien bieten soll.

### 3 Problembeschreibung

Es soll eine Konsolenapplikationen entwickelt werden, die medizinische Texte analysiert und das darin enthalten Wissen strukturiert in einer XML-Datei dokumentiert oder mittels eines RDF-Modells visualisiert werden. Die Konsolenapplikation soll in Python programmiert werden, wobei für die Textanalyse Klassen und Methoden der Open-Source-Bibliothek spaCy 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 dann, 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 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.

Neben der grammatikalischen Analyse besteht eine wesentliche Aufgabe von spaCy darin, Schlüsselbegriffe zu finden und zu kategorisieren. Diese Art der Analyse wird allgemein als Named Entity Recognition (NER) bezeichnet. Bei spaCy übernimmt diese Aufgabe der regelbasierte EntityRuler oder der auf statistischen Modellen basierende EntityRecognizer. Eine weitere Komponente ist der EntityLinker, mit dem Begriffe eindeutig den in einer Wissensbasis gespeicherten Entitäten zugeordnet werden können. Auch der EntityLinker 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-Komponenten von spaCy verwendet werden kann.

Die auf dem Campus des National Institute of Health (NIH) im US-Bundestaat Maryland angesiedelte National Library of Medicine 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 der vom National Library of Medicine entwickelte Medical Subject Heading (MeSH)-Thesaurus. 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.

Basierend auf den gefundenen Entitäten soll das Python-Programm in der Lage sein, Begriffe wie Krankheiten und Symptome richtig zuzuordnen. Hierzu müssen Aussagesätze, Fragen, Aufzählungen und Überschriften richtig gedeutet werden.

- Medizinische Begriffe, die mithilfe von Named Entity Recognition gefunden werden, müssen zueinander in Zusammenhang gebracht werden.
- Es gibt keine vorgegebene / einheitliche Struktur für die Wissensrepräsentation von Entity Linking.

## 4 Forschungsfragen und Ziele

**Frage 1** *Wie kann die MESH Datenbank in spaCy importiert und von der NER-Komponente genutzt werden?*

**Frage 2** *Welche Entities werden mit der MESH Datenbank gefunden? Wie unterstützen diese Entities den Aufbau der Wissensrepräsentation?*

**Frage 3** *Was muss in der Wissensdatenbank der Entity-Linker-Komponente von spaCy gespeichert werden?*

**Frage 4** *Wie kann der Text in einzelne Sätze gegliedert werden, so dass auch Überschriften, Leerzeilen und Aufzählungen berücksichtigt werden?*

**Frage 5** *Wie können Überschriften erkannt werden?*

**Frage 6** *Wie können Aufzählungen erkannt und einem Thema zugeordnet werden?*

**Frage 7** *Wie sieht eine sinnvolle Wissensrepräsentation aus?*

## 5 Lösungsansatz

Arbeitspakete:

- Ein medizinisches Vokabular muss eingebunden werden um Named Entities identifizieren zu können (MEsH).
- Eine Named Entity Recognition muss für den zu analysierenden Text ausgeführt werden.
- (Eine Knowledge Base für das Entity Linking muss erstellt werden.)
- Eine manuelle Wissensrepräsentation basierend auf dem zu analysierenden Text muss erstellt werden.
- Die automatisch erstellte Wissensrepräsentation muss mit der manuell erstellten Wissensrepräsentation verglichen werden.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

## 6 Stand der Wissenschaft

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.



## 7 Arbeits- und Zeitplan

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.